Francesco Bellotti ·
Miltos D. Grammatikakis · Ali Mansour ·
Massimo Ruo Roch · Ralf Seepold ·
Agusti Solanas · Riccardo Berta   *Editors*

# Applications in Electronics Pervading Industry, Environment and Society

## APPLEPIES 2023

Springer

# Lecture Notes in Electrical Engineering 1110

## Series Editors

The book series *Lecture Notes in Electrical Engineering* (LNEE) publishes the latest developments in Electrical Engineering—quickly, informally and in high quality. While original research reported in proceedings and monographs has traditionally formed the core of LNEE, we also encourage authors to submit books devoted to supporting student education and professional training in the various fields and applications areas of electrical engineering. The series cover classical and emerging topics concerning:

- Communication Engineering, Information Theory and Networks
- Electronics Engineering and Microelectronics
- Signal, Image and Speech Processing
- Wireless and Mobile Communication
- Circuits and Systems
- Energy Systems, Power Electronics and Electrical Machines
- Electro-optical Engineering
- Instrumentation Engineering
- Avionics Engineering
- Control Systems
- Internet-of-Things and Cybersecurity
- Biomedical Devices, MEMS and NEMS

For general information about this book series, comments or suggestions, please contact leontina.dicecco@springer.com.

To submit a proposal or request further information, please contact the Publishing Editor in your country:

**China**

Jasmine Dou, Editor (jasmine.dou@springer.com)

**India, Japan, Rest of Asia**

Swati Meherishi, Editorial Director (Swati.Meherishi@springer.com)

**Southeast Asia, Australia, New Zealand**

Ramesh Nath Premnath, Editor (ramesh.premnath@springernature.com)

**USA, Canada**

Michael Luby, Senior Editor (michael.luby@springer.com)

**All other Countries**

Leontina Di Cecco, Senior Editor (leontina.dicecco@springer.com)

**\*\* This series is indexed by EI Compendex and Scopus databases. \*\***

Francesco Bellotti · Miltos D. Grammatikakis ·
Ali Mansour · Massimo Ruo Roch ·
Ralf Seepold · Agusti Solanas · Riccardo Berta
Editors

# Applications in Electronics Pervading Industry, Environment and Society

APPLEPIES 2023

Springer

*Editors*
Francesco Bellotti 
Department of Electrical, Electronic,
Telecommunication Engineering and Naval
Architecture (DITEN)
University of Genoa
Genoa, Italy

Ali Mansour 
Lab-STICC, CNRS UMR 6285
École nationale supérieure de techniques
avancées Bretagne
Brest, France

Ralf Seepold 
Department of Computer Science
HTWG Konstanz—University of Applied
Sciences
Konstanz, Germany

Riccardo Berta
Department of Electrical, Electronic,
Telecommunication Engineering and Naval
Architecture (DITEN)
University of Genoa
Genoa, Italy

Miltos D. Grammatikakis
Department of Electrical and Computer
Engineering
Hellenic Mediterranean University
Heraklion, Greece

Massimo Ruo Roch 
Department of Electronics
and Telecommunications
Politecnico di Torino
Turin, Italy

Agusti Solanas 
Department of Computer Engineering
and Mathematics
Universitat Rovira i Virgili
Tarragona, Catalonia, Spain

# Introduction

The 2023 edition of the conference on Applications in Electronics Pervading Industry, Environment and Society (Applepies) is held in Genova, Italy, on September 28–29, at the Polytechnic School (Aula Benvenuto, Department of Architecture and Design). It is the eleventh edition of the conference, the first one without the presence of its founder, Prof. Alessandro De Gloria, who passed away this spring. The first article of these proceedings is devoted to retracing his outstanding human and professional profile.

After a thorough blind-review selection process, 34 full papers and 39 short presentations have been accepted, covering edge computing, embedded systems, machine learning, hardware acceleration, computer architecture, etc. The application domains include transportation, energy, security, health, etc. Moreover, a call for special sessions has been launched this year. Besides the conference's main track, there will be for the first time a special track on "In-Home Measurement and Analysis of Health-Related Physiological Parameters."

All these topics show the pervasiveness of today's electronic engineering, which calls for multifaceted knowledge and vision to exploit state-of-the-art methods and tools in information processing, storage, and networking.

Applepies 2023 offers academicians and industry practitioners the possibility of discussing the design, prototyping, and testing of various electronics-enabled systems and applications characterized by innovation, high performance, real-time operations, and requirement compliance (production time, cost, device size, weight, power consumption, etc.).

| | |
|---|---|
| Genoa, Italy | Francesco Bellotti |
| Heraklion, Greece | Miltos D. Grammatikakis |
| Brest, France | Ali Mansour |
| Turin, Italy | Massimo Ruo Roch |
| Konstanz, Germany | Ralf Seepold |
| Tarragona, Spain | Agusti Solanas |
| Genoa, Italy | Riccardo Berta |

# Contents

## Machine Learning

## Edge Computing and Sensing

## Short Contributions

**Alessandro De Gloria**

# Alessandro De Gloria, a Pioneer in Electronic Engineering Applications

Francesco Bellotti[1]([✉]) [iD], Elisa Bricco[1], Agostino Bruzzone[1], Daniele Caviglia[1], Ermanno Di Zitti[1], Paolo Gastaldo[1], Daniele Grosso[1], Lauro Magnani[1], Mauro Olivieri[2], Marco Raggio[1], Maurizio Valle[1], Alessandro Verri[1], and Riccardo Berta[1] [iD]

[1] Università degli Studi di Genova, Genova, Italy
{francesco.bellotti,riccardo.berta}@unige.it
[2] Università La Sapienza, 00184 Roma, Italy

**Abstract.** This article aims to sketch the figure of Alessandro De Gloria, a professor who dedicated his entire life, with generosity and enthusiasm, to engineering and scientific research. He designed one of the first chips in an Italian university and then a set of digital system architectures that contributed to shaping the then-fledgling field of microprocessors. Also, by founding the Applepies conference (International Conference on Applications in Electronics Pervading Industry, Environment and Society), he promoted the intuition of the value for the whole population of the applications of electronic systems. He pioneered the field of mobile apps and started a fruitful dialogue with the humanistic culture, particularly exploiting virtual reality technologies. In the last few years, he also contributed to machine learning and big data, particularly in the domain of automated driving. He founded an MSc program on strategy, intending to shape education in that field from a rigorous engineering point of view. We believe that his generous and enthusiastic academic life can be an outstanding example in the face of today's challenges of Electronic Engineering and higher education in the broader societal context.

**Keywords:** Electronic engineering · Higher education · Electronic system applications · Embedded systems · Cultural heritage · Virtual reality · Machine learning

## 1 Introduction

The 2023 edition of the International Conference on Applications in Electronics Pervading Industry, Environment and Society (Applepies) is the first without Alessandro De Gloria. He founded the conference in 2013 (Rome) and chaired all the editions until 2017. Then he served as emeritus chair—not honorary role, in his interpretation -. His last public presence was at Applepies 2022, in Villa Cambiaso, the seat of the Aula magna of the Polytechnic School of the University of Genova, for the commemoration of his teacher, Alessandro Chiabrera. Heavily hit by a long-lasting illness, he had carefully organized that event, besides overviewing the whole organization, from the keynote

speeches to the industry-academy dialogue panels. He would have then continued his academic life, even from remote, and died on March 20th, 2023, aged 68, surrounded by family affection. Given Alessandro's outstanding figure in the academic landscape, the authors—professors who collaborated with him at the University of Genova—thought it would be appropriate to recall him to pass on his memory to new generations and to those who did not know him personally (Fig. 1).



**Fig. 1.** Prof. Alessandro De Gloria

Alessandro De Gloria was a Full Professor of Electronic Engineering at the Electrical, Electronics and Telecommunication Engineering and Naval Architecture Department (DITEN) of the University of Genova, where he was teaching at the Electronic Engineering MSc, Engineering Technology for Strategy MSc., Computer Science MSc. And Art History MA.

He was passionate about Electronic Engineering, particularly digital system design. He was very attentive to the applications of electronics, as the impactful, "visible" end of an extremely powerful technology that had rapidly become pervasive, but also "hidden", in the world. To promote the study of electronic engineering applications and the research dialogue with the industries, he founded the Applepies conference and served as responsible for the "Electronic Systems and Applications" scientific area within the Società Italiana di Elettronica (SIE), for which he also seated as national councilor. This was an insightful intuition, as, in the current shortage of Engineering students, we have realized the value of the applications as a privileged way to illustrate the fascination and potential of Electronic Engineering to young people.

Alessandro was a person with a remarkable breadth of vision and strongly believed in the dialogue among people and disciplines. He was long-sighted, full of ideas and

gifted in planning, but wanted to share his ideas and plans with a wide community. As mentioned, he served in the SIE, was president of the Electronic Engineering degree course, and founded the Serious Games Society. He strongly believed in the importance of the mentality of the Electronic Engineering (able to manage and design hardware and software systems in an integrated and harmonized view) in the current pervasiveness of information and communication technologies (ICT).

This article aims to sketch the figure of a professor who dedicated his entire life, with generosity and enthusiasm, to University, Electronic Engineering and the broader development of project ideas and knowledge at local, national and international levels. Having he preferred to be coherent with his intuitions rather than following easy success, his academic life was challenging and met not infrequent disappointments. The authors argue that recalling him and learning from his experience can benefit the academic and R&D community in Electronic Engineering and far beyond.

## 2 Designing One of the First Chips in an Italian University

Alessandro enrolled in the Electronic Engineering degree course at the Engineering Faculty of the University of Genoa in 1974. In 1978, under the guidance of Prof. Alessandro Chiabrera and Giancarlo Parodi, he began the first studies that would lead him, together with his colleague Daniele Caviglia, to the realization of the degree thesis project. The aim was to create an arbitrary waveform generator designed to study the reactions of living cells (primarily) "in vitro" when subject to electromagnetic fields. The generated signal would then have been suitably amplified to drive the coils surrounding the cellular samples. The degree exam was held on January 29, 1980, and Alessandro passed it with honors. An extension of the thesis project, completed the following year, was published in IEEE Tr. on Instrumentation and Measurement [1].

The collaboration with the Institute of Electrical Engineering was the natural continuation of the thesis activity. Since Ph.D. courses were not yet active in Italy in those years, Alessandro enrolled, in the same year, in the newly established "School of Specialization in Computer Engineering" of the University of Genoa, chaired by Prof. Arrigo Frisiani. It was a very stimulating path, which also, in this case, led to an innovative thesis work: it was, in fact, under the supervision of Prof. Joy Marino, the design of a chip (probably one of the first ones designed and realized in an Italian University) in Negatively Doped Metal Oxide Semiconductor (NMOS) 5m technology, of a Bus Arbiter for a Multiprocessor Systems [2]. The chip was fabricated, with the support of Prof. Paolo Antognetti, in one of the very first "runs" offered by the newborn Metal Oxide Semiconductor Implementation System (MOSIS) service. The specialization exam was held in 1982.

This was an important turning point for Alessandro's career, which introduced him to the world of Very large-scale integration (VLSI) design. At the same time, he undertook a collaboration with Ansaldo Impianti, and in particular with its Automation Division, regarding the implementation of the monitoring systems for the PEC (Prova Elementi Combustibile) experimental nuclear power plant, which was under construction in those years. This activity led to the foundation of a start-up, namely Cybertek, which he carried on until November 1983, when he joined the Institute of Electrical Engineering of the University of Genoa as a Research Fellow.

## 3   Digital System Architectures

In the mid of the 1980s, Alessandro targeted the design and development of a 32-bit microprocessor, a very challenging research objective at that time. The project had two main purposes: (i) to study in depth the VLSI design process; and (ii) to design the core hardware for the development of processors for in-house applications. The P32 was characterized by orthogonality and composability of data types, instructions and addressing modes; these features, in conjunction with regularity in the implementation of the instructions, allowed an easy implementation of compilers and a more efficient code generation. The P32 provided a large uniform address space (4 gigabytes), a virtual memory capability, and basic support of Operating Systems (OS) functions. Four privilege levels, privileged instructions and support of memory access privileges provided the basic mechanisms to guarantee the security of the computer system. The instruction set supported general-purpose instructions and could be extended by a co-processor capability. The microarchitecture was designed to be expandable. Alessandro led, with passionate enthusiasm and clear research vision, a group of master thesis students who, under his guidance, achieved remarkable results [3].

Starting from 1990 and for almost a decade, Alessandro continuously supported and contributed to a research field that, at the time, was substantially unprecedented in the Italian academic scenario, i.e. microprocessor architecture design. The first efforts were devoted to the application of instruction level parallelism to Prolog processing [4], followed by a general-purpose innovative study based on static scheduling on Very Long Instruction Word (VLIW) microarchitectures [5, 6]. In this period, Alessandro organized in Santa Margherita Ligure (Genova) the 1995 edition of the IEEE/ACM International Symposium on Computer Architecture (ISCA). For the first time in the history of that outstanding conference, it was held in Italy.

Further research directions addressed the implementation of self-timed microarchitectures, for which a complete library of dedicated standard cells was designed [7] as well as innovative arithmetic units [8], and the implementation of co-processors dedicated to fuzzy logic processing [9–11] and for the Boltzmann Machine [12, 13], which anticipated the advent of computation hardware acceleration units that is mainstream nowadays, after 25 years.

## 4   Dialoguing with Humanistic Culture

Beside continuing the work in the field of electronic systems and architectures, Alessandro (with the support of his group, the ELIOS Lab) pioneered a research dialogue with the humanistic culture, particularly developing two main threads: cultural heritage and languages.

In the city of Genova, the 1990s had been characterized by a renewed focus on the historical and artistic heritage in the wake of a major exhibition entitled "Genoa in the Baroque Age", which emphasized the city's historical role in the circuit of European power during the ancien régime, while at the same time pointing to the theme of local cultural heritage as a great potential for a contemporary city. At the end of that decade, Alessandro, who had realized the importance of that line far beyond the art-historical

context, got in touch with the organizers of the exhibition, Giovanna Terminiello and Ezia Gavazza, who introduced him to Prof. Lauro Magnani, knowing his interest in new experiments. This began a dialogue and friendship that continued for more than twenty years. The early stages were very intense and exciting: alongside the emerging potential of integrating engineering and humanistic knowledge, there was the need to find a common language. A first major achievement consisted of a set of 3D renderings of "cubist" figures by Luca Cambiaso, a renowned Genoise Renaissance painter. After several experiments, the work, together with the relevant virtual reality movie, was taken to Germany on the occasion of an exhibition dedicated to the artist in Osnabrueck. The reading and presentation to the public of the urban artifact in its historical transformations analyzed through the sources, verified on the ground and visualized in a 3D reconstruction, was another of the focuses of the research collaboration between the humanist and engineering research groups.

Along the same line, the Strada Nuova project realized a 3D Virtual Reality video clip illustrating the history of the Renaissance "Via Aurea" road that was built as the representative district of the, at that time, world-powerful Genoise bank families. The clip was exhibited for years in the Museo Rosso, and Carlo Azeglio Ciampi, the Italian President of the Republic, watched it with interest with 3D glasses.

A particularly lively field of experimentation and debate has been the museum field. The discussion on the use of the virtual as a 'recontextualisation' focused on issues that had emerged during the layout design of the Museo Diocesano in Genova and then on the possibility of working on "virtual" museums [14] and experimenting interactive forms with the public [15].

The second thread of interest in the humanistic area concerned natural languages. In 2003, Alessandro contacted the F@rum research group and, in particular, Prof. Elisa Bricco, of the Faculty of Languages and Foreign Literatures of the University of Genova. The group was working on applying new Information and Communication Technologies (ICT) to research and language teaching, and Alessandro wanted to create a multimedia object for language teaching with a playful perspective. The term "edutainment" he proposed made his colleagues skeptical, but Alessandro was far-sighted and had already understood the importance of marketing and communication, even for educational products.

His proposed project was appreciated, and a collaboration was established to create a series of educational paths. These mini-courses were based on accounts from writers who had passed through Genoa during the "Grand Tour", describing places and customs. Thanks to funding from the Liguria Region, the "Scuolagiocando lingue" (School playing with languages) was created, which was available on the Region's portal for several years. The final product consisted of three videos with animations. Starting from a chance encounter between a young modern traveler and the writer (Alexandre Dumas for the French language, Charles Dickens for English and Miguel de Cervantes for Spanish), the users could follow them on a discovery journey that mirrored what the writers had described in their works. The illustrious traveller accompanied the young person, engaging in conversations in the foreign language and sharing overall impressions, sensations, and appreciation for the food and Genoese hospitality. Each course presented a range of glottodidactics tools for lexical and syntactic insights.

## 5   Games and Simulations

Towards the end of the millennium, Alessandro had the intuition that applications of electronic systems would soon have had a significant impact on the daily life of a large part of people. Particularly, he committed to studying and developing the application of two, at that time, early-emerging technologies: touch screen and 3D virtual reality. Serious Games represented a challenging domain for both these technologies, intending to build potentially great tools for education and training.

In the previous section, we presented his dialogue with the humanistic culture that spurred from these premises. Here we recall some results he and his group achieved (also by leading and/or participating in European research projects, such as e-Tour, games@large, ELU, Chi:Kho, VitalMind), particularly the design of applications that now look like early prototypes of smartphones' apps. The tour guide for Genoa's Aquarium [16] and the VeGame territorial game for Venice [17] were developed as pioneer demonstrators of such new and widely spreading technologies.

Based on these and similar experiences, he led the Games and Learning Alliance (GaLA) Network of Excellence (NoE) (2010–2014). To continue the aims of the GALA NoE, Alessandro founded in July 2012 the Serious Games Society (SGS), of which he was president until 2017. From the SGS, two main initiatives blossomed: the International Journal of Serious Games (IJSG) and the Games and Learning Alliance (GALA) Conference (12 editions until 2023). The IJSG, which Alessandro firmly wanted to be a gold open-access publication, is indexed in Scopus and ESCI Web of Science and has reached the tenth year of publications. Alessandro was its Editor in Chief since the foundation until his death. His achievements in serious games are outstanding and briefly summarized in [18].

Anticipating the current popularity of the topic, Alessandro's talent in finding new applications of electronics in useful sectors could not miss to address health and medicine. Particularly, his contribution concerned simulation-based training, following the 'serious game' approach [19]. As an evolution of the pre-existing medical simulation laboratory, Alessandro contributed to creating the University of Genova's Center for Simulation and Advanced Training (SIMAV), of which Scientific Council he was a member from 2015 to 2020. He significantly contributed to the mission of the Center by developing simulation-based training and assessment tools for medicine, health and accessibility [20]. He also proposed and contributed to developing a brand-new approach in simulation-based medical serious games. The approach follows the Agent-Based Model Simulation (ABMS) paradigm, exploiting 3D VR devices. The approach led to an easy-to-use, immersive, real-time simulation, providing a realistic experience to the user [21].

## 6   Automotive Applications, Machine Learning and Big Data Management

The Active FP5 project, which developed the first digital dashboards in automotive cockpits—in 1998, when the liquid crystal displays (LCDs) were still an expensive rarity-, marked the beginning of a long-time collaboration with Centro Ricerche Fiat

(CRF) and several other original equipment manufacturers (OEMs). In the years, the collaboration involved several research fields and projects, such as Human-Computer Interaction (Comunicar and Aide), machine learning (Edel), vehicle-to-X cooperation (Safespot), and collaborative mobility (Team). In the last period, Alessandro and his group have been deeply involved in piloting automated driving functions (ADFs) for SAE level 3 and 4 vehicle automation, particularly by developing the big data management architecture [22], based on the Measurify open-source framework [23], and developing machine learning models for perception (e.g., [24]) and decision making (e.g., [25]). Also, in this area, a key appreciated peculiarity was the multidisciplinary approach, which meant designing applications and systems able to effectively exploit the underlying hardware in a holistic, application-requirement-centred view.

## 7   Higher Education

Alessandro was keen to combine his enthusiasm for research with a passion for higher education, to which he dedicated relentlessly, also launching and leading several new initiatives. From the late '90s to 2016, he served as the President of the Degree Course in Electronic Engineering at the University of Genoa. He managed the subdivision of the 5-year course into a BSc and an MSc course. He invented the role of the didactic manager, who acted as an appreciated interface between teachers and students. The realization of a new website for the courses facilitated the management of formal procedures (stages, thesis selection, etc.) and helped students access practical information and educational material (lecture notes for all the teachings were also printed and distributed to the students). He introduced in the syllabus new courses that only recently have become mainstream, like "Soft Skills", "Entrepreneurship", and "Orientation workshops", to introduce students to the labour market. By proposing projects for his courses' exams, he promoted the design and development of educational games and also organized game-based events with students to challenge and stimulate them in unusual contexts.

Alessandro also felt the need to translate into didactics the ongoing experimental research dialogue between information science and technology and the humanistic culture: this is how the teaching of "Art & Image Narrative for Virtual Worlds" was born in the Electronics Engineering MSc and then "3D visualizations for the analysis of the artistic and architectural heritage" in the MA in History of Art and Enhancement of the Artistic Heritage. Alessandro's work for the DIRAAS courses allowed many students with a humanistic background to learn and practice 3D reconstruction techniques. The numerous theses of which De Gloria was supervisor or co-rapporteur were very useful gyms for combining the art history vocation with digital applications to the point of being translated into scientific publications, doctoral experiences, research grants, or even specialized work activities.

Alessandro's creativity and visionary leadership in conceiving new generations of engineers also led to the creation of the MSc program in "Engineering Technology for Strategy (and Security)" in 2018, after some years in which he did not shy away from academic battles conducted with his proverbial stubbornness and strength. The program combines modelling, simulation, machine learning, big data analysis and social, economic and political studies to prepare the figure of a professional able to manage

strategic decision-making processes based on a solid engineering foundation. It is now one of the most attended MSc programs of the University of Genova, attracting students from all over the world.

## 8    Conclusions

As an outstanding academic and charismatic leader, Alessandro has been a highly valued supervisor and mentor for many generations of students and doctoral candidates and an essential support for his colleagues at the University of Genova, with not a few of whom had an intense dialogue.

He always demonstrated remarkable open-mindedness and a genuine curiosity towards the whole reality, without preclusions and with a visionary yet project-oriented approach. He used to make to students and colleagues challenging proposals, with no discounts. They were difficult to implement but long-sighted and impactful. He never surrendered to the difficulties of various types that he met on his paths. Even from his bed of death—that he never referred to as such—he continued to discuss with his colleagues and make plans for the future. You thought you had to bring him words of comfort; he gave them to you instead, together with ideas for collaborations.

He felt an unquenchable push to look at and shape the future. We are confident that he has now reached the plenty of the life that he so much desired.

## References

1. Caviglia D, De Gloria A, Donzellini G, Parodi G, Ponta D (2023) Design and construction of an arbitrary waveform generator. IEEE Trans Instrum Meas 32(3). https://ieeexplore.ieee.org/document/4315094
2. Caviglia D, De Gloria A (1982) Bus arbiter for multiprocesor systems—specialization thesis. Università di Genova
3. De Gloria A (1985) External and internal architecture of the P32, a 32 bit microprocessor. In: Antognetti P, Anceau F, Vuillemin J (eds) NATO ASI series. Springer Netherlands, Dordrecht, pp 127–146. https://doi.org/10.1007/978-94-009-5143-3_5
4. De Gloria A, Faraboschi P (1992) Instruction-level parallelism in Prolog: analysis and architectural support. ACM SIGARCH Comput Archit News 20(2):224–233. https://doi.org/10.1145/146628.139730
5. Costa A, De Gloria A, Faraboschi P, Olivieri M (1993) An analysis of dynamic scheduling techniques for symbolic applications. In: Proceedings of the 26th annual international symposium on microarchitecture, pp 185–191. https://doi.org/10.1109/MICRO.1993.282754
6. De Gloria A, Faraboschi P, Olivieri M (1992) A non-deterministic scheduler for a software pipelining compiler. ACM SIGMICRO Newslett 23(1–2):41–44. https://doi.org/10.1145/144965.144995
7. De Gloria A, Faraboschi P, Olivieri M (1994) Design and characterization of a standard cell set for delay insensitive VLSI design. IEEE Trans Circ Syst II Analog Digit Signal Process 41(6):410–415. https://doi.org/10.1109/82.300201
8. De Gloria A, Olivieri M (1996) Statistical carry lookahead adders. IEEE Trans Comput 45(3):340–347. https://doi.org/10.1109/12.485572
9. Costa A, De Gloria A, Giudici F, Olivieri M (1997) Fuzzy logic microcontroller. IEEE Micro 17(1):66–74. https://doi.org/10.1109/40.566209

10. Costa A, De Gloria A, Olivieri M (1996) Hardware design of asynchronous fuzzy controllers. IEEE Trans Fuzzy Syst 4(3):328–338. https://doi.org/10.1109/91.531774
11. Costa A, De Gloria A, Faraboschi P, Pagni A, Rizzotto G (1995) Hardware solutions for fuzzy control. Proc IEEE 83(3):422–434. https://doi.org/10.1109/5.364488
12. De Gloria A, Olivieri M (1996) An asynchronous distributed architecture model for the Boltzmann machine control mechanism. IEEE Trans Neural Netw 7(6):1538–1541. https://doi.org/10.1109/72.548186
13. De Gloria A, Faraboschi P, Olivieri M (1993) Clustered Boltzmann machines: massively parallel architectures for constrained optimization problems. Parallel Comput 19(2):163–175. https://doi.org/10.1016/0167-8191(93)90046-N
14. Angeloni I, Bisio F, De Gloria A, Mori D, Capurro C, Magnani L (2012) A virtual museum for Flemish artworks. A digital reconstruction of Genoese collections. In: 2012 18th international conference on virtual systems and multimedia, pp 607–610. https://doi.org/10.1109/VSMM.2012.6365989
15. Mori D, Berta R, De Gloria A, Fiore V, Magnani L (2013) An easy to author dialogue management system for serious games. J Comput Cult Herit 6(2):1–15. https://doi.org/10.1145/2460376.2460381
16. Bellotti F, Berta R, De Gloria A, Margarone M (2002) User testing a hypermedia tour guide. IEEE Pervasive Comput 1(2):33–41. https://doi.org/10.1109/MPRV.2002.1012335
17. Bellotti F, Berta R, Gloria AD, Ferretti E, Margarone M (2003) VeGame: exploring art and history in Venice. Computer 36(09):48–55. https://doi.org/10.1109/MC.2003.1231194
18. Bellotti F et al (2023) Alessandro De Gloria: 1955–2023. Int J Serious Games 10(1), Art no 1. https://doi.org/10.17083/ijsg.v10i1.612
19. Gongsook P et al (2014) Designing a serious game as a diagnostic tool. In: Revised selected papers from the third international conference on games and learning alliance, vol 9221, in GALA 2014. Springer-Verlag, Berlin, Heidelberg, pp 63–72
20. Pescosolido L, Berta R, Scalise L, Revel GM, De Gloria A, Orlandi G (2016) An IoT-inspired cloud-based web service architecture for e-health applications. In: 2016 IEEE international smart cities conference (ISC2). IEEE, Trento, Italy, pp 1–4. https://doi.org/10.1109/ISC2.2016.7580759
21. De Gloria A et al (2019) 3D software simulator for primary care training. In: Proceedings of the 8th international workshop on innovative simulation for healthcare (IWISH 2019), CAL-TEK srl, pp 114–119. https://doi.org/10.46354/i3m.2019.iwish.019
22. Bellotti F et al (2020) Managing big data for addressing research questions in a collaborative project on automated driving impact assessment. Sensors 20(23), 23. https://doi.org/10.3390/s20236773
23. Berta R, Kobeissi A, Bellotti F, De Gloria A (2021) Atmosphere, an open source measurement-oriented data framework for IoT. IEEE Trans Ind Inform 17(3):1927–1936. https://doi.org/10.1109/TII.2020.2994414
24. Cossu M, Berta R, Capello A, De Gloria A, Lazzaroni L, Bellotti F (2023) Developing a toolchain for synthetic driving scenario datasets. In: Berta R, De Gloria A (eds) Applications in electronics pervading industry, environment and society, Lecture notes in electrical engineering. Springer Nature Switzerland, Cham, pp 222–228. https://doi.org/10.1007/978-3-031-30333-3_29
25. Bellotti F, Lazzaroni L, Capello A, Cossu M, De Gloria A, Berta R (2023) Explaining a deep reinforcement learning (DRL)-based automated driving agent in highway simulations. IEEE Access 11:28522–28550. https://doi.org/10.1109/ACCESS.2023.3259544

# System Architecture and Hardware Acceleration

# Heterogeneous Tightly-Coupled Dual Core Architecture Against Single Event Effects

Marcello Barbirotta[(✉)], Francesco Menichelli, Antonio Mastrandrea, Abdallah Cheikh, Marco Angioli, Saeid Jamili, and Mauro Olivieri

Department of Information Engineering, Electronics and Telecommunications (DIET), "La Sapienza" University of Rome, Via Eudossiana, 18, 00184 Rome, Italy
Marcello.Barbirotta@uniroma1.it

**Abstract.** Among all the fault tolerance (FT) techniques developed over the years, Dual-core lock-step techniques have emerged as an effective approach to enhance the fault tolerance capabilities of these systems. However, they cannot withstand Hard Errors occurring in the architecture, and they have certain drawbacks for checkpoints and restore methodologies necessary to save and restore the correct state of the core. This paper shows an execution paradigm within a new architectural approach to overcome the disadvantages related to the long checkpointing/restoring procedures that affect the classical lock-step architectures, also improving the hard, destructive faults resilience, leveraging the advantages of the already published dynamic-TMR technique inside the Klessydra-dfT03 architecture.

**Keywords:** Fault tolerance · Dual core lock-step · Fault resilience

## 1 Introduction

The probability of faults in digital electronic devices has increased with technology scaling, voltage margin reduction and statistical process variations [1]. In the case of radiation-induced faults caused by high energy particles present in space, there can be several errors caused by Single Event Effects (SEE) [2] that can be divided into Soft Errors (recoverable) and Hard Errors (non-recoverable), which in Digital architectures are such as Single-Event Latch-up (SEL), Single-Event Burnout (SEB) or Single-Event Gate Rupture (SEGR). All of them can lead to permanent damage if not properly detected and solved with power cycling [3]. Comparing the probabilities of Hard and Soft Errors can be challenging, as different factors and circumstances influence them. Soft Errors are generally more likely than Hard Errors, especially in systems operating in environments with

higher radiation levels or electrical interference. Regardless of the application, having a system able to fix Soft Errors quickly without going through the long checkpoint/restore routines or being able to solve Hard Errors while maintaining the program instruction flow is the bottleneck of all lock-step systems.

## 2   Related Works

Lock-stepped architectures involve two identical processor cores running in parallel, executing the same instructions simultaneously, and comparing the results at predefined checkpoints. The authors in [4] build a lock-step DMR system with rollback for error recovery, periodically issuing an interrupt request for checkpoints. In [5], the dual-core lock-step architecture with two CPUs in an ARM Cortex-A9 processor is combined with a FreeRTOS to handle the checkpoint operations and rollback. The classic Dual-Core lock-step configuration is present in ARM Cortex-M7 processors [6], as well as Cortex-R where the execution between the two cores is cycle by cycle compared and the recovery mechanism is done resetting the cores or returning to some previous checkpoints. Authors in [7] provide the Triple Core Lockstep Architecture with three identical Cortex-R5 lock-step cores handled by a central Assist Unit that vote the generated outputs. When there is a mismatch, some assembly routines are launched to save the state inside a specific ECC-protected Stack composed by 113 registers, which are restored after a global reset between all the three architectures in a total of 2351 clock cycles (save + restore time). Authors in [8] propose an FPGA methodology creating a heterogeneous architecture through Dual-Core Lock-step (DCLS) technique at ISA-level. Each core receives the same input, and the system can stop, restart or restore from a checkpoint. Another heterogeneous architecture is described in [9], where two ARM cores and a MicroBlaze TMR core simultaneously runs in a TCLS fashion with verification points in the code used to write the execution state inside a BRAM memory and compare the registers. All the dual-core lock-step architecture techniques in the literature suffer the overhead due to the implementation of checkpoints and restore techniques, involving the periodic interruption of the running program to activate the processor state-saving routine.

## 3   Architecture

The multi-core architecture used as a basis for our work is made of two single-core processors with a shared TMR Register file. The two cores are based on an open-source RISC-V softcore family, namely Klessydra-T, which interleaves two or more hardware threads in a round-robin fashion on a four-stage in-order pipeline [10,11]. The first *Head* core on the right in Figure 1 is the Klessydra-dfT03 microarchitecture, built on the Dynamic-TMR principle described in [12], where two threads work in a DMR way, enabling the third one only in case of discrepancy. With a shared IMT pipeline structure, this core can only react against Soft Errors. Therefore, it needs another decoupled core architecture to detect and

react to Hard Errors, affecting the common pipeline units, which could, other-wise, create multiple consecutive failures. The *Tail* core on the left in Fig. 1 is the simplest version of the Klessydra family cores in a single-thread implementation called Klessydra-S0. Looking at Fig. 1, the main characteristics of this multi-core architecture are the Register-file sharing, able to implement the automatic processor state saving, improving and lightening the checkpointing/restoring phases, and the Error reporting mechanism, able to give fault-hardening features against destructive faults.



**Fig. 1.** Dynamic heterogeneous tightly-coupled dual core klessydra-dHTCDC microarchitecture. *Head* and *Tail* core on the right and left respectively. blue arrows: normal mode; black arrows: restore mode

The core can work in three modes:

– **Normal execution mode**: The *Head* core works in interleaved DMR mode, with two threads active (Red and Orange in Fig. 1), which execute the same instructions comparing and saving the Program Counter values inside the third Program Counter register (Green in Fig. 1) and the results coming from the critical units in a restore block mechanism (Black arrows in Fig. 1). During this phase, the *Tail* core is clock-gated, remaining at the reset state to save power.

– **Restore mode**: When the comparing logic gives a negative result due to a fault, the core enters the recovery mode activating the third green thread. Since a fault is always detected before the Register File is updated with a wrong result using the faulted instruction [12], the new green thread can Fetch, Decode and Execute the previously successfully executed instruction having the previous Core state saved inside the shared Register File. During

this phase, the *Tail* core is prepared to exit the clock-gating state, enter the reset state, and delete the possible accumulated errors.

– **End of Restore mode**: Once the *Head* core finishes the previous corrected instruction, the produced results wait inside the Error Reporting blocks, and the *Tail* core start Fetching Decoding and Executing the same instruction already executed by the green *Head* core thread, taking the same PC value and reading from the shared redundant Register-file. Once the result from the *Tail* core reaches the Error Reporting blocks from WB logic or LS logic, it is compared with the previous result by the *Head* core:

  • If they are equal (a common event), the received error was a Soft Error, the produced value is committed to the Register file, and the *Tail* core returns to the clock-gating state re-establishing the normal execution mode.
  • If they are different, the *Head* core got a Hard Error (not a common event). In that case, the *Head* core can be reset to possibly eliminate the error. Nevertheless, if the error is permanent, the execution can continue unprotected with the *Tail* core, thanks to the shared protected Register file which maintains the correct state of the core.

## 4    Results and Comparisons

Thanks to this mode and depending on the current instruction, the correct execution can be restored in a few tens of clock cycles, in contrast to the normal lock-step architecture, as shown in [7], in which the restoring procedure can consume as many as 1180 clock cycles. The presented architecture also does not need any checkpoint software procedures, which can be particularly time and hardware-consuming, especially if executed periodically [13] as happens in [7] with 1171 cycles. Klessydra-dHTCDC can also automatically preserve the correct state inside the shared protected Register file without adding additional main memory locations to save the data produced by the redundant instructions, never interrupting the program execution. It can also present better power consumption results than a normal lock-step architecture due to the clock-gating of the *Tail* core. At the same time, as observed from Table 1, the hardware occupation could be slightly greater due to the multi-thread environment if compared to a classical lock-step architecture obtained with two Klessydra-S0 cores. Unlike normal dual-core lock-step architectures, the proposed scheme can detect and react against Hard Errors, with efficiency in fault detection almost comparable to TCLS architectures, having less power consumption reached without triplicating the architecture. The various fault scenarios in Table 2 demonstrate that all Soft Errors can be easily detected. However, one challenging Hard Error remains undetectable: when a fault strikes the *Head* core pipeline during the initial thread execution (depicted as Orange in Fig. 1). Due to the fixed IMT order in Klessydra cores, Thread 2 will always be the first to fetch further instructions. Consequently, any hard fault in the shared pipeline units occurring during its execution will also affect the following threads, creating two wrong results

bypassing the voting systems and the activation of the *Tail* core, resulting in an error that cannot be identified.

**Table 1.** Microprocessor synthesis results on a Xilinx Kintex-7 FPGA for the classical dual Core lock-step Klessydra-S0 and the dynamic heterogeneous tightly-coupled dual core Klessydra-dHTCDC.

|  | Dual Core lock-step | Klessydra-dHTCDC |
|---|---|---|
| *LUTs* | 7656 | 10864 |
| *FFs* | 4030 | 7002 |

**Table 2.** Fault scenarios and their outcome with shared or no-shared pipeline units referring to the *Head* core. All soft errors affecting Thread 2 or Thread 1 can be easily detected for shared or no-shared pipeline units. In the hard error case, the failure appears when the hard error hits the shared pipeline units during the Thread 2 execution.

| | *Head* core | | | | | | *Tail* core | | | |
| | Thread 2 (Orange) | | Thread 1 (Red) | | Thread 0 (Green) | | Thread (Blue) | | Results | |
| Pipeline units | Shared | No-shared | Shared | No-shared | Shared | No-shared | Shared | No-shared | Shared | No-shared |
|---|---|---|---|---|---|---|---|---|---|---|
| *Soft error* | × | × | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | *Correct* | *Correct* |
| *Soft error* | ✓ | ✓ | × | × | ✓ | ✓ | ✓ | ✓ | *Correct* | *Correct* |
| *No error* | ✓ | ✓ | ✓ | ✓ | – | – | – | – | *Correct* | *Correct* |
| *Hard error* | ✓ | ✓ | × | × | × | ✓ | ✓ | ✓ | *Correct* | *Correct* |
| *Hard error* | × | × | × | ✓ | × | ✓ | – | ✓ | *Error* | *Correct* |

As demonstrated by [14], Ion strikes are the major source of Hard Errors in SRAMs technology. The hard-error probabilities caused by single Ion Strike impacts with 1MeV of energy on GEO orbits can vary from $3 \times 10^{-5}$ Single Hard Error (SHE) per device per day, own to $3 \times 10^{-7}$ SHE for earth-orbiting satellites. Considering all the mentioned advantages related to the proposed design, without assuming the incremental Hard Error rate due to aging effects [15], the above hard-error rate is further reduced by the fact that only one of the two threads can lead to failures in case of Hard Errors. The error rate is also reduced if we consider the limited percentage of area occupied by common pipeline units within the *Head* core, leading to an estimated reduction of an order of magnitude. Overall, the case of a permanent event on the first thread results extremely unlikely, even considering the worst application scenario of $3 \times 10^{-5}$ SHE per dvice per day. It is worth underlying that most Hard Errors inside Digital architectures can be solved by power cycling the entire architecture [3], which in our case, is made possible in the *Head* core during the End of the Restore mode phase when the *Tail* core carries out the program execution. Obviously, this option is available only if the two cores operate with different power lines.

## 5   Conclusions

Thanks to the dynamic TMR principle, this work demonstrates how it is possible to obtain a high performance multi-core architecture compared to the classical lock-step schemes in the literature, overcoming the disadvantages due to the heavy checkpointing and restore procedures and allowing a higher level in terms of fault protection against hard errors. Despite the disadvantages related to the case of a Hard Error occurring during the execution of Thread 2 inside common pipeline units, the architecture is usable for fault-tolerance applications, considering the probability of this event is very rare.

## References

1. Barbirotta M, Mastrandrea A, Cheikh A, Menichelli F, Olivieri M (2022) Improving set fault resilience by exploiting buffered DMR microarchitecture. In: Annual meeting of the Italian electronics society. Springer, pp 233–238
2. Barbirotta M, Mastrandrea A, Menichelli F, Vigli F, Blasi L, Cheikh A, Sordillo S, Di Gennaro F, Olivieri M (2020) Fault resilience analysis of a risc-v microprocessor design through a dedicated UVM environment. In: 2020 IEEE international symposium on defect and fault tolerance in VLSI and nanotechnology systems (DFT). IEEE, pp 1–6
3. Secretariat E (2008) Space product assurance
4. Violante M, Meinhardt C, Reis R, Reorda MS (2011) A low-cost solution for deploying processor cores in harsh environments. IEEE Trans Ind Electron 58(7):2617–2626
5. de Oliveira ÁB, Rodrigues GS, Kastensmidt FL (2017) Analyzing lockstep dual-core arm cortex-a9 soft error mitigation in Freertos applications. In: Proceedings of the 30th symposium on integrated circuits and systems design: chip on the sands, pp 84–89
6. Yiu J (2015) Design of SOC for high reliability systems with embedded processors. In: Embedded World Conference
7. Iturbe X, Venu B, Ozer E, Poupat JL, Gimenez G, Zurek HU (2019) The arm triple core lock-step (TCLS) processor. ACM Trans Comput Syst (TOCS) 36(3):1–30
8. Marques I, Rodrigues C, Tavares A, Pinto S, Gomes T (2021) Lock-V: a heterogeneous fault tolerance architecture based on Arm and RISC-V. Microelectron Reliabil 120:114120
9. Kasap S, Wächter EW, Zhai X, Ehsan S, McDonald-Maier KD (2021) Novel lockstep-based fault mitigation approach for SOCs with roll-back and roll-forward recovery. Microelectron Reliabil 124:114297
10. Cheikh A, Sordillo S, Mastrandrea A, Menichelli F, Scotti G, Olivieri M (2021) Klessydra-t: designing vector coprocessors for multithreaded edge-computing cores. IEEE Micro 41(2):64–71
11. Barbirotta M, Cheikh A, Mastrandrea A, Menichelli F, Olivieri M (2022) Design and evaluation of buffered triple modular redundancy in interleaved-multithreading processors. IEEE Access 10:126074–126088
12. Barbirotta M, Cheikh A, Mastrandrea A, Menichelli F, Ottavi M, Olivieri M (2023) Evaluation of dynamic triple modular redundancy in an interleaved-multithreading RISC-V core. J Low Power Electron Appl 13(1). https://www.mdpi.com/2079-9268/13/1/2

13. Liu K, Li Y, Ouyang L (2022) A survey of fault tolerance hardware architecture. In: 2021 international conference on advanced computing and endogenous security. IEEE, pp 01–06
14. Poivey C, Carriere T, Beaucour J, Oldham T (1994) Characterization of single hard errors (SHE) in 1 m-bit SRAMS from single ion. IEEE Trans Nucl Sci 41(6):2235–2239
15. Amrouch H, van Santen VM, Ebi T, Wenzel V, Henkel J (2014) Towards interdependencies of aging mechanisms. In: 2014 IEEE/ACM international conference on computer-aided design (ICCAD). IEEE, pp 478–485

# Wire Bonding: Limitations and Opportunities for High-Speed Serial Communications

Gabriele Ciarpi[(✉)], Marco Mestice, Daniele Rossi, and Sergio Saponara

Department of Information Engineering, University of Pisa, 56122 Pisa, Italy
`gabriele.ciarpi@unipi.it`

**Abstract.** This paper explores the impact of wire bonding solutions on high-speed serial links, focusing on the co-design between on-chip and off-chip parameters to achieve a boost in system bandwidth. Despite traditional limitations, wire bonding can outperform flip-chip techniques when properly designed. This study investigates various cases with different in-band ripples and describes how to optimize the on-chip and bonding parameters accordingly. Additionally, the adoption of shunt peaking techniques is examined, which allows us to achieve a bandwidth boost of up to 3.82 times. While on-chip parameters can be well controlled, the limited mechanical precision of wire bonding may lead to suboptimal wire lengths. An analysis of the variation in bandwidth boost and ripple is presented, indicating that wire length variations below 10% result in minimal impact on boosting effects and ripple. In conclusion, wire bonding offers advantages for high-speed serial links, and co-design optimizes system performance. Properly designed, wire bonding is an attractive choice for high-speed applications, outperforming flip-chip techniques in bandwidth enhancement.

**Keywords:** High-Speed Links · Serial Communications · Wire Bonding

## 1 Introduction

Thanks to continuous improvement over the past decades in semiconductor technologies, on-chip links are today able to reach several tens of GHz [1–3]. However, the advancement of packaging has lagged behind that of semiconductor technology, and the interconnects between chips or chips and boards are presenting significant barriers to the creation of broadband communication links [4]. The fundamental cause is that scaling criteria for on-chip and off-chip structures differ.

To mitigate the reduced frequency performance of off-chip structures, parallel data lines are employed. However, the large form factor and high pin count of parallel connections result in high system costs. Instead, an accurate design of high-speed differential serial links allows for increased system bandwidth with a low pin count [5].

Currently, the primary bonding techniques for both chip-to-chip and chip-to-board connections are wire bonding and flip-chip techniques. The former is the most commonly used due to its adaptability and lower cost. The latter solution allows for a greater pin

count. In this case, the connections are placed not only on the chip's border but also in the internal region, enabling the management of a high number of parallel communication links.

Regarding the bandwidth limitation of off-chip structures for radio-frequency (RF) or millimeter-wave (MMW) applications, which operate in a narrow band around the carrier frequency, the limitation imposed by off-chip structures can be reduced by making them resonate at the carrier frequency [4]. However, for high-speed digital communications requiring very broadband characteristics, the solution adopted for RF and MMW applications cannot be adopted. Therefore, an accurate co-design of the on-chip and off-chip structures is required. This paper aims to explore the limitations and advantages of using cost-effective wire bonding solutions for high-speed digital links.

In particular, Sect. 2 presents the modelling of a generic high-speed serial interface. The system optimization and its results are described in Sect. 3 for several cases. The limitations of the wire bonding solution are discussed in Sect. 4. Conclusions are drawn in Sect. 5.

## 2   Wire Bonding System Modeling

The first step in analyzing the limits of the wire bonding solution for high-speed links is to model a generic output driver and bonding wire connections. The most commonly used output architectures for high-speed serial drivers are the current-mode logic architecture (CML), low voltage differential signaling (LVDS), or emitter-coupled logic (ECL) [6–10]. These drivers consist of a differential couple with a tail current, along with pull-up resistors (or transistors biased in resistive regions). Hence, in the simplified single-ended model shown in Fig. 1, these output drivers are generically represented as a current generator $I_{in}$ plus a parallel resistor $R$. To achieve broadband output matching, in this work, the value of the output resistor is set to 50 $\Omega$.

One of the main limiting factors that reduce the bandwidth of broadband serial links is the presence of parasitic elements, such as on-chip pads and electrostatic discharge (ESD) protections. In the model depicted in Fig. 1, their effects are accounted for by the capacitance $C_1$. Additionally, $C_1$ considers the non-negligible output capacitance of the driver. The inductor $L_b$ represents the primary parasitic element of a wire bond. For the simplified model, the resistivity part of the bonding wire is neglected. The capacitor $C_2$ serves as a simplified model of the second chip or board pad.

Considering that the driver is supposed to drive a transmission line, a resistive termination R of 50 $\Omega$ is included in the model.

An additional inductor is placed in series with the pull-up resistance of the driver to implement a shunt peaking technique. This solution is widely adopted in the literature to further extend the bandwidth of high-speed serial links [11]. The typical value for the inductor that enables bandwidth extension while providing the maximally flat gain is $L = R^2 C / 0.4$. However, this value is derived for a shunt peaking technique with a pull-up resistance R and loaded with only a capacitor C. Instead, the effects of wire bonding should be considered because they could generate overdamping or underdamping in the frequency response of the systems, either magnifying or eliminating the widening bandwidth effect introduced by the shunt peaking solution.

**Fig. 1.** Schematic of the system model

## 3    System Optimization

In order to find the optimum values of the circuit parameters for maximum bandwidth extension, a new set of more significant parameters than the circuital one is defined and presented in Eq. (1). Here, $\omega_0$ represents the pulsation that corresponds to the bandwidth of a basic RC output driver, which serves as a reference. This value can be associated with the bandwidth of a driver in a flip-chip configuration (near zero-length bonding) without the shunt peaking technique. Parameters $Q_b^2$ and $Q_s^2$ are the square of the quality factor of an RLC network for the bonding and shunt inductors, respectively. Additionally, $k_C$ accounts for the capacitance distribution between $C_1$ and $C_2$.

Considering the number of variables that need to be optimized to maximize the system bandwidth, a numerical approach is adopted instead of an analytical one. The optimization method used combines the Quasi-Newton and Random algorithms [12]. Compared to the Brute-Force algorithm, the Quasi-Newton method allows us to find the minimum of the cost function in a few steps, while the Random algorithm generates a new initial guess to avoid local minimums.

$$\omega_0 = \frac{2}{R(C_1 + C_2)}; Q_b^2 = \frac{R^2(C_1 + C_2)}{4L_b}; Q_s^2 = \frac{R^2(C_1 + C_2)}{4L_s}; k_C = \frac{C_1}{(C_1 + C_2)} \quad (1)$$

Even though the bandwidth of the system can be enhanced thanks to the shunt peaking technique and the bonding wire, the flatness of the frequency response is degraded due to the generation of a ripple. For a given ripple amplitude, the bandwidth has a maximum for a particular combination of the parameters defined in (1). In Table 1, the optimal values of the system parameters obtained by the algorithm are reported for 0 dB, 1 dB, and 2 dB ripple, along with the corresponding bandwidth boosting (BWBS) level. The BWBS is defined as the ratio between the bandwidth of the optimized solution and the reference bandwidth associated with a flip-chip solution without shunt peaking techniques ($\omega_0$ in (1)).

As shown in Table 1, different combinations of parameters are required to achieve the maximum bandwidth boost for a specific ripple level. When the bonding wires are properly designed, taking into account the other system parameters defined in (1), they can introduce a boost in the system bandwidth of up to 3.82.

The bandwidth boosting effect of the bonding wire is also illustrated in Fig. 2, displaying the frequency behavior of the system with the optimized parameter values

**Table 1.** Optimized system values for different cases

| Bond | Shunt | Ripple (dB) | $k_C$ | $Q_s^2$ | $Q_b^2$ | BWBS |
|------|-------|-------------|-------|---------|---------|------|
|      |       | 0           |       |         |         | 1.00 |
|      | x     | 0           |       | 0.64    |         | 1.52 |
|      | x     | 1           |       | 0.43    |         | 2.21 |
|      | x     | 2           |       | 0.38    |         | 2.89 |
| x    |       | 0           | 0.5   |         | 0.51    | 1.47 |
| x    |       | 1           | 0.5   |         | 1.02    | 1.48 |
| x    |       | 2           | 0.5   |         | 1.73    | 1.48 |
| x    | x     | 0           | 0.75  | 0.91    | 1.08    | 2.47 |
| x    | x     | 1           | 0.69  | 0.54    | 1.97    | 3.31 |
| x    | x     | 2           | 0.64  | 0.49    | 2.71    | 3.82 |

for the case with a 1 dB ripple as a function of the frequency normalized for the reference bandwidth (bandwidth for flip-chip solution without shunt peaking). Since the values of the optimized parameters do not depend on $\omega_0$ they are applicable to any $\omega_0$ (which has been verified with simulations), leading to an optimized design for every high-speed link operating at different bit rates.



**Fig. 2.** Module of the transfer functions of the system for different optimized cases reported in Table 1 for 1 dB ripple condition as a function of the normalized frequency.

## 4  Wire Bonding Limitation

Although the optimized parameters found in the previous section are applicable to any frequency, there are some limitations on the boosting effect of wire bonding techniques at high frequencies due to physical and technological constraints.

The system bandwidth is primarily determined by $\omega_0$, which is limited by the parasitic capacitors of on-chip pads and ESD protections. However, similar pad and ESD structures are used for other bonding techniques, such as flip-chip.

Regarding the boosting effect of $Q_s^2$, a wide range of inductor values from a few pH to a few nH can be integrated on a chip. Since the inductors are series-connected to pull-up resistances in the targeted application using CML, LVDS, or ECL drivers, $Q_s^2$ depends mostly on the pull-up resistance and the inductor value, rather than the quality factor of the inductor itself. Therefore, an inductor high-quality factor is not necessary.

Considering $Q_b^2$, the main limitation in the boosting level is related to the bonding length. The bonding wire length is associated with its parasitic inductance following the relation of 1 nH/mm. While long bonding wires might impact packaging costs, they are generally not a technological problem. On the other hand, the realization of short bonding wires can be challenging due to the minimum distance of the pads. From the literature, it is found that with proper chips placement and alignment, a minimum bonding wire length of about 200 µm can be achieved [13].

However, as explained in the previous sections, minimum bonding length does not necessarily lead to the maximum boost of the system bandwidth. Instead, an optimized length should be used. For example, if a bandwidth of 5 GHz is required and no shunt peaking technique is adopted to save chip area, the wire bonding length that allows reaching the target frequency while reducing the constraints of on-chip pads is equal to 1.15 mm (calculated with optimized data in Table 1).

The minimum bonding wire length can be used to determine the frequency limit of the boosting techniques. For an output driver with a shunt peaking technique and a 1 dB ripple in bandwidth, the wire bonding technique can provide its maximum boosting effect up to $\omega_0$ approximately equal to 63.3 Grad/s. This indicates that the system bandwidth, inclusive of the boosting effect, could be close to 33.4 GHz. Considering the rule of thumb that the system bandwidth should be at least 2/3 of the bit rate to have a good eye diagram [14], the data rate enabled in this system configuration is about 50 Gb/s.

Although bonding wire technology is improving its ability to create precise wire loop structures, uncertainties can still occur due to the mechanical alignment of the chips.

Figure 3 illustrates the variation of the BWBS and the ripple as a function of the relative variation of the wire bonding length over the optimal case (driver with shunt peaking technique and 1 dB ripple). As can be observed, if the bonding wires are longer than the optimal case, the BWBS increases, but so does the ripple. An increase of over 55% in the bonding wire length results in a 3 dB ripple. Consequently, considering the 3 dB bandwidth definition, the system bandwidth is drastically reduced. Conversely, if the wire bonding length is lower than the optimized case, the BWBS is reduced, and the ripple is increased. If the bonding wire length is lower than $-70\%$ of the optimized length, the BWBS continues to decrease, and the ripple increases. In this scenario, the bonding wire length tends to zero, losing its boosting effect on the bandwidth.

## 5  Conclusions

The proposed work aims to investigate the impact of wire bonding solutions on high-speed serial links. Conventionally, bonding wires are known to impose limitations on broadband system bandwidth. However, the paper demonstrates that through a proper

**Fig. 3.** Bandwidth boost and ripple value as a function of the relative variation of the bonding wire length over the optimized case.

co-design between on-chip and off-chip parameters, it is possible to achieve a boost in the system bandwidth. As a result, for serial link applications, using wire bonding techniques instead of flip-chip can lead to performance improvement. The paper examines several cases with different in-band ripples and optimizes them accordingly. Additionally, the adoption of shunt peaking techniques is analyzed for systems requiring an additional bandwidth boost of up to 3.82 times.

While on-chip parameters such as pull-up resistors, shunt peaking inductors, and parasitic capacitance of pads and ESD protections can be well controlled, the limited mechanical precision of wire bonding techniques can result in suboptimal wire lengths. Therefore, an analysis of the variation in bandwidth boost and ripple is also presented, highlighting that for wire length variations lower than 10%, the boosting effects change at a maximum of about 5%, and the ripple increases at a maximum of 0.38 dB. Since both negative and positive variations in the bonding wire length increase the ripple, a slightly longer wire length is preferable for increasing BWBS, as opposed to the case with a shorter wire.

# References

1. Sun L, et al (2022) SI/PI co-simulation analysis of high-speed I/O link. In: 2022 ICEPT, Dalian, China, pp 1–5
2. Monda D et al (2020) Analysis and comparison of rad-hard ring and LC-tank controlled oscillators in 65 nm for SpaceFibre applications. Sensors 20(16):4612
3. Monda D, et al (2021) Design and verification of a 6.25 GHz LC-tank VCO integrated in 65 nm CMOS technology operating up to 1 Grad TID. IEEE TNS 68(10):2524–2532
4. Božani´c M, et al (2019) Systems-level packaging for millimeter-wave transceivers. Springer, Vienna, Austria

5. Chang PH, et al (2021) Signal and power integrity analysis of a 0.38 pJ/bit 12.8 Gb/s parallel interface for die-to-die link applications. In: 2021 ECTC, CA, USA, pp 1264–1269
6. Ciarpi G et al (2022) Design and characterization of 10 Gb/s and 1 Grad TID-tolerant optical modulator driver. In: IEEE TCAS-I, vol 69, no 8, pp 3177–3189
7. Ciarpi G, et al (2019) Design, implementation, and experimental verification of 5 Gbps, 800 Mrad TID and SEU-tolerant optical modulators drivers. In: IEEE TCAS-I, vol 67, no 3, pp 829–838
8. Palla F et al (2019) Design of a high radiation-hard driver for Mach-Zehnder Modulators based high-speed links for hadron collider applications. Nucl Instrum Meth Phys Res, Sect A 936:303–304
9. Mestice M, et al (2020) Analysis and design of integrated blocks for a 6.25 GHz spacefibre PLL. Sensors 20(14):4013
10. Ciarpi G, et al (2023) A 10 Gb/s line driver in 65 nm CMOS technology for radiation-pervaded and high-temperature applications. IEEE Access
11. Kim J (2009) Design optimization of on-chip inductive peaking structures for 0.13-$\mu$m CMOS 40-Gb/s transmitter circuits. IEEE TCAS-I 56(12):2544–2555
12. Mokhtari A, et al (2020) Stochastic Quasi-Newton methods. Proc IEEE 108(11):1906–1922
13. Tao Z, et al (2017) High-gain low-cost broadband 60 GHz differential integrated patch array antennas with wire-bonding packaging and on-board compensation network. IET Microwaves Antennas Propag 11.7:971–975
14. Voinigescu S (2013) High-frequency integrated circuits. Cambridge University Press

# LOKI Low-Latency Open-Source Kyber-Accelerator IPs

Alessandra Dolmeta[(✉)], Mattia Mirigaldi, Maurizio Martina, and Guido Masera

Politecnico di Torino, Torino, Italy
{Alessandra.Dolmeta,Mattia.Mirigaldi,Maurizio.Martina, Guido.Masera}@polito.it

**Abstract.** CRYSTALS-Kyber is a lattice-based key encapsulation mechanism (KEM) recognized as one of the finalist algorithms in NIST's post-quantum cryptography (PQC) standardization process. Polynomial multiplications and hash functions, as essential operations in lattice-based PQC schemes, pose a significant time consumption challenge with respect to nowadays cryptographic protocols. This work addresses these computational efforts by incorporating LOKI, an accelerator, into a RISC-V microcontroller. By leveraging the accelerator, the performance can be enhanced, contributing to the overall efficiency of Kyber in the fundamental tasks of key generation, encryption, and decryption operations. Through empirical evaluations and benchmarking, the effectiveness and practicality of the proposed hardware architectures are demonstrated, highlighting their potential to advance the field of post-quantum cryptography.

**Keywords:** Hardware design · Accelerator · Security · Post-quantum cryptography · CRYSTALS-Kyber · Keccak · NTT · RISC-V

## 1 Introduction

Since its discovery in 1994, the Shor algorithm has garnered significant interest in the field of cryptography. In actual fact, it has the potential to break widely-used public-key encryption schemes. Hence, its realization has raised significant concerns about the security of modern cryptographic systems, leading to the urgent need for **post-quantum cryptography** (PQC) in today's digital landscape. PQC refers to cryptographic schemes specifically designed to withstand attacks from quantum computers. The field of PQC encompasses multiple mathematical approaches and primitives, including lattice-based cryptography, code-based cryptography, hash-based signatures, and more. The current goal of researchers and organizations worldwide is to identify and standardize quantum-resistant cryptographic schemes that can seamlessly replace the vulnerable algorithms currently in use. Hence, in 2012, the National Institute of Standards

and Technology (NIST) launched a public evaluation process to standardize quantum-resistant public key algorithms. After three rounds of solicitations, in 2022, CRYSTALS-Kyber has been the first algorithm to be selected for standardization.[1] However, these algorithms have a non-negligible computational burden if compared to their pre-quantum counterparts. When running them on a microcontroller, performance is not optimal, resulting in slow and energy-consuming execution. To overcome these challenges, **accelerators** are often used. In fact, by offloading tasks to dedicated hardware, microcontrollers can handle complex operations effectively, improving overall performance and energy efficiency. Another important feature of an accelerator is versatility, namely the property of being easily integrable and testable within different microcontrollers. Being open-source, **RISC-V** has been selected as the target. In fact, the main aim of RISC-V is to provide a free and open ISA suitable for processor designs, without imposing a particular implementation technology.

## 2    Preliminaries

CRYSTALS-Kyber is a **lattice-based cryptosystem** based on the hardness of solving the Module Learning-with-Error (Module-LWE) problem. It implements a KEM, a probabilistic algorithm that produces a random symmetric key and encrypts it. Compared to traditional protocols, KEM adds a shared secret to the process, encapsulating and decapsulating the key using the public and the secret keys . In principle, it is a triplet of algorithms: a **key generation algorithm**, an **encapsulation algorithm**, and a **decapsulation algorithm**. Kyber includes **three parameter sets**, corresponding to three security levels of NIST: Kyber512 (I-level), Kyber768 (III-level), and Kyber1024 (V-level). Independently of the levels of security, Kyber always uses the same polynomial arithmetic for all variants. In particular, it works on rings of integer polynomials modulo prime $q$, denoted as $Z_q[X]$. Polynomials modulo both $q$ and $X^n + 1$ compose the ring $R_q = Z_q[X]/(X^n + 1)$. The only difference is given by the number of polynomial vectors required. A more detailed description of the algorithm can be found in the official documentation [1]. In fact, the aim of this study is to speed up the most onerous part of the algorithm. Indeed, we are focusing our attention on Numeric Theoretic Transform (NTT) and hash functions.

**Numeric Theoretic Transform**. NTT is a peculiar case of DFT which is conducted in the finite field $Z_q[X]$. It is generally exploited to speed up polynomial multiplication and reduce the complexity of multiplying two n-terms polynomials. NTT is commonly executed through a butterfly unit in Gentleman-Sande (GS) or Cooley-Tukey (CT) configuration. In our case, there is only one unified butterfly. Twiddle Factors are precomputed and stored in tables. Moreover, butterfly computations are executed modulo $q$ employing Barrett and Montgomery reductions [8].

**Keccak**. Kyber algorithm involves several cryptographic primitives from the Secure Hash Algorithm (**SHA 3**) standard. It lists four specific instances of

---

SHA-3 and two extendable output functions. All instances are based on a 1600-bit state, which is the parallelism of the most important part of these primitives: the Keccak permutation function. In each of its 24 rounds, it calls the *f-1600* function, which is characterized by the five consecutive steps: $\theta, \rho, \pi, \chi$, and $\iota$ [2].

## 3   Implementation

In this section, LOKI is presented. The proposed hardware architecture and its main blocks are explained, starting from the implementation and ending with their integration into the PULPissimo microcontroller (RISC-V microcontroller from the open-source PULP[2]). In this work, we concentrate more on the design of the blocks conceived for NTT, INTT, and PWM accelerators. Therefore, for an exhaustive description of the Keccak blocks, refer to [3]. For the implementation, the official c-code reference model provided by NIST was used, to ensure the execution of the algorithm to be compliant with the standards.[3]



**Fig. 1.** Architecture's block diagram

Figure 1 represents the block diagram of the architecture. A unified butterfly unit (BU) is used, designed to work both for NTT and INTT. Although the mathematical calculations involved in both algorithms are comparable, their primary distinctions lie in how they store and retrieve data and in the execution of the butterfly operation. The butterfly operation encompasses multiplying by a specific constant, called Twiddle Factor, followed by addition and subtraction within $Z_q[X]$. Twiddle Factors are powers of $\omega_n \in Z_q[X]$, which is the $n$-th root of unity. These values are precomputed and stored in the BROM memory shown in Fig. 1. Then, there are four dual-port BRAMs: the first two are

---

[2] https://github.com/pulp-platform.

[3] Reference code is taken from PQClean: https://github.com/PQClean/PQClean.git.

used to store the first input polynomial and output polynomial, while the other
two are used to store the second input polynomial [12]. Before any operation,
input polynomials are stored in the BRAMs using input multiplexers. Subse-
quently, the control unit initiates its operation based on start signals, and the
resulting polynomial is retrieved using output multiplexers after the operation.
As previously stated, LOKI can perform NTT, INTT, and PWM. Both NTT
and INTT can be divided into 7 stages, each composed of 128 butterfly oper-
ations. NTT and INTT make use of different computational structures: NTT
required CT butterfly, while INTT makes use of GS structure. The first follows
the decimation-in-time approach, i.e., $a + b \cdot \omega_n (\bmod\ q)$ and $a - b \cdot \omega_n (\bmod\ q)$,
while the latter follows the decimation-in-frequency approach, i.e., $a + b(\bmod\ q)$
and $(a-b) \cdot \omega_n(\bmod\ q)$. With respect to the reference code mentioned above, the
accelerator is able to execute the functions `poly_ntt` and `poly_invntt_tomont`
respectively. The architecture of the BU is shown in the left part of Fig. 2. There
are one integer multiplier, one modular adder, one modular subtractor, and two
reduction units (Montgomery and Barrett, outlined in the right part of Fig. 2).



**Fig. 2.** Block diagram of the butterfly unit

Our BU can also be adapted to carry out PWM operations. This operation
corresponds to `basemul` function in the reference model, where 5 multiplication
and 2 addition operations are executed. The single butterfly unit takes 906, 906,
and 649 clock cycles to execute NTT, INTT, and PMW operations.

**Integration**. The accelerator is driven in a memory-mapped fashion style
and integrated through a robust hardware/software interface. Its versatility
is achieved thanks to a *generic register interface*, which is simple, faster, and
equipped with protocol adapters from APB, AXI-Lite, and AXI. This makes it
easy to attach to most microcontrollers. We use the PULP RISC-V Toolchain[4]
with the optimization flag 'O3' to compile the code and increase the stack's mem-
ory size. The `randombytes` file is modified as done by [5], generating a pseudo-
random sequence of bytes. Considering the latency due to the interchange of
data on the bus, the overall architecture takes now 2856, 2856, and 3373 clock
cycles to execute NTT, INTT, and PMW operations.

---

[4] https://github.com/pulp-platform/pulp-riscv-gnu-toolchain.

# 4    Results and Comparison

In this section, we report our implementation results and their comparison with other prior work in the state-of-art.

First, we evaluate the efficiency of our design, we conduct several tests running Kyber with and without the accelerator. Speed-up factors obtained are **4,49/5,08/4,15** for the I-level, **4,35/4,63/4,09** for the III-level, and **4.39/4,81/4,21** for the V-level. Each value represents the speed-up factor respectively for key generation, encapsulation, and decapsulation.

Then, cycle counts and area results of our and related works are reported in Table 1. LOKI outperforms most of the prior works reported. In fact, it shows up to $2\times/3\times$ better performance for all the three levels of security examined with respect to all the software optimizations [4,7,11]. Both in [5] and

**Table 1.** Implementation results and comparison to prior works. [*] indicates occupation results referring to NTT/INTT/PWM hardware only.

| Ref. | Ver. | Device | LUT/REG/DSP/BRAM | KeyGen | Encaps. | Decaps. |
|------|------|--------|------------------|--------|---------|---------|
| [5] | I | RISC-V (Pulpino) | 2908/179/9/- | 150,106 | 193,076 | 204,843 |
|  | III |  |  | 273,370 | 325,888 | 340,418 |
|  | V |  |  | 349,673 | 405,477 | 424,682 |
| [6] | I | RISC-V (VexRiscv) | 1842/1634/5/34* | 710,00 | 971,000 | 870,00 |
|  | III |  |  | – | – | – |
|  | V |  |  | 2,203,000 | 2,619,000 | 2,429,000 |
| [7] | I | ARM Cortex-M4 | –/–/–/– | 455,191 | 586,334 | 543,500 |
|  | III |  |  | 864,008 | 1,032,540 | 969,867 |
|  | V |  |  | 1,404,695 | 1,605,707 | 1,525,805 |
| [4] | I | ARM Cortex-M4 | –/–/–/– | 514,291 | 652,69 | 621,245 |
|  | III |  |  | 976,757 | 1,146,556 | 1,094,849 |
|  | V |  |  | 1,575,052 | 1,779,848 | 1,709,348 |
| [8] | I | CVA6 | 178/0/5/0* | 419,597 | 438,280 | 100,796 |
|  | III |  |  | 694,504 | 731,597 | 130,348 |
|  | V |  |  | 1,090,458 | 1,116,462 | 159,639 |
| [11] | I | ARM Cortex-M4 | –/–/–/– | 499,000 | 634,000 | 597,000 |
|  | III |  |  | 947,000 | 1,113,000 | 1,059,000 |
|  | V |  |  | 1,525,000 | 1,732,000 | 1,653,000 |
| **OUR** | I | **RISC-V (Pulpissimo)** | **938/539/10/2.5*** | **245,130** | **282,272** | **344,802** |
|  | III |  |  | **407,204** | **492,221** | **552,033** |
|  | V |  |  | **630,158** | **703,284** | **796,591** |
| [10] | I | Artix-7 | 25674/3137/64/6 | 9,400 | 19,000 | 43,000 |
|  | III |  |  | 14,200 | 26,2000 | 59,100 |
|  | V |  |  | 18,500 | 33,700 | 77,500 |
| [9] | I | Artix-7 | 7412/4644/2/3 | 3,800 | 5,100 | 6,700 |
|  | III |  |  | 6,300 | 7,900 | 10,000 |
|  | V |  |  | 9,400 | 11,300 | 13,900 |

[8], they proposed a dedicated Post-Quantum Arithmetic Logic Unit in the pipeline of the RISC-V processor. With respect to [5], ours are $0.6/0.68/0.59\times$, $0.67/0.66/0.62\times$, and $0.55/0.58/0.55\times$ slower for Kyber512, Kyber768, and Kyber1024 respectively. However, LOKI was not included in the pipeline. Therefore, in contrast to [5,8], the higher overhead due to data interchange between core and co-processor must be taken into account. This leads to worse performance than in [5], but better accelerator flexibility. This is precisely the major strength of LOKI, which, having a generic interface, can easily be connected to any microcontroller. Compared to [8] instead, performance is still better, although the results shown in Table 1 refer to encryption and decryption, and not encapsulation and decapsulation.

References [9,10] have been reported for the sake of completeness, despite they rely on a very different approach. In the case of [10], performances are clearly better, but the area occupied is 27 times larger. Similar reasoning can be made with [9], which achieves better performance by implementing the entire algorithm in hardware, but with significantly higher design effort.

## 5  Conclusion

This study establishes an ideal starting point for the optimization of the PQC algorithm in today's microcontroller. The hardware design was done by following as closely as possible the code provided by CRYSTALS in order to obtain a version fully compliant with the software implementation. With respect to the other hardware-only solutions, LOKI provides a considerable increase in performance. In future works, we would like to study the weaknesses of the algorithm, while analyzing it on the hardware level through electromagnetic and instantaneous power consumption analyses. In this way, we will be able to understand which points are vulnerable to side-channel attacks and modify both the code and the hardware accordingly.

## References

1. Soni D et al (2019) Power, Area, Speed, and Security (PASS) Trade-offs of NIST PQC signature candidates using a C to ASIC design flow. In: 2019 IEEE 37th international conference on computer design (ICCD). Abu Dhabi, United Arab Emirates, pp 337–340. https://doi.org/10.1109/ICCD46524.2019.00054
2. Bertoni G et al (2015) FIPS PUB 202—SHA-3 standard: permutation-based hash and extendable-output functions. https://keccak.team/hardware.html
3. Dolmeta A et al (2023) Implementation and integration of Keccak accelerator on RISC-V for CRYSTALS Kyber. In: 20th ACM international conference on computing frontiers, CF23-OSHW. https://doi.org/10.1145/3587135.3591432
4. Matthias J et al (2019) pqm4: Testing and Benchmarking NIST PQC on ARM Cortex-M4, Cryptology ePrint Archive, Paper 2019/844. https://ia.cr/2019/844
5. Fritzmann T et al (2020) RISQ-V: Tightly Coupled RISC-V accelerators for post-quantum cryptography. IACR TCHES 2020(4):239–280. https://doi.org/10.13154/tches.v2020.i4.239-280

6. Alkim E et al (2020) ISA extensions for finite field arithmetic: accelerating kyber and newhope on RISC-V. IACR TCHES 219–242. https://doi.org/10.46586/tches.v2020.i3.219-242

7. Alkim E et al (2020) Cortex-M4 optimizations for R,M LWE schemes. In: IACR TCHES (3):336–357. https://doi.org/10.13154/tches.v2020.i3.336-357

8. Nannipieri P et al (2021) A RISC-V post quantum cryptography instruction set extension for number theoretic transform to speed-up crystals algorithms. IEEE Access 9:150798–150808. https://doi.org/10.1109/ACCESS.2021.3126208

9. Xing Y et al (2021) A compact hardware implementation of CCA-secure key exchange mechanism CRYSTALS-KYBER on FPGA. IACR TCHES (2):328–356. https://doi.org/10.46586/tches.v2021.i2.328-356

10. Zhao Y et al (2022) A high-performance domain-specific processor with matrix extension of RISC-V for module-LWE applications. IEEE Trans Circuits Syst I: Regular Papers 69(7):2871–2884. https://doi.org/10.1109/TCSI.2022.3162593

11. Leon Botros L et al (2019) Memory-efficient high-speed implementation of kyber on cortex-M4. Cryptology ePrint Archive, Paper 2019/489 (2019). https://eprint.iacr.org/2019/489

12. Yaman F et al (2021) A hardware accelerator for polynomial multiplication operation of CRYSTALS-KYBER PQC scheme. In: 2021 design, automation and test in Europe conference and exhibition (DATE). Grenoble, France, pp 1020–1025. https://doi.org/10.23919/DATE51398.2021.9474139

# A Universal Hardware Emulator for Verification IPs on FPGA: A Novel and Low-Cost Approach

Saeid Jamili[✉], Antonio Mastrandrea, Abdallah Cheikh, Marcello Barbirotta,
Francesco Menichelli, Marco Angioli, and Mauro Olivieri

Sapienza Universitá di Roma, Via Eudossiana, 18, 00184 Roma, RM, Italy
{saeid.jamili,antonio.mastrandrea,abdallah.cheikh,
marcello.barbirotta,francesco.menichelli,
marco.angioli,mauro.olivieri}@uniroma1.it

**Abstract.** Efficient and cost-effective functional verification strategies are more and more essential in digital integrated system design. This paper presents a low-cost approach to meet this challenge, introducing a universal hardware emulator designed for verifying various Intellectual Property (IP) cores on FPGA. We demonstrate the practicality and performance of this emulator through a use-case of verifying an advanced Integer Arithmetic Logic Unit (ALU) for the RISC-V ISA vector extension. The process of integration and the results obtained from the verification are presented and discussed. Preliminary findings indicate that the emulator can perform more than 1.1 million tests per second. This research contributes to the advancement of hardware verification techniques, providing researchers, universities, and small businesses with an accessible and effective solution.

**Keywords:** Hardware emulator · Verification · Co-simulation · FPGAs · UVM

## 1 Introduction

The escalating complexity of digital integrated systems compels us to adopt robust and comprehensive functional verification strategies, which are critical in assuring the dependability of the final product [1]. Recent studies have highlighted that verification efforts consume about 50% of development time [2,3]. Functional design verification is typically done via logic simulation, pre-silicon hardware emulation, or ultimately silicon prototyping. Logic simulation offers low cost, flexibility, and precision but may be slow for big designs. Emulation facilitates very extensive test patterns thanks to a much higher speed, yet the implementation of an ad-hoc emulation environment from scratch is time-

consuming and error prone. Silicon prototyping is superior in terms of speed for almost final designs, but its implementation is costly and offers limited debugging capabilities [4, 5]

Hardware emulation has garnered significant attention due to its ability to accelerate verification by modelling hardware designs on a configurable devices hardware. FPGAs provide a suitable platform for the emulation of various IPs, enabling the rapid testing of new designs. Nevertheless, building a versatile and cost-effective hardware emulator on FPGA, which is able to accommodate a wide range of IPs, is highly demanding.

This work proposes a universal hardware emulator that can verify a variety of IPs on FPGA, also offering different verification modes that may involve hardware test-benches generated via High-Level Synthesis (HLS), or software test-benches running on a built-in processor on the FPGA chip, or even co-emulation with test-benches concurrently simulated on a host PC. To illustrate the utility and performance of our proposed universal emulator, we report the use-case of verifying an advanced arithmetic unit for vector processors. We demonstrate how the IP is integrated into the emulation system and present the results obtained from this exercise.

The paper concludes by elucidating the future directions for enhancing the capabilities of our emulator and exploring potential applications and improvements.

## 2   Background and Related Work

Numerous methodologies, including logic simulation, hardware emulation, prototyping, and formal verification, have been devised over the years to confront the challenge of ensuring the functional correctness of digital chip designs. Each of these approaches exhibits its unique set of advantages and limitations [4, 5]. Here we refer to the works related to hardware emulation of systems or individual IP blocks.

Existing literature documents several endeavours aimed at establishing specialized emulators on FPGA platforms, geared towards specific IPs [6–8]. However, these solutions inherently suffer from a lack of versatility and universality, as they are constrained by their tailored design.

In the commercial arena, universal hardware emulators such as those belonging to Cadence Palladium series, Synopsys ZeBu and Mentor Graphics Veloce series represent the pinnacle of performance and versatility. However, these high-end emulators come with substantial financial requisites, rendering them inaccessible for many researchers, educational institutions, and small-scale enterprises.

Addressing the apparent gaps in the existing body of work, this paper seeks to design and implement a universal hardware emulator that provides a cost-effective, yet high-performance solution for verifying a diverse array of IPs on FPGA.

# 3   Proposed Universal Hardware Emulator Design

## 3.1   Design Principles

The proposed approach aims at maximizing flexibility, enabling it to universally support a diverse range of Devices Under Test (DUTs) at a low cost. The system architecture incorporates three distinct environments that we refer to as the Personal Computer (PC), Programmable Logic (PL), and Soft Processor (SP), facilitating the implementation of both the verification process and the reference model. Moreover, the environment permits the parallel verification of multiple instances of the same IP blocks. A significant portion of the system has been implemented using HLS.

## 3.2   Architecture of the Emulator

Figure 1 presents the architecture of the proposed hardware emulator. The three environments PC, PL, SP are suitably connected via high-performance interfaces to facilitate simulation, co-simulation, and hardware emulation. The system includes a Universal Verification Methodology (UVM) environment, with a reference model that can be implemented in C++. Thanks to HLS technology, this model can be placed in any of the three environments. Running the reference model on FPGA PL environment offers speed and the potential for use in system prototyping. However, one must consider limitations for large models in terms of resource requirements and library usage constraints. In the PL environment, besides running the synthesized reference model, we may also use an already verified IP as a reference. This is particularly useful when one needs to update or modify existing IPs and may run a regression test comparing with the previously verified version.

A crucial part of the system is a constrained random generator module for generating stimuli. We utilized the Linear Feedback Shift Register (LFSR) method for random generation due to its simplicity and speed. Through the control software, initial values or seeds are updated or reseeded at regular intervals to maximize randomness and prevent output signal repetition [9]. In addition, this module can generate all the signals in parallel.

A clock generator is included in the architecture. It is designed to control the clock of the DUTs, thereby supporting cycle-accurate verification. Additionally, transaction-level modeling aids in efficiently checking functionality, and an operational mode that utilizes an array of stimuli can expedite the verification process. In this work, we utilized the VCU108 Xilinx platform, and for enhancing the ease of system configuration, we developed a user-friendly Graphical User Interface (GUI) using C++ and C# for the implementation.

As a result, the presented hardware emulator integrates multi-environment verification, synthesizable reference model implementation, and parallale random stimuli generation. By leveraging HLS technology, the emulator allows for reference model deployment across these environments, further enhancing testing flexibility.

**Fig. 1.** Architecture of the proposed universal hardware emulator

## 4   Use Case: Verification of an Advanced Arithmetic Unit for Vector Processors

### 4.1   Overview

In this section, we present a use case that demonstrates the application of our proposed universal hardware emulator for the verification of an advanced Arithmetic and Logic Unit (ALU) conceived to be integrated in vector processors. The ALU was specifically designed at the Barcelona Supercomputing Center (BSC) to support the RISC-V "V" vector extension v1.0 [10]. The ALU, expressed in SystemVerilog, supports a wide range of operations including addition, subtraction, multiplication, comparison, bitwise operations, shifting, and division on data widths of 64, 32, 16 and 8-bits, operating in packed-SIMD fashion.

### 4.2   Integration of ALU Unit into Emulation Process

The integration of the ALU i nthe verification system started with the design of a behavioral reference model, implemented in C++, whose block diagram is illustrated in Fig. 2. Using the HLS tool, the C++ model was converted into RTL to be synthesized on FPGA. The verification process was carried out within the PL environment, enabling high-speed execution. The integration and setup of the system consisted of several distinct steps:

1. Defining Input/Output Ports of the DUTs.
2. Defining Test Scenarios and Signal Constraint for generating stimuli.
3. Creating and integrating the Behaviornal Reference Model.
4. Defining the Verification Process: this step consists of programming the test sequences.
5. Running, Debugging, and Reporting.

### 4.3   Results

Our experimental evaluation of the proposed emulation process involved conducting more than 37 billion tests across 25 distinct scenarios to examine the

**Fig. 2.** Schematic of the advanced integer ALU unit used for verification

functional operations and control signals of the ALU, as shown in Table 1. The emulator executed approximately 1.1 million tests per second, showcasing its considerable speed. In contrast, traditional simulation methods managed about 33 thousand tests per second, showing the advantage of our emulation approach.

**Table 1.** Test results of ALU unit for 37 billion tests with 25 scenarios

| Tested Func. | Key features | Ctrl. sig | Number of tests (M) | Number of failed | Execution Time (min) |
|---|---|---|---|---|---|
| ADD/SUB | Average: × | × | 2 K | 0 | 29 |
| | | × | 4 K | 0 | 58 |
| | | ✓ | 500 | 0 | 7 |
| Multiplier | fi mode: × | × | 2 K | 0 | 29 |
| | | × | 2 K | 0 | 29 |
| | | ✓ | 500 | 0 | 7 |
| FMA | Add/sub: ✓ | × | 2 K | 0 | 34 |
| | Add/sub: ✓ | ✓ | 500 | 0 | 8 |
| Bitwise | | × | 2 K | 0 | 28 |
| | | ✓ | 500 | 0 | 7 |
| Unary | VPOPCNT, VMFIRST | × | 2 K | 0 | 33 |
| | VPOPCNT VMFIRST   \| Input src1: −1 | × | 2 K | 0 | 34 |
| | VMSIF, VMSIF, VMSIF | × | 2 K | 0 | 33 |
| | VMSIF, VMSIF, VMSIF \| Input src1: 0 | × | 2 K | 0 | 34 |
| | VPOPCNT, VMFIRST | ✓ | 500 | 1.2 M | 7 |
| | VMSIF, VMSIF, VMSIF | ✓ | 500 | 0 | 7 |
| | VMSIF, VMSIF, VMSIF \| Input src1: 0 | ✓ | 500 | 750 K | 7 |
| Comp. | | × | 2 K | 0 | 31 |
| | | ✓ | 500 | 0 | 8 |
| Shifter | Left \| Round: ×\| Arith.: ×\| Out. Mode: BASE, WIDEN | × | 2 K | 0 | 27 |
| | Left \| Round: ×\| Arith.: ×\| Out. Mode: MASK | × | 2 K | 500 | 29 |
| | Out. Mode: NARROW | × | 500 | 0 | 7 |
| | Right \| Round \| Out. Mode: BASE, WIDEN | × | 2 K | 0 | 28 |
| | Right \| Round \| Out. Mode: MASK | × | 2 K | 0 | 28 |
| | Right \| Round \| Out. Mode: MASK | ✓ | 500 | 0 | 7 |

*Note* ✓ = parameter enabled; ×= disabled

# 5    Conclusion and Future Work

In this study, we introduced a hardware emulator designed to test various IPs on FPGAs. Our initial results indicated that our emulator exhibited exceptional

speed. However, our testing was restricted and did not include a comparison with other emulators. For future work, We intend to evaluate our emulator using a wider range of IPs, including floating-point units, artificial intelligence IPs, and reconfigurable HW accelerators [11]. This comprehensive analysis will enable us to better understand the strengths of our emulator, as well as areas that may require improvement.

# References

1. Vasudevan S (2017) Still a fight to get it right: verification in the era of machine learning. IEEE Int Conf Rebooting Comput (ICRC) 2017:1–8
2. Siemens (2022) Part 3: The 2022 wilson research group functional verification study [Online]. Available: https://blogs.sw.siemens.com/verificationhorizons/2022/10/30/part-3-the-2022-wilson-research-group-functional-verification-study/
3. Akram A, Sawalha L (2019) A survey of computer architecture simulation techniques and tools. IEEE Access 7:78120–78145
4. Scheffer L, Lavagno L, Martin G (2006) EDA for IC system design, verification, and testing (Electronic Design Automation for Integrated Circuits Handbook). CRC Press, Inc.
5. Mehta AB (2018) ASIC/SoC functional design verification. Springer, Berlin
6. Shi K et al (2023) ENCORE: efficient architecture verification framework with FPGA acceleration. In: Proceedings of the 2023 ACM/SIGDA international symposium on field programmable gate arrays (FPGA 2023). ACM, Monterey, CA, USA, pp 209–219
7. Lagos-Benites J et al (2011) An FPGA-emulation-based platform for characterization of digital baseband communication systems. IEEE Int Symp Defect Fault Toler VLSI Nanotechnol Syst 391–398
8. Tomas BJ, Jiang Y, Yang M (2014) SoC scan-chain verification utilizing FPGA-based emulation platform and SCE-MI interface. In: 2014 27th IEEE international system-on-chip conference (SOCC), pp 398–403
9. Goss C (2011) Improving verification productivity with the dynamic load and reseed methodology. Cadence Des Syst
10. Minervini F et al (2023) An area-efficient RISC-V decoupled vector coprocessor for high performance computing applications. ACM Trans Arch Code Optim 20(2):1–25
11. Jamili S et al (2022) Implementation of dynamic acceleration unit exchange on a RISC-V soft-processor. In: *International conference on applications in electronics pervading industry, environment and society.* Springer, Berlin, pp 300–306

# Single Event Transient Reliability Analysis on a Fault-Tolerant RISC-V Microprocessor Design

Marcello Barbirotta[✉], Marco Angioli, Antonio Mastrandrea,
Abdallah Cheikh, Saeid Jamili, Francesco Menichelli, and Mauro Olivieri

Department of Information Engineering, Electronics and Telecommunications
(DIET), "La Sapienza" University of Rome, Via Eudossiana, 18, 00184 Rome, Italy
{Marcello.Barbirotta,Marco.Angioli,Antonio.Mastrandrea,Abdallah.Cheikh,
Saeid.Jamili,Francesco.Menichelli,Mauro.Olivieri}@uniroma1.it

**Abstract.** The miniaturization of electronic devices and the improved operating speeds increase the likelihood of single event faults. Differently from Single Event Upset (SEU) faults, Single Event Transient (SET) faults generally affect combinational logic, making all voting systems vulnerable to errors. The proposed work uses an ad-hoc fault-simulation campaign employing signal glitching to identify SET vulnerabilities inside a RISC-V core already equipped with resilience logic against Single Event Upset (SEU) faults. The faults target the majority voting logic structures, highlighting how they can be susceptible to faults depending on the width of the injected pulses, and showing how the use of Buffered Triple Modular Redundancy (BTMR) allows decreasing the total failure probability due to erroneous majority voters.

**Keywords:** Fault-resilience · Single event transient · Fault-injection

## 1 Introduction

Single Event Transient (SET) faults refer to the transient electrical effect that can occur in microelectronics. Differently from Single Event Upset (SEU), SET do not intrinsically induce a state change in sequential components, unless the transient effect is sampled by a memory element in the circuit. The impacts of SETs range from minor effects, such as a temporary signal glitch, to more serious consequences, such as permanent failure of the device. Due to device geometry scaling, multiple transistors fit into the impact area of a single ion strike. For this reason, SETs may represent an issue not only in fields where electronics are subjected to high levels of ionizing radiation - such as aerospace, military, and nuclear power - but also in consumer, automotive, and industrial electronics. This work aims to perform a resilience analysis of SET-induced failures on a RISC-V-compliant processor core that is SEU-fault-tolerant by design, to evaluate the impact of the transient pulse width on the failure probability, and the possible improvements derived from the use of the Buffered TMR technique [1].

## 2    Background and Related Works

Single Event phenomena occur when a single energetic particle hits a target device, generating electron-hole pairs through ionization processes. The generated current movement in the presence of an electric field causes a glitch defined as a Single Event Transient. In digital electronics, sufficient transient currents in OFF elements can temporally upset the stable condition, generating a SET fault. In contrast, when a charge is added to an element that is already ON, it does not affect the logic behavior. When a SET pulse propagates through the device nodes may be incorrectly latched by a register as data to be stored. SET propagation happens if the duration of the SET pulse is longer or comparable with a logic cell delay, so that for more advanced technology nodes SET fault propagation is more likely [2]. Moreover, SET latching happens if the transient pulse reaches a register and the clock edge occurs before the transient pulse has finished, so that for higher clock frequency SET latching is more likely [2].

In the literature, many works about SET fault resilience base their interest on eliminating the SET pulse from the circuit or eliminating the effects of the SET pulse from the circuit. The former approach lead to custom designs, with the increase in Area and Power consumption [3]. For this reason, in recent years, more attention has been made to the search for standard cell solutions capable of providing an adequate degree of fault tolerance to designs without increasing costs or production time. Authors in [3] present different TMR standard-cell Flip Flop where the robustness is achieved by an integrated transient filter on the data path obtained with low/high drive inverter. Other works like [4] use the self-voter approach to create the Self-Voting DMR FF, but they add delays in the architecture. Authors in [5] propose a multi-bit FF able to check errors with bit-even parity and clock skew, filtering up to 142, 150, and 165 ps pulse widths. However, for all these approaches, it is not always possible to definitively eliminate the probability that the fault SETs could damage the circuit. In many cases, the leaf nodes (outputs) of the proposed circuits that may interface with control inputs of a flip-flop are still vulnerable to SETs [6]. Generally speaking, any SET on the output multiplexer that drives the final output can result in a glitch, which could eventually propagate to data inputs of subsequent stage FFs [5].

## 3    *Klessydra-fT03* Fault Tolerant Microarchitecture

Klessydra-fT03 (Fig. 1) is a fault-tolerant RISC-V processor core based on an open-source RISC-V softcore family named Klessydra-T [7], developed at the Digital System Lab group at the Sapienza University of Rome [8]. It uses the Buffered TMR paradigm to achieve SEU Radiation Hardening features [1,8–10] with three different threads executing the same code with clock cycle granularity for temporal redundancy and with the corresponding hardware support for spatial redundancy.

**Fig. 1.** Klessydra-fT03 microarchitecture with shared pipeline units and spatial hardware support for three different threads represented by *orange*, *red* and *green* colours, with five majority voting units handling 1171 flip-flops.

## 4   Single Event Transient Analisys

SET pulses depend significantly on operating conditions such as the supply voltage, temperature, load, and driving strength [11]. To take into account all these characteristics, physical mixed-mode modelling and simulation modelling should be used, however, these techniques are time-consuming, resource-hungry, and require full design structure details, resulting not applicable in practice [12]. For these reasons, other techniques have been introduced, and in this work, we will rely on the RTL fault injection simulation using a Universal Verification Methodology (UVM) fault-injection environment [13], to firstly evaluate the vulnerability to SET faults occurring within the combinational logic of voting circuits inside a fault-tolerant core previously tested for SEU-type faults. From a fault simulation point of view, current pulses were simulated with different pulse widths inside the voter output. Existing experimental results indicate that the pulse widths vary from approximately 900 ps to over 3 ns [14], decreasing with scaling of approximately 100–700 ps [15,16]. For these reasons, we decided to reproduce different width patterns described in Fig. 2, also increasing the failure rate from zero to 2000 multiple consecutive faults, with a 5% error margin in a Monte Carlo scenario for a *Coremark* benchmark. We observed that the failure rates grow as the pulse width increases because it increases the probability that the sequential logic elements capture these faults. As a proof of concept, according to [17], the probability of latching the fault is:

$$P = \frac{\delta}{T_{CLK}} \tag{1}$$

and results as the ratio between the pulse width and clock period, as reflected in Fig. 2. Decreasing the pulse width while maintaining a constant frequency

decreases the latching probability and thus reduces the failure probability. If the flip-flop data input is not stable during the lockout window (–tSETUP, +tHOLD) due to the transient fault, the real flip-flop will latch "1" or "0", but using a logic simulation, we have no way of telling which [17]. To avoid this uncertainty, since RTL simulations have no information about logical delays, all faults within the performed simulation do not occur during the rising edge of the clock, but only before that, so as to be within setup and hold times.



**Fig. 2.** 5% error margin Montecarlo analysis for SET random Fault Injection inside Klessydra-fT03 pipeline voters. Comparisons between different Single Event Rate with different SET-width on Coremark benchmark

## 5   Performance Evaluations

The bottleneck of majority voting systems is voting itself, despite the development of several costly solutions (reported in Sect. 3), able to reduce the total failure probability. Another way to see the reduction of this probability could be to reduce the number of voting systems that appears inside the entire design. As can be observed from the pipeline shown in Fig. 1, unlike basic TMR or full TMR fault-tolerant systems, Klessydra-fT03 has a very low number of voting structures inside. For the Program Counter (PC), Register File (RF), Load Store Unit (LSU), and Write Back (WB), there are a total of 3513 bits with 1171 voters. To better clarify the hardware overhead of a voting block from FPGA synthesis evaluation, we found that a 32-bit voter with no SET protection mechanism inside occupies around 32 LUT in Xilinx Kintex-7 FPGA. From that, it is possible to roughly consider 1 LUT for each 1 TMR voter. Wanting to theoretically evaluate the SET fault tolerance performances of the Buffered TMR core linked to the voting structures, we have to consider its single-thread version Klessydra-S1 and its hardware occupation reported in Table 1 [1]. A basic-TMR version of Klessydra-S1 would replicate the hardware three times, with 200% hardware overhead in flops and 55% in LUTs. In contrast, a full-TMR version would triplicate flops, combinational logic and voter logic, adding extra voters for leaf nodes. Considering the total size of the Klessydra-fT03 core,

compared with the basic-TMR version, there would be approximately 16.06% less hardware occupation for FFs while there would be 21,76% more LUTs. On the other hand, as stated before, considering voter occupation, we have 66.52% fewer voters than basic-TMR and 400% fewer voters than full-TMR systems. Assuming these values, we could roughly state that the probability that a voter is hit by a fault inside the Buffered-TMR architecture is around 66% less than the one inside a basic-TMR architecture and 400% time less than that one in a full-TMR system, proving the way of reducing failures by reducing the number of voting systems inside the architecture, without increasing power, area and time overhead adopting expensive SET tolerance solutions.

**Table 1.** Microprocessor Synthesis results on a Xilinx Kintex-7 FPGA for the Buffered TMR Klessydra-fT03 and single thread Klessydra-S1. Basic-TMR and full-TMR values are estimated considering 3x improvement in flops and logic [18].

|      | S1 no-redundancy | fT03 buffered-TMR | Estimated basic-TMR | Estimated full-TMR |
|------|------------------|-------------------|---------------------|--------------------|
| LUTs | 3528             | 6670              | 5478                | 16434              |
| FFs  | 1950             | 4910              | 5850                | 5850               |

## 6    Conclusions

This work presented a Single Event Transient analysis of a Buffered TMR microprocessor architecture already resilient to SEU faults. We performed SET fault-injection using a custom UVM environment with different fault rates and SET widths ranging from 600 ps to 7 ns, confirming that, in the hypothesis of a single fault not overlapping with the rising edge of the clock, the failure levels depend on the fault width as a function of the probability of latching a fault. We also theoretically proved that the Buffered TMR technique reduces the total failure probability because it implies a low percentage of voters inside, compared with standard approaches like basic-TMR or full-TMR. Moreover, due to its structure, Klessydra-fT03 intrinsically has good resilience to SET-type faults since the combinatorial logic of the pipeline is temporally shared among the various threads, and any transient faults will result in incorrect values within each thread's spatially redundant sequential logic, possibly deleted by voting systems.

# References

1. Barbirotta M, Cheikh A, Mastrandrea A, Menichelli F, Olivieri M (2022) Design and evaluation of buffered triple modular redundancy in interleaved-multi-threading processors. IEEE Access 10:126074–126088

2. Narasimham B, Bhuva BL, Holman WT, Schrimpf RD, Massengill LW, Witulski AF, Robinson WH (2006) The effect of negative feedback on single event transient propagation in digital circuits. IEEE Trans Nucl Sci 53(6):3285–3290

3. Schrape O, Andjelković M, Breitenreiter A, Zeidler S, Balashov A, Krstić M (2021) Design and evaluation of radiation-hardened standard cell flip-flops. IEEE Trans Circuits Syst I: Regul Pap 68(11):4796–4809

4. Teifel J (2008) Self-voting dual-modular-redundancy circuits for single-event-transient mitigation. IEEE Trans Nucl Sci 55(6):3435–3439

5. Jain A, Veggetti AM, Crippa D, Benfante A, Gerardin S, Bagatin M (2022) Radiation tolerant multi-bit flip-flop system with embedded timing pre-error sensing. IEEE J Solid-State Circuits 57(9):2878–2890

6. Sajjade FM, Goyal NK, Varaprasad B (2021) Single event transient (set) mitigation circuits with immune leaf nodes. IEEE Trans Device Mater Reliab 21(1):70–78

7. Cheikh A, Sordillo S, Mastrandrea A, Menichelli F, Scotti G, Olivieri M (2021) Klessydra-t: designing vector coprocessors for multithreaded edge-computing cores. IEEE Micro 41(2):64–71

8. Barbirotta M, Cheikh A, Mastrandrea A, Menichelli F, Ottavi M, Olivieri M (2022) Evaluation of dynamic triple modular redundancy in an interleaved-multi-threading risc-v core. J Low Power Electron Appl 13(2)

9. Barbirotta M, Cheikh A, Mastrandrea A, Menichelli F, Angioli M, Jamili S, Olivieri M (2023) Fault-tolerant hardware acceleration for high-performance edge-computing nodes. Electronics 12(17). https://www.mdpi.com/2079-9292/12/17/3574

10. Barbirotta M, Cheikh A, Mastrandrea A, Menichelli F, Olivieri M (2022) Analysis of a fault tolerant edge-computing microarchitecture exploiting vector acceleration. In: 2022 17th conference on Ph. D research in microelectronics and electronics (PRIME). IEEE, pp 237–240

11. Veeravalli VS, Polzer T, Schmid U, Steininger A, Hofbauer M, Schweiger K, Dietrich H, Schneider-Hornstein K, Zimmermann H, Voss KO et al (2013) An infrastructure for accurate characterization of single-event transients in digital circuits. Microprocess Microsyst 37(8):772–791

12. Hamad GB, Mohamed OA, Savaria Y (2016) Towards formal abstraction, modeling, and analysis of single event transients at rtl. In: 2016 IEEE international symposium on circuits and systems (ISCAS). IEEE, pp 2166–2169

13. Barbirotta M, Mastrandrea A, Menichelli F, Vigli F, Blasi L, Cheikh A, Sordillo S, Gennaro FD, Olivieri M (2020) Fault resilience analysis of a risc-v microprocessor design through a dedicated uvm environment. In: 33rd IEEE international symposium on defect and fault tolerance in VLSI and nanotechnology systems, DFT 2020. Institute of Electrical and Electronics Engineers Inc.

14. Narasimham B, Gadlage MJ, Bhuva BL, Schrimpf RD, Massengill LW, Holman WT, Witulski AF, Zhu X, Balasubramanian A, Wender SA (2008) Neutron and alpha particle-induced transients in 90 nm technology. In: 2008 IEEE international reliability physics symposium. IEEE, pp 478–481

15. Ferlet-Cavrois V, Massengill LW, Gouker P (2013) Single event transients in digital cmos-a review. IEEE Trans Nucl Sci 60(3):1767–1790

16. Chen Q, Liu Y, Wu Z, Liao J (2021) A single event effect simulation method for risc-v processor. In: 2021 IEEE 15th international conference on anti-counterfeiting, security, and identification (ASID). IEEE, pp 106–110
17. Alexandrescu D, Anghel L, Nicolaidis M (2002) New methods for evaluating the impact of single event transients in vdsm ics. In: 17th IEEE international symposium on defect and fault tolerance in VLSI systems. DFT 2002. Proceedings. IEEE, pp 99–107
18. Petrovic V, Krstic M (2015) Design flow for radhard tmr flip-flops. In: 2015 IEEE 18th international symposium on design and diagnostics of electronic circuits and systems. IEEE, pp 203–208

# Secure Data Authentication in Space Communications by High-Efficient AES-CMAC Core in Space-Grade FPGA

Luca Crocetti[1]([✉]), Francesco Falaschi[2], Sergio Saponara[1], and Luca Fanucci[1]

[1] Department of Information Engineering, University of Pisa, Via G. Caruso 16, 56122 Pisa, Italy
{luca.crocetti,sergio.saponara,luca.fanucci}@unipi.it
[2] IngeniArs S.r.l., Via Ponte a Piglieri 8, 56121 Pisa, Italy
francesco.falaschi@ingeniars.com

**Abstract.** Latest technological improvements and investments from government agencies and private companies pushed to the limits the requirements related to both data rate speed and security of the communication links in space applications. The high volume of data and the continuous integration of services opened the path to hackers for new and increasingly diffused cyberattacks. Governmental agencies are attempting to stem this problem by issuing and updating accordingly a series of reports and standards through the Consultative Committee for Space Data Systems (CCSDS). In this work, we present the implementation of an Advanced Encryption Standard—Cipher-based Message Authentication Code (AES-CMAC) core on space-grade FPGAs, that is compliant with the latest CSSDS security standards and outperforms the state-of-the-art in terms of resource efficiency.

**Keywords:** Space security · Cybersecurity · Hardware security · CCSDS · 352.0-B-2 · 355.0-B-2 · SDLS · Protocol · AES · CMAC · Data · Integrity · Authentication · Space-grade · FPGA

## 1 Introduction

Over time, highly technological fields of application are always characterized by the constant growth of volume and communication speed of the data transfers, and continuous operations required to support innovative and even more complex services and interconnected systems. In recent years, this trend especially invested fields such as automotive, industry 4.0, Internet-of-Things (IoT), and space, which actually reached an average data rate per satellite in the range of 2.04–22.74 Gbps [1]. As a drawback, the increased data traffic, often confidential, constitutes an augmented surface of attacks for hackers, in the other sectors as well as in the space one [2,3]. Indeed [4], reports a growing number of incidents in satellites due to cyberattacks and violations of security (Fig. 1).

**Fig. 1.** Space sector [4]: number of operational satellites (red line) versus number of security incidents (blue bars).

To counteract this phenomenon, the Consultative Committee for Space Data Systems (CCSDS) occupies of issuing and maintaining a series of reports, and standards to meet the requirements of space applications, including also security. These requirements are applicable to both the satellite-to-ground (*downlink*) communications, mainly employed to transmit confidential scientific payload data, and ground-to-satellite (*uplink*) connections, used to send remote controls that need protection against unauthorized tampering. The CCSDS reports 350.0-G [5] and 350.1-G [6] provide an overview of possible threats and introduce multiple security mechanisms. In particular [5], presents the space link reference model that finds correspondence with the OSI model layers. According to [7], we focused our work on the implementation of a hardware accelerator to address the space security requirements at the Data Link layer, by generating a solution that is both high-speed and efficient in terms of supported data rate per used logic resources. The equivalent layer in the CCSDS space link reference model is the Space Data Link Security (SDLS) layer.

## 2  Integrity and Authentication for Secure Space Data Links

The SDLS layer is regulated by the CCSDS standards 352.0-B [8] and 355.0-B [9], and it introduces a frame format integrating a Security Header and an optional Security Trailer to envelope and protect the payload (i.e. the Frame Data). For this purpose [8,9], specify the usage of 4 different algorithms for two different security services. The first service is authenticated encryption, that provides at the same time confidentiality, integrity, and authentication of data. It is realized

by employing the Advanced Encryption Standard—Galois/Counter Mode (AES-GCM) [10], which is used to encrypt the Frame Data and transmit a Message Authentication Code (MAC) as Security Trailer, providing authenticity of both Frame Data and the Security Header. The second service offers only integrity and authentication, and it can be employed in contexts in which confidentiality is not needed, lowering the processing time of data. This service can be realized with any of the remaining 3 security schemes, i.e. Cipher-based MAC (CMAC) using the AES algorithm [11], the Hash-based MAC (HMAC) using the SHA2 algorithm [12], and the Galois MAC (GMAC) [10] that exploits multipliers over a Galois finite field.

In previous works [13–15] we concentrated on the optimization of a hardware AES core for high-performance computing that was used as the main building block of a hardware accelerator for to the AES-GCM algorithm and to be used in the MACsec scheme [16] for securing the Automotive Ethernet links [17]. For this reason, in this work, we decided to concentrate on the implementation of an IP-core for integrity and authentication through MACs. According to [18], the usage of an AES-based core has to be preferred with respect to a SHA2-based core for both data rate and efficiency reasons because the solution based on the AES algorithm leads to hardware modules that achieve higher throughput and consume fewer logic resources. In addition, the GMAC scheme exploits the same multiplier integrated into the AES-GCM algorithm, and from [13] it emerges that the cost in terms of logic resources is higher for the Galois multiplier with respect to the AES core, without advantages in terms of maximum frequency and throughput. Therefore, we opted for the implementation of an AES-CMAC hardware accelerator on space-grade FPGAs.

## 3    Design of the AES-CMAC Hardware Accelerator

The AES-CMAC algorithm [19] is a chaining mode of operation of the AES cipher, which exploits only the AES encryption algorithm to generate MACs. With reference to Fig. 2(a), the input message is split into $N$ blocks ($M_1$, $M_2$ to $M_N$), and each block is firstly XORed with the output of the previous block and then processed through the encryption algorithm AES with the key $K$; finally, the output MAC is generated by processing the last input block $M_N$ that is XORed also with the sub-key $K_1$ or $K_2$ (derived from the primary key $K$), according to the bit length of $M_N$.

Because of its nature, parallelization or pipelining approaches of the internal architecture are useless because they would not give any advantage in terms of speed or data rate. For this reason, we aimed to optimize at the most efficiency by implementing a core that constitutes the best trade-off between the utilization of logic resources and critical delay. Exploiting the design architecture and the optimization strategies proposed in [13], we developed an AES-CMAC core in SystemVerilog that is mainly composed of an AES encryption core and a control unit to handle the chaining process (CCU, Chaining Control Unit, in Fig. 2(b)).

**Fig. 2.** AES-CMAC core: algorithm scheme (**a**), and architecture outline (**b**).

## 4    Implementation Results on Space-Grade FPGAs

We implemented the AES-CMAC core on different space-grade FPGAs, in particular on a Kintex Ultrascale KU060 by Xilinx/AMD, and an RTG4 by Microsemi. The implementation results shown in Table 1 refer to the proposed AES-CMAC core configured to support both 128-bit and 256-bit keys.

**Table 1.** Implementation results on space-grade FPGAs for 128-bit/256-bit keys.

| Space-qualified FPGA | Frequency (MHz) | Resource utilization | Throughput (Gbps) | Power consumption (mW) |
|---|---|---|---|---|
| KU060 | 224.37 | 434 CLBs | 2.871/2.051 | 138 |
| RTG4 | 81.59 | 3842 LEs | 1.044/0.746 | 149 |

CLBs (Configurable Logic Blocks) and LEs (Logic Elements) in Table 1 are the programmable logic resources of the corresponding FPGA device; the column Frequency indicates the maximum frequency; the double value in column Throughput indicates the data rate for 128-bit/256-bit keys respectively. In both cases, the throughput achieved on the FPGA KU060 is able to support at least the minimum average data rate per satellite (2.04 Gbps).

The state-of-the-art related to the implementation of AES-CMAC cores compliant with CCSDS security standards [8,9] offers few (recent) works [20,21], and none of them employs space-qualified FPGAs. For this reason, we synthesized our core also on the Virtex-5 FPGA device used by both [20,21], and, accordingly, we configured our core to support only 128-bit keys. The results of the comparison are shown in Table 2.

Referring to the results in Table 2, although our core is characterized by the lowest frequency and the lowest throughput, it shows the best efficiency with respect to other solutions, because it is able to support data rates of 1.96 Gbps occupying only 462 CLBs, respectively, about one-third and one-half the logic resources required by [20,21].

**Table 2.** Comparison with the state-of-the-art.

|          | Fmax (MHz) | Resource utilization (CLBs) | Throughput (Gbps) | Resources efficiency (Mbps/CLB) |
|----------|-----------|------------------------------|-------------------|---------------------------------|
| [20]     | 320.84    | 1355                         | 3.80              | 2.80                            |
| [21]     | 268.85    | 844                          | 2.87              | 3.40                            |
| This work| 152.98    | 462                          | 1.96              | 4.24                            |

## 5  Conclusions

In this work, we presented the development of an AES-CMAC core in SystemVerilog to provide data authentication in space communications according to the SDLS specifications. Our core has been implemented on space-grade FPGAs obtaining results that confirm its employment in modern satellite applications. Moreover, the comparison with the state-of-the-art highlights that our core is about 1.25 to 1.5 times more efficient than the other existing solutions compliant with the CCSDS security standards 352.0-B/355.0-B, and it occupies about 50–67% less logic resources than other cores. Such a result implies that the proposed core can be replicated in parallel 2 to 3 times occupying the same logic resources and offering data rates much higher than other solutions.

Future works will include the realization of a merged hardware accelerator for both AES-CMAC and AES-GCM algorithms by combining the results of the work inhere presented and the previous works about the AES-GCM core [13–15], and the analysis of Side-Channel vulnerabilities [22] with the implementation of related countermeasures.

## References

1. Del Portillo I, Cameron BG, Crawley EF (2019) A technical comparison of three low earth orbit satellite constellation systems to provide global broadband. Acta Astron 159:123–135
2. Naja G, Mathieu C (2015) Space and Security in Europe. In: Handbook of space security. Springer, pp 371–383
3. Orlova A, Nogueira R, Chimenti P (2020) The present and future of the space sector: a business ecosystem approach. Space Policy 52:101374
4. Manulis M, Bridges CP, Harrison R, Sekar V, Davis A (2021) Cyber security in new space: analysis of threats, key enabling technologies and challenges. Int J Inf Secur 20(3):287–311 (Springer)

5. CCSDS (2019) The application of security to CCSDS protocols. Inf Rep CCSDS 350.0-G-3
6. CCSDS (2022) Security threats against space missions. Inf Rep CCSDS 350.1-G-3
7. Baldanzi L, Crocetti L, Di Matteo S, Fanucci L, Saponara S, Hameau P (2019) Crypto accelerators for power-efficient and real-time on-chip implementation of secure algorithms. In: 26th IEEE international conference on electronics, circuits and systems (ICECS). IEEE, pp 775–778
8. CCSDS (2019) CCSDS cryptographic algorithms. Recomm Stand CCSDS 352.0-B-2
9. CCSDS (2022) Space data link security protocol. Recomm Stand CCSDS 355.0-B-2
10. NIST (2007) Recommendation for block cipher modes of operation: galois/counter mode (GCM) and GMAC. Special Publication (SP) 800-38D
11. NIST (2001) Advanced encryption standard (AES). Federal Information Processing Standards (FIPS) publication 197
12. NIST (2015) Secure hash standard. Federal Information Processing Standards (FIPS) publication 180-4
13. Nannipieri P, Baldanzi L, Crocetti L, Di Matteo S, Falaschi F, Fanucci L, Saponara S (2022) CRFlex: a flexible and configurable cryptographic hardware accelerator for AES block cipher modes. In: Applications in electronics pervading industry, environment and society (APPLEPIES 2021). Springer, pp 31–38
14. Nannipieri P, Di Matteo S, Baldanzi L, Crocetti L, Zulberti L, Saponara S, Fanucci L (2022) VLSI design of advanced-features AES cryptoprocessor in the framework of the European processor initiative. IEEE Trans Very Large Scale Integr (VLSI) Syst 30(2):177–186. https://ieeexplore.ieee.org/document/9631958
15. Nannipieri P, Crocetti L, Di Matteo S, Fanucci L, Saponara S (2023) Hardware design of an advanced-feature cryptographic tile within the European processor initiative. IEEE Trans Comput
16. IEEE (2018) IEEE standard for local and metropolitan area networks-media access control (MAC) security. Standard IEEE 802.1AE-2018
17. Carnevale B, Falaschi F, Crocetti L, Hunjan H, Bisase S, Fanucci L (2015) An implementation of the 802.1AE MAC security standard for in-car networks. In: 2015 IEEE 2nd world forum on internet of things (WF-IoT). IEEE, pp 24–28
18. Baldanzi L, Crocetti L, Falaschi F, Belli J, Fanucci L, Saponara S (2020) Digital random number generator hardware accelerator IP-core for security applications. In: Applications in electronics pervading industry, environment and society (APPLEPIES). Springer, pp 117–123
19. NIST (2006) Recommendation for block cipher modes of operation: the CMAC mode for authentication. Special Publication (SP) 800-38B
20. Dhaou IB, Gia TN, Liljeberg P, Tenhunen H (2017) Low-latency hardware architecture for cipher-based message authentication code. In: IEEE international symposium on circuits and systems (ISCAS). IEEE, pp 1–4
21. Pirzada SJH, Murtaza A, Hasan MN, Xu T, Jianwei L (2019) The implementation of AES-CMAC authenticated encryption algorithm on FPGA. In: IEEE 2nd international conference on computer and communication engineering technology (CCET). IEEE, pp 193–197
22. Crocetti L, Baldanzi L, Bertolucci M, Sarti L, Carnevale B, Fanucci L (2019) A simulated approach to evaluate side-channel attack countermeasures for the advanced encryption standard. Integration 68:80–86

# On the Usage of Isomorphic Fields in Hardware AES Modules for Optimizing the Efficiency

Luca Crocetti$^{(\boxtimes)}$ and Sergio Saponara

Department of Information Engineering, University of Pisa, Via G. Caruso 16, 56122 Pisa, Italy
{luca.crocetti,sergio.saponara}@unipi.it

**Abstract.** The Advanced Encryption Standard (AES) is widely accepted as the de-facto standard for symmetric-key encryption, and it is going to be used in the coming decades because of its resistance against Post-Quantum Cryptography. For this reason, it is the subject of many research works, and almost all converge on the usage of composite/tower fields for the hardware implementation of the S-box, the most expensive circuit in terms of both area and critical delay. Anyway, the debate is still open on applying isomorphic fields also to the other AES algorithm operations. In the attempt to give an answer, it is analyzed the application of the two approaches to the most recent and performing solutions from the state-of-the-art with the synthesis of the corresponding circuits on a 7 nm standard-cell technology. In addition, the presented work constitutes also a guideline for implementing hardware AES modules that execute all operations over composite/tower fields.

**Keywords:** AES · Round · S-box · Composite · Tower · Field · Galois · Encryption · High performance · Hardware accelerator

## 1 Introduction

The Advanced Encryption Standard (AES) [1] was released by the National Institute of Standards and Technology (NIST) and represents the de-facto standard for symmetric-key encryption, also because of its efficiency and performance [2]. Indeed, it is employed in several application fields such as High-Performance Computing [3,4] and Automotive Security [5], and it is going to be used in the coming decades because of its resistance against Post-Quantum Cryptography [6]. For this reason, a high volume of works focusing on its optimization can be found in the literature. Concerning hardware implementations, almost all of the works converge on the usage of composite (or tower) fields to reduce the complexity of the S-box circuit [7–13], which consists in the calculation of the multiplicative inverse of a byte and an affine transformation. This approach consists in mapping the AES native field $GF(2^8)$ to an isomorphic field $GF((2^4)^2)$

(*composite* field) or $GF(((2^2)^2)^2)$ (*tower* field) [11], reducing the problem of the multiplicative inversion on a 4-bit vector, $GF(2^4)$, or on a 2-bit vector, $GF(2^2)$. In both cases, the basis used to represent the bytes on the isomorphic field can be a *Normal Basis* (NB), a *Polynomial Basis* (PB) [12,13], a *Redundant Representation Basis* (RRB) [10,11], or a *Mixed Basis* (MB) [7–9]. Whatever the composite/tower field and the basis used, the S-box circuit has always the same structure that is illustrated in Fig. 1.



**Fig. 1.** Outline of the composite/tower field S-box circuit.

In Fig. 1, the blocks Map and Map$^{-1}$ represent respectively the isomorphic mapping ($M$) and the inverse isomorphic mapping ($M^{-1}$): this last to map back the multiplicative inverse on $GF(2^8)$ before the affine function (block Affine Trans.). Because both the inverse mapping and the affine transformation correspond to a matrix-vector multiplication, respectively, with the matrix $M^{-1}$ and the matrix $A$, these two operations are merged into a unique matrix-vector multiplication ($A \cdot M^{-1}$), reducing both the gate count and the gate delay. Anyway, some works analyzed also the effects of extending the isomorphism to other AES operations, and the researchers have not found a common conclusion on this aspect. Some works propose architectures that extend the isomorphism to the other AES transformations [7], some others explicitly declare that this approach has no advantages, rather it worse the efficiency both in terms of area and maximum frequency [8,9], or they do not consider this aspect at all, making the implicit assumption that the isomorphism should involve only the multiplicative inverse [10–13]. To investigate the effects of isomorphism on the other AES operations, we selected the most recent and performing works from the literature that use different composite/tower fields, and both approaches were implemented for each of them using SystemVerilog. The circuits were synthesized on a 7 nm standard-cell technology and characterized by maximum frequency, area, and efficiency (expressed as frequency per area). The analysis and the results that follow refer to the encryption algorithm of AES.

## 2   Application of the Isomorphism Only to the AES S-Box

The AES encryption algorithm iteratively processes 16-byte data blocks arranged in a $4 \times 4$ matrix according to a certain number of rounds. In the case of 128-bit keys, the encryption process performs 10 (main) rounds composed by the operations *SubBytes* (the substitution of each byte by means of the S-box, Fig. 1), *ShiftRows* (a byte re-ordering), *MixColumns* (a linear transformation), and *AddRoundKey* (a 128 bitwise XOR between the data block and a round-key), as shown in Fig. 2. The illustrated architecture includes also an additional

XOR between the input data block (data in) and the input key (key in) for the preliminary round, and a multiplexer after the *MixColumns* block to skip this operation in the last round. The round-keys are derived from the input key using operations similar to the ones of the AES round, such as the *SubWord*, the substitution through the S-box of a 32-bit word (4 bytes of the key). The only different operation is the XOR with the *Rcon* constant, anyway, the overall area and timing complexity of the key expansion circuit (Fig. 2) is lower than the one of the AES round circuit [3,6], which contains the critical, and it is:

$$t_R = \underbrace{t_{Map} + t_{MultInv} + t_{InvMap||Aff}}_{t_{S\text{-}box}} + t_{MixCol} + 2 \cdot t_{MUX} + t_{XOR} \qquad (1)$$

In Eq. 1, $t_{Map}$, $t_{MultInv}$, $t_{InvMap||Aff}$, $t_{MixCol}$, $t_{MUX}$, $t_{XOR}$ are, respectively, the propagation delays of isomorphic mapping, multiplicative inverse, inverse isomorphic mapping merged with the affine transformation, the *MixColumns*, a multiplexer, and an XOR gate (*AddRoundKey*).



**Fig. 2.** Architecture of the AES encryption round using the composite/tower field for the S-box (SubBytes) and the native Galois field for the remaining (linear) operations. The internal architecture of each S-box corresponds to the one shown in Fig. 1.

## 3   Application of the Isomorphism Also to Other AES Encryption Operations

To apply the isomorphism of the composite/tower field to the other encryption round operations, the *MixColumns* transformation can be expressed as:

$$b_o = 2 \cdot (b_{i_0} \oplus b_{i_1}) \oplus 1 \cdot (b_{i_1} \oplus b_{i_2} \oplus b_{i_3}) \qquad (2)$$

In Eq. 2, the output byte $b_o$ is generated by multiplying ($\cdot$) the inputs bytes $b_{i_j}$ by the coefficients 1 and 2 that corresponds, respectively, to the identity function, and the shift on the left by one position plus the XOR with *00011011* if the most significant bit of the multiplied byte is *1*. Both coefficients can be

expressed in the matrix form, respectively, as:

$$C_1 = \begin{bmatrix} 1\,0\,0\,0\,0\,0\,0\,0 \\ 0\,1\,0\,0\,0\,0\,0\,0 \\ 0\,0\,1\,0\,0\,0\,0\,0 \\ 0\,0\,0\,1\,0\,0\,0\,0 \\ 0\,0\,0\,0\,1\,0\,0\,0 \\ 0\,0\,0\,0\,0\,1\,0\,0 \\ 0\,0\,0\,0\,0\,0\,1\,0 \\ 0\,0\,0\,0\,0\,0\,0\,1 \end{bmatrix} \quad C_2 = \begin{bmatrix} 0\,1\,0\,0\,0\,0\,0\,0 \\ 0\,0\,1\,0\,0\,0\,0\,0 \\ 0\,0\,0\,1\,0\,0\,0\,0 \\ 1\,0\,0\,0\,1\,0\,0\,0 \\ 1\,0\,0\,0\,0\,1\,0\,0 \\ 0\,0\,0\,0\,0\,0\,1\,0 \\ 1\,0\,0\,0\,0\,0\,0\,1 \\ 1\,0\,0\,0\,0\,0\,0\,0 \end{bmatrix} \tag{3}$$

Exploiting Eqs. 3 and 2 can be reformulated as:

$$b_o = C_2 \cdot (b_{i_0} \oplus b_{i_1}) \oplus C_1 \cdot (b_{i_1} \oplus b_{i_2} \oplus b_{i_3}) \tag{4}$$

From a mathematical point of view, the bytes $b_{i_j}$ in Eqs. 2 and 4 can be expressed as $b_{i_j} = (A \cdot M^{-1}) \cdot b'_{i_j}$, where $b'_{i_j}$ is the isomorphic multiplicative inverse. Hence, the isomorphic byte at the output of the *MixColumns* can be computed as:

$$\begin{aligned} b'_o = M \cdot b_o &= M \cdot C_2 \cdot (b_{i_0} \oplus b_{i_1}) \oplus M \cdot C_1 \cdot (b_{i_1} \oplus b_{i_2} \oplus b_{i_3}) \\ &= \underbrace{M \cdot C_2 \cdot A \cdot M^{-1}}_{C'_2} \cdot (b'_{i_0} \oplus b'_{i_1}) \oplus \underbrace{M \cdot C_1 \cdot A \cdot M^{-1}}_{C'_1} \cdot (b'_{i_1} \oplus b'_{i_2} \oplus b'_{i_3}) \end{aligned} \tag{5}$$

According to Eq. 5, the isomorphic output bytes of *MixColumns*, $b'_o$, can be directly obtained from $b'_{i_j}$ by implementing the matrix-vector multiplications with $C'_1$ and $C'_2$. The other round operations do not require modification because the *ShiftRows* is just a byte re-ordering, and the *AddRoundKey* (i.e. an XOR) is invariant with respect to the isomorphism. This approach consents to eliminate the timing cost of the isomorphic mappings in Eq. 1, according to the architecture shown in Fig. 3. For this purpose, also the key derivation process has to be modified accordingly by using the isomorphic values of the *Rcon* constant, i.e. $Rcon(r)' = M \cdot Rcon(r)$, for $r = 1, 2, \ldots, 10$. Hence, the critical delay of the fully isomorphic AES round is:

$$t'_R = t_{MultInv} + t_{LinOp} + 2 \cdot t_{MUX} + t_{XOR} \tag{6}$$

If $t_{LinOp}$, the delay of merged linear operations (Eq. 5), is such that $t_{LinOp} < t_{Map} + t_{InvMap\|Aff} + t_{MixCol}$, then $t'_R < t_R$ (i.e. the frequency increases).

## 4   Results

The two architectural approaches described in Sects. 2 and 3 were implemented in SystemVerilog reproducing the state-of-the-art works that use composite/tower fields, i.e. [7,10,12,13]. It is to be noted that [7] already proposes the usage of isomorphism for linear operations, hence the counterpart with isomorphic mapping applied only to the S-box has been derived, and vice versa for

**Fig. 3.** AES round using composite/tower field for all the operations. The S-box in Fig. 2 is substituted by the only multiplicative inverse (block 1/x in Fig. 1), while the affine transformation (block Affine Trans. in Fig. 1) is merged with the MixColumns inside the block Linear Operations. The isomorphic mapping (Map) and inverse mapping (Map$^{-1}$) blocks are moved, respectively, to the beginning of the input data paths (key in, data in) and to the end of the output data path (data out).

**Table 1.** Approaches comparison for [7,10,12,13]. The second area value in column Area (*) is the area for the same maximum synthesis frequency of the corresponding architecture that uses the AES native field for linear operations.

| Refs. | Field for linear operations | Maximum frequency (GHz) | Area (Gate equivalent, GE) (kGE) | Area efficiency (GHz/kGE) |
|---|---|---|---|---|
| [7] | Native AES field | 2.83 | 14.07 | 0.201 |
| | Isomorphic field | 3.06 | 14.78 (13.27*) | 0.207 |
| [10] | Native AES field | 2.83 | 13.79 | 0.206 |
| | Isomorphic field | 3.10 | 14.75 (12.88*) | 0.210 |
| [12] | Native AES field | 2.69 | 16.61 | 0.162 |
| | Isomorphic field | 3.08 | 18.12 (15.28*) | 0.170 |
| [13] | Native AES field | 2.80 | 17.67 | 0.158 |
| | Isomorphic field | 2.94 | 18.01 (16.87*) | 0.163 |

[10,12,13]. The hardware circuits were synthesized on a 7nm standard-cell technology included in the logic product kit SCH300MCPP64 from TSMC process CLN07FF41001, in the PVT corner slow process, 0.90 V and 125 °C. Table 1 reports the synthesis results.

Referring to results in Table 1, the usage of isomorphic mapping also for the linear operations of the AES encryption round is advantageous. In each case, this approach allows reducing the area for the same synthesis frequency and increasing the maximum supported frequency with a low area cost, which however leads to an overall improvement of the efficiency. For instance, focusing on the experiment based on the work proposed in [10], the usage of the isomorphic mapping only in the composite/tower field S-box (case Native AES field) achieves a maximum frequency of 2.83 GHz at the cost of 13.79 kGE, i.e. an efficiency in terms of frequency per area of 0.206 GHz/kGE. If extending the same isomorphic mapping also to the remaining (linear) operations of the AES algorithm (case Isomorphic field), one effect is that for the same synthesis fre-

quency of 2.83 GHz the area consumption is reduced to 12.88 kGE (about the 6.6%), i.e. the area value indicated between the round brackets and the symbol *. This would correspond to an improved area efficiency of $\frac{2.83 \text{ GHz}}{12.88 \text{ kGE}} \approx 0.220$ GHz/kGE. On the other hand, another effect is the reduction of the critical path of the round according to Fig. 3 and Eq. 6, therefore the possibility to increase the maximum synthesis frequency which rises to 3.10 GHz at the cost of 14.75 kGE. This gives again an improved efficiency of 0.210 GHz/kGE.

The same effects can be found in every experiment conducted, hence the overall results suggest that the implementation of all the AES operations on isomorphic fields leads to more efficient hardware solutions, improving both the area consumption and the maximum supported frequency.

## 5    Conclusions

This work presents the investigation of the usage of composite/tower fields in the AES algorithm. To the best of our knowledge, this is the first work that analyzes this aspect in a systematic fashion clearly pointing out how to implement the linear AES round operations on isomorphic fields. Indeed, the provided mathematical analysis and the highlighted correspondence with the hardware architectures constitute also a guideline for hardware designers of AES modules. In addition, the systematic approach used in this work allows us to easily extend the analysis (and the hardware implications) to the inverse transformations of the AES decryption algorithm for implementing decryption-only modules or encryption/decryption modules, and to AES modules supporting also (or only) 192-bit keys (12 rounds) and (or) 256-bit keys (14 rounds).

Future works will include the evaluation of the effects on the resistance to the Side-Channel attacks due to the application of the isomorphic mapping to the linear operations of AES, according to the methodology presented in [14].

## References

1. NIST (2001) Advanced encryption standard (AES). Federal Information Processing Standards (FIPS) publication 197
2. Singha TB, Palathinkal RP, Ahamed SR (2023) Securing AES designs against power analysis attacks: a survey. IEEE Internet Things J
3. Nannipieri P, Di Matteo S, Baldanzi L, Crocetti L, Zulberti L, Saponara S, Fanucci L (2022) VLSI design of advanced-features AES cryptoprocessor in the framework of the european processor initiative. IEEE Trans Very Large Scale Integr (VLSI) Syst 30(2):177–186. https://ieeexplore.ieee.org/document/9631958

4. Nannipieri P, Crocetti L, Di Matteo S, Fanucci L, Saponara S (2023) Hardware design of an advanced-feature cryptographic tile within the european processor initiative. IEEE Trans Comput

5. Carnevale B, Falaschi F, Crocetti L, Hunjan H, Bisase S, Fanucci L (2015) An implementation of the 802.1AE MAC security standard for in-car networks. In: 2nd IEEE world forum on internet of things (WF-IoT). IEEE, pp 24–28

6. Nannipieri P, Baldanzi L, Crocetti L, Di Matteo S, Falaschi F, Fanucci L, Saponara S (2022) CRFlex: a flexible and configurable cryptographic hardware accelerator for AES block cipher modes. In: Applications in electronics pervading industry, environment and society (APPLEPIES 2021). Springer, pp 31–38

7. Ueno R, Morioka S, Miura N, Matsuda K, Nagata M, Bhasin S et al (2019) High throughput/gate AES hardware architectures based on datapath compression. IEEE Trans Comput 69(4):534–548

8. Ueno R, Morioka S, Homma N, Aoki T (2016) A high throughput/gate AES hardware architecture by compressing encryption and decryption datapaths—Toward efficient CBC-mode implementation. In: 18th international conference on cryptographic hardware and embedded systems (CHES). Springer, pp 538–558

9. Ueno R, Homma N, Sugawara Y, Nogami Y, Aoki T (2015) Highly efficient GF $(2^8)$ inversion circuit based on redundant GF arithmetic and its application to AES design. In: 17th international conference on cryptographic hardware and embedded systems (CHES). Springer, pp 63–80

10. Reyhani-Masoleh A, Taha M, Ashmawy D (2018) Smashing the implementation records of AES S-box. IACR Trans Cryptogr Hardw Embedd Syst 298–336

11. Reyhani-Masoleh A, Taha M, Ashmawy D (2019) New low-area designs for the AES forward, inverse and combined S-boxes. IEEE Trans Comput 69(12):1757–1773

12. Gaded SV, Deshpande A (2019) Composite field arithematic based S-Box For AES algorithm. In: 3rd international conference on electronics, communication and aerospace technology (ICECA). IEEE, pp 1209–1213

13. Kumar TM, Reddy KS, Rinaldi S, Parameshachari BD, Arunachalam K (2023) A low area high speed FPGA implementation of AES architecture for cryptography application. Electronics 10(16)

14. Crocetti L, Baldanzi L, Bertolucci M, Sarti L, Carnevale B, Fanucci L (2019) A simulated approach to evaluate side-channel attack countermeasures for the advanced encryption standard. Integration 68:80–86

# Speeding Up Non-archimedean Numerical Computations Using AVX-512 SIMD Instructions

Lorenzo Fiaschi, Federico Rossi[(✉)], Marco Cococcioni, and Sergio Saponara

University of Pisa, Largo Lucio Lazzarino 1, 56122 Pisa, Italy
{lorenzo.fiaschi,federico.rossi}@ing.unipi.it,
{marco.cococcioni,sergio.saponara}@unipi.it

**Abstract.** This work presents the acceleration of a Bounded Algorithmic Number (BAN) library exploiting vector instructions in general-purpose processors. With the use of this encoding, it is possible to represent non-Archimedean numbers that are not only finite (like real numbers) but also infinite or infinitesimal. The tremendous growth in non-Archimedean numerical computations over the past 20 years and the resulting applications spurred this study's development. Enabling acceleration of BANs processing can significantly increase the throughput of non-Archimedean numerical computations, enlarging the spectrum of possible applications to industrial and real-time ones.

**Keywords:** Non-archimedean fields · Alpha theory · Bounded Algorithmic Number (BAN) · Single Instruction Multiple Data (SIMD) · AVX-512 instruction set

## 1 Introduction

Numerical non-Archimedean computations have been pioneered by Sergeyev and his Grossone Methodology [1], and allow for the use of infinitely large and infinitely small numbers in machines, other than finite ones as usual. From their advent, numerous applications benefited, especially in the domain of multi-objective optimisation, e.g., linear programming [2,3], quadratic programming [4], evolutionary algorithms [5], game theory [6], artificial intelligence [7,8], etc.

The reference framework of this study is the non-Archimedean model built upon the Alpha Theory [9], which introduces the set of Euclidean numbers and their associated numerical encoding called Bounded Algorithmic Number (BAN) format. The BAN encoding is a fixed length representation, guaranteeing that any operation involving two BANs outputs a result that occupies the same memory as the operands, exactly as happens with computations between 32-bit IEEE 754 floats, where the result of addition, subtraction, addition, and division is again a 32-bit float. However, since the Euclidean numbers form a superset of the

real ones, their numerical representation is heavier than the one of floats, making the processing of BANs cumbersome. CPU vectorization can be exploited to optimise computations since BANs encoding can be implemented through multiple fixed-length coefficients, thus fitting them inside vector registers and efficiently computing operations in a single clock cycle.

Using vector instructions had already proved to be a good solution for handling non-native data types that do not have hardware acceleration [10,11]. In this paper, we present an optimisation of a C++ BAN library called *BANcpp* [12]. The goal is to optimise the library to leverage CPU vectorization when dealing with BAN coefficients. In particular, we focused on BAN numbers with eight 64-bit coefficients and 512-bit vector instruction sets (e.g., Intel AVX-512). We also present a benchmark application that consists of several iterations of a non-Archimedean optimisation problem. We evaluate the goodness of automatic vectorization as a baseline model and then we enhance the automatic vectorization whenever the compiler fails to automatically optimise the code.

The paper is organised as follows: (i) Sect. 2 briefly introduces Alpha Theory and the Euclidean numbers; (ii) Sect. 3 details the BAN format; (iii) Sect. 4 explains the choices made to implement the enhanced vectorization of the library; (iv) Sect. 5 shows the application benchmark and the results obtained in terms of timing performance and throughput.

## 2    Alpha Theory and The Euclidean Numbers

Alpha Theory reference set of non-Archimedean numbers is indicated with the symbol $\mathbb{E}$ and it is called the set of $\alpha$-Euclidean numbers, or, in brief, just Euclidean numbers. The peculiar name of the theory comes from the definition of a reference infinite value within $\mathbb{E}$ and it is indicated by the symbol $\alpha$. Then, any Euclidean number can be represented as a function of $\alpha$, and only functions of $\alpha$ are numbers in $\mathbb{E}$, which guarantees that the Euclidean numbers and the mathematical operations among them behave according to their counterparts in $\mathbb{R}$, i.e., commutative operations continue to be commutative, differentiable functions are still differentiable, etc. For instance, the following are all Euclidean numbers:

$$\alpha^3, \qquad \frac{1}{\alpha}, \qquad \frac{1}{\alpha^2} - e^\alpha, \qquad -\frac{1}{2^\alpha}, \qquad -\ln\left(\frac{1}{\alpha}\right). \qquad (1)$$

As opposed to Archimedean Mathematics the concepts of infinite and infinitesimal numbers are sharply defined rather than vague concepts.

**Definition 1.** Given $\xi \in \mathbb{E}$, then

– $\xi$ is infinite $\iff \forall\, n \in \mathbb{N}, |\xi| > n$
– $\xi$ is finite $\iff \exists n \in \mathbb{N}, \frac{1}{n} < |\xi| < n$
– $\xi$ is infinitesimal $\iff \forall n \in \mathbb{N}, |\xi| < \frac{1}{n}$.

Therefore, in (1) the first and fifth numbers are positive and infinite, the third is negative and infinite, the second is positive and infinitesimal, while the fourth is negative and infinitesimal. A more detailed presentation of Alpha Theory and the set $\mathbb{E}$ can be found in [13].

## 3    The BAN Format

The BAN encoding consists in a finite length representation for Euclidean numbers; however, as for IEEE 754 floating point numbers, it cannot represent the whole set $\mathbb{E}$ because it would require infinitely many binary possible representations, i.e., an infinite computer memory. Any Euclidean number compliant with the following representation can be represented as a BAN:

$$\xi = \sum_{i=1}^{L} r_i \alpha^{p-i},$$

where $L \in \mathbb{N}$ is the encoding length, $r_i \in \mathbb{R}$ and $p \in \mathbb{Z}$. Changing perspective, one can define a BAN as a Euclidean number that can be represented by a linear combination of $L$ subsequent integer powers of $\alpha$.

From the very first glimpse, one may notice that the BAN representation of a Euclidean number is very similar to the one of a polynomial, suggesting how cumbersome can be to execute algebraic operations between BANs. An example of addition and multiplication between BANs with $L = 3$ follows:

$$(3.2\alpha^2 - 0.5\alpha + 1.4) + (0.2\alpha + 1 - 1.5\alpha^{-1}) = 3.2\alpha^2 - 0.3\alpha + 2.4 - 1.5\alpha^{-1}$$

$$(3.2\alpha^2 - 0.5\alpha + 1.4) \times (0.2\alpha + 1 - 1.5\alpha^{-1}) = 0.64\alpha^3 + 3.1\alpha^2 - 5.02\alpha - 2.15 - 2.1\alpha^{-1}$$

Both computations output a result that is not a BAN since it requires more than three consecutive powers of $\alpha$ to be represented. Therefore, there is the need to approximate them considering only the first three highest powers of $\alpha$, the most significant ones somehow, that is executing a truncation of the result. Below we report the numerical execution of the previous two operations, along with the division, realized by a software simulator of a BAN Processing Unit (BPU, whose hardware design has been recently proposed in [12]). The latter manipulates and outputs BANs in the normal form [13], i.e., in the standardized format which guarantees the uniqueness of the representation.

```
Operands: α^2(3.2 - 0.5η^1 + 1.4η^2) and α^1(0.2 + 1η^1 - 1.5η^2)

Sum: α^2(3.2 - 0.3η^1 - 2.4η^2)
Product: α^3(0.64 - 3.1η^1 - 5.02η^2)
Division: α^1(16 - 82.5η^1 + 539.5η^2)
```

## 4    Vectorization of BANcpp Library

We vectorized the *BANcpp* library by mixing two approaches: (i) leveraging the automatic vectorization offered by the compiler; (ii) enhancing vectorization manually whenever the automatic optimisation of the compiler fails. In particular, we needed to implement manual vectorization in the following cases:

– when a for-loop contains control-flow instructions on the BAN coefficients: due to possible branches and FPU exceptions the compiler refuses to insert vector instructions for these loops. The solution is to provide vectorization manually exploiting masked instructions based on comparisons. This is the case of comparison between BANs or checks on BAN values.
– when we have an outer and inner for-loop whose indexes are depending on each other. The solution is to implement vectorization manually leveraging the "geometry" of the problem. This is the case of multiplication between two BANs, which ends up in a one-dimensional convolution.

**Listing 1.** Manually vectorized $8 \times 8$ 1-D convolution for AVX512 SIMD.

```
void convmul(const double* a, const double* b, double* dst) {
    __m512 va0 = _mm512_set1_pd(a[0]);
    ...
    __m512 va7 = _mm512_set1_pd(a[7]);

    __m512 vb07 = _mm512_loadu_pd(&b[0]);
    __m512 vb17 = _mm512_mul_pd(vb07,masked1);
    ...
    __m512 vb77 = _mm512_mul_pd(vb07,masked7);

    __m512 va0b = _mm512_mul_pd(va0,vb07);
    __m512 va1b = _mm512_slide(_mm512_mul_pd(va1,vb17),1);
    __m512 va2b = _mm512_slide(_mm512_mul_pd(va2,vb27),2);
    ...
    __m512 va7b = _mm512_slide(_mm512_mul_pd(va7,vb77),7);

    __m512 accr = _mm512_add_pd(va0b,va1b);
    accr = _mm512_add_pd(accr,va2b);
    ...
    accr = _mm512_add_pd(accr,va7b);
    _mm512_storeu_pd(&dst[0],accr);
}
```

## 5   Benchmark Application and Results

The problem used as a benchmark in this study is one of the first ever used for testing and showing the efficacy of non-Archimedean numerical computations, namely Kite [3]. It consists of a bi-objective lexicographic linear programming problem, i.e., an optimisation problem of two linear functions, ordered by strict priority, over a linearly defined domain. To solve the problem, we adopted a Simplex-like non-Archimedean algorithm [3], precisely tailored to this type of task. To make it more realistic, we wrapped the problem within the I-Big-M framework [14], which adds a third objective to generalize the optimisation to the case of unknown starting feasible basis.

We ran the benchmark application for $10^5$ steps with both the auto-vectorized (namely, *baseline*) and the enhanced-auto-vectorized (namely *enhanced*) versions of the BAN library, collecting the time spent for each iteration. We smoothed the data by a 200-steps moving average window and computed a least square fit to plot the metric trends for the average time spent during an iteration and the average throughput (in terms of iterations per second). The benchmark was run on an Intel Xeon Gold 6238R processor running at 2.2 GHz with eight 64-bit BAN coefficients. Figure 1 shows the comparison between the two versions in terms of average time spent per iteration and overall throughput (iterations per second). Mean value and standard deviation of the two are reported in Table 1.



**Fig. 1.** Comparison between time spent for each iteration (**left**) and throughput (iterations per second, **right**) in the two different versions of the BAN library with the associated fitted curve.

**Table 1.** Mean value and standard deviation over $10^5$ iterations of the benchmark application, 8 64-bit BAN Coefficient with AVX-512.

|                          | Time (ns, $\times 10^4$) | Throughput (iter/s, $\times 10^5$) |
| ------------------------ | ------------------------ | ---------------------------------- |
| Baseline (Non-vector)    | $9.44 \pm 1.11$          | $1.07 \pm 0.87$                    |
| Baseline (Vector)        | $5.18 \pm 0.21$          | $1.93 \pm 0.04$                    |
| Enhanced                 | $4.68 \pm 0.19$          | $2.14 \pm 0.05$                    |

## 6    Conclusions

In this work, we presented the acceleration of a C++ library for Bounded Algorithmic Numbers (BAN) exploiting vector instructions, testing it on a non-Archimedean optimisation benchmark. The results showed how manually enhancing the automatic vectorization produced by the compiler can improve the performance of such applications even without complete hardware support for BANs. We have found that the performance of compiler automatic vectorization is significantly inferior to that achieved by manually optimising which intrinsics to use (and in which order) and how to load the information (again, in which order and according to which scheme). This can be helpful for the community of compiler developers too since it means that there is room for improving the compilers for handling the specific use case tackled in this work.

# References

1. Sergeyev YD (2017) Numerical infinities and infinitesimals: Methodology, applications, and repercussions on two Hilbert problems. EMS Surv Math Sci 4(2):219–320
2. De Cosmis S, De Leone R (2012) The use of grossone in mathematical programming and operations research. Appl Math Comput 218:8029–8038
3. Cococcioni M, Pappalardo M, Sergeyev YD (2018) Lexicographic multi-objective linear programming using grossone methodology: theory and algorithm. Appl Math Comput 318:298–311
4. Fiaschi L, Cococcioni M (2022) A non-archimedean interior point method and its application to the lexicographic multi-objective quadratic programming. Mathematics 10(23):4536
5. Lai L, Fiaschi L, Cococcioni M, Deb K (2021) Solving mixed pareto-lexicographic many-objective optimization problems: the case of priority levels. IEEE Trans Evol Comput 25:971–985
6. Cococcioni M, Fiaschi L, Lambertini L (2021) Non-Archimedean Zero Sum Games. J Comput Appl Math 393:113483
7. Astorino A, Fuduli A (2020) Spherical separation with infinitely far center. Soft Comput 24(23): 17 751–17 759
8. Cavoretto R, De Rossi A, Mukhametzhanov MS, Sergeyev YD (2021) On the search of the shape parameter in radial basis functions using univariate global optimization methods. J Global Optim 79(2):305–327
9. Benci V, Di Nasso M (2018) How to measure the infinite: mathematics with infinite and infinitesimal numbers. World Scientific, Singapore
10. Cococcioni M, Rossi F, Ruffaldi E, Saponara S (2020) Fast deep neural networks for image processing using posits and ARM scalable vector extension. J R-Time Image Process 17(3):759–771 Jun
11. Cococcioni M, Rossi F, Ruffaldi E, Saponara S (2021) Faster deep neural network image processing by using vectorized posit operations on a RISC-V processor. In: Kehtarnavaz N, Carlsohn MF (eds)Real-time image processing and deep learning 2021, vol 11736. International Society for Optics and Photonics. SPIE, p 1173604
12. Rossi F, Fiaschi L, Cococcioni M, Saponara S (2023) Design and FPGA synthesis of BAN processing unit for non-archimedean number crunching. In: Berta R, De Gloria A (eds) Applications in electronics pervading industry, environment and society. Springer Nature Switzerland, Cham, pp 320–325
13. Benci V, Cococcioni M, Fiaschi L (2022) Non-standard analysis revisited: an easy axiomatic presentation oriented towards numerical applications. Appl Math Comput 32(1):65–80
14. Cococcioni M, Fiaschi L (2021) The big-M method using the infinite numerical M. Optim Lett 15:2455–2468

# A 0.94 V Dynamic Bias Double Tail Comparator for High-Speed Applications in 5 nm Technology

Valentina Marazzi[1](✉), Enrico Monaco[2], Claudio Nani[2], and Danilo Manstretta[1]

[1] University of Pavia, Pavia, Italy
valentina.marazzi02@universitadipavia.it
[2] Marvell, Viale Della Repubblica 38, Pavia, Italy

**Abstract.** The necessity of efficient comparators in modern ADCs calls for new two-stages topologies that overcome the Strongarm limitations in terms of common-mode offset and gain dependency. The latest fashion is represented by dynamic bias integrators coupled with low power and low noise latches. The dynamic bias significantly reduces the overall power consumption, a feature which the second stage has to maintain; also, a CMOS implementation is usually adopted to gain more robustness. Starting from the comparator with dynamic floating inverter amplifier [1], a new structure has been derived and explained. The architecture has been simulated in 5 nm FinFET technology and compared to the state-of-the-art for very stringent power, noise and speed targets. Nonetheless, the proposed topology matches various applications thanks to its multiple degrees of freedom which allow the designer plenty of room for further improvements. The continuously scaling technology will favour this CMOS dynamic bias implementation even more.

**Keywords:** StrongARM · Dynamic Bias · Comparator

## 1 Introduction

### 1.1 Analog-to-Digital Converters (ADCs) for Serial Links

In modern communications, speed and efficiency are the key features and the motivations behind the present work. Nowadays, serial data are transferred at 56 Gb/s [2] or even higher rates. This requires extensive equalization that can only be achieved in the digital domain. As a result, high-speed ADCs are key components of serial data receivers. Such great speed may be obtained only through at least 64 time-interleaved Successive Approximation Register (SAR) converters. The single SAR should be carefully designed as its power, speed and resolution greatly affects the overall performance. In a typical SAR the most energy-hungry block is the comparator [3]. The aim of this article is to present a power optimized comparator for high-speed applications.

**Table 1.** Comparator targets.

| Targets | | |
|---|---|---|
| Decision Time | $t_{DEC}$ | $\approx 10ps$ |
| Charge in a period | Q | $<2fC$ |
| Noise | $N_i$ | $<2$ mV |
| Integrator Gain | $A_V$ | $\in [6, 9]$ |



**Fig. 1.** **a** Dynamic Floating Inverter Amplifier, **b** Strongarm, **c** Tang Comparator Outputs

## 1.2 Comparator Targets

As explained in Sect. 1.1, analog-digital converters represent a bottleneck regarding the overall performance of the entire receivers. This scenario imposes stringent targets in terms of speed [4], noise and energy consumption. For this work, the target specifications of Table 1 have been adopted. Combining these three metrics a Figure of Merit (1) can be defined, i.e. a single number which allows to rank different topologies based on their efficiency at first glance. This FoM [5] fairly evaluates the time needed to perform a comparison, the charge consumed in a clock period (T) and the noise produced in the process [1]. The smaller FoM will guarantee the higher efficiency.

$$\text{FoM} = t_{DEC} * Q * N_i^2 \tag{1}$$

$$t_{DEC} = t \,(\textit{differential output} = 0.5 * V_{DD}) \tag{2}$$

$$Q = I_{SUPPLY}\,dt \tag{3}$$

$$Nt = \sqrt{N_{o,integrator} + \frac{No, latch}{A^2_{V,integrator}}} \qquad (4)$$

**Technology**. It is important to remark that 5 nm technology has been used for this work. Such scaled technology node is necessary due to successive digital processing required for signal equalization. Topology considerations are linked to the adopted technology. Complementary (CMOS) structures are more suited for Finfet [6] with respect to much older technology nodes, e.g. 90nm or similar, where it would be preferable to adopt nmos implementations, such as the dynamic bias integrator of [7].

## 2 Comparator Topologies Overview

### 2.1 Strongarm (SA)

Strongarm [8] is one of the most popular comparator implementations thanks to its extreme simplicity and efficiency. As shown in Fig. 1b, the integrator ($M_{1,2}$) and the latch ($M_{3-6}$) are implemented with the minimum number of transistors and the whole architecture has only one stage. This compactness guarantees both velocity and reduced energy consumption, while the noise is further reduced by the cascode enhanced gain. Unfortunately, the one-stage approach presents major drawbacks as both the offset and the gain heavily depend on the input common mode. These flaws may be overcome by construction with a two-stage comparator [1].



**Fig. 2.** Integrator FoM for different sizing of $C_{RES}$.

## 2.2 Dynamic Floating Inverter Amplifier Comparator (Tang)

A two-stage approach allows to overcome the common mode dependency, whereas a CMOS implementation is preferable in the adopted technology.

The need for low power comparators requires dynamic biasing, a technique for reducing the energy consumption by improving the transistors transconductance over current ratio [7]. These three features point to the Tang comparator, shown in Fig. 1.

**Dynamic Floating Inverter Amplifier**. This first stage—Fig. 1a—is a dynamic bias CMOS integrator providing all the advantages of an independent power domain [1] thanks to the capacitive supply offered by the reservoir capacitor ($C_{RES}$). This capacitance is charged between ground and $V_{DD}$ during a reset phase, then it supplies the inverters ($M_{1p,1n}$ and $M_{2p,2n}$) through a charge sharing mechanism. Once the transistors size has been chosen, the sizing of the reservoir allows to control speed, power and gain [1]. The bigger the reservoir, the greater the gain and the charge consumption, whereas the time will be shortened. A simple way out of these trade-offs is to select a certain gain range and then to look for the minimum integrator figure of merit—Fig. 2, which accounts for the integration time, the charge and the noise of the first stage.

**Strongarm as Second Stage**. The Strongarm presents not only a latch, but also an additional amplifier ($M_{1,2}$), therefore Tang comparator relies on two integrators and a latch in only two stages. As stated in 2.1, this topology has excellent characteristics in terms of noise and power thanks to the gradual turning on of its transistors. SA performance as second stage may be further optimized by adding a small delay at the tail ($M_7$) so to allow the first integrator to produce a bigger gain. An effective delay should increase the integrator gain without an excessive degradation of the decision time, so as to improve the FoM. Hence, the required delay should be extremely low, in the order of fraction of ps. In this work such delay has been implemented by changing the threshold of the tail transistors. If an older technology is employed, the same delay can be realized through some logic ports on the clock path or resetting the drain of $M_7$ at $V_{DD}$.

## 3 Proposed Comparator

### 3.1 Double Tail Comparator

In this work a further optimization of Tang comparator is proposed, the Double Tail, Fig. 3. Once understood that delaying the Strongarm turning on brings significant advantages in terms of power, gain and noise, a delayed Tang comparator may seem the most straightforward approach. Nonetheless, the SA presents a major flaw when used as second stage: this circuit is reset at $V_{DD}$, but its input voltages ($O_1$, $O_2$) are expected to be around $0.5*V_{DD}$. Consequently, designing a reset network at half the supply seems more convenient in terms of noise and power.

A second improvement is the additional tail ($M_8$) which cuts the power consumption of $M_{5,6}$. Also this additional tail is delayed through the threshold voltage. Such delay should match the delay of the lower tail ($M_{TAIL}$) to maximize the power efficiency of this circuit, which is named "double tail" after these two fets turning on the structure.

**Fig. 3.** **a** Dynamic Floating Inverter Amplifier, **b** Double Tail Latch, **c** Outputs of Double Tail Comparator.

The main strength of the Double Tail is maintaining the SA advantages while improving the noise and power efficiency at the cost of a slight delay in the overall decision time, as shown in Fig. 3c. As previously explained, the delay should be kept reasonable compared to the total decision time. The easiest way to estimate the optimal delay is to look at the minimum FoM in a certain decision time range.

## 4   Conclusion

### 4.1   State of the Art Comparison

In order to validate the Double Tail topology in a fair comparison, the state of the art has been considered, redesigned in 5nm and simulated through Cadence Virtuoso to meet the given targets. The results are reported in Table 2. Note that the SA is still one of the most performing comparators thanks to its one-stage implementation.

### 4.2   Final Considerations

The Double Tail comparator stands out as the most efficient topology for the given targets as it employs a low power dynamic bias integrator coupled with a latch which is equally optimized for energy minimization.

This topology suits different applications thanks to the two main degrees of freedom, i.e. the reservoir capacitor in the first stage and the tunable delays in the second. The gain and noise performance may be further improved by stacking the integrator transistors or the tails.

**Table 2.** Comparator topologies simulated in 5 nm TSMC technology.

| Topologies | Decision time charge noise | FoM |
|---|---|---|
| Strongarm [8] | 9.63ps 1.596fC 2.58 mV | $102.31*10^{-33}CsV^2$ |
| Elzakker [9] | 10.89ps 2.411fC 2.37 mV | $147.48*10^{-33}CsV^2$ |
| Bindra [7] | 13.40ps 1.249fC 2.52 mV | $106.28*10^{-33}CsV^2$ |
| Tang [1] | 10.47ps 1.631fC 2.93 mV | $146.60*10^{-33}CsV^2$ |
| Delayed tang | 11.89ps 1.726fC 2.34 mV | $112.37*10^{-33}CsV^2$ |
| Double tail | 10.77ps 1.588fC 2.22 mV | $84.29*10^{-33}CsV^2$ |

## Future Work

A further clarification should be added: this work was meant to be built inside a given system where two clock signals (Ck and Nck) and a half supply reference ($0.5*V_{DD}$) had been already employed, therefore the double tail topology required no additional circuitry. Whenever such conditions are not met, the nmos solution mentioned in the "technology" paragraph would represent the most advisable choice.

## References

1. Tang X et al (2020) An energy-efficient comparator with dynamic floating inverter amplifier. IEEE J Solid-State Circ 55(4):1011–1022. https://doi.org/10.1109/JSSC.2019.2960485
2. Pisati M et al (2020) A 243-mW 1.25–56-Gb/s continuous range PAM-4 42.5-dB IL ADC/DAC-based transceiver in 7-nm FinFET. IEEE J Solid-State Circ 55(1):6–18. https://doi.org/10.1109/JSSC.2019.2936307
3. Li S, et al (2020) Low power SAR ADC design with digital background calibration algorithm. Symmetry 12(11):1757. https://doi.org/10.3390/sym12111757
4. Kull L, et al (2013) A 3.1mW 8b 1.2GS/s single-channel asynchronous SAR ADC with alternate comparators for enhanced speed in 32 nm digital SOI CMOS. In: 2013 IEEE international solid-state circuits conference digest of technical papers, San Francisco, CA, USA, pp 468–469. https://doi.org/10.1109/ISSCC.2013.6487818
5. Bindra HS et al (2022) A 174μVRMS input noise, 1 GS/s comparator in 22nm FDSOI with a dynamic-bias preamplifier using tail charge pump and capacitive neutralization across the latch. IEEE Int Solid-State Circ Conf (ISSCC) 2022:1–3. https://doi.org/10.1109/ISSCC42614.2022.9731728
6. Shang E, et al (2020) The factors that influence the effective mobility in 5 NM PMOS Finfet design. In: 2020 China semiconductor technology international conference (CSTIC), Shanghai, China, pp 1–3. https://doi.org/10.1109/CSTIC49141.2020.9282475
7. Bindra HS et al (2018) A 1.2-V dynamic bias latch-type comparator in 65-nm CMOS with 0.4-mV input noise. IEEE J Solid-State Circ 53(7):1902–1912. https://doi.org/10.1109/JSSC.2018.2820147
8. Kobayashi T, et al (1992) A current-mode latch sense amplifier and a static power saving input buffer for low-power architecture. In: 1992 symposium on VLSI circuits digest of technical papers, Seattle, WA, USA, pp 28–29. https://doi.org/10.1109/VLSIC.1992.229252
9. van Elzakker M et al (2010) A 10-bit charge-redistribution ADC consuming 1.9 μW at 1 MS/s. IEEE J Solid-State Circ 45(5):1007–1015. https://doi.org/10.1109/JSSC.2010.2043893

# A RISC-V Hardware Accelerator
# for Q-Learning Algorithm

Damiano Angeloni, Lorenzo Canese, Gian Carlo Cardarilli, Luca Di Nunzio,
Marco Re, and Sergio Spanò(✉)

Tor Vergata University of Rome, Via del Politecnico 1, 00133 Rome, Italy
spano@ing.uniroma2.it

**Abstract.** We propose a Q-Learning hardware accelerator for a RISC-V platform. In particular, our work focuses on the Klessydra processor. To the best of our knowledge, this is the first work in the literature that addresses this topic. We implemented the system on an AMD-Xilinx Zed-Board development board using a small amount of hardware resources and requiring a limited dynamic power of 1.528 W. The data we obtained are compatible with the future implementation of more accelerators on the same device to enhance the capabilities of the system. Compared to a standard software version of the algorithm, our accelerator allows a speed-up of ×36 in convergence time and an energy saving of ×34. The results obtained prove how our proposed system is suitable for high-speed and low-energy applications like Edge Machine Learning and embedded IoT systems.

**Keywords:** RISC-V · Hardware acceleration · Edge machine learning · Reinforcement learning · Q-learning

## 1 Introduction

In recent years, RISC-V processors have gained significant interest due to their open-source architecture, allowing for cost-effective customized solutions [1]. Unlike CISC structures, these cores combine a reduced and highly optimized instruction set with a simple instruction processing pipeline. These characteristics make the components more energy efficient and better suited for processing large amounts of data like in Edge Machine Learning applications [2–5].

The VLSI design research community has shown considerable interest in the hardware acceleration of Reinforcement Learning, as evidenced by its increasing popularity [6]. Among the notable algorithms in the field, Q-Learning [7] stands out with its simple tabular approach that relies on a Q-Matrix approach. In recent years, numerous hardware implementations of this algorithm have been proposed [8–11]. The latest survey [6] reveals that the work of Spanò et al. [11] is the standard de facto architecture for the implementation of Q-Learning.

---

Moreover, the authors proposed an automatic tool to generate the hardware system [12] that can be fruitfully used as an off-the-shelf IP.

Having regard to the above, we propose a Q-Learning hardware accelerator for a RISC-V platform. To our knowledge, this is the first work in the literature that addresses the topic.

## 2 Hardware Architecture

We introduce some information about the chosen RISC-V processor and the hardware architectures of the Q-Learning accelerator.

### 2.1 RISC-V Processor

Our work is based on the Klessydra processor [13]. The Klessydra family is fully compatible with the RISC-V ISA, specifically RV32I, and with the technical support of the PULP research group, it has been developed to seamlessly integrate with the open-source PULPino SoC. Compliance with PULPino is achieved at both the hardware and software levels. In fact, in addition to being a pin-to-pin alternative to one of the SoC cores (RI5CY), they have a memory map that aligns perfectly with the PULPino specifications [14].

Figure 1 shows the block diagram of the PULPino SoC. The system consists of a processor, two separate 32kB RAMs for data and instructions (following the Harvard architecture), and a 512B Boot ROM. For 32-bit data exchange with peripherals, the SoC utilizes an AXI bus connected to APB through a bridge. Additionally, there is a JTAG debug interface and an SPI slave port that operates on the communication bus as an AXI Master for loading the compiled code into memory.



**Fig. 1.** Top-level view of PULPino processor.

## 2.2    Hardware Accelerator

The update rule for the Q-Learning algorithm [7] can be described as follows:

$$Q_{new}(s_t, a_t) = (1 - \alpha)Q(s_t, a_t) + \alpha[r_t + \gamma_t \max_A Q(s_{t+1}, A)] \tag{1}$$

In the given equation, $Q$ denotes the value in the Q-Matrix, $s$ and $a$ refer to the states and actions respectively, $r$ represents the reward, $\alpha$ represents the learning rate, and $\gamma$ represents the discount factor. The variable $t$ represents the iteration index of the algorithm. The capital $A$ signifies the entire row of the Q-Matrix associated with a specific state.

The accelerator architecture is depicted in Fig. 2, with bold text and thick arrows used to represent vector signals. A detailed description of the architecture can be found in [11,15].



**Fig. 2.** Hardware architecture of the Q-Learning accelerator.

The input of the system consists of the aforementioned parameters, while the output corresponds to the entire row of the Q-Matrix that is associated with the current iteration $t$. This functionality is designed to enable the accelerator to align with the input requirements of a conventional Reinforcement Learning action-policy generator [16].

## 3    System Implementation

PULPino can be synthesized and executed on an FPGA, and although this type of implementation may not be as optimized as an ASIC implementation, it serves as a valid emulation environment. Within the PULPino-Klessydra project on GitHub [17], everything necessary for compilation on a Xilinx-AMD ZedBoard is already available, so we chose such a platform for our work.

To incorporate the design of the Q-learning algorithm into the PULPino-Klessydra architecture, it is necessary to package the VHDL code into an IP with

an AXI-Lite interface. However, since the PULPino architecture is equipped with only an AXI4-Full communication bus, it is necessary to convert the protocol to accommodate the IP within the system using a Xilinx AXI Protocol Converter IP. This approach allowed us to overcome compatibility issues without complicating the basic interface designed for the block. The choice is also supported by the fact that the peripheral does not require burst memory accesses, which is a feature of the full protocol.

## 3.1   Hardware Implementation Results

Figure 3 shows the utilization of resources of the system deployed to the target device. The evaluated dynamic power of the design is 1.528 W, considering a vectorless approach at 20 MHz clock. The data obtained are compatible for the future implementation of more accelerators on the same device to enhance the capabilities of the system.



**Fig. 3.** FPGA resources utilization.

## 4   Experimental Results

We tested the system considering a $4 \times 4$ Q-Matrix and considered 10 training processes. We ran the same experiments on the RISC-V using a software-only Q-Learning and a hardware-accelerated one. The policy generation is done for both cases in software. Since the two approaches give the same convergence values, in Fig. 4 we only compare the convergence time of all our tests.   Thanks to the hardware acceleration, the algorithm has an average speed-up rate of about $\times 36$ than the software approach. Thanks to the accelerator speed-up, the hardware approach allows for a $\times 34$ reduction of energy.

Figure 5 shows the average energy required to reach convergence.

## 5   Conclusions

We proposed a Q-Learning hardware accelerator for a RISC-V platform. In particular, our work was focused on the Klessydra processor. To the best of our knowledge, this is the first work in the literature that addresses this topic.

**Fig. 4.** Convergency time for software and hardware accelerated approaches.



**Fig. 5.** Average energy required to reach convergence for software and hardware accelerated approaches.

The system was implemented on an AMD-Xilinx ZedBoard development using a small amount of hardware resources and requiring a limited dynamic power of 1.528 W. The data obtained are compatible with the future implementation of more accelerators on the same device to enhance the capabilities of the system. Compared to a standard software version of the algorithm, our accelerator allows a speed-up of ×36 in convergence time and an energy saving of ×34. The results obtained prove how the proposed system is suitable for high-speed and low-energy applications like Edge Machine Learning and embedded IoT systems.

# References

1. Dörflinger A, Albers M, Kleinbeck B, Guan Y, Michalik H, Klink R, Blochwitz C, Nechi A, Berekovic M (2021) A comparative survey of open-source application-class risc-v processor implementations. In: Proceedings of the 18th ACM international conference on computing frontiers, pp 12–20

2. Ramírez C, Castelló A, Quintana-Orti ES (2022) A blis-like matrix multiplication for machine learning in the risc-v isa-based gap8 processor. J Supercomput 78(16):18051–18060
3. Kovačević N, Mišeljić D, Stojković A (2022) Risc-v vector processor for acceleration of machine learning algorithms. In: 2022 30th Telecommunications Forum (TELFOR). IEEE, pp 1–4
4. Ottavi G, Garofalo A, Tagliavini G, Conti F, Benini L, Rossi D (2020) A mixed-precision risc-v processor for extreme-edge dnn inference. In: 2020 IEEE computer society annual symposium on VLSI (ISVLSI). IEEE, pp 512–517
5. Ciccarella G, Giuliano R, Mazzenga F, Vatalaro F, Vizzarri A (2019) Edge cloud computing in telecommunications: case studies on performance improvement and tco saving. In: 2019 fourth international conference on fog and mobile edge computing (FMEC). IEEE, pp 113–120
6. Rothmann M, Porrmann M (2022) A survey of domain-specific architectures for reinforcement learning. IEEE Access 10:13753–13767
7. Watkins CJ, Dayan P (1992) Q-learning. Mach Learn 8(3):279–292
8. Liu X, Diao J, Li N (2022) A fpga-based accelerator implementation for path planning using q_learning algorithm. J Phys: Conf Ser 2245. IOP Publishing
9. Sahoo SS, Baranwal AR, Ullah S, Kumar A (2021) Memorel: a memory-oriented optimization approach to reinforcement learning on fpga-based embedded systems. In: Proceedings of the 2021 on Great Lakes Symposium on VLSI, pp 339–346
10. Meng Y, Kuppannagari S, Rajat R, Srivastava A, Kannan R, Prasanna V (2020) Qtaccel: a generic fpga based design for q-table based reinforcement learning accelerators. In: 2020 IEEE international parallel and distributed processing symposium workshops (IPDPSW). IEEE, pp 107–114
11. Spanò S, Cardarilli GC, Di Nunzio L, Fazzolari R, Giardino D, Matta M, Nannarelli A, Re M (2019) An efficient hardware implementation of reinforcement learning: the q-learning algorithm. Ieee Access 7:186340–186351
12. Canese L, Cardarilli GC, Di Nunzio L, Fazzolari R, Re M, Spanó S (2022) Automatic ip core generator for fpga-based q-learning hardware accelerators. In: International conference on applications in electronics pervading industry, environment and society. Springer, Berlin, pp 242–247
13. Cheikh A, Sordillo S, Mastrandrea A, Menichelli F, Scotti G, Olivieri M (2021) Klessydra-t: designing vector coprocessors for multithreaded edge-computing cores. IEEE Micro 41(2):64–71
14. Gautschi M, Schiavone PD, Traber A, Loi I, Pullini A, Rossi D, Flamand E, Gürkaynak FK, Benini L (2017) Near-threshold risc-v core with dsp extensions for scalable iot endpoint devices. IEEE Trans Very Large Scale Integr (VLSI) Syst 25(10):2700–2713 (2017)
15. Cardarilli GC, Di Nunzio L, Fazzolari R, Giardino D, Re M, Ricci A, Spano S (2022) An fpga-based multi-agent reinforcement learning timing synchronizer. Comput Electr Eng 99:107749
16. Cardarilli GC, Di Nunzio L, Fazzolari R, Giardino D, Matta M, Re M, Spanò S (2020) An action-selection policy generator for reinforcement learning hardware accelerators. In: International conference on applications in electronics pervading industry, environment and society. Springer, Berlin, pp 267–272
17. Klessydra: Klessydra/pulpino-klessydra: an open-source microcontroller system based on risc-v. https://github.com/klessydra/pulpino-klessydra

# Efficient Optimization of SFQ-Based Logic Circuits: Introducing a Novel Methodology for Performance and Design Enhancement

Laura Di Marino[1(✉)] , Francesco Fienga[1] , Vincenzo Romano Marrazzo[1] ,
Alessandro Borghese[1] , Giovanni Breglio[1] , Andrea Irace[1] ,
Federico Vittorio Lupo[2] , Oleg Mukhanov[3] , Marco Arzeo[2] ,
and Michele Riccio[1]

[1] Department of Electrical Engineering and Information Technologies, University of Naples
"Federico II", Naples, Italy
laura.dimarino@unina.it
[2] SeeQC-EU Srl, Naples, Italy
[3] SeeQC, Inc. Elmsford, New York, USA

**Abstract.** Single-Flux-Quantum (SFQ) logic is a digital electronic technology known for its very low-power consumption (nW-µW) and high operating frequency (up to 100 GHz). Like any other device, SFQ-based logic circuits suffer from manufacturing process issues, specifically concerning variations in the determined values of individual components such as the critical current of a Josephson junction and inductances. This leads to the need for a deep understanding of the circuit performances, its tolerance range and, furthermore, an optimization tool to improve it achieving a certain margin for each component. In this regard, the present article delves into the techniques and the development of a new design parameter optimization algorithm, whose main goal is to increase the critical margin of the circuit. By using such a simple and efficient technique, failures due to the fabrication are avoided and performance enhancement is achieved.

**Keywords:** SFQ logic · PSCAN2 · Optimization tool · Superconducting circuit

## 1 Introduction

In the last decades, SFQ technology has been widely studied thanks to its characteristics that make it promising in the field of computing and information processing. Even though SFQ is a binary logic, the way in which information is stored and transmitted changes: information is stored in the form of magnetic flux quantum ($\Phi_0$) and transmitted as SFQ voltage pulses with quantized areas [1]. The main component of such a technology is the Josephson junction (JJ), a two terminal device made of two superconducting electrodes separated by a thin non-superconducting layer [2]. Due to the presence of superconductive materials, the use of dedicated simulators specifically designed for superconducting circuits is necessary. In this study, simulators such as PSCAN2 or JSIM

[3] were employed to simulate and verify the behavior of SFQ-based logic circuits. To this purpose, simulations, SFQ-based circuits tests and their related optimization are discussed. At the state of the art, optimization algorithms are commonly used to bring SFQ circuit to work into their operating region. The approaches are multiple and very different from each other: some use the Monte-Carlo method to scatter the circuit parameter around the initial value [4], other use SQIR (small quantum intermediate representation) in Coq to describe a quantum circuit as a program [5]. While the latter uses a low-level quantum language, in this work, the optimization algorithm is developed with MATLAB and Python programming language. Moreover, this study looks at the problem mostly from a numerical point of view, by comparing simulations results with the desired ones. Hence, this approach provides an open, non-proprietary solution. In the following section simulations are presented, and the necessity for dimensioning optimization are explained. In Sect. 2, optimization algorithm methods are shown and finally, in Sect. 3, results and conclusions will follow.

### 1.1 SFQ-Based Circuit Design and Simulations

Every SFQ-based circuit design is based on the elementary concept of loop. A superconductive loop can merely be made with a superconductive wire, which being superconductive brings flux quantization in multiple of $\Phi_0$. However, in such a loop suppression and restoration of superconductivity takes too much time, but if the loop is interrupted with a JJ, switching is made way faster [1]. During the circuit's dimensioning process, the user makes a decision based on the desired behavior for the circuit itself. One can choose to implement a quantizing loop to store and retain information within it or a non-quantizing loop to allow the information transmission [6]. The simplest SFQ-based cell is the Josephson Transmission Line, used to transmit data between two cells, shown in Fig. 1 [7]. In this case, the loop containing J1-L2-J2 is dimensioned to be non-quantizing, since the main goal is to let the information pass through [8]. Once the dimensioning is done, it can be tested and simulated in the PSCAN2 simulator to verify the performance of the device under study, in terms of output phase or voltage.



**Fig. 1.** Josephson Transmission Line circuit

It is important to highlight that any manufacturing process suffers from fabrication tolerances thus leading to actual component whose values deviate from the designed ones [9]. In these terms, the margin of a circuit component refers to the difference between the actual operating point and the desired one. Therefore, it becomes vital to deeply

investigate the lower and upper margins of all the circuit elements (i.e., inductances, JJs and dc-current generators), and verify whether the target requirements are met. A component margin defines an interval in which the value of the component itself can vary not compromising the correct functioning of the circuit. Consequently, a margin analysis is an iterative process whose result are the margins of each circuit component. During this process the circuit is simulated several times, and each time margins are calculated for all the circuital elements. In the single simulation, one parameter is changed within a chosen interval, while the others are fixed. This method conducts simulations at a test point slightly above the initial value, as well as at lower and upper estimation points, subsequently identifying the operational region of the circuit [3].

## 2 Methods

The purpose of the iterative optimization algorithm is to bring all the circuital element composing the circuit, to a certain value of margin. To this end, a margin analysis is executed within the PSCAN2 simulator, and its result is saved in the form of a text file. Then, the algorithm proceeds with the optimization in terms of minimization of a multivariable function, $f(x)$, where $x$ is the project variables vector (i.e., inductances, critical current of JJs and bias currents). During each simulation, a margin for every single parameter is calculated. Then, starting from this information a vector *error* is created, whose elements represent the displacement of each margin component to the desired one. Once margins have been estimated the $f(x)$ function value is determined as follows:

$$f(x) = \|\frac{m_o - g(x)}{m_o}\|$$ (1)

where $m_o$ is the objective margin, while $g(x)$ is the margin obtained from the i-th simulation for each circuit component. The optimization function goes from a space $\mathbb{R}^N$ to $\mathbb{R}$, since it has $N$-project variables as input and returns a scalar as output, which is the $f(x)$ value (1). The iterative algorithm repeats itself until the $f(x)$ is minimized. As a case study, a MATLAB script has been developed. Every single iteration of the optimization is handled with the *fminsearch* function [10], which finds the minimum of unconstrained multivariable function using a derivative-free method. A vector of initial points, $x_0$, defined by user with the design parameters, is used to run the first simulation as a starting point for the fminsearch to minimize the function $f(x)$. After the first simulation, the $x_0$ vector will be updated by the fminsearch in this way:

- *Step* 1: several simulations will be executed changing one parameter at a time;
- *Step* 2: more parameters will be changed at the same time, and this process will be iterated until $f(x)$ is minimized.

The algorithm convergence strongly depends on the initial points vector which, being the result of an analytic solution, is very close to the circuit solution and so the convergence is guaranteed. In the following section, the workflow of the optimization algorithm is described.

## 2.1 Case Study in MATLAB Environment

The MATLAB script contains the definition of the $x_0$ vector, that will be used in the first simulation and as a starting point for the *fminsearch* function, whose job is to minimize, and hence find the zero of the $f(x)$, calculated using Eq. (1) [10]. The algorithm is an iterative one, which means that it has the following workflow, also reported in Fig. 2: anytime the fminsearch calls over the optimization function, a simulation is carried out with the values of the $x_0$ vector. Then, a margin analysis is done and from its result an estimation of the error is calculated. After that, the $f(x)$ value is returned and the whole process repeats itself. By setting the *fminsearch* options, user is able to define in which conditions the parameter optimization must end. To this end, TolX and TolFun [10] options have been used to give a termination tolerance respectively to the project variables vector and the current function value. Since in the manufacturing process there is a certain resolution for the values of currents and inductances, TolX should be properly set in relation to this. However, TolFun and TolX boundary are defined by the following inequalities:

$$|x_i - x_{i+1}| < TolX * (1 + |x_i|) \tag{2}$$

$$|f(x_i) - f(x_{i+1})| < TolFun * (1 + |f(x_i)|) \tag{3}$$

When both (2) and (3) are satisfied, the fminsearch algorithm stops.



**Fig. 2.** Block-diagram of the optimizer workflow

## 3 Results

The iterative optimization algorithm described in Sect. 2 was tested over the SFQ-based merger circuit (see Fig. 3) to achieve a 25% margin of tolerance. To perform the algorithm, a computer with an Intel-Core I5-8250U processor and 8 GB of RAM has been employed. By executing simulations of such a cell, it was found by the margin analysis that the desired margin was not accomplished, thus requiring an optimization cycle of the circuit. As a primary goal, all the circuit components need to fall within the required tolerance (functionality specification). As a second objective, such ranges should be identified as fast as possible. To achieve the best possible result, multiple optimizations

were performed by adjusting the options of the *fminsearch* function: *TolFun* and *TolX*. These options establish lower bounds for the current function value (which is the $f(x)$) and the step size for each $x_0$ element (12 elements in total, corresponding to 12 parameters to optimize).



**Fig. 3.** SFQ-based merger circuit.

Starting with an initial error value of 2.27, it was successfully reduced resulting in an error in the order of $10^{-16}$. Table 1 presents the margins of the merger components before and after the optimization process: in the "Pre-optimization Margins (%)" column, margins already to 25% are reported in italic.

**Table 1.** Merger circuit components margins before and after optimization: currents and inductances values are respectively in μA and pH.

| Parameters | Pre-optimization values | Pre-optimization margins | Optimized values | Optimized margins |
|---|---|---|---|---|
| J1 | 7 | *(12.41; 25)* | 7.42 | (50; 49.9) |
| J2 | 14 | *(25; 4.5)* | 12.72 | (38.2; 33.4) |
| J3 | 14 | *(25; 3.05)* | 12.76 | (38.2; 33.4) |
| J4 | 10 | *(6.2; 25)* | 10.9 | (50; 49.9) |
| J5 | 7 | *(10.9; 25)* | 7.54 | (50; 49.9) |
| J6 | 10 | *(7.46; 25)* | 10.99 | (50; 50) |
| L1 | 0.16 | *(25; 8.72)* | 0.154 | (50; 50) |
| L2 | 0.08 | *(25; 25)* | 0.064 | (50; 50) |
| L3 | 0.08 | *(4.5; 25)* | 0.091 | (39.8; 49.9) |
| L4 | 0.16 | *(25; 7.46)* | 0.147 | (50; 50) |
| L5 | 0.16 | *(25; 7.46)* | 0.149 | (50; 50) |
| I1 | 21 | *(21.6; 6.2)* | 22.7 | (28.9; 50) |

Once the optimizer performance is tested over the merger and has succeeded, to verify its robustness, other tests have been conducted over different SFQ-based cells

[1]. The equivalent results are reported in Table 2, along with the number of parameters optimized for each circuit, and the simulation time, which multiplied by the number of iterations returns the time it takes to optimize all of them to the desired margin value. Since its circuital complexity and the number of parameters growing, the counter circuit optimization time is bigger than the others, leading the fminsearch, and the optimization, to be time consuming. Of course, this aspect can be improved in future works.

**Table 2.** Simulation and optimization time of SFQ-based logic cells.

| Cells | Number of parameters | Simulation time (s) | Optimization time |
|---|---|---|---|
| JTL | 5 | 0.04 | 2 min |
| Splitter | 13 | 0.05 | 4 min |
| Merger | 12 | 0.07 | 8 min |
| T-flip flop | 21 | 0.46 | 11 min |
| Counter | 27 | 1.05 | 6 h |

## 4   Conclusions

In this paper, an automated iterative routine for the optimization of SFQ-based logics has been described. SFQ-based logic circuits can be simulated in the PSCAN2 simulator, but due to convergence problems affecting the software, there is the need for an optimization tool. So, in this work a novel method for design optimization of SFQ-based logic circuits has been introduced and tested, by showing a case study in the MATLAB environment and PSCAN2 simulations via Python. For all the considered test-cases, the optimization routine has allowed to find component values that met the circuit requirements, hence proving its usefulness in assisting the design of SFQ-based circuits. In future works, the objective is to speed up the optimization process even when used over larger circuits and bring it to a new release, in which the algorithm is fully developed in Python programming language.

## References

1. Paul Bunyk et al (2001), RSFQ technology: physics and devices. Int J High Speed Electron Syst 11(1):257–305
2. Likharev KK (1986) Dynamics of Josephson junctions and circuits. Gordon and Breach, New York
3. Fang ES et al. A Josephson integrated circuit simulator (JSIM) for superconductive electronics application. Ext Abstr 1989 Int Supercond Electron Conf ISEC 89:407–410, 198
4. Mori N et al (2001) A new optimization procedure for single flux quantum circuits. Phys C 357–360:1557–1560
5. Hietala K et al (2020) A verified optimizer for quantum circuits. arXiv:1912.02250v3
6. Van Duzer T et al (1981) Principles of superconducting circuits. Elsevier, New York

7. Likharev KK (1993) Rapid single-flux-quantum logic. In: Weinstock H, Ralston RW (eds) The new superconducting electronics. NATO ASI Series, vol 251. Springer, Dordrecht. https://doi.org/10.1007/978-94-011-1918-4_14

8. Semenov VK et al (1997) Digital-to-analog converter based on processing of SFQ pulses. Ext Abstr Int Supercond Electron Conf, PTB, Berlin, 320–322

9. Tinkham M (1996) Introduction to superconductivity, 2nd edn. McGraw-Hill, New York

10. Lagarias et al (1998) Convergence properties of the Nelder-Mead simplex method in low dimensions. SIAM J Opt 9(1):112–147

# A PUF-Based Secure Boot for RISC-V Architectures

Stefano Di Matteo[(✉)] , Luca Zulberti , Federico Cosimo Lapenna,
Pietro Nannipieri , Luca Crocetti , Luca Fanucci , and Sergio Saponara

Department of Information Engineering, University of Pisa, Via G. Caruso 16, Pisa,
Italy
stefano.dimatteo@dii.unipi.it

**Abstract.** Recently, there has been a growing interest in Physically
Unclonable Functions (PUFs). These electronic circuits possess several
key characteristics such as unpredictability and uniqueness that make
them particularly attractive for security applications. PUFs offer an
appealing solution for secure boot applications, providing a hardware-
based mechanism for generating unique cryptographic keys. These keys
can be used to encrypt the bootloader and operating system, thereby
enhancing security. In this paper, we propose an innovative, secure boot
scheme that leverages the functionality and characteristics of a PUF. Our
approach eliminates the need for physical storage of the encryption key
of the boot code, which enhances security and provides the possibility
of securely updating the firmware. We will present an architecture that
comprises essential components, along with a demo board on FPGA. The
demo board features a general-purpose 64-bit RISC-V-based system that
leverages the proposed PUF-based secure architecture, enabling secure
boot and firmware update functionalities.

**Keywords:** PUF · RISC-V · Secure boot · Hardware · FPGA

## 1 Introduction

In today's digital age, security has become an increasingly critical concern for
individuals, organizations, and governments. With the growing complexity and
sophistication of cyber threats, it is crucial to implement robust security mea-
sures to protect computer systems and sensitive data. Secure boot technology is
one of such security measure that is integrated into modern computer systems
[5,7,10]. Its primary aim is to verify the digital signatures of the bootloader and
operating system before loading, ensuring that only trusted software is executed
during the boot process. This is achieved by creating a chain of trust that begins
with the system firmware and continues through the bootloader and operating
system [8]. If any of these components are compromised, the secure boot process
will fail, and the system will not boot, ensuring that only reliable software is
executed.

Dedicated hardware can adopted to enhance the security of the chain of trust [8,9] and the performance of the entire system [2,3]. Physically Unclonable Functions (PUFs) are electronic circuits that have gained increasing interest in recent years due to their unique properties, making them ideal for security applications [11,12,14]. PUFs can generate unpredictable and unique responses to input stimuli, which is advantageous for creating cryptographic keys that are difficult to replicate or guess. The most attractive properties of a PUF, in fact, are the unpredictability of its response and the uniqueness of the response to each instance of the circuit. These properties result from the inherent variations in the manufacturing process fabricating the circuit. The response generated by the PUF can serve as a cryptographic key for encrypting the boot code, firmware, and operating system. Furthermore, another advantage of using the PUF response as a key is that it eliminates the need to preserve the secret key within the device.

Our paper presents a novel approach for enhancing security during boot-up, utilizing a security subsystem that leverages the functionality and unique characteristics of a PUF. By implementing this approach, we aim to address certain security concerns associated with secure boot, while also enabling secure firmware updates when necessary or in the event of a successful attack on a chip. This is crucial in preventing potential security breaches that could impact other chips of the same type. The paper is organized as follows: Chap. 2 explains our PUF-based secure boot concept and the hardware components needed to implement it. Chapter 3 presents a complete RISC-V-based system prototype composed of a general-purpose subsystem enhanced by a secure subsystem equipped with an FPGA-based Ring-Oscillator PUF to perform secure boot and secure firmware updates. Chapter 4 will draw the conclusions of this work.

## 2   Proposed PUF-Based Secure Boot

One of the classical approaches to verify the authenticity and integrity of the boot code is to include Original Equipment Manufacturer (OEM)-related information in the system. For instance, the public key of the system owner can be used to verify the ownership of the boot code at the startup; in this way, the owner's public key shall be embedded in a Read Only Memory (ROM), and it shall be used at the startup of the system to verify the signature of the boot code that has been previously signed with the corresponding private key. This approach implies that all devices use the same public key. As reported in [1], "storage of the public key for the root of trust can be problematic; embedding it in the on-SoC ROM implies that all devices use the same public key. This makes them vulnerable to class-break attacks if the private key is stolen or successfully reverse-engineered". The same concept can be applied to symmetric keys in case they are used to encrypt the boot code. Starting from these issues, the main goals of our secure boot concept are:

– Avoid the physical storage of secret keys: this eliminates the possibility of disclosure of secret keys stored in memories.

– Unique secret key for each physical device: if the secret key of a device is discovered, it cannot affect the security of the other devices as each key generated by the PUFs is unique.
– Possibility to modify OEM public key stored in the device: in case of disclosure of the OEM private key, the OEM should be able to update its public key stored in the device.
– Unpredictability of the secret keys for each physical device: this resolves issues related to the trustiness of the electronic manufacturers and supply chain.



**Fig. 1.** Concept architecture of the secure system. NVM is a Non-Volatile Memory. OTP is a One-Time Programmable memory. VPK is the Vendor's Public Key. FSB is the First-Stage Bootloader of the CPU.

The system we propose is illustrated in Fig. 1, and it is composed of three main parts:

– Secure Environment (S-Env): provides PUF-based encryption and authentication (hardware fingerprint) and can reset the Non-Secure part. We expect the need of an OTP memory to store the PUF challenge plus auxiliary information to retrieve the corresponding response (e.g. redundancy of correction codes, if needed).
– Non-Secure Environment (NS-Env): is a standard CPU-based SoC whose boot sequence is regulated by the Secure part.
– A Non-Volatile Memory (NVM): contains the encrypted Vendor Public Key (VPK) and the encrypted First-Stage Bootloader (FSB) for the CPU.

The S-Env exposes the following Secure Application Program Interface (SAPI) to the NS-Env, based on the status of the OTP memory:

– WriteNVM: this operation initializes the S-Env. It can be issued only if the OTP memory is not written.
– UpdateNVM: the NS-Env can request an update of the NVM employing Publick Key Authentication using the VPK. It can be issued only if the OTP memory has been written.

The boot sequence of the system depends on the S-Env status. As long as the S-Env is not initialised (fresh OTP memory), the Secure Element (SE) copies the content of the NVM into the Secure Memory (SMem) and wakes up the CPU. No verification is performed in this procedure, and the boot is not secure.



**Fig. 2.** Secure environment initialisation procedure

*Secure Environment Initialisation.* The vendor of the system can initialise the S-Env by performing the following procedure, illustrated in Fig. 2:

1. Vendor prepares the NVM with the VPK-FSB couple and a bootable code executed by the CPU to perform the initialisation sequence.
2. At system boot, the SE will perform a non-secure boot sequence (the OTP is fresh) by coping the NVM content at the boot address of the SMem and waking up the CPU.
3. The code executed by the CPU prepares the SMem with the VPK-FSB couple and the Boot Challenge (BC). The latter can be hard-coded into the initialisation code or determined at run-time. Then, the CPU issues a WriteNVM request to the SE.
4. After halting the CPU, the SE retrieve the Boot Key (BK) from the PUF by providing the BC.
5. The SE encrypts the VPK-FSB couple found in the SMem with the BK using an encryption algorithm (e.g. Advanced Encryption Standard) and stores it in the NVM with a hash digest used for integrity check during secure boot. Depending on the computational capabilities of the SE, the vendor may employ a Digital Signature Algorithm (DSA) or a Hash-based Message Authentication Code (HMAC) to ensure the integrity and authenticity of the FSB code.
6. Finally, the SE saves the BC into the OTP (with any other required parameters related to PUF implementation), and the S-Env can be considered initialised.

7. From now, the WriteNVM operation cannot be issued any more, and the vendor can ship the system with Secure-Boot enabled.



**Fig. 3.** Flow diagram describing the secure-boot sequence

*Secure Boot Sequence.* When the S-Env is initialised, the boot sequence includes the PUF-based verification of the NVM, illustrated in Fig. 3:

1. At system boot, the SE reads the BC from the OTP and retrieves the BK by challenging the PUF.
2. The SE reads the encrypted VPK-FSB couple from the NVM, decrypts it with the BK using an encryption algorithm (e.g. Advanced Encryption Standard) and verifies its integrity and authenticity using DSA or HMAC.
3. Finally, the VPK and the FSB are copied into the SMem, and the CPU is woken up.



**Fig. 4.** Flow diagram for updating the NVM

*Update NVM Sequence.* During normal operation, the CPU can issue an UpdateNVM operation to update the VPK or the FSB contained in the NVM, as illustrated in Fig. 4:

1. While the system is running, the CPU prepares a new VPK-FSB couple and signs it with the active private key.

2. The CPU store the signed content into the SMem and issues the UpdateNVM operation.
3. The SE halts the CPU and verifies the signature in the SMem using the active VPK.
4. If successful, the SE encrypts the new VPK-FSB couple and stores it in the NVM.
5. Finally, the CPU is woken up, and the new Secure Boot process is performed.

The proposed PUF-based Secure Boot procedure and SAPI avoid the physical storage of secret keys and allow to have unpredictable and unique secret keys for each physical device. In addition, the update sequence allows modifying the OEM VPK in case of private-key class-break attacks to limit the number of compromised devices.

## 3  Implementation on FPGA

To prototype the proposed PUF-based Secure Boot procedure, we implemented a System-on-Chip (SoC) on a Xilinx ZCU106 FPGA Board that emulates the architecture of Fig. 1.



**Fig. 5.** SoC implemented on a Zynq Ultrascale+ FPGA.

As illustrated in Fig. 5, the system is divided into a Non-secure Environment and a Secure Environment. The former (that emulates a general-purpose application system) includes a RISC-V CPU (i.e., the CVA6 64-bit processor [13]) connected through an AXI4 interconnect to its Main Memory (DDR4 controller) and the UART JTAG peripheral provided as Xilinx IPs. The latter implements the SE using a smaller RISC-V CPU (i.e., the CV32E40P 32-bit processor [4]) with its Program Memory. On-chip memory is used to emulate the SMem shared between the two environments and the OTP of Fig. 1. The VPK-FSB couple is stored encrypted on the external SD memory (NVM in Fig. 1). The PUF we implemented is a Ring Oscillator PUF (RO-PUF), generating a response of 325 bits using a challenge of 220 bits. To improve the reliability of the responses,

the output is corrected using an Error Correcting Code (ECC) based on the Reed-Muller code, and we used the VHDL implementation published in [6]. The Red-Muller ECC needs a helper generated during the initialisation procedure and stored on the OTP memory. To generate the AES-CBC-128 key, the SHA256 algorithm is applied to the PUF response and truncated to 128 bits. While to ensure the integrity and authenticity of the VPK-FSB couple, the ECDSA-256 algorithm is used.

## 4    Conclusions

Our research paper presents a novel implementation of a Secure Boot process, which leverages PUF technology to eliminate the requirement of storing secret keys in physical storage. PUF technology enables the generation of unique and unpredictable secret keys for each individual device. Furthermore, our updating process streamlines the ability to modify VPK in order to mitigate the effect of class-break attacks in case of disclosure of the vendor's private key. As a proof-of-concept, we implemented the proposed PUF-based Secure Boot procedure on a Xilinx ZCU106 FPGA board, which includes a 32-bit RISC-V-based secure environment exposing a Secure API for the 64-bit RISC-V application-class processor.

## References

1. ARM security technology—Building a secure system using Trust Zone technology (2005–2009). https://developer.arm.com/documentation/PRD29-GENC-009492/c
2. Baldanzi L, Crocetti L, Di Matteo S, Fanucci L, Saponara S, Hameau P (2019) Crypto accelerators for power-efficient and real-time on-chip implementation of secure algorithms. In: 2019 26th IEEE international conference on electronics, circuits and systems (ICECS), pp 775–778. https://doi.org/10.1109/ICECS46596.2019.8964731
3. Di Matteo S, Lo Gerfo M, Saponara S (2023) Vlsi design and fpga implementation of an ntt hardware accelerator for homomorphic seal-embedded library. IEEE Access 11:72498–72508. https://doi.org/10.1109/ACCESS.2023.3295245
4. Gautschi MEA (2017) Near-threshold risc-v core with dsp extensions for scalable iot endpoint devices. IEEE Trans Very Large Scale Integr (VLSI) Syst 25(10):2700–2713. https://doi.org/10.1109/TVLSI.2017.2654506
5. Haj-Yahya J, Wong MM, Pudi V, Bhasin S, Chattopadhyay A (2019) Lightweight secure-boot architecture for risc-v system-on-chip, pp 216–223 (2019). https://doi.org/10.1109/ISQED.2019.8697657, https://www.scopus.com/inward/record.uri?eid=2-2.0-85065164064&doi=10.1109%2fISQED.2019.8697657&partnerID=40&md5=dbd40305c40b567a1f6089561f3d8863 cited by: 17; All Open Access, Green Open Access

6. Liguori P, Reed-muller-decoder. https://github.com/piliguori/Reed-Muller-Decoder

7. Liu, Y., Briones, J., Zhou, R., Magotra, N.: Study of secure boot with a fpga-based iot device. vol. 2017-August, p. 1053 - 1056 (2017). https://doi.org/10.1109/MWSCAS.2017.8053108, https://www.scopus.com/inward/record.uri?eid=2-s2.0-85034084065\&doi=10.1109\%2fMWSCAS.2017.8053108\&partnerID=40\&md5=9cba3d9340f4b2807acc0ab560f9da3fcitedby:24

8. Nannipieri P, Crocetti L, Di Matteo S, Fanucci L, Saponara S (2023) Hardware design of an advanced-feature cryptographic tile within the european processor initiative. IEEE Trans Comput 1–14. https://doi.org/10.1109/TC.2023.3278536

9. Nannipieri P, Matteo SD, Baldanzi L, Crocetti L, Zulberti L, Saponara S, Fanucci L (2022) Vlsi design of advanced-features aes cryptoprocessor in the framework of the european processor initiative. IEEE Trans Very Large Scale Integr (VLSI) Syst 30(2):177–186. https://doi.org/10.1109/TVLSI.2021.3129107

10. Sabt M, Achemlal M, Bouabdallah A (2015)Trusted execution environment: What it is, and what it is not 1:57–64. https://doi.org/10.1109/Trustcom.2015.357, https://www.scopus.com/inward/record.uri?eid=2-s2.0-84967164163&doi=10.1109%2fTrustcom.2015.357&partnerID=40&md5=358ca116c45a82c981b7654c287b8ef6, cited by: 275; All Open Access, Green Open Access Access, Green Open Access

11. Shamsoshoara A, Korenda A, Afghah F, Zeadally S (2020) A survey on physical unclonable function (puf)-based security solutions for internet of things. Comput Netw 183:107593. https://doi.org/10.1016/j.comnet.2020.107593, https://www.sciencedirect.com/science/article/pii/S1389128620312275

12. Wang W, Chen Q, Yin Z, Srivastava G, Gadekallu TR, Alsolami F, Su C (2022) Blockchain and puf-based lightweight authentication protocol for wireless medical sensor networks. IEEE Internet Things J 9(11):8883–8891. https://doi.org/10.1109/JIOT.2021.3117762

13. Zaruba F, Benini L (2019) The cost of application-class processing: Energy and performance analysis of a linux-ready 1.7-ghz 64-bit risc-v core in 22-nm fdsoi technology. IEEE Trans Very Large Scale Integr (VLSI) Syst 27(11):2629–2640. https://doi.org/10.1109/TVLSI.2019.2926114

14. Zhang J, Qu G (2020) Physical unclonable function-based key sharing via machine learning for iot security. IEEE Trans Ind Electron 67(8):7025–7033. https://doi.org/10.1109/TIE.2019.2938462

# Machine Learning

# Neural Architecture for Tennis Shot Classification on Embedded System

Ali Dabbous[✉] , Matteo Fresta , Francesco Bellotti , and Riccardo Berta

Department of Electrical, Electronic and Telecommunication Engineering (DITEN), University of Genoa, Via Opera Pia 11a, 16145 Genova, Italy
{ali.dabbous,matteo.fresta}@edu.unige.it, {francesco.bellotti, riccardo.berta}@unige.it

**Abstract.** Data analysis has become a common practice in professional and amateur sport activities, to monitor the player state and enhance performance. In tennis, performance analysis requires detecting and recognizing the different types of shots. With the advances in microcontrollers and machine learning algorithms, this topic becomes ever more considerable. We propose a 1-D convolutional neural network (CNN) model and an embedded system based on Arduino-Nano system for real-time shot classification. The network is trained through a dataset composed of three different tennis shot types, with 6 features recorded by an inertial device placed on the racket. Results demonstrate that the proposed model is able to discriminate the tennis shots with high accuracy, also generalizing to different users. The network has been deployed on a low-cost Arduino nano 33 IoT model, with an inference time of 65 ms.

**Keywords:** 1-D CNN · Arduino Nano · Embedded device · Tennis shot

## 1 Introduction

Thanks to developments in microelectronics, wearable technology [1–3] has attracted a lot of attention for research and commercial applications, also in the sport field, where field data collection and processing is essential for identifying flaws and enhance skills. In tennis, this requires detecting the main shots, such as forehand, backhand and serve.

This paper focuses on the classification of various tennis shots. A lot of research work has been carried out in the field of tennis sport activity detection. Several studies have used various machine learning methods to tackle this classification problem [4–7]. For instance, Büthe et al. [8] created a pipeline to detect and categorize leg and arm movement, as well as execute gesture recognition for the firing arm using the longest common subsequence. In [9], Hazem and Al-Sadek employed an artificial neural network (ANN) to determine stroke type and forecast how accurate the player will play in subsequent plays based on statistics.

In all these works, data is recorded in the field and processed in the cloud or on the desktop. However, it would be important to process information online by using a low-cost, low-power embedded device capable of handling AI models, in order to have

real-time feedback to the athlete. Here the literature has much less examples, because the processing is not lightweight, while edge devices have strong limitations in terms of computational capabilities and memory availability.

The main goal of this paper is to propose a new model of neural architecture for real-time classification of tennis shots on a micro-controller directly attachable to a tennis racket. The system processes six inertial signal timeseries collected from an IMU sensor that represents the three axes of accelerometers and gyroscopes.

In literature, timeseries are frequently processed through feature engineering in order to provide suitable input to a machine learning model [10]. Deep-learning techniques, on the other hand, tend to extract features directly from the raw data (e.g., image pixels); however, this capability comes at the cost of significantly increase the number of model parameters, increasing memory needs, computing burden, and energy consumption [11]. This increase is critical for low-power edge devices. Therefore, we targeted the development of a deep learning architecture able to fulfill the strict hardware and cost limitations of edge devices.

Particularly, we propose a deep network featuring 1-D convolutional layers to classify tennis shots. 1-D convolutional neural networks (CNNs) can process 1-D time series. 1-D CNN can handle multi-dimensional input tensors, as in our case, by applying the convolution to each single dimension, independent of the other, producing a feature map for each channel. 1-D convolutions just involve scalar multiplications and additions, and have already proved to be particularly effective also in complex classification tasks (e.g., touch modalities [10, 12]). Thus, we consider them particularly promising for real-time classification of tennis shots.

The contributions of this paper are summarized as follow:

- Adopting 1-D CNNs to classify different sport tennis shots with high accuracy.
- Deploying the 1-D CNN classifier on a low-cost Arduino nano IOT microcontroller.
- Analyzing overall performance of the embedded system.

## 2   Experimental Setup

### 2.1   Dataset

The dataset employed in this work is made up of 6-axis time-series maps for three distinct tennis shots: forehand, backhand, and serve. To capture these shots, we asked a amateur tennis player to make several repetitions of the three types of shots, using his racket, to which an Arduino Nano IoT [13] device was connected. During each trial, the time-series data from the IMU (LSM6DS3) sensor was recorded. This time-series has six axes: x, y, and z accelerometers and x, y, and z gyroscopes. The total number of samples collected for all shots is 213 divided into 80% for training and 20% for validation. Additionally, 45 separate time-series samples were recorded for the three shots for testing purposes.

### 2.2   1-D Classifier Network Architecture

Figure 1 shows the architecture of the proposed neural network, which consists of two 1-D convolutional layers, interleaved by two max pooling layers, followed by a dense

layer, and a three neuron softmax output layer. The convolutional layers aim at extracting the classification features, while the max pool layers, with non-overlapping receptive field, aim at reducing dimensionality. We discuss the actual hyperparameter values (e.g., number of neurons and of filters) later in the next section describing the training, which is where the actual values were set.



**Fig. 1.** Tennis shot architecture. (Left) 1-D convolutional neural network. (Right) Arduino Nano 33 IoT for model deployment and real time classification.

## 2.3  Training Strategy

The first main challenge for the training consisted in maintaining consistency in the input shape of the 1-D CNN model, given the different lengths of the shot sequences. To this end, we employed Dynamic Time Warping (DTW) [14] techniques on the samples of each shot. This allowed us to clean up the dataset by identifying and excluding anomalous or poorly created samples. After applying DTW, we calculated the maximum number of time steps across all the samples, which was 20 in our case. For samples with fewer than 20 time-steps, we used padding techniques to maintain homogeneity. This entailed adding more time steps to ensure that all samples had the proper length. As a result, our 1-D CNN model's input shape was specified as (20, 6), where 6 denoted the number of features recorded by each time step.

The actual training procedure then starts with fine-tuning the classifier architecture's hyperparameters, through cross-validation on the training set. A grid of potential candidates was investigated, with a focus on solutions with few parameters and high computing efficiency. The process finally selected a network architecture with two 1-D convolutional layers. The first layer has 8 filters, whereas the second 16. The kernel size in both layers is 3. The Dense layer comprised 32 units. The Rectified Linear Unit (ReLU) activation function is applied to both the convolutional and dense layers. Nevertheless, the total number of trainable parameters was equal to 2,219.

The training strategy's efficacy was assessed by measuring the model's accuracy on the testing set, after 50 epochs.

## 2.4  Embedded System for Real Time Classification

We deployed the 1-D CNN on an Arduino Nano 33 IoT, with the goal to classify shots in real-time. The Arduino Nano 33 IoT is a cost-effective device that can be easily installed into any tennis racket. It is powered by a SAMD21 Cortex®-M0+32-bit low power ARM MCU with a 48 MHz CPU. The device includes 256 KB of flash memory and 32 kB of SRAM, and it operates at a voltage of 3.3 V.

We were able to export the trained model in the float 32 format suitable with the Arduino Nano 33 IoT after finishing the training of the 1-D CNN using Python and the TensorFlow library. The model was converted to the TFLite format, which is suitable for deployment on low-resource devices, including the Arduino Nano microcontroller. The embedded memory footprint of the proposed neural network is of 14 KB.

# 3  Experimental Results

To evaluate the model's performance, we utilized a test set in which each one of the 45 samples represents an individual shot: 15 forehand strokes (red colour in Fig. 2), 14 backhand (blue), and 16 serves (green). The trained model could correctly classify each shot in the test set.



**Fig. 2.** Tennis shots classification results. Red represents the Forehand shot, blue represents the Backhand shot and green represents the serve shot.

As a field-validation step, the model deployed on the Arduino Nano 33 IoT was evaluated with three distinct players, each of whom repeated three separate shots 15 trials. The complete system also included a threshold approach on the accelerometer measurements to assess whether there was movement and thus compute the classification.

Also in this case, the embedded model could detect and classify all the shots, as shown in Table 1. The inference time was 65 ms, which is suited for real-time performance. These findings show that the proposed network is effective in terms of accuracy, and, advancing the state-of-the-art work [8, 9], has been deployed on an embedded device, executing real time inference.

**Table 1.** Real time tennis shots classification.

|          | Forehand | Backhand | Serve | Inference time |
|----------|----------|----------|-------|----------------|
| Forehand | 45       | 0        | 0     | 65 ms          |
| Backhand | 0        | 44       | 1     |                |
| Serve    | 0        | 1        | 44    |                |

## 4  Conclusion

This paper presented the development of a 1-D CNN to distinguish between different types of tennis shots. The network is capable of discriminating among three types tennis shots, with high accuracy. Advancing the state of art (e.g., [8, 9]) the trained model, converted into a micro model using TFlite Libraries, works in real-time on an embedded device in the field. The memory footprint of the micro model is of 14 KB, which is acceptable for the vast majority of current microcontrollers. The inference time for detecting and classifying one shot after applying it by the athletes' in real time was equal to 65 ms, which is suited for real time performance, which can be important for a athlete's training. This preliminary work calls for further studies on time series analysis for sport activities. We will keep on resorting computational primitives that have a corresponding hardware implementation using low-power microelectronic devices, with the aim of designing compact and efficient end-to-end systems that can analyze sports activity signals in the field.

## References

1. Sakr F, Bellotti F, Berta R, De Gloria A (2020) Machine learning on mainstream microcontrollers. Sensors 20(9). https://doi.org/10.3390/s20092638
2. Coyle S, Morris D, Lau K-T, Diamond D, Moyna N (2009) Textile-based wearable sensors for assisting sports performance. In: 2009 Sixth International Workshop on Wearable and Implantable Body Sensor Networks, pp 307–311. https://doi.org/10.1109/BSN.2009.57
3. Patel S, Park H, Bonato P, Chan L, Rodgers M (2012) A review of wearable sensors and systems with application in rehabilitation. J NeuroEng Rehabil 9(1):21. https://doi.org/10.1186/1743-0003-9-21
4. Petkovic M, Jonker W, Zivkovic Z. Recognizing strokes in tennis videos using hidden Markov models
5. Kelly P, O'Connor NE (2012) Visualisation of tennis swings for coaching. In: 2012 13th international workshop on image analysis for multimedia interactive services, May 2012, pp 1–4. https://doi.org/10.1109/WIAMIS.2012.6226750
6. Connaghan D, Ó Conaire C, Kelly P, O'Connor NE (2010) Recognition of tennis strokes using key postures. In: IET Irish signals and systems conference (ISSC 2010), pp 245–248. https://doi.org/10.1049/cp.2010.0520
7. Shah H, Chokalingam P, Paluri B, Pradeep N, Raman B (2007) automated stroke classification in tennis. In: Kamel M, Campilho A (eds) Image analysis and recognition. Lecture notes in computer science, vol 4633. Springer Berlin Heidelberg, Berlin, Heidelberg, pp 1128–1137. https://doi.org/10.1007/978-3-540-74260-9_100

8. Büthe L, Blanke U, Capkevics H, Tröster G (2016) A wearable sensing system for timing analysis in tennis. In: 2016 IEEE 13th International Conference on Wearable and Implantable Body Sensor Networks (BSN), pp 43–48. https://doi.org/10.1109/BSN.2016.7516230

9. Hazem O, Al-Sadek AF (2022) Detection of tennis strokes using wearable sensor. In: 2022 international conference on software, telecommunications and computer networks (SoftCOM), pp 1–6. https://doi.org/10.23919/SoftCOM55329.2022.9911405

10. Gianoglio C, Ragusa E, Zunino R, Valle M (2021) 1-D convolutional neural networks for touch modalities classification. In: 2021 28th IEEE international conference on electronics, circuits, and systems (ICECS), pp 1–6. https://doi.org/10.1109/ICECS53924.2021.9665576

11. Ragusa E, Gianoglio C, Zunino R, Gastaldo P (2020) A design strategy for the efficient implementation of random basis neural networks on resource-constrained devices. Neural Process Lett 51(2):1611–1629. https://doi.org/10.1007/s11063-019-10165-y

12. Sakr F, Younes H, Doyle J, Bellotti F, De Gloria A, Berta R (2022) A tiny CNN for embedded electronic skin systems. In: Advances in system-integrated intelligence: proceedings of the 6th international conference on system-integrated intelligence (SysInt 2022), September 7–9, 2022, Genova, Italy, Springer, pp 564–573

13. Nano 33 IoT | Arduino Documentation. https://docs.arduino.cc/hardware/nano-33-iot. Accessed 19 Jun 2023

14. Ma R, Ahmadzadeh A, Boubrahimi SF, Angryk RA (2019) A scalable segmented dynamic time warping for time series classification. In: Rutkowski L, Scherer R, Korytkowski M, Pedrycz W, Tadeusiewicz R, Zurada JM (eds) Artificial intelligence and soft computing. Lecture notes in computer science. Springer International Publishing, Cham, pp 407–419. https://doi.org/10.1007/978-3-030-20915-5_37

# Evaluating the Effect of Intrinsic Sensor Noise for Vibration Diagnostic in the Compressed Domain Using Convolutional Neural Networks

Federica Zonzini[1(✉)], Edoardo Ragusa[2], Luca De Marchi[1], and Paolo Gastaldo[2]

[1] Department of Electrical, Electronic, and Information Engineering (DEI), University of Bologna, Bologna, Italy
`federica.zonzini@unibo.it`

[2] Department of Electrical, Electronic, Telecommunications Engineering and Naval Architecture (DITEN), University of Genoa, Genoa, Italy

**Abstract.** Machine learning allows designing intelligent sensing networks capable to perform automatic inferences about the integrity of technical facilities. Compression techniques decrease significantly energy requirements of the sensing networks proving essential when sensing nodes are not supported by constant power sources. Existing schemes pass through the reconstruction of the original time series data before moving to the diagnosis phase. However, this passage can be avoided, i.e., inference can be performed directly in the compressed domain, by exploiting the specific information retained in the compressed patterns. This paper fulfills the goal above in the context of vibration-based structural health monitoring by proving, from an empirical perspective, that Convolutional Neural Networks (CNNs) can be used to predict the structural health status directly in the compressed domain when properly combined with adapted Compressed Sensing mechanisms. Importantly, the study analyses the effect of the intrinsic noise that affects digital accelerometer sensors. Results confirm that CNNs can mine information in the compressed domain even in presence of strong noise components, i.e., accuracy remains above 94% even for ultra-low-cost solutions featuring a signal-to-noise-ratio below 20 dB.

**Keywords:** CNNs · Compressed sensing · Intrinsic accelerometer noise · Structural health monitoring · Vibration analysis

## 1   Introduction

Structural Health Monitoring (SHM) implemented via dense and distributed sensor networks has become a dominant paradigm in modern inspection systems aimed at ensuring the serviceability of civil and industrial appliances. SHM has

largely benefited from the progress achieved by the Information and Embedded system community, paving the way to the permanent inspection of large facilities via low-cost and low-power electronics offering sensor-near processing functionalities [1]. Indeed, implementing algorithms at a sensor level has significant advantages, such as the possibility to extract damage-sensitive features at the sensing stage, opening new scenarios for the compression of long time series into a reduced vector of meaningful parameters, which could be passed directly to the diagnostic unit for the retrieval of the health status. In this framework, data reduction techniques, such as those based on the Compressed Sensing (CS) theory, are among the most effective for vibration-based SHM [2]. The latter is a non-destructive testing technique which can be used to inspect structures dominated by an oscillatory behaviour: accordingly, the analysis is usually performed by observing the evolution over time of spectral features [3].

The existing approaches for damage detection with CS must face a burdensome decoding procedure: after compression at the sensor level, data transmitted to the aggregation unit are reconstructed in the time domain with onerous iterative procedures before the meaningful information is extracted, and then passed to the diagnostic process [4]. However, since the reconstruction stage works by applying some optimization metrics to the residual data [5], this means that, in principle, the residual information itself could suffice to discriminate healthy from damaged patterns without the need to reconstruct the global properties of the structure. In this work, we explore the feasibility of this decompression-free approach from a Machine Learning (ML) perspective, also encompassing the effect of intrinsic sensor noise typical of digital low-cost accelerometers. The motivation is that this unavoidable source of noise can easily contaminate the quality of vibration data under noise ambient excitation, i.e., scenarios in which the vibration response is induced by very faint forces (e.g., environmental agents or vehicles) and, hence, it can significantly compromise the discrimination capabilities of ML models.

## 2   Processing Framework

The monitoring network consists of a distributed set of extreme edge vibration nodes (i.e., sensing units), each of them equipped with a Microelectro-Mechanical System (MEMS) accelerometer mastered by a microprocessor with enabled digital signal processing functionalities. Accordingly, individual nodes acquire and compress onboard vibration signals, thus outsourcing the compressed information to a central aggregator where the inference phase is performed in a supervised manner by means of a tiny Convolutional Neural Network (CNN). It is worth noting that the merging step is mandatory to provide a global and robust judgment of the health status, hence preventing the relative position of single devices to bias the results [3].

## 2.1 Extreme Edge Level: CS-Based Compression

CS works by compressing data via a matrix-vector multiplication: in case an $N$-long time series $x_i$ is collected by a sensor $S_i$, the reduced version $\hat{x}_i = \Theta x_i$ is computed by exploiting the so-called *compression matrix* $\Theta \in \mathbb{R}^{M \times N}$ having $M << N$. A proper design of $\Theta$ is crucial to boost the compression performance and maximize the content retained at the reduction stage. To this end, adaptive methods have demonstrated superior capabilities thanks to their physics-informed nature, i.e., they can allow deriving the compression matrix as a function of the *a priori* known spectral properties of the structure under analysis. Adapted approaches show two additional advantages: (i) they can provide conservative solutions, namely, the resulting sensing scheme does not over adapt to the statistics used during the design process, therefore, (ii) they accurately capture variations even in the presence of damages, without needing $\Theta$ to be updated over time. This is a desirable feature in sensor-near computing frameworks, in which the limited storage and processing capabilities of the microprocessor hamper a streaming estimation of the compression matrix.

## 2.2 Aggregation Level: CNN-Based Inference

In the proposed sensing scheme, a sensing node sends a set of compressed measures $\hat{X} = [\hat{x}_1, \ldots, \hat{x}_R] \in \mathbb{R}^{M \times R}$ collected over $R$ successive intervals. The aggregating unit then receives $N_s$ matrices, one for each sensor, that are further aggregated into a tensor $\mathcal{X} \in R^{M \times R \times N_s}$ such that it can be fed as input of a ML detector for structural diagnostics. Hence, classification can take place in the compressed domain: the key idea investigated in this work is to seek changes in the configuration and magnitude of the compressed coefficients, that are likely to be caused by changes in the vibration pattern due to potential defects. The literature presented many approaches to mining and classifying information based on the resources available on the sensing nodes [6,7]. Among the possible options, CNN based on 1D filters through time can represent an appealing solution to detect such abnormalities because these models prove effective in retrieving defective recurrent patterns or changes in the overall behavior, while maintaining limited computing requirements [7]. The same approach can easily be extended to multi-class classification problems by properly defining the output layer of the convolutional architecture with minimal increment in the overall computational burden [8].

Accordingly, the devised approach consists of training CNNs with a filter of size $M \times K_s \times N_s$, where $K_s$ is the kernel size, that are convoluted with the original input. The overall structure of the proposed architecture is schematized in Fig. 1. Since the information is embedded in the specific configuration of the coefficients, multiple levels of stratification are not involved and a relatively shallow, single-branch architecture can instead be adopted: a 1D convolutional block (Conv-1D) with $N_{f1}$ filters processes the input followed by a max pooling layer. A second Conv-1D layer with $N_{f2}$ filters is stacked, followed by a Global Average Pooling for the extraction of a second set of features; a fully connected layer predicting the (binary) class of the input datum completes the model.

**Fig. 1.** Summary of the proposed architectures

## 3 Experimental Validation

**Dataset**: The benchmark use case of the Z24 bridge has been exploited to showcase the performance of the proposed monitoring approach. The dataset consists of 5651 labelled time series collected over one year of monitoring of a highway viaduct between Bern and Zürich, which was then interrupted in 1999 for modernization. The infrastructure has been artificially damaged via a rigorous experimental protocol, introducing slight to very severe flaws [9]. The first 4922 coincide with healthy conditions, while the remaining ones pertain to the structure under progressive damage tests.

**Noise injection**: High-sensitivity piezolectric accelerometers were used in the original sensor installation. However, their cost, size, and electrical features are far from the modern low-cost and low-power MEMS counterparts integrated in embedded devices, as those dominating in state-of-the-art SHM architectures. To assess the robustness of the devised strategy against these additional sources of noise, data were perturbed by adding the intrinsic noise characterizing two cut-of-the-shelf MEMS accelerometers widely applied for vibration-related applications: the LSM6DSL from STMicroelectronics [10] (noise density 80 $\mu$g/$\sqrt{\text{Hz}}$) belonging to the class of ultra-low-cost sensors, and the ADXL355 accelerometer fabricated by Analog Devices [11] and characterized by better noise properties (noise density 24.5 $\mu$g/$\sqrt{\text{Hz}}$). By simulating the characteristics of the former accelerometer, a Signal-to-Noise Ratio (SNR) equal to 20 dB was measured, while the same quantities increases to 30 dB for the second and more performative sensor.

**CS setting and CNN architectures**: Data were compressed using the Model-assisted Rakeness-based (MRak-CS) strategy in [12] because (i) it demonstrated better robustness against the perturbations induced by environmental and operational parameters and (ii) it allows to attain high compression level with negligible loss in the accuracy of the reconstructed damage-sensitive features. A constant compression ratio $CR = 8$ has been imposed, a quantity which is comparatively higher with respect to conventional solutions found in the literature, where the compression depth hardly exceeds 4–5 [13]. From a ML viewpoint, two CNNs architectures were designed: the first one, termed as *small*

*CNN*, used 5 and 3 filters for the first and second layer, respectively. The second architecture, called *large CNN*, disposed of 30 and 10 filters in the convolutional layers. For both configurations, $K_s = 3$ was set. Models were validated using a standard 5-fold cross-validation procedure. A validation set was extracted from the training set selecting a random subset of 20% of the training data. Then, the testing fold was not involved in any parameter tuning. Keras and Tensorflow supported the implementation of the ML models.

**Results**: Table 1 summarizes the results in terms of standard classification metrics (accuracy, precision, recall, F1) including the model type, the *CR*, the MEMS-related corrupting noise. Moreover, scores obtained by the baseline one-class classifier neural network (OCCNN) model in [14] are reported to corroborate the quality of the attained performances in comparison with state-of-the-art supervised alternatives dealing with the same noise settings, but with lower $CR = 6$. Best scores for each setting have been magnified in bold font.

**Table 1.** Classification results of the proposed workflow by considering various MEMS-related noise levels. Better values for each noise case are highlighted in bold

| Model | $CR$ | Sensor type | Accuracy | Precision | Recall | F1 |
|---|---|---|---|---|---|---|
| Small CNN | 8 | LSM6DSL | 0.933 | 0.952 | 0.972 | 0.962 |
| Large CNN | 8 | LSM6DSL | **0.945** | **0.961** | 0.976 | **0.969** |
| OCCNN [14] | 6 | LSM6DSL | 0.901 | 0.893 | **1.000** | 0.940 |
| Small CNN | 8 | ADXL355 | 0.949 | 0.966 | 0.976 | 0.971 |
| Large CNN | 8 | ADXL355 | **0.956** | **0.973** | **0.977** | **0.975** |
| OCCNN [14] | 6 | ADXL355 | 0.935 | 0.927 | 0.994 | 0.962 |
| Small CNN | 8 | N/A | 0.957 | 0.974 | 0.977 | 0.975 |
| Large CNN | 8 | N/A | **0.967** | **0.983** | **0.979** | **0.981** |
| OCCNN [14] | 6 | N/A | 0.930 | 0.940 | 0.950 | 0.910 |

As can be observed, CNNs can extract valuable information in the compressed domain as proved by classification scores always above 94%, even when a very low-cost and high noise density (i.e., LSM6DSL) sensor is considered. Among the two ML implementations, the largest CNN model attained the best performances, with an average increment of 1% in comparison to the tinier version. This outcome is reasonably due to the increased number of filters which allows a better understanding of the damage pattern contained within compressed data. Beside, it is paramount to emphasize that the proposed workflow always outperforms, apart from one isolated case, the prediction capabilities characterizing the previous OCCNN implementation: proof is the fact that the metrics increase from a minimum of 2% to a maximum of 7%, additionally working with a compression level which is 1.33x higher than the one tested in the previous configuration. Finally, compared with compression-free scenarios (N/A

label), a negligible decrement in all the scores has been reported, with a worst case reduction of 2.5% for the LSM6DSL sensor.

## 4    Conclusion

This paper presented an analysis of the interaction between intrinsic sensor noise typical of MEMS accelerometers, compressed sensing, and structural classification implemented via CNNs in the compressed domain. The empirical results confirm that CNNs can master noisy data yielding high-quality prediction even when noise signal ratio deteriorates significantly, improving up to 7% the inference metrics of state-of-the-art alternatives in the field.

## References

1. Ragusa E, Zonzini F, De Marchi L, Gastaldo P (2023) Vibration monitoring in the compressed domain with energy-efficient sensor networks. IEEE Sens Lett
2. Hartmann M, Hashmi US, Imran A (2022) Edge computing in smart health care systems: review, challenges, and research directions. Trans Emerg Telecommun Technol 33(3):e3710
3. Yang Y, Zhang Y, Tan X (2021) Review on vibration-based structural health monitoring techniques and technical codes. Symmetry 13(11):1998
4. Ganesan V, Das T, Rahnavard N, Kauffman JL (2017) Vibration-based monitoring and diagnostics using compressive sensing. J Sound Vib 394:612–630
5. Li L, Fang Y, Liu L, Peng H, Kurths J, Yang Y (2020) Overview of compressed sensing: sensing model, reconstruction algorithm, and its applications. Appl Sci 10(17):5909
6. Ragusa E, Gastaldo P, Zunino R, Cambria E (2020) Balancing computational complexity and generalization ability: a novel design for elm. Neurocomputing 401:405–417
7. Gianoglio C, Ragusa E, Gastaldo P, Valle M (2023) Trade-off between accuracy and computational cost with neural architecture search: a novel strategy for tactile sensing design. IEEE Sens Lett
8. Lamrini M, Chkouri MY, Touhafi A (2023) Evaluating the performance of pre-trained convolutional neural network for audio classification on embedded systems for anomaly detection in smart cities. Sensors 23(13):6227
9. Peeters B, De Roeck G (2001) One-year monitoring of the Z24-bridge: environmental effects versus damage events. Earthq Eng Struct Dyn 30(2):149–171
10. ST Microelectronics: LSM6DSL. iNEMO inertial module: always-on 3D accelerometer and 3D gyroscope (2017)
11. Analog Devices: ADXL355. Low Noise, Low Drift, Low Power, 3-Axis MEMS Accelerometers (2016)
12. Zonzini F, Zauli M, Mangia M, Testoni N, De Marchi L (2021) Model-assisted compressed sensing for vibration-based structural health monitoring. IEEE Trans Ind Inform 17(11):7338–7347
13. Premanand B, Sheeba V (2020) Compressed encoding of vibration signals using extremum sampling. SN Appl Sci 2:1–10
14. Zonzini F, Carbone A, Romano F, Zauli M, De Marchi L (2022) Machine learning meets compressed sensing in vibration-based monitoring. Sensors 22(6):2229

# Cooperative Driver Assistance
# for Electric Wheelchair

Federico Pacini$^{(\boxtimes)}$ ⓘ, Pierpaolo Dini ⓘ, and Luca Fanucci ⓘ

Department of Information Engineering, University of Pisa, Pisa, Italy
`federico.pacini@phd.unipi.it`

**Abstract.** Driving a motorized wheelchair is not without risk and requires high cognitive effort to obtain a good environmental perception to complete tasks such as obstacle avoidance and path planning. Nonetheless, completely autonomous wheelchairs release the users from any cognitive and muscular effort, resulting in a progressive degradation of their residual capacities. Thus, we developed a collaborative driver assistance system that engages selectively, specifically when it detects a positive or negative height object within a defined safety zone. The system first reduces linear speed to assist and prioritize the user in making a decision and ultimately calculates a steering direction that avoids obstacles. An extensible mathematical model is devised by relying on lidar and ultrasonic sensors. Simulations have been carried out, confirming the effectiveness of this proof-of-concept.

**Keywords:** Wheelchair · Cooperative driving assistance · Obstacles avoidance

## 1 Introduction

The World Health Organization estimates that 1.3 billion people, or 16% of the global population, experience significant disability [1]. Among them, over 200 million individuals have mobility challenges. The Convention for the Rights of Persons with Disabilities (CRPD) [2] emphasizes personal autonomy in disability policies. Limited mobility can lead to a decrease in social life and well-being [3]. Power wheelchairs are a solution for those with motor skill impairments, but they come with risks, as more than 54% of wheelchair users have accidents [4]. Robotic assistive solutions can enhance safety [5–7]. Autonomous technology has made progress, but completely autonomous wheelchairs can lead to a decline in users' capabilities. Therefore, a cooperative driver assistance system is proposed. This system activates when necessary, reducing accidents and allowing users to handle navigation tasks independently. The goal is to design a cooperative driving assistance system that works indoors and outdoors and activates in dangerous situations, such as approaching obstacles. The system utilizes 2D lidar scanners and ultrasonic sensors for obstacle detection and can be easily extended with other sensors.

## 2    Cooperative Control System

To achieve a robust obstacle avoidance control system that prevents the user from crashing into obstacles or falling down from pavement edges, while still granting the user control over the wheelchair, we propose a cooperative control system algorithm. The process begins by fusing data from a 2D lidar scan and an ultrasonic sensor. This fusion enables the creation of a map that distinguishes between positive height obstacles (those above the ground level of the wheelchair) and negative ones, such as steps or downward slopes, all represented on a single plane. Both the obstacles map and the user's steering direction are fed into the VFH+ [8] algorithm. This algorithm then computes a steering direction that balances smoothness and safety, taking into account the distribution of obstacles. In addition, to facilitate the user to identify obstacles, we added a speed control system that activates progressively as approaching a dangerous situation.

### 2.1    Modelling

As shown schematically in Fig. 1, we use the following nomenclature: $L$ represents the origin of lidar rays; $G$ denotes the mass-centre of the wheelchair; $s^{ith}$ is the ith ultrasonic sensor; $\beta^{ith}$ represent the orientation of the ith ultrasonic sensor; $d^{ith}$ denotes the distance of the ith ultrasonic sensor from R; $\alpha^{ith}$ is the inclination of the ith ultrasonic sensor; $h^{ith}$ is the position height of the ith ultrasonic sensor; $y^{ith}$ represent the distance between the ith ultrasonic sensor ray and the first reflection object.



**Fig. 1.** Model of the wheelchair with lidar and ultrasonic sensors from top view (LEFT); model of the wheelchair with details on the $k$th ultrasonic sensor from lateral view (RIGHT).

If the wheelchair is standing on flat terrain, we define:

$$\overline{y^{ith}} = \frac{h^{ith}}{\cos \alpha^{ith}} \tag{1}$$

Then, given a measurement $y_k^{ith}$, we define a boolean function which states if the $i^{th}$ ultrasonic sensor is detecting an edge as:

$$edge^{ith}(y_k^{ith}) = \begin{cases} 1, & \text{if}(y_k^{ith} < \overline{y^{ith}} - \epsilon) \vee (y_k^{ith} > \overline{y^{ith}} + \epsilon) \\ 0, & \text{if}(y_k^{ith} > \overline{y^{ith}} - \epsilon) \wedge (y_k^{ith} < \overline{y^{ith}} + \epsilon) \end{cases} \tag{2}$$

where $\epsilon$ is a small value taking into consideration measurement errors by the sensor. If a edge is detected, taking as reference origin $L$, we assume that there is a edge $e^{ith}$ which has the Cartesian coordinates:

$$e_k^{ith} = (d^{ith} + y_k^{ith} \cos \beta^{ith}; y_k^{ith} \sin \beta^{ith}) \tag{3}$$

When considering N sensors, we can write the edges matrix at $k$th time instant as:

$$E_k = \begin{bmatrix} e_k^{1st} & e_k^{2nd} & \dots & e_k^{ith} & \dots & e_k^{Nth} \end{bmatrix} \tag{4}$$

where $e_k^{jth}$ is a valid point if $edges^{jth}(y_k^{jth}) = 1$. By applying linear interpolation between two consecutive valid elements, we define the shape of the obstacles as:

$$ob_k^{ith,(i+1)th}(x) = e_{y_k}^{ith} + \frac{e_{y_k}^{(i+1)th} - e_{y_k}^{ith}}{e_{x_k}^{(i+1)th} - e_{x_k}^{ith}}(x - e_{x_k}^{ith}) \tag{5}$$

where $x \in [e_{x_k}^{(i+1)th}; e_{x_k}^{ith}]$.

Due to the decision to align the origin of the 2D ultrasonic data plane with the origin of the lidar rays, we can effortlessly represent both data on the same plane without any additional computing. The result is an obstacles map that contains both positive and negative height obstacles, coming respectively from lidar and ultrasonic sensors.

## 2.2   Steering Direction

To prevent the user from crashing into positive height obstacles or falling into negative ones, the cooperative control system calculates an obstacle-free steering direction taking into consideration user input and the obstacles map. According to the VFH+ algorithm [8], the obstacle map is converted into a polar histogram. Each sector in the polar histogram contains a value representing the polar obstacle density in that direction. The algorithm selects the most suitable direction based on the masked polar histogram and a cost function, which takes into consideration the user's target direction, the current wheelchair direction and the previously selected direction—see Fig. 2. To ensure the algorithm remains inactive in the absence of dangerous situations, it is triggered only when a safety distance threshold is surpassed.

**Fig. 2.** Example of selected steering direction given an obstacle map and a user target direction.

## 2.3   Speed Regulation

To prioritize the user to make a decision, we design a speed control in order to gradually decrease the forward speed. This choice has the dual effect of allowing the user, who may have a slower reaction time than average, to realise that a dangerous situation may arise and to avoid an awkward turn due to the selected steering direction. Depending on the minimum distance $2d^*$ between an obstacle that lays in a collision cone of amplitude $2\gamma$, a forward speed reduction multiplier is calculated (see Fig. 3).



$$m = \frac{1}{1 + e^{-k\,x + (d^* - 1)}}$$

**Fig. 3.** Collision cone centred in lidar origin and safety distance $d^*$ (LEFT); examples of forward speed reduction multiplier varying the k parameter. As k increases, later and stronger speed reduction will be applied (RIGHT).

# 3    Simulations and Results

The simulation has been carried out making use of Gazebo. Gazebo is an open-source well-established simulator in the field of robotics that provides a robust physics engine, and convenient programmatic and graphical interfaces. The algorithms related to the logic of the wheelchair control have been implemented making use of Robotic Operating System (ROS) version 2. ROS2 is a well-established set of software libraries and tools for building robot applications.

A schematic representing the relationships between the User interface, ROS2 nodes and the Gazebo simulator is represented in Fig. 4.

The User Interface gathered input from the user, specifically the desired steering direction and forward speed.



**Fig. 4.** Schematic representation of the virtual environment.

The Gazebo simulator simulated the behaviour of the lidar, ultrasonic sensors, and wheelchair dynamics. In the end, a group of ROS2 nodes processed the data from the lidar and ultrasonic sensors to generate an obstacle map.

They also computed the appropriate steering direction, forward speed, and wheel torques to ensure that the wheelchair followed them.

Simulations have been carried out on a virtual path that included sharp corners, rounded obstacles and edges. To be as close as possible to commercially available wheelchairs, we set a maximum forward speed of 1.3 m/s and a maximum rotational speed of 0.4 rad/s. The task was to reach the end of the path by controlling the wheelchair with a joystick. An example of a trial is represented in Fig. 5 which includes both positive height obstacles such as walls and geometrical volumes and negative height ones such as curbs. As we can see, when the wheelchair is approaching an obstacle, the speed control is progressively reducing the forward speed to let the user have time to react. For the sake of the trial, we did not turn on purpose. For this reason, eventually, the cooperative control system helps the user by calculating a steering direction that is obstacle-free.

**Fig. 5.** Example of wheelchair trajectory on simulation environment. Speed varies according to the speed regulation algorithm. Speed variations are represented by a coloured bar.



**Fig. 6.** Subset of reference steering direction from the user (input) and from the cooperative control system (output), related to first obstacle avoidance.

An example of cooperative control system action is shown in Fig. 6. As depicted in the figure, even though the user is not steering, the wheelchair is requested to steer at roughly 24°. On purpose, we tried to steer toward the obstacles by requesting a steering direction of roughly –45°. However, the system partially accomplished our request by slightly steering but not completely because we would have crashed into the obstacle. As soon as the system realizes that the dangerous situation is over, the wheelchair follows the user's request.

## 4     Discussion and Conclusions

Driving a motorized wheelchair is a very demanding task. People with low vision, visual field reduction, spasticity, tremors or cognitive deficits have severe difficulties in carrying out this task in a safe way and this has a negative effect on their autonomy. The possibility of moving in an autonomous way gives individuals a remarkable physical and psychological sense of well-being. For this reason, this paper has presented a cooperative control system that aims to help, but not replace, the user in driving their wheelchair safely. A mathematical model has been derived and tested in a simulation environment. Results demonstrated the effectiveness of this proof-of-concept. Further development of this research project involves building a demonstration prototype, which includes a mechatronic design comprising electric drives, control and modulation algorithms for power and control electronics. These components will be optimized for seamless integration with the navigation sensor system, enabling the safe execution of the cooperative-driving algorithm developed and tested in a virtual environment, as shown in these articles [9–11]. Another aspect of the development will focus on utilizing the AI-based model for two main purposes. Firstly, it will be employed to enhance the autonomous driving system, making it even more effective. Secondly, the model will be utilized to expand the range of scenarios covered, all with the ultimate goal of ensuring user safety [12,13].

## References

1. Organization WH, Bank W (2011) World report on disability 2011
2. Un general assembly, convention on the rights of people with disabilities (2007)
3. Rosso AL, Taylor JA, Tabb LP, Michael YL (2013) Mobility, disability, and social engagement in older adults. J Aging Health 25(4):617–637
4. Chen W-Y, Jang Y, Wang J-D, Huang W-N, Chang C-C, Mao H-F, Wang Y-H (2011) Wheelchair-related accidents: relationship with wheelchair-using behavior in active community wheelchair users. Arch Phys Med Rehabil 92(6):892–898
5. Dini P, Saponara S (2022) Processor-in-the-loop validation of a gradient descent-based model predictive control for assisted driving and obstacles avoidance applications. IEEE Access 10:67958–67975
6. Cosimi F, Dini P, Giannetti S, Petrelli M, Saponara S (2021) Analysis and design of a non-linear MPC algorithm for vehicle trajectory tracking and obstacle avoidance. In: Applications in electronics pervading industry, environment and society: APPLEPIES 2020, vol. 8. Springer, pp 229–234

7. Bernardeschi C, Dini P, Domenici A, Mouhagir A, Palmieri M, Saponara S, Sassolas T, Zaourar L (2021) Co-simulation of a model predictive control system for automotive applications. In: International conference on software engineering and formal methods. Springer, pp 204–220
8. Ulrich I, Borenstein J (1998) Vfh+: reliable obstacle avoidance for fast mobile robots. In: Proceedings. 1998 IEEE international conference on robotics and automation (Cat. No.98CH36146), vol 2, pp 1572–1577
9. Dini P, Saponara S (2019) Cogging torque reduction in brushless motors by a nonlinear control technique. Energies 12(11):2224
10. Dini P, Saponara S (2020) Design of adaptive controller exploiting learning concepts applied to a BLDC-based drive system. Energies 13(10):2512
11. Bernardeschi C, Dini P, Domenici A, Palmieri M, Saponara S (2020) Formal verification and co-simulation in the design of a synchronous motor control algorithm. Energies 13(16):4057
12. Dini P, Begni A, Ciavarella S, De Paoli E, Fiorelli G, Silvestro C, Saponara S (2022) Design and testing novel one-class classifier based on polynomial interpolation with application to networking security. IEEE Access 10:67910–67924
13. Begni A, Dini P, Saponara S (2022) Design and test of an LSTM-based algorithm for Li-Ion batteries remaining useful life estimation. In: International conference on applications in electronics pervading industry, environment and society, Springer, pp 373–379

# Tiny Neural Deep Clustering: An Unsupervised Approach for Continual Machine Learning on the Edge

Giovanni Poletti, Andrea Albanese$^{(\boxtimes)}$, Matteo Nardello, and Davide Brunelli

University of Trento, Department of Industrial Engineering, Via Sommarive 9, 38123 Trento, Italy
{giovanni.poletti,andrea.albanese,matteo.nardello,
davide.brunelli}@unitn.it

**Abstract.** Traditional tiny machine learning systems are widely employed because of their limited energy consumption, fast execution, and easy deployment. However, such systems have limited access to labelled data and need periodic maintenance due to the evolution of data distribution (i.e., context drift). Continual machine learning algorithms can enable continuous learning on embedded systems by updating their parameters, addressing context drift, and allowing neural networks to learn new categories over time. However, the availability of labelled data is scarce, limiting such algorithms in supervised settings. This paper overcomes this limitation with an alternative approach which combines supervised deep learning with unsupervised clustering to enable unsupervised continual machine learning on the edge. Tiny Neural Deep Clustering (TinyNDC) is deployed in an OpenMV Cam H7 Plus and tested with the MNIST dataset reaching a classification accuracy of 92.3% and a frame rate of 44 FPS.

**Keywords:** Continual machine learning · On-device training · TinyML · Unsupervised learning · Deep learning

## 1 Introduction

Traditional tiny Machine Learning (tinyML) systems are trained on a fixed dataset and assume that the data distribution does not change over time. In many real-world applications, the data distribution can evolve, leading to a problem known as context drift. It occurs when the relationship between input data and output labels changes over time, causing degradation of ML model performance. The main reasons are changes in the environment, changes in user behavior, or changes in the underlying data distribution. Context drift can be a challenging problem for ML models deployed in real-world applications, where the input data may change over time. If a model cannot adapt to the changing context, its performance can suffer, resulting in incorrect predictions or decisions. A solution to overcome context drift is Continual Machine

Learning (CML). CML is a subfield of machine learning that focuses on developing algorithms and techniques that can continuously learn and adapt to new data and tasks over time, without forgetting previously learned information. However, this approach is computationally expensive and is usually relegated only to high performance cloud computing. An example is the "teacher-student" paradigm where the knowledge of a large and accurate neural network (teacher) trained using high performance computing, is transferred to a more efficient model (student), making it feasible to run complex machine learning tasks on resource-constrained devices [1]. A different approach consists of applying CML to resource-constrained devices (e.g., MCUs) creating intelligent self-updated IoT systems with high flexibility requiring low maintenance in challenging scenarios [2]. Common state-of-the-art CML strategies, such as CWR [3], LWF [4], and tinyOL [5], can mitigate the effects of context drift; however, they need a ground truth to carry out the backpropagation for weights and biases updates. Those approaches are usually referred to as supervised learning. Thus, the incoming data must be labeled, making such systems unfeasible for real-world applications because of the unavailability of ground truths. Other works, such as Online Deep Clustering [6] and Unsupervised Continual Learning for Gradually Varying Domains [7] avoid the need for labeled data. In the former case, the authors use a combination of clustering and deep learning algorithms to implement a semi-unsupervised continual learning system. In the latter, authors use clustering as a classification system to adapt to gradually varying domains. Even though these works show successful solutions for unsupervised CML, they do not include embedded implementation, which is challenging for real-world applications. So far, only the work presented in [8] presents a solution for on-device training with a low-memory budget. However, this solution still works in a supervised setting, and during the on-device retraining, the data need the ground truth. Thus, an embedded solution for implementing unsupervised CML to get closer to real-world exploitation is still missing making this work a useful contribution to prove the effectiveness of unsupervised CML on the edge. Multiple the domains that could benefit from unsupervised learning on the edge, in the industrial visual inspection [9], for surveillance tasks [10,11], smart agriculture [12,13], and even in the medical field (e.g., histopathology) to enhance the period of use without maintenance [14,15]. In this paper, we present Tiny Neural Deep Clustering (TinyNDC), an innovative approach that combines supervised learning with unsupervised clustering to enable unsupervised continual online machine learning on resource-constrained devices. To evaluate the proposed algorithm, the MNIST dataset was used. The first 6 digits were used to initialize the model, while the remaining 5 were used to test the online learning algorithm. Results show that TinyNDC can achieve 44 FPS and a learning accuracy of 92.3% on average among the 10 MNIST digits. The paper is organized as follows: Sect. 2 presents the design of TinyNDC, while Sect. 3 presents the experimental setup. Finally, Sect. 4 discusses the experimental results, and Sect. 5 concludes the paper.

## 2 System Design

Combining clustering and deep learning systems can be proven a successful solution. It mixes the flexibility of clustering in adapting to different domains with the stability of DL in avoiding drift from the current domain. We propose an algorithm called **TinyNDC** that has been tailored for running in an ARM Cortex M7 processor.



**Fig. 1.** TinyNDC architecture. Data are fed to the *Frozen Model* that returns a set of features. Online deep clustering clusters the features, calculating the closest centroid, while the *Active Model* carries out the classification. Then the active model prediction is used to update the centroids of online deep clustering, and the pseudo label produced by the clustering algorithm is used to update the weights and biases of the active model.

Figure 1 presents the architecture of TinyNDC and is composed of three main components:

– **Frozen model**. The static part of the model which is not updated during runtime. It can be any NN truncated in the last layer (e.g., the Softmax classifier).
– **Active model**. The dynamic part of the model which is updated during runtime. It is composed of two sub-modules, namely the *consolidated layer* and the *training layer*. The first acts as a memory to mitigate catastrophic forgetting, while the latter is used for domain adaptation. The active model can update its weights and biases and add new classes if the current data is not associated with an already known class.
– **Online deep clustering**. It is fed with the output of the frozen model and produces the pseudo-label estimation. This label is fundamental for updating the active model. Furthermore, the clustering can update its centroids by using the prediction of the active model. The clusterings centroids need to be initialized with a small amount of data (e.g., a hundred samples).

## 3    Experimental Setup

To evaluate the performance of the proposed TinyNDC algorithm, we used the MNIST dataset and an OpenMV Cam H7 Plus (OMV)[1] device. The frozen model was trained to classify digits from 0 to 5, while the remaining digits were used to test the CML capabilities in an unsupervised scenario. Specifically, the dataset consists of a total of 60000 samples. 35000 samples are used for training the frozen model, 100 for cluster initialization (10 per class), 7000 for the online learning phase, and 1000 for assessing the CML performances. Three different online learning algorithms, namely CWR [3], LWF [4], and tinyOL [5], were tested as update mechanisms for the active model. Accuracy was selected as the performance metric. CWR, LWF, and tinyOL reach an accuracy of 92.3, 91.3 and 91.7%, respectively. As CWR was proved to be the best-performing algorithm it was chosen as the target update rule for performing further tests.

The experimental setup consists of correctly identifying data shown to the OMV through a PC screen.

### 3.1    Frozen Model Training

TinyNDC is pre-trained in a supervised fashion on a PC to recognize digits from 0 to 5. The frozen model generation and training were carried out with the TensorFlow library and Python on the PC. We used a CNN specifically designed for image processing composed of four 2D convolutional layers. For the training phase, Adam was used as an optimizer, categorical cross entropy as a loss function, with a total of 40 epochs, and a batch size of 32. After the training phase, the frozen model shows an accuracy of 99.64% on average among the first 6 digits.

### 3.2    Centroids Inizialization

After the initialization of the frozen model, we moved to the centroid initialization, needed as a preliminary step. We used 10 samples per category, thus we extracted 100 samples from the 60000 training images. The OMV snaps a photo of each digit that shows up on the screen, feeds it to the trained frozen model for feature extraction, and sends the resulting array to the PC. Finally, the PC collects the 100 arrays, where each image feature is represented by a vector of 512 samples. The brightness of the screen is fixed. The clusters are defined by the sample labels and each centroid is calculated through the mean of all the respective cluster samples.

## 4    Experimental Results

The proposed system is compared with the same system composed of the frozen and active model, *TinyCML*, in a supervised setting (i.e., without the clustering

---

[1] https://openmv.io/products/openmv-cam-h7-plus.

algorithm for estimating the pseudo-label) to monitor performance degradation in moving from supervised to unsupervised CML. Furthermore, a comparison is also conducted with the clustering algorithm *TinyODC*, which acts as an unsupervised CML system to ensure the stability and feasibility of the proposed system. In this case, the clustering algorithm can classify data and update its centroids during runtime. The comparison, presented in Table 1, shows that even though our system works in a more challenging setting (i.e., unsupervised), it reaches good performance almost comparable with the supervised system. Furthermore, the clustering algorithm alone performs worse than the *TinyNDC*, confirming the effectiveness of the combination of clustering and DL. Figure 2a presents the learning curve comparison considering the accuracy. It shows that *TinyNDC* starts learning with a low accuracy level compared to *TinyCML*, but, after 8000 samples, it almost reaches the learning level of *TinyCML*. This further confirms the stability of clustering and DL combination, which can reach the performance of the associated supervised system. It can be argued that the best performing algorithm, as shown in Table 1, is *TinyCML*. However, it uses a supervised learning approach that is not feasible in real-world scenarios. On the contrary, *TinyNDC* trades a mere 2% accuracy loss with the ability to work in an unsupervised fashion, making it feasible to be implemented in real-world applications.

**Table 1.** Scores comparison between supervised CML (*TinyCML*), unsupervised CML with only clustering (*TinyODC*), and the proposed solution *TinyNDC*. *TinyCML* uses an initial learning rate of 0.8, a decay rate of 2, a learning rate step size of 500, and a batch size of 16. Tests were conducted at 160 lux using 7000 samples for the training and 1000 for the validation.

| Algorithm | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|
| TinyCML | 94.3% ± 1.2% | 95.0% ± 1.2% | 94.0% ± 1.2% | 94.0% ± 1.2% |
| TinyODC | 90.7% ± 0.6% | 91.0% ± 0.6% | 91.0% ± 0.6% | 91.0% ± 0.6% |
| TinyNDC | 92.3% ± 1.2% | 93.0% ± 1.2% | 92.0% ± 1.2% | 92.0% ± 1.2% |

Finally, the proposed system is characterized by measuring the execution time of each task to ensure real-time capability and possible employment in real-world applications. The *TinyNDC* tasks execution times are compared with supervised CML (*TinyCML*) and unsupervised CML with clustering (*TinyODC*). As shown in Table 2, *TinyNDC* needs more time than the other strategies to process one sample and perform the model update. On the other hand, as shown in Fig. 2b, most of the processing time is needed by the frozen model, which requires 16.39 ms to produce the inference result, while the active model and the clustering need only 6.12 ms for inference and update. Thus, even though *TinyNDC* uses a combination of clustering and DL, the processing complexity is not increased considerably, confirming the feasibility of using *TinyNDC* for real-time applications.

(a)                                                    (b)

**Fig. 2.** Comparison of *TinyCML*, *TinyODC*, and *TinyNDC* for **a** accuracy learning curves, and **b** tasks' execution time.

**Table 2.** Execution times of different sections of supervised CML (*TinyCML*), unsupervised CML with clustering (*TinyODC*), and *TinyNDC*.

| Section | TinyCML [ms] | TinyODC [ms] | TinyNDC [ms] |
|---|---|---|---|
| Frozen model | 16.39 | 16.39 | 16.39 |
| Clustering | – | 2.18 | 2.18 |
| Active model | 1.38 | – | 1.38 |
| Clustering update | – | 0.22 | 0.22 |
| Active model update | 2.34 | – | 2.34 |
| Total | 20.10 | 18.79 | 22.50 |

## 5    Conclusion

Continual Machine Learning is an active research area that aims to provide self-updating systems. This work presents TinyNDC, an initial approach to unsupervised learning on the edge by combining clustering with deep learning. The algorithm was tested on an MCU-based camera, namely OpenMV Cam H7 Plus, and evaluated against the MNIST dataset. Results show that the system can achieve a frame rate of 44 FPS and a learning accuracy of 92.3%. Future work will employ the usage of a more deep update by considering the inner layers and benchmarking with real case studies.

# References

1. Cai H, Gan C, Wang T, Zhang Z, Han S (2020) Once for all: Train one network and specialize it for efficient deployment. In: International conference on learning representations
2. Pilati F, Sbaragli A (2023) Learning human-process interaction in manual manufacturing job shops through indoor positioning systems. Comput Ind 151:103984
3. Maltoni D, Lomonaco V (2019) Continuous learning in single-incremental-task scenarios. Neural Netw 116:56–73
4. Zenke F, Poole B, Ganguli S (2017) Continual learning through synaptic intelligence. In: International conference on machine learning. PMLR, pp 3987–3995
5. Ren H, Anicic D, Runkler TA (2021) Tinyol: Tinyml with online-learning on microcontrollers. In: 2021 international joint conference on neural networks (IJCNN). IEEE, pp 1–8
6. Zhan X, Xie J, Liu Z, Ong Y-S, Loy CC (2020)Online deep clustering for unsupervised representation learning. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. IEEE, pp 6688–6697
7. Taufique AMN, Jahan CS, Savakis A (2022) Unsupervised continual learning for gradually varying domains. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp 3740–3750
8. Lin J, Zhu L, Chen W-M, Wang W-C, Gan C, Han S (2022) On-device training under 256kb memory. Adv Neural Inf Process Syst 35:22 941–22 954
9. Albanese A, Nardello M, Fiacco G, Brunelli D (2023) Tiny machine learning for high accuracy product quality inspection. IEEE Sens J 23(2):1575–1583
10. Nardello M, Desai H, Brunelli D, Lucia B (2019) Camaroptera: a batteryless long-range remote visual sensing system. In: Proceedings of the 7th international workshop on energy harvesting and energy-neutral sensing systems, ser. ENSsys'19. Association for Computing Machinery, New York, NY, USA, pp 8–14
11. Desai H, Nardello M, Brunelli D, Lucia B (2022) Camaroptera: a long-range image sensor with local inference for remote sensing applications. ACM Trans Embed Comput Syst 21(3)
12. Albanese A, Nardello M, Brunelli D (2021) Automated pest detection with DNN on the edge for precision agriculture. IEEE J Emerg Sel Top Circuits Syst 11(3):458–467
13. Brunelli D, Albanese A, d'Acunto D, Nardello M (2019) Energy neutral machine learning based iot device for pest detection in precision agriculture. IEEE Internet Things Mag 2(4):10–13
14. Muhammad K, Hussain T, Del Ser J, Palade V, de Albuquerque VHC (2020) Deepres: A deep learning-based video summarization strategy for resource-constrained industrial surveillance scenarios. IEEE Trans Ind Inform 16(9):5938–5947
15. Thandiackal K, Piccinelli L, Pati P, Goksel O (2023) Multi-scale feature alignment for continual learning of unlabeled domains. arXiv:2302.01287

# Investigating Adversarial Policy Learning for Robust Agents in Automated Driving Highway Simulations

Alessandro Pighetti[1]([✉]) [ID], Francesco Bellotti[1] [ID], Changjae Oh[2], Luca Lazzaroni[1] [ID], Luca Forneris[1] [ID], Matteo Fresta[1] [ID], and Riccardo Berta[1] [ID]

[1] Department of Electrical, Electronic and Telecommunication Engineering (DITEN), University of Genoa, Via Opera Pia 11a, 16145 Genoa, Italy
{alessandro.pighetti,luca.lazzaroni,luca.forneris, matteo.fresta}@edu.unige.it, {francesco.bellotti, riccardo.berta}@unige.it

[2] School of Electronic Engineering and Computer Science, Queen Mary University of London, London, England
c.oh@qmul.ac.uk

**Abstract.** This research explores an emerging approach, the adversarial policy learning paradigm, that aims to increase safety and robustness in deep reinforcement learning models for automated driving. We propose an iterative procedure to train an adversarial agent acting in a highway-simulated environment to attack a victim agent that is to be improved. Each training iteration consists of two phases. The adversarial agent is first trained to disrupt the victim-agent policy. The victim model is then trained to overcome the defects observed by the attack from the adversarial agent. The experimental results demonstrate that the victim agent trained with adversarial attacks outperforms the original agent.

**Keywords:** Automated driving · Adversarial policies · Autonomous agents · Proximal policy optimization · Reinforcement learning · Highway driving · Decision making · Highway-env simulator

## 1 Introduction

Deep reinforcement learning (DRL) has emerged as a powerful approach to train online agent policies, enabling them to learn complex decision-making behaviors in dynamic environments, which is being used also in automotive applications (e.g., [1–3]). However, ensuring the safety and robustness of the learned policy remains a critical challenge. Different approaches are being studied to address safety and robustness in deep learning models. The adversarial paradigm [4] represents one of such emerging approaches in various research fields, such as computer vision and reinforcement learning. It is defined as the procedure of purposefully attacking the model under development, also known as the victim, through different methods depending on its purpose and area of deployment. In particular, DRL policies are vulnerable to adversarial attacks involving perturbations

to the models state, action or transition dynamics [4] or by applying destabilizing forces on the trained system [5].

The existing work in [6] shows that, in competitive two-player games, it is possible to build an adversarial policy for an agent physically present in the target environment to get the victim model in a failure state. Following this approach, we propose a new method to increase robustness of an a highway-driving DRL agent by implementing the adversarial training paradigm in the well-established highway-env [7] DRL environment. The method relies on adapting the iterative training procedure of the victim and adversarial agents, which is at the core of the adversarial learning (e.g., [8]).

The remainder of the paper is organized as follows. Section 2 describes the environment for the proposed training methodology. Section 3 highlights the core ideas behind our two-phase adversarial training process. Section 4 presents the experimental settings, while Sect. 5 discusses the results and Sect. 6 draws the conclusions.

## 2   Environment

We present a modified version of the highway-env [7] open-source simulator to develop a DRL training framework, given the simplicity and ease of customization which allows for quick prototyping and testing of different mechanics and details as well as a direct comparison to previous work [9, 10].

Figure 1 depicts the environment we used for developing our method. The scene contains a customizable number of heuristic Non-Player Vehicles (NPVs) navigating the scene alongside two vehicles playing as the part of the victim (VV) and the adversarial vehicles (AV), that are guided by two different DRL policies. Two sets of actions are provided for the proposed experiment: the lower-level action space already present in the original framework (faster, slower, lane change right, lane change left) and a high-level action space (keep lane, go to right-most lane, overtake) illustrated in [10]. The observations for both agents are named "kinematic", as per the original highway-env's definition and consist of presence (whether a vehicle is present or not in the scene), longitudinal and lateral position and velocity. These observations are taken for the ego vehicle and the $N = 7$ (for our experiments) nearest vehicles.



**Fig. 1.** Snapshot of the highway environment. The victim vehicle is colored in light-green, the adversary in purple, while the NPVs in light-blue.

## 3   Training Cycle Definition

For training the victim and adversarial agents, respectively driving VV and AV, we defined a two-phase cycle, which is depicted in Fig. 2. We begin the process by freezing a pre-trained victim model (particularly, the one illustrated in [10]), and let it act in pure

exploitation as a part of the highway environment. The goal of VV is to drive safely (i.e., no collisions with other vehicles) for as many kilometers as possible. In this first phase, the learner is AV acting upon the scene alongside with VV and the NPVs. The goal of AV is to hinder VV. To this end, the AV observation includes also the VV last observation and action. Unlike VV, which uses the aforementioned high-level action space designed to reach better "normal" driving performance, AV uses the original lower-level action space, which allows a greater degree of freedom in driving, which is needed to perform abrupt maneuvers to challenge VV. The idea behind the first phase of training is to build a strong adversarial policy and put the pre-trained vehicle in a danger/failure state.



**Fig. 2.** Outline of the proposed two-phase adversarial training cycle.

As the main effect of the first training phase, we expect VV to collide significantly more frequently than in the normal situation (i.e., without trained AV) and maintain a lower mean cruising speed. We conclude the first training phase when the metrics indicate a clear degradation of VV performance and proceed to the second phase. Here, AV is frozen and VV becomes the subject of the training. Thus, we resume the victim's model training through a continued learning procedure without any changes to the original reward function or state and fine-tune it in the adversarial scenario. The goal of this second phase is to make VV policy robust to the newly introduced adversarial attack, thus enhancing VV capability to deal with unexpected dangerous behaviors by other vehicles, that would previously induce a failure state of the agent. This two-phase process is then repeated a certain number of times, until the predefined stopping criteria is met (e.g., in terms of VV policy robustness) or results suggest using different hyper-parameter values.

## 4   Experiment

For the implementation, we opted for the Gymnasium Python framework (formerly known as OpenAI Gym [11]), which provides an intuitive API for DRL model training and environment configuration. The DRL model deployed for both victim and adversarial agents is the Proximal Policy Optimization (PPO) algorithm [12]. By repeatedly updating the policy parameters through a clipped surrogate objective function, PPO balances the exploration and exploitation learning stages, assuring stability and dependable

performance. The neural network at the core of the algorithm is a 2-layer multi-layer perceptron, with 64 neurons per layer for both analyzed models. The environment policy frequency (i.e., decision rate) is set to 1 Hz., to give the agents sufficient time to observe the effect of their previous action. We set the maximum training episode duration to 60 s. In both training phases an episode is terminated whenever a collision occurs. Also, when training AV, the episode is terminated if AV is overtaken by VV, since AV did not succeed in making VV fail. The number of NPVs is set to a randomly chosen value from 10 to 15 for each episode in both training phases.

The reward function (RF) is the key to the success of any DRL agent training, since RF numerically shapes the model behavior. For the first training phase (in which the VV policy is frozen), RF aims to reward the AV's ability to make VV fail. Thus, RF of AV is simply the opposite of RF with which VV was trained. That RF included rewards for high speed and right-most lane, and huge penalty for collision, as detailed in [10]. It is important to stress that the considered quantities for the AV training (i.e., speed, lane collision) are referred to the VV. However, the RF of AV additionally includes a penalty for its own collision with NPVs (a collision would prevent AV from continuing its disturbance task), while a rear-front collision with VV is not penalized (we do not reward it to prevent AV from learning to really chase for VV, which would be an unrealistic behavior). Finally, if VV manages to overtake AV, the episode is terminated, with a training penalty. In the second phase, the AV policy is frozen, while VV is fine-tuned with its "normal" RF, which is described in [10]. It is important to highlight that, in each phase, only one agent is trained, so to prevent non-stationary conditions that would hinder the training.

## 5   Results

The initial victim policy was trained for a total of 600k steps before being frozen for the first cycle of adversarial training. Phases 1 and 2 of the process took 300k steps each, for a total of 600k steps for the execution of the whole first cycle. For the preliminary experiment presented in this paper, we considered a single cycle iteration.

We evaluate the victim's policy before and after adversarial re-training. Particularly, the increase in robustness of the VV policy is evaluated by comparing the episode metrics reported in Table 1 averaging 1000 test episodes. Tests were conducted in the highway 3-lane environment in two conditions: with and without AV in the scene, while NPVs are always present. In the first phase, the emergent behavior adopted by our best AV model can be described as a continuous, yet not random, swerving procedure. As soon as AV detects the presence of VV behind, it visibly tries to move up or down the lanes to block VV. If possible, AV tends to steer at the very last possible moment to cut right in front of VV and thus have a chance of causing a "good" collision (rear-front). As shown in Table 1, this adversarial behavior leads VV to increase the collision rate by one order of magnitude (14.7% vs. 1.4% of the episodes) and lower its speed significantly when compared to the case without the AV. The VV model also shows a reduction in average deceleration and an increase in average acceleration. We argue that this is due to the fact that VV drives more slowly, due to the behavior of AV (less average deceleration), then it tries to abruptly accelerate to quickly perform the overtaking.

The second phase aims to improve the pre-trained policy of VV, which is the actual goal of the overall process, through a continued learning procedure. VV has now to learn to deal with AV, which hinders the normal behavior of VV through selective lane and speed changes. Comparing the third and first column of Table 1, we see that, the VV policy at the end of this phase is still negatively influenced by the AV presence, mainly in terms of collision number and average kilometers travelled, compared to the original VV without AV. The average left-lane change number is also increased, suggesting that the VV has learned to overtake vehicles more often in the adversarial scenario. On the other hand, comparing the second and third column, we see with no surprise that, in the adversarial scenario, the re-trained agent outperforms the original one in all the metrics.

**Table 1.** Performance over 1000 episodes of the victim models

| Considered metrics | Baseline | Baseline versus AV | Re-trained versus AV | Re-trained |
|---|---|---|---|---|
| Mean episode duration (s) | 59 | 53 | 57 | **60** |
| Number of collisions | 14 | 147 | 65 | **2** |
| Average travelled kilometers | 1.84 | 1.5 | 1.82 | **1.95** |
| Average speed (km/h) | 111 | 89 | **118** | **118** |
| Average deceleration (m/s$^2$) | $-5.38$ | $-4.52$ | $-5.47$ | **$-5.57$** |
| Average acceleration (m/s$^2$) | 1.62 | 1.84 | **1.28** | 1.34 |
| Average left-lane changes | 0.32 | 0.38 | 0.9 | **0.96** |
| Average right-lane changes | **1.33** | 1 | 1.09 | 1.27 |

Finally, testing the agent in a normal, non-adversarial scenario (fourth and first column), we see that the re-trained model collides less frequently, while traveling more kilometers at a slightly higher speed and accelerates less abruptly. This indicates a significant agent improvement in terms of safety and overall robustness.

## 6   Conclusions and Future Work

We explored the implementation of an adversarial learning procedure to increase robustness of a DRL agent for automated driving in a 3-lane highway-env simulated environment. Experimental results show that the re-trained model's policy has proven to be more robust and safer, outperforming the original one in all the metrics. These initial results

suggest that research should be conducted on the reward function, improving the realism of the AV disturbances, and loss function customization for the employed learning algorithm to stabilize the training phase. Moreover, the method may be transferred to more complex driving scenarios and simulators.

## References

1. Folkers A, Rick M, Buskens C (2019) Controlling an autonomous vehicle with deep reinforcement learning. In: 2019 IEEE intelligent vehicles symposium (IV), pp 2025–2031. IEEE, Paris, France
2. Bellotti F, Lazzaroni L, Capello A, Cossu M, De Gloria A, Berta R (2023) Explaining a deep reinforcement learning (DRL)-based automated driving agent in highway simulations. IEEE Access. 11:28522–28550. https://doi.org/10.1109/ACCESS.2023.3259544
3. Lazzaroni L, Bellotti F, Capello A, Cossu M, De Gloria A, Berta R (2023) Deep reinforcement learning for automated car parking. In: Berta R, De Gloria A (eds) Applications in electronics pervading industry, environment and Society. Springer Nature Switzerland, Cham, pp 125–130
4. Zhang H, Chen H, Xiao C, Li B, Liu M, Boning D, Hsieh C-J (2021) Robust deep reinforcement learning against adversarial perturbations on state observations. http://arxiv.org/abs/2003.08938
5. Pinto L, Davidson J, Sukthankar R, Gupta A (2017) Robust adversarial reinforcement learning. http://arxiv.org/abs/1703.02702
6. Gleave A, Dennis M, Wild C, Kant N, Levine S, Russell S (2021) Adversarial policies: attacking deep reinforcement learning. http://arxiv.org/abs/1905.10615
7. Leurent E (2018) An environment for autonomous driving decision-making. https://github.com/eleurent/highway-env
8. Goodfellow I et al (2020) Generative adversarial networks. Commun ACM 63:139–144. https://doi.org/10.1145/3422622
9. Campodonico G et al (2021) Adapting autonomous agents for automotive driving games. In: De Rosa F, Marfisi Schottman I, Baalsrud Hauge J, Bellotti F, Dondio P, Romero M (eds) Games and learning alliance. Springer International Publishing, Cham, pp 101–110
10. Pighetti A et al (2022) High-level decision-making non-player vehicles. In: Kiili K, Antti K, de Rosa F, Dindar M, Kickmeier-Rust M, Bellotti F (eds) Games and learning alliance. Springer International Publishing, Cham, pp 223–233
11. Brockman G, Cheung V, Pettersson L, Schneider J, Schulman J, Tang J, Zaremba W (2016) OpenAI Gym. http://arxiv.org/abs/1606.01540
12. Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O (2017) Proximal policy optimization algorithms. http://arxiv.org/abs/1707.06347

# Evaluation of AI and Video Computing Applications on Multiple Heterogeneous Architectures

Federico Rossi[1(✉)], Giacomo Mugnaini[1], Sergio Saponara[1], Carlo Cavazzoni[2], and Antonio Sciarappa[2]

[1] Department of Information Engineering, University of Pisa, Pisa, Italy
{Federico.Rossi,Giacomo.Mugnaini,Sergio.Saponara}@unipi.it
[2] Leonardo S.p.A., via Pieragostini 80, Genova 16149, Italy
{Carlo.Cavazzoni,Antonio.Sciarappa}@unipi.it

**Abstract.** This paper evaluates an AI video surveillance application on diverse high-performance computing (HPC) architectures. AI-powered video surveillance has emerged as a vital tool for security and monitoring, relying on hardware infrastructure for efficient processing. We present a benchmark of an AI application based on the YOLO object dection framework to track downed pepole in critical scenarios. This study investigates the impact of different architectural designs, including CPUs and GPUs on video analysis performance. Evaluation metrics encompass computational speed, power consumption and resource utilisation.

**Keywords:** hpc · Video processing · Object detection

## 1 Introduction

Video surveillance has grown significantly as a result of recent developments in artificial intelligence (AI), which have changed many areas. Applications for video surveillance that are AI-powered have become effective tools for boosting security and monitoring in public areas, vital infrastructures, and commercial buildings [1]. These programs make use of the skills of computer vision and machine learning algorithms to automatically scan video streams, spot anomalies, identify objects and people, and send operators real-time alerts.

However, the underlying hardware infrastructure has a significant impact on the efficiency and functionality of AI video surveillance applications [2,3]. The successful implementation of AI surveillance systems depends on efficient and accurate video processing, which is made possible by high-performance computing (HPC) architectures [4–6]. With their vast computational capacity and parallel processing capabilities, HPC systems may greatly improve the speed and accuracy of video analysis operations, allowing for real-time monitoring and prompt reaction to possible security risks [7].

The purpose of this article is to assess an AI video surveillance application using various HPC architectures. We'll investigate how various architectural layouts, such as conventional CPU-based systems and Graphics Processing Units (GPUs), affect the effectiveness and speed of video processing activities. We will evaluate these designs' advantages, disadvantages, and applicability for various surveillance scenarios by comparing them to a set of agreed-upon evaluation metrics.

To conduct this study, we leverage state-of-the-art deep learning frameworks and libraries. We present a video surveillance application based on the well-known YOLO object detection framework [8], paired with the DeepSort [9] tracking algorithm to detect and identify downed people in critical scenarios on the Pytorch framework. Through an exhaustive benchmarking procedure, we assess important performance indicators including processor power, resource use, and energy efficiency. In the context of AI video processing, findings from prominent computing systems like x86 and ARM are compared.

The paper is organised as follows: in Sect. 2 we present and detail the people down application, highlighting its components; in Sect. 3 we describe the benchmark setup; in Sect. 4 we report the benchmark metrics for the people tracking application.

## 2   People Down Tracking Application

### 2.1   Base Application

The man down application is designed to automatize the task of identifying people lying on the ground and classifying them as "resting" (i.e. not in a dangerous situation) or "man down" (i.e. injured or otherwise needing assistance). The application is based on the YOLOv5 [8] (an acronym for You Only Look Once) object detection + classification model, one of the most popular one-stage object detector algorithms, trained on a custom labelled man down dataset; the size of the model is the "small" one (YOLOv5s) and the commit tags considered have been the v5.0 (for the basic model) and the v6.1 (for TensorRT and ONNX extensions [10]) ones. The man down application workload can thus be divided into three main components:

– Pre-processing
– Inference
– Non-Maximum Suppression or NMS.

The pre-processing step performs simple preliminary operations when necessary, such as data type conversion and data normalisation; the inference step, which typically is the most computational expensive one, detects all possible instances of people lying on the ground and labels them as "resting" or "man down"; finally, the NMS step selects only the most probable instances and bounding boxes identified by the inference, thus eliminating all other less likely predictions.

Fig. 1. Computation flow for the YOLOv5 component

## 2.2 Tracking Application

The people down tracking application is an extended and more complex version of the previously described man down application which can also track and, if needed, count people lying on the ground in a video sequence, again based on YOLOv5 (size 'x', release 6.2). Its workload can be divided into five main components:

– Image pre-processing
– People detection
– Non-Maximum Suppression or NMS
– Man down classification
– Tracking.

Differently from before, people detection and man down classification have been separated into two different steps: the man down classifier, once received the YOLOv5 people detection output (bounding boxes, scores and class IDs), examines all the bounding boxes and selects all those who have an aspect ratio less than 1. The final step of tracking all the detected man down across different frames is instead implemented using DeepSORT [9,11] (an acronym for Deep Simple Online and Realtime Tracking), which takes advantage of an OSNet x1.0 [12] (a convolutional neural network used for person re-identification) in order to identify the same person in the video stream.

## 3    Benchmark Architectures

We evaluated the aforementioned deep learning application through Pytorch on several architectures. We employed architectures with and without GPUs, to evaluate also the performance of vector units inside the CPUs. The tested CPU-only architectures are:

– ARM Cortex-A72 CPU with four cores running at 1.5 GHz, with support for the 128-bit Neon SIMD extension and 4 Gbyte of RAM.
– HiFive Unmatched board powered by a SiFive U740 SoC running at 1.2 GHz for the scalar RISC-V unit.
– HPE Apollo80 system, powered by a Fujitsu A64FX, running at 2.2 GHz with support for the ARM 512-bit Scalable Vector Extension (SVE) SIMD unit.

– ARM Neoverse N1 CPU running at a maximum frequency of 3.3 GHz with support to 128-bit NEON SIMD instructions.
– Cascade Lake-based Intel Xeon Silver CPU running at 2.2 GHz with support for the 512-bit AVX SIMD instructions.

The tested GPU-based architectures are:

– NVIDIA Jetson Orin SoC with an Ampere-based GPU
– NVIDIA Tesla T4
– NVIDIA A100
– NVIDIA GTX1650Ti.

## 4   Benchmark Results

In this section we report the benchmark results for the platforms listed in Sect. 3. In Table 1 we report average number of frame processed per second for each platform on the people down application. Furthermore we detail, for each architecture, the timing performance for the three different components of the application. In Table 2 we report on the average resource utilisation of the architectures during the execution of the people down application and, for the architectures that allow it, we report in Fig. 2 the average performance-per-watt value for the different setups.

**Table 1.** Average FPS and time performance for the people tracking application with different architectures.

| Platform | FPS | YOLO (ms) | Man down classifier (ms) | DeepSORT (ms) |
|---|---|---|---|---|
| Cortex-A72 | 0, 1 | 7493 | 1, 2 | 1107 |
| Neoverse N1 | 0, 8 | 758 | 0, 6 | 503 |
| A64FX | 0, 4 | 1292 | 1,1 | 1054 |
| Arriesgado (Scalar) | 0,006 | 148987 | 2, 6 | 13835 |
| AGX Orin | 0, 05 | 2084 | 0, 5 | 18155 |
| + Ampere | 7, 7 | 38,8 | 0, 5 | 54,9 |
| i7-10750H | 1 | 807 | 0, 2 | 207 |
| + GeForce 1650Ti | 7, 3 | 85, 4 | 0, 3 | 11, 9 |
| Xeon | 1, 8 | 335 | 0, 3 | 197 |
| + Tesla T4 | 12, 8 | 37, 8 | 0, 3 | 15, 1 |
| Xeon | 3, 8 | 153 | 0, 4 | 97,5 |
| + Tesla A100 | 21, 3 | 10, 9 | 0, 3 | 13, 9 |

**Table 2.** Average resource utilisation for the people down benchmark for CPU and/or GPU in different architectures.

| Platform | CPU (%) | CPU (C) | CPU (W) | GPU (%) | GPU (C) | GPU (W) |
|---|---|---|---|---|---|---|
| Cortex-A72 | 96, 6 | 83 | – | – | – | – |
| Neoverse N1 | 5, 8 | 51, 7 | – | – | – | – |
| A64FX | 39, 5 | – | – | – | – | – |
| Sifive U740 | 79, 4 | 42, 5 | – | – | – | – |
| AGX Orin | 98, 3 | 58, 2 | 16, 8 | – | – | – |
| + Ampere | 45, 7 | 52, 2 | 6, 7 | 34, 1 | 47 | 16 |
| i7-10750H | 40, 8 | 91, 8 | 43, 2 | – | – | – |
| + GeForce 1650Ti | 32, 1 | 93, 5 | 36, 2 | 65, 8 | 78, 7 | 41, 2 |
| Xeon | 93, 5 | – | – | – | – | – |
| + NVIDIA Tesla T4 | 49 | – | – | 54 | 43, 7 | 62, 5 |
| Xeon | 50, 8 | 40, 3 | – | – | – | – |
| + NVIDIA Tesla A100 | 25, 9 | 32, 3 | – | 30, 2 | 38, 9 | 51, 9 |



**Fig. 2.** Performance per Watt for different architectures based on available metrics

## 5   Conclusions

In this paper, we provided a thorough assessment of an AI video surveillance application on various high-performance computing (HPC) architectures in this research. Our study emphasized the significance of hardware infrastructure for effective processing in AI-powered video surveillance by focusing on the application of the YOLO object identification framework to locate humans in dangerous situations. We looked into the effects of different architectural designs, including CPUs and GPUs, on the performance of the program for video analy-

sis. Important evaluation parameters like processing speed, power consumption, and resource use were included in our evaluation.

# References

1. Sharma V et al (2021) Video processing using deep learning techniques: a systematic literature review. IEEE Access 9:139489–139507. https://doi.org/10.1109/ACCESS.2021.3118541
2. Wang X (2013) Intelligent multi-camera video surveillance: a review. Pattern Recognit Lett 34(1). Extracting Semantics from Multi-Spectrum Video, pp. 3–19. issn: 0167–8655. https://doi.org/10.1016/j.patrec.2012.07.005, https://www.sciencedirect.com/science/article/pii/S016786551200219X
3. Dilshad N et al (2020) Applications and challenges in video surveillance via drone: a brief survey. In: 2020 international conference on information and communication technology convergence (ICTC), pp 728–732. https://doi.org/10.1109/ICTC49870.2020.9289536
4. Webb JA (1994) High performance computing in image processing and computer vision. In: Proceedings of the 12th IAPR international conference on pattern recognition, vol. 2—conference B: computer vision and image processing (Cat. No.94CH3440-5), vol 3, pp 218–222. ICPR.1994.577165. https://doi.org/10.1109/ICPR.1994.577165
5. Cococcioni M et al (2020) Fast deep neural networks for image processing using posits and ARM scalable vector extension. J R-Time Image Process 17(3):759–771. issn: 1861-8219. https://doi.org/10.1007/s11554020-00984-x
6. Cococcioni M et al (2021) Faster deep neural network image processing by using vectorized posit operations on a RISC-V processor. In: Kehtarnavaz N, Carlsohn MF (eds) Real-time image processing and deep learning, vol 11736. International Society for Optics and Photonics. SPIE, p 1173604. https://doi.org/10.1117/12.2586565
7. Yi S et al (2012) The model of face recognition in video surveillance based on cloud computing. In: Jin D, Lin S (eds) Advances in computer science and in-formation engineering. Springer, Berlin, pp 105–111. isbn: 978-3-642-30126-1
8. Jocher G et al, ultralytics/yolov5: v7.0 - YOLOv5 SOTA Realtime Instance Segmentation. Version v7.0. Nov. 2022.5281/zenodo.7347926. https://doi.org/10.5281/zenodo.7347926
9. Wojke N, Bewley A (2018) Deep cosine metric learning for person re-identification. In: 2018 IEEE winter conference on applications of computer vision (WACV). IEEE, pp 748–756. https://doi.org/10.1109/WACV.2018.00087

10. Bai J, Lu F, Zhang K et al (2019) ONNX: open neural network ex-change. github.com/onnx/onnx
11. Wojke N, Bewley A, Paulus D (2017) Simple online and realtime tracking with a deep association metric. In: 2017 IEEE international conference on image processing (ICIP). IEEE, pp 3645–3649. https://doi.org/10.1109/ICIP.2017.8296962
12. Ploco A, Rodriguez AM, Geradts Z (2020) Spatial-temporal omni-scale feature learning for person re-identification. In: 2020 8th international workshop on biometrics and forensics (IWBF), pp 1–5. https://doi.org/10.1109/IWBF49977.2020.9107966

# Edge Computing and Sensing

# Preliminary Frequency Response Analysis of a Contact Force Measurement System for Rail Applications

Giovanni Bellacci[1(✉)], Francesco Neri[1], Luca Pugi[1], Andrea Giachetti[2], Enzo Barlacchi[2], and Niccolò Baldanzini[1]

[1] University of Florence (DIEF), Florence, Italy
`giovanni.bellacci@unifi.it`
[2] University of Florence (DICEA), Florence, Italy

**Abstract.** This paper presents a finite element (FE) modeling and frequency response analysis of an instrumented wheelset for rail applications. The wheelset was equipped with strain gauges to estimate contact forces at the interface between the wheels and the rails. The objective is to showcase the simulation results of the wheelset state-space reduced model, which represents the dynamic behavior of the FE wheelset. The frequency response results of the FE and state-space models are compared. The paper demonstrates how the reduced model accurately represents the behavior of the FE system, resulting in significantly improved simulation and calculation speed, with a great reduction in memory usage. Furthermore, the state-space representation provides the flexibility to augment and manage data points for enhanced post-processing manipulation.

**Keywords:** Railway · Instrumented Wheelset · Modal Analysis · Frequency Response Analysis · Finite Element Method · State-Space

## 1 Introduction

To ensure a proper diagnosis of the under-investigation system state (e.g., train or infrastructure), two main acquisition techniques can be employed: wayside monitoring systems and onboard monitoring systems. In the case of onboard monitoring, portable, non-invasive devices, such as GNSS, accelerometers, Wheatstone bridges, etc., can be installed on the vehicle to gather data on both track and train. This enables a reliable diagnosis of the system's state. When such devices are equipped on the rolling stock, they form an instrumented wheelset that functions as a sensing system for the train. These technologies are used to measure wheel-rail contact forces for several purposes, such as monitoring of rolling stock to evaluate ride quality of the tested vehicle. Conversely some measurement system can be used to locate defects and identify the state of the line contributing to a faster and safer management of maintenance activities [1]. It is well-known from numerous studies [2, 3] that signals originating from wheel/rail irregularities and defects, such as squats, roughness, wheel flats, etc., can be captured. These

signals provide information to estimate their position along the track and ensure accurate and targeted predictive maintenance of the infrastructure or rolling stocks. To achieve this, understanding the frequency response of the wheelset plays a crucial role in inverse-forces identification from in-service high-speed trains. Conventionally, a bandwidth of 20–40 Hz limits data acquisition to low-speed information for running dynamic tests [4], while no bandwidth is still specified for diagnostic purposes on the norms. Improvements can be made to enable the consideration the high-frequency content related to defects. While various sensors can be installed on wheelsets, this paper focuses on a specific case where an existing wheelset equipped with strain gauges is replicated and simulated using the finite element method (FEM). The model and its frequency analysis results are then compared with a reduced state-space model of the same wheelset, extracted from the FEM model, using MATLAB (Fig. 1).



**Fig. 1.** Rendering of the instrumented wheelset.

## 2    Wheelset: Strain Gauges Layout

The instrumented wheelset was equipped with 16 strain gauges placed on the wheel web field sides. This configuration included a total of 2 Wheatstone bridges for each side of the axle shaft [5, 6]. Figure 2 illustrates the positioning of the strain gauges on the wheel web. The arrangement was as follows: 8 strain gauges (4 + 4), denoted as quadrature bridge "Q_x_s", symmetrically positioned relative to the vertical and horizontal axes of the face, where x indicates the orientation (vertical or horizontal), and s denotes the side (A or B for right and left wheel, respectively).

Referring to the scheme in Fig. 3a, strain signals originating from the wheel-rail interaction are treated as follows from Eq. (1)

$$U = \frac{VG}{4}(\varepsilon_a - \varepsilon_b + \varepsilon_c - \varepsilon_d) \tag{1}$$

where U is the Wheatstone bridge's output signal, V denotes the input tension, G is the gauge factor, and ε_(a-d) are the strains of the respective strain gauges (or equivalent strain in case of multiple series/parallel configurations of strain gauges in the same leg circuit leg). The root sum square (RSS) of the vertical and horizontal quadrature bridge

Fig. 2. Schematic representation of the instrumented wheelset strain gauges layout.

signals is calculated to obtain a continuous, nearly flat signal (otherwise sinusoidal-like, due to the alternating nature of the acquisition), that is independent of the wheel rotation angle [7]. Equation (4) provides an example of the combined signal corresponding to the Wheatstone bridges.

$$Q\_ve = \frac{VG}{4}(\varepsilon_1 - \varepsilon_4 + \varepsilon_8 - \varepsilon_5) \tag{2}$$

$$Q\_ho = \frac{VG}{4}(\varepsilon_3 - \varepsilon_6 + \varepsilon_2 - \varepsilon_7) \tag{3}$$

$$Q = \sqrt{Q\_ve^2 + Q\_ho^2} \tag{4}$$



Fig. 3. Wheatstone bridge schematic representation: a full-bridge configuration [6], b quadrature vertical (left) and horizontal (right) Wheatstone bridges.

## 3   Wheelset FE Simulation Frequency Response Analysis

A finite element model was developed to represent and study the wheelset under different load cases. It was determined that the chosen mesh sizing adequately approximated some reference results and no further refinements were necessary to enhance the precision and accuracy of calculated outcomes. A finer mesh was then employed for the instrumented components of the assembly such as the rolling stocks and shaft, while a coarser mesh was used for the non-instrumented parts, including the brakes sub-assembly. The deformation of the strain gauges was modeled as linear monoaxial strain according to Eq. (5), as the

distance difference between two points along the same line of the sensor patch, which represents the effective action line of the strain gauge divided for the undeformed initial distance. This was due to the impossibility of directly extracting deformations as outputs from the frequency analysis.

$$\varepsilon = \frac{l - l_0}{l_0} \tag{5}$$

Given the significant complexity of the actual system mounted on a bogie during operation, a preliminary study aimed to approximate the behavior of a free-free system was undertaken. A 'free-free condition' denotes a state where the system is entirely devoid of external mechanical constraints or boundaries. This allows it to move and behave naturally without any imposed limitations or restrictions. Such an ideal condition permits the observation of the natural frequencies of the studied object without disturbances arising from interactions with mechanical constraints. The focus was specifically on evaluating the performance of a single Wheatstone bridge, namely Q_ve_A quadrature vertical Wheatstone bridge, which is bonded to the wheel web field side, located on wheel "A". To begin the study, a modal analysis was conducted to determine the eigenvectors and eigenvalues (mode shapes and natural frequencies) of the model. Subsequently, a frequency response analysis was performed to examine the behavior of the wheelset in the frequency range from 0 to 1500 Hz, under various load cases and with different positions of input forces [8, 9]. The FE model employed in this study yielded a total of 32 outputs (2 points for each strain gauge) and 15 inputs, representing the forces applied in each direction (X, Y, and Q: longitudinal, lateral, and vertical, respectively) at five different locations on the wheel tread (Fig. 4).



**Fig. 4.** FE model of the instrumented wheelset.

## 4   Wheelset FE Simulation Frequency Response Analysis

Referring to the general theory of state-space (SS) reduced models, the n second-order equations representing and approximating the dynamic system, can be transformed into 2n decoupled first-order equations:

$$\{x\} = \left\{ \begin{array}{c} \{z\} \\ \{\dot{z}\} \end{array} \right\} \tag{6}$$

$$\begin{cases} \dot{x} = Ax + Bu \\ y = Cx + Du \end{cases} \tag{7}$$

where x represents the state vector, $z$ and $\dot{z}$ are the modal coordinates and velocities, A, B, C, and D correspond to the system (or state), input, output, and feedthrough matrices, respectively, and u is the input vector. The matrices and vectors were derived from simulations conducted in FE software and processed using MATLAB. Output signals y, which indicate Wheatstone bridges deformation, were obtained by solving the model equations. Transfer functions between the input u and output signals y were calculated for each Wheatstone bridge and input force. To analyze the Q_ve_A bridge, the initial MIMO (Multiple-Input Multiple-Output) setup was decomposed into a set of SISO (Single-Input Single-Output) state-space models. This involved considering one input force at a time for each of the four Wheatstone bridges and performing linear matrix operations on the state-space model. A transfer function was then obtained for the Q_ve_A bridge and compared with the results of the FEM simulation.



**Fig. 5.** Comparison of the state-space (blue) and FEM (red) transfer functions.

The two models are in accordance for low frequency response, up to approximately 250 Hz. The state-space calculated transfer function accurately depicts the first three represented natural frequencies in Fig. 5 (corresponding to 78.5 Hz, 121.6 Hz and 185.3 Hz respectively), relative to the bending deformation of the wheel web, and maintaining a correspondence with the fifth and sixth natural frequencies. For the Nyquist–Shannon sampling theorem, the FE model response was trimmed at 750 Hz. The amplitude is of minor interest due to the use of a first try damping factor of 0.001. It can be seen from Fig. 5 that the FEM-calculated transfer function trims resonance peaks due to resolution limitations, necessary to maintain the lowest computational cost and simulation time, while the state-space representation offers higher resolution, depending on a user input parameter on the MATLAB script, fully highlighting resonance peaks and anti-resonances.

## 5  Conclusions and Future Developments

This paper presents a comparison between the frequency response analysis results of a state-space model derived from a FEM-simulated instrumented wheelset and the FEM-simulated frequency response itself. The real system's approximate linearity enables its linear state-space representation and the subsequent application of superposition techniques. By leveraging linearity and modal superposition, displacements are evaluated as a linear combination of eigenmodes [10]. Output signals are treated as engineering strain due to the impossibility of directly extracting strains from the FEM simulation, which may affect result accuracy. A strong dependence was observed between the output nodes positions on the FE model, corresponding to the strain gauges, and their respective output signals. A strategic selection of output nodes can enhance strain values and, consequently, the simulation's accuracy. The study demonstrates the state-space model's capability to represent the dynamic behavior of the FE wheelset under free-free constrained conditions, significantly reducing computational time. The calculation time has been reduced from 16 h using FEM to just 45 s using the state-space model in MATLAB. Additionally, it improves memory usage, reducing it from 80 Gb in the low-resolution FEM simulation to just 79 kb in the MATLAB state-space script. Furthermore, the resolution of the results has been significantly enhanced, allowing users to manage the calculated transfer function with a user-defined number of points for more effective post-processing. The damping factor of 0.001 utilized in the FE model will be further calibrated in further validation activities. Future improvements of the model will focus on directly extracting strain output and developing a refined non-physics state-space model in the MATLAB environment. This will involve using the FE-obtained state-space as initial conditions for the definitive state-space identification. Subsequently, different constrained configurations will be modeled to replicate more realistic operating conditions, and the study will be extended to include all Wheatstone bridges. Moreover, impulse tests will be conducted using an instrumented hammer to validate the obtained models with experimental data.

## References

1. Bracciali A, Cavaliere F, Macherelli M (2014) Review of instrumented wheelset technology and applications. In: Pombo J (ed) Proceedings of the second international conference on railway technology: research, development and maintenance
2. Nielsen JCO (2008) High-frequency vertical wheel-rail contact forces-validation of a prediction model by field testing. Wear 2008(265):1465–1471
3. Gullers P, Dreik P, Nielsen JCO, Ekberg A, Andersson L (2011) Track condition analyser: identification of rail rolling surface defects, likely to generate fatigue damage in wheels, Using instrumented wheelset measurements. In: Proceedings of the institution of mechanical engineers, vol 225. Part F: J Rail Rapid Transit, pp 1–13
4. Licciardello RV, Broggiato GB, Bruner M, Corazza GR, Malavasi G (2006) A comparative study of the results of field tests carried out with different contact force measurement methods. In: Proceeding of the 7th world congress on railway research (WCRR). Montreal, Canada, pp 1–10
5. Alessandria M, Dotta B, Licciardello RV (2011) Long-term contact force measurements with the CML method. Scienza e Tecnica, Ingegneria Ferroviaria, pp 929–947, 11

6. Broggiato GB, Cosciotti E (2006) Ottimizzazione numerica della disposizione di ponti esten-simetrici sulle sale ferroviarie utilizzate nella misura delle forze di contatto ruota-rotaia. Ingegneria Ferroviaria 11:879–889
7. Gómez E, Giménez JG, Alonso A (2011) Method for the reduction of measurement errors associated to the wheel rotation in railway dynamometric wheelsets. Mech Syst Signal Process 25(8):3062–3077
8. Bracciali A, Rissone P (1994) Analisi modale sperimentale di una sala portante del convoglio ETR500. Ingegneria Ferroviaria 7(8):394–407
9. Ewins DJ (2009) Modal testing: theory, practice and application, 2nd edn. Mech Eng Res Stud Eng Dyn Ser
10. Pugi L, Abati A (2020) Design optimization of a planar piezo-electric actuation stage for vibration control of rotating machinery. Meccanica 55(3):581–596

# Assembly of Solder Beads with a Surface Mount Technology Resistor with Optoelectronic Tweezers and Freezing-Drying Techniques

Abdussalam Elhanashi[1(✉)], Sergio Saponara[1], Pierpaolo Dini[1], Qinghe Zheng[2], Abdurazak Saide[3], Weizhen Li[4], and Steven Neale[4]

[1] Department of Information Engineering, University of Pisa, Pisa, Italy
`a.elhanashi@studenti.unipi.it`
[2] School of Intelligence Engineering, Shandong Management University, Dingxiang Road 250357, Jinan, China
[3] University of Dayton, Corresponding Email:300 College Park, Dayton, USA
[4] University of Glasgow, James Watt School of Engineering, Glasgow, Scotland

**Abstract.** This research explores the use of optoelectronic tweezers (OET) as a method of manipulating microparticles. It focuses on how the dielectrophoretic (DEP) force can be used to trap solder beads with different concentrations of solution in an OET device, achieving velocities of up to 1350 µm/s. This study analyzes new techniques for creating complex patterns and examines how they are affected by various parameters, including light intensity, pattern size, and solution conductivity, to manipulate particles into desired structures. The freezing-drying technique has been used to fix solder beads with a surface mount technology (SMT) resistor to construct an electronic circuit, and we found that the freezing stage was sufficient when using Peltier to fix the assembled circuit. The paper discusses the temperature and time required for the heating process to connect the circuit by partially melting the solder and measures and analyzes the I-V characteristics to determine the manufactured device's resistance compared to the installed resistor. Overall, this study presents a promising result for microparticle manipulation and demonstrates its potential application in micro-electronic circuit construction or reworking.

**Keywords:** Optoelectronic tweezers · Microparticles · Solder beads · OET device · Freezing drying

## 1 Introduction

Micromanipulation technologies refer to a set of techniques used to manipulate and control microscale objects, such as cells and particles. Optical tweezers [1–3] use focused laser beams to trap and move particles, while acoustic tweezers [4, 5] use ultrasonic waves to create standing waves that trap particles. Dielectrophoresis (DEP) employs electric fields to control the movement of particles [6]. These technologies have wide-ranging potential applications in various fields, including biology, medicine, and material

science. Currently these techniques have had the greatest impact in bioscience where they have provided enhanced cell handling capabilities. In particular Optoelectronic Tweezers (OET) has the capability of trapping 15,000 traps simultaneously [7, 8], see Fig. 1. Touchless force is a methodology that can move objects without physically touching them. This work aims to assemble microelectronic components into a circuit using light-patterned electrical forces in OET [9]. The OET can be utilized to manipulate micro- and nanoscale particles, such as cells, nanowires, and carbon nanotubes. OET makes use of light-induced dielectrophoresis to enable this optically controlled manipulation. Dielectrophoresis is an electrokinetic force prompted by particles in a non-uniform electric field. OET combines the flexibility and control of optical manipulation with the parallel manipulation and sorting capabilities of dielectrophoresis [10]. The OET device is made up of two electrodes. The upper electrode comprises a thin layer of transparent conductive material and indium tin oxide (ITO) on a glass substrate. The lower electrode includes thin layers of ITO and amorphous silicon (a-Si) on a glass substrate. The amorphous silicon layer alters its electrical resistance in response to light, allowing the optical control of dielectrophoresis in the OET device. An aqueous solution containing the particles under manipulation is added between the electrodes. An electric field is created in the device by incorporating an AC bias across the bottom and top electrodes. Optical patterns are focused onto the photosensitive surface to develop the dielectrophoretic force [11]. Two electrodes with a spacer form a sample chamber. Layers of a-Si:H, ITO, and glass are at the bottom. A fluid buffer with microparticles, DI water, and Tween 20 (0.05%) is between electrodes. Impedance is high initially, but imaging reduces it, creating non-uniform fields. Polarizable particles move based on the electric field's polarity, allowing precise control of particle position with light in the OET device. [12].



**Fig. 1.** Schematic of the OET device, the light pattern creates the force to move and trap the metallic beads

The main contributions of this research are as the following:-

1. The proposed approach allows for the manipulation of a wide range of component sizes, shapes, and materials, enabling the assembly of complex microelectronics

structures that may have been difficult or impossible to assemble using traditional techniques.
2. The use of freezing-drying techniques during micro-electronics assembly can help to reduce the potential for damage to delicate components, particularly those that may be sensitive to high temperatures or physical stress.
3. The experiment aims to assess the I-V characteristics of micro-assembled electronic components and to determine the resistance of the manufactured device by comparing it to the original resistor value. This will involve measuring and analyzing the I-V characteristics of the components. A thorough evaluation of the experimental data will be conducted to verify the device's resistance.

In this research, we used OET to manipulate electrical components on an amorphous silicon substrate by assembling $40\,\mu m$ diameter Sn62Pb36Ag2 solder microspheres with a resistor. We employed a technique involving freeze-drying to remove the liquid medium after assembly. This involved freezing the liquid medium in the OET device and reducing the surrounding pressure, causing the frozen medium to sublimate directly from its solid to its gas phase. Hereafter, the paper is organized as follows: after the introduction in Sect. 1, Sect. 2 presents a related work for the proposed approach. Section 3 shows the methodology. Section 4 describes the discussion and experimental results. Conclusions are drawn in Sect. 5.

## 2 Related Work

Assembly of micro-electronics components with optoelectronic tweezers and freezing-drying techniques is a cutting-edge technology that enables precise and efficient manipulation of small-scale materials [13]. This technology combines light-based optoelectronic tweezers and freezing-drying methods for precise assembly of microelectronics, with applications in sensors, medical devices, and electronics. It offers non-contact, high-precision assembly, potentially revolutionizing microelectronics production for complex structures. Recent years have seen increased focus on light-based micro-assembly techniques for applications like sensing, imaging, and communication. This paper reviews key works showcasing light's role in manufacturing optoelectronic microstructures. In this paper, we review the most relevant works that demonstrate the use of light to manufacture optoelectronic microstructures. Assembly of multi-component structures" by Melzer et al. [14]. This recent study shows how optical tweezers, which use focused laser beams to trap and manipulate microscopic particles can be used for micro-assembly techniques utilizing light. Kristin et al. [15] demonstrated laser-induced forward transfer use to assemble microstructures with high precision and accuracy. Du et al. [16] presented an optical projection tomography technique for constructing microstructures, allowing for three-dimensional visualization and manipulation of the assembly process. Optical Alignment and Assembly of Nanoparticles" by Toh [17]. This work explored optical trapping and alignment to assemble nanoparticles into ordered arrays. Optoelectronic Tweezers have also been previously used to assemble silicon chips [18].

## 3   Methodology

In this experiment, 40 μm diameter Sn62Pb36Ag2 solder microspheres were assembled on an a-SI substrate. Various solutions were prepared with different conductivities, and when an AC voltage was applied, the particles were observed under microscopy. The velocities of the particles were recorded to determine how they can be trapped using different voltages and frequencies. The speed is an important parameter when using an OET unit. The unit was modeled as a lumped equivalent circuit containing impedance components representing the device and the liquid volume. An increase in the moderate conductivity reduced the impedance of this layer and led to a more significant voltage drop across the OET device. The stimulation benefits of the voltage drop throughout the center of an illuminated aSi layer were observed for moderate conductivities of 6, 11, and 16 mS/m with an AC signal of 30 V at 20, 18, and 12 kHz and light intensity of 2 W/cm$^2$. The voltage dropped primarily throughout the liquid layer for conductivities of 6 mS/m, resulting in the sample's best non-uniform electrical area. When conductivities were raised to 11 and 16 mS/m, the voltage primarily dropped with the OET, creating a low electrical area magnitude in the liquid. The OET device's functionality was tested whilst altering liquid conductivities using the above technique. The frequency and the complex permittivity of the medium and particle determine the bead manipulation. For smaller beads, their surface conductance and the electrical double layer affects the contaminants' conductivity, hence the DEP effect. Velocities of 1350 μm/s were achieved with a conductivity of 6 mS/m and a frequency of 18 kHz. The OET functionality in various liquid conductivities was observed, and it is essential when offering a sample suspended in a specific solution with high conductivities. The conductivity stability for a solution was recorded over time. (Fig. 2 illustrates the schematic diagram for the OET device).



**Fig. 2.** The schematic diagram for the experiment.

# 4  Discussion and Experimental Results

## 4.1  Assembly of the Resistor with Solder Beads

An a-Si substrate was produced by depositing an electrode and creating a gap of approximately 200 μm to connect an SMT component (a 10 MΩ resistor). The resistor was placed on the Si substrate using tweezers and solder beads with a diameter of 40 μm made of Sn62Pb36Ag2. An OET device was then assembled using ITO on the top and a-Si on the bottom to create a chamber connected to a power supply function generator and placed on a Peltier for freezing. A projector delivered light to create a DEP force, and a filter was used to reduce emission. A microscope was used to monitor the experiment live with magnifications of 4× and 10×, and images and videos were recorded for future analysis. The function generator was set to a voltage of 28 V p.p and a frequency of 18 kHz, and the projector was turned on to manipulate and trap the beads. The resistor was placed near the electrode using a manual tweezer and glass magnifier, OET trapping was used to assemble beads on both sides of the resistor to the electrodes. As the trapping was unable to be achieved with a resistor on the Si substrate, two beads were trapped and connected on the top of the resistor, and seven beads were trapped at the bottom of the resistor to the electrode. Finally, the Peltier was used to cool the substrate, which froze the OET device for approximately 35 min. The device was then transferred to a Nitrogen freezer, which was kept at a temperature of −79 °C for two hours to further increase the device's freezing. The OET device was then moved to the Eclipse dryer and kept for approximately an hour for sublimation. After the sublimation process, the ITO electrode was removed and inspected, and the manufactured structure was successfully created.

The device was then baked at a temperature of 179 °C for around four minutes to melt and connect the solder beads. The parameters and materials used to manufacture the device have been described in Table 1. The study demonstrates that the freezing-drying method can assemble SMT components (resistors) by manipulating particles and connecting them to deposited electrodes with solder beads in the OET device. The technique enables the creation of parallel and straight lines to form electronic circuits. A resistor of 10 MΩ is used as an SMT component in the circuit, and it requires delicate positioning using a glass magnifier and a tweezer. The experiment used a microelectronic resistor of 10 MΩ placed with solder beads in the device using the gap between ITO and Si as 100 μm. However, the resistor's thickness was higher than the gap between the electrodes, and it broke when placed on the ITO surface. Thus, the gap was increased to 200 μm by adding a double layer of tape on both sides of the a-Si substrate. The conductivity of the solution used in the experiment was DI water + 0.05% Tween 20. The experiment used maximum DEP by adjusting the function generator at 30 V and frequency at 18 kHz. However, the study found that the existing setup using the light from the projector and voltage supply of 30 V is insufficient to manipulate the resistor, as the DEP force is not strong enough with these conditions. The analysis suggests that DEP force, especially when the resistor is placed in the OET device with a 200 μm gap is insufficient to move such large particles.

**Table 1.** Parameters and material specifications used for manufactured devices.

| Method | Specification |
|---|---|
| Concertation of solution | DI water + 0.005% of Tween 20 |
| Substrate structure | ITO & substrate and gap of 200 µm |
| Resist type for a-Si layer | S1818 |
| Dimension for the device | L 2.5 cm, W 1.5 cm, T 1 mm |
| Freezing time in Peltier | 35 min |
| Solder Beads | 40 µm diameter Sn62Pb36Ag2 |
| Electrode Type | Ti/Au |
| Time for freezing | 120 min |
| Time for drying | 60 min |
| Temp for melting the beads | 179–180 °C |
| Time to melt the beads | 4 min |

## 4.2 Freezing and Drying Techniques

The process of removing ice from a substrate without letting it go through a liquid state, known as freeze drying, was utilized in this study. A Peltier device connected to a DC power supply was used as a freezer in the OET setup to freeze the liquid medium of the fabricated device after constructing the desired circuit. The liquid medium used was a solution of DI water and 0.005% Tween 20, which supported the manipulation of particles with less tension on the surface of a substrate and reduced the displacement of solder beads during the freezing process. In addition, the criticality of the freezing temperature was demonstrated, our initial experiments used the lowest value of $-7$ °C, resulting in bead displacement after the dryer. In this study, the Peltier device was used for up to 35 min. The substrate was then placed in a nitrogen refrigerator at $-79$ °C for 2 h, increasing the efficiency of fixing the beads and resistors without changing their positions.

## 4.3 The I-V Characteristic Measurement

We performed I-V measurements on the substrate using a Casca-de-Microtech MPS150 probe station at various stages of analysis. We aimed to identify opportunities for process improvement based on our findings as shown in Fig. 3, we used the gradient to calculate the resistance value of the resistor from the I-V curve, which yielded a value of 3.3 MΩ when the device was uncleaned, and the beads were not melted. After the beads were baked and melted, the resistance value decreased to 1.58 MΩ. When the device was cleaned with water, the resistance value dropped to 0.874 MΩ. We also measured the resistance of the a-Si substrate, which was 1 MΩ. However, since the 10 MΩ resistor was connected in parallel with the substrate resistance, we used the following formula

for calculating total parallel resistance, see Eq. (1).

$$\frac{1}{R} = \frac{1}{R1} + \frac{1}{R2} \tag{1}$$

$$\frac{1}{R} = \frac{1}{1} + \frac{1}{10} = 0.9\mathbf{Mohms}$$



(a)



(b)



(c)



(d)

**Fig. 3 a** Measurement of resistance before melting beads. **b** Measurement of resistance after melting beads. **c** Measurement of resistance after cleaning the device. **d** Measurement of resistance of the aSi substrate without the resistor

In measuring accuracy, percent error, also known as fractional difference, is used to quantify the discrepancy between the measured value E and the actual value A. The percent error can be calculated using the following formula, see Eq. (2):

$$\% \text{ Error} = \frac{|E - A|}{A} \tag{2}$$

$$(0.874 - 0.9)/0.9 = 2.8\%$$

## 5  Conclusions

This research introduced a novel technique for manipulating solder beads at the particle level using optoelectronic tweezers and freeze-drying. It showcased the precision of this method in crafting small-scale electronic components through Dielectrophoresis forces. Multiple solder beads could be handled concurrently, aiding precise assembly with other electronic parts. This marks a significant leap in using OET for circuit manufacturing over conventional methods. The study used ITO and amorphous silicon structures, conducted experiments, and addressed ion concentration-related issues. It also investigated the I-V characteristics of the device post freeze-drying, revealing the influence of solder bead melting and cleaning on resistance measurements, eventually reaching the anticipated values only after thorough cleaning.

## References

1.  A Ashkin J Dziedzic T Yamane 1987 Optical trapping and manipulation of single cells using infrared laser beams Nature 330 769 771 https://doi.org/10.1038/330769a0
2.  D Grier 2003 A revolution in optical manipulation Nature 424 810 816
3.  R Dasgupta SK Mohanty PK Gupta 2003 Controlled rotation of biological microscopic objects using optical line tweezers Biotechnol Lett 25 1625 1628
4.  Xiaoyun Ding Zhangli Peng Sz-Chin Lin Michela Geri Sixing Li Peng Li Yuchao Chen Ming Dao Subra Suresh Tony Huang 2014 Cell separation using tilted-angle standing surface acoustic waves Proc Natl Acad Sci USA 111 https://doi.org/10.1073/pnas.1413325111
5.  Feng Guo Zhangming Mao Yuchao Chen Zhiwei Xie James Lata Peng Li Liqiang Ren Jiayang Liu Jian Yang Ming Dao Subra Suresh Tony Huang 2016 Three-dimensional manipulation of single cells using surface acoustic waves Proc Natl Acad Sci 113 201524813 https://doi.org/10.1073/pnas.1524813113
6.  Pei-Yu Chiou Aaron Ohta Ming Wu 2005 Massively parallel manipulation of single cells and microparticles using optical images Nature 436 370 2 https://doi.org/10.1038/nature03831
7.  HA Pohl K Pollock JS Crane 1978 Dielectrophoretic force: a comparison of theory and experiment J Biol Phys 6 133 160 https://doi.org/10.1007/BF02328936
8.  S Neale A Ohta H Hsu J Valley A Jamshidi M Wu 2009 Trap profiles of projector-based optoelectronic tweezers (OET) with HeLa cells Optics Express 17 7 5231 5239 https://doi.org/10.1364/OE.17.005231
9.  R Pethig 2010 Review article—dielectrophoresis: status of the theory, technology, and applications Biomicrofluidics 4 022811
10. Loire S (2008) Manipulation of microsize to nanosize particles with AC electrokinetic forces (Order No. 3342025). Available from SciTech Premium Collection. (230793349). https://www.proquest.com/dissertations-theses/manipulation-microsize-nanosize-particles-with-ac/docview/230793349/se-2
11. MC Zhong XB Wei JH Zhou ZQ Wang YM Li 2013 Trapping red blood cells in living animals using optical tweezers Nat Commun 4 1768
12. Zhang S, Juvert J, Cooper JM, Neale SL (2016) Manipulating and assembling metallic beads with Optoelectronic Tweezers. Sci Rep

13. S Zhang Y Liu Y Qian W Li J Juvert P Tian 2017 Manufacturing with light-micro-assembly of opto-electronic microstructures Opt Express 25 28838 50 https://doi.org/10.1364/oe.25.028838

14. JE Melzer E McLeod 2021 Assembly of multicomponent structures from hundreds of micron-scale building blocks using optical tweezers Microsyst Nanoeng 7 45 https://doi.org/10.1038/s41378-021-00272-z

15. Charipar KM, Díaz-Rivera RE, Charipar NA, Piqué A (2018) Laser-induced forward transfer (LIFT) of 3D microstructures. In: Proceedings of SPIE 10523, Laser 3D Manufacturing V, 105230R. https://doi.org/10.1117/12.2294578

16. W Du C Fei J Liu Y Li Z Liu X Zhao J Fang 2020 Optical projection tomography using a commercial microfluidic system Micromachines 11 293

17. MC Toh WK Phua EH Khoo 2020 Optical alignment of achiral nanoparticles via induced chiral force SN Appl Sci 2 1428

18. Plochowietz A et al (2019) Programmable micro-object assembly with transfer. In: 2019 20th international conference on solid-state sensors, actuators and microsystems & Eurosensors XXXIII (Transducers & Eurosensors XXXIII), Berlin, Germany, pp 390–393. https://doi.org/10.1109/TRANSDUCERS.2019.8808800

# A Low Cost Open Platform
# for Development and Performance
# Evaluation of IoT and IIoT Systems

Massimo Ruo Roch$^{(\boxtimes)}$ and Maurizio Martina

Politecnico di Torino, Dipartimento di Elettronica e Telecomunicazioni, Torino, Italy
`massimo.ruoroch@polito.it`

**Abstract.** The Internet of Things (IoT) paradigm is nowadays pervasive in a variety of applications. Distributed computing, i.e., local processing of data inside IoT nodes is mandatory to reduce power consumption and data communication cost, too. An evaluation of the needed computing power and of the hardware requirements is then as much useful as it is conducted in an early stage of the product development. The possibility to perform on-the-field experiments at an early stage is a desirable feature, too. In this paper, a low cost and open system for IoT evaluation is described. It contains basic blocks of an IoT/IIoT system, and software and hardware structures aimed to assess its performances with negligible overhead. Power consumption and timing characteristics of firmware (MCU) and hardware (FPGA) IoT nodes components can be collected and analyzed, to obtain real system requirements.

**Keywords:** IoT · IIoT · Ubiquitous computing · Artificial intelligence · BLE · EtherCAT

## 1 Introduction

The design of a new Internet of Things system is a complex problem, due to the overwhelming number of degree of freedom. In fact, this kind of application is based mainly on computation and communication tasks, which must be distributed between a variety of objects of different types [1,2]. As an example, a possible deployment is built up by a multitude of smart nodes, communicating through a low power protocol with a smart gateway. The latter, in turn, exchange data with servers located in remote data centers. This scenario is depicted in Fig. 1.

The represented topology is just an example, as intermediate layers can be added, or even removed. Moreover, power consumption and computing power can vary according to environmental factors (energy sources availability) and strongly influences the location of computing tasks.

Real life technology constraints must be considered, too, as the usage of off-the-shelf modules, typical in this field, can introduce far then optimal performances at the interface level between building blocks. As an example, better

**Fig. 1.** Typical structure of an IoT system: Servers, on the left, are connected through internet to a gateway (black box) which in turn communicates with end-nodes

described later, a commonly used EtherCAT slave IC has nominal 80 Mb/s QSPI interface, but real performances could drop off to less than 1 Mb/s, due to synchronization delays related to the internal chip architecture.

An early evaluation of the performances of the used modules, and of the related interfaces (wired or wireless) is then mandatory to be able to design the overall system, and to assess its feasibility for the desired system specifications. In fact, the design an IoT system can be a high risk process, as deeply discussed in [3].

Performance assessment is typically conducted in a lab, through the usage of common measurements instruments (MSOs, Logic state analyzers) on early prototypes of the nodes. But this approach requires money, and even worse, time, as it requires the physical realization of the nodes and of the measurement bench.

A possible alternative is the usage of so-called software models able to mimic the behaviour of parts of the system. Anyway, the unavailability of suitable models for some components is a not avoidable drawback of this methodology. Precision of results related to the interactions between different parts of the system is a concern, too. Problem complexity can be somewhat reduced introducing *standard* architectures, as described in [4], but a complete solution is currently missing.

Last, both approaches can not be easily adapted to true on-the-field tests, leading to a relative uncertainty in the obtained results.

Industrial IoT is a quite recent development of the IoT paradigm, sharing similar characteristics, but adding specific requirements from the point of view of reliability, safety, and robustness. In the following, the term IoT is used to indicate IIoT, too, as the proposed solution is applicable to both fields without changes. In fact, hardware components used to build IIoT devices are exactly the same used for consumer IoT.

In this paper, a new methodology to early assess real system performances is presented. It is based on the adaptation of the VirtLAB board [6], initially developed by the authors for teaching tasks, to optimize it for IoT performance data collection and analysis. This board has the peculiarity to host both software defined measurement instruments, and the hardware on which the application

runs. It means that with a reduced cost, and in a reduced space, it is already present everything needed to operate an IoT node or gateway, and to collect performance data. If needed, data analysis can be performed locally, too.

Both timing and power performances can be evaluated with the proposed system, in lab, and on-the-field. Last, in the context of Open Science and Open Innovation [5], the system is based on an open-source approach for both the hardware and the software, i.e., the system is freely available for research and development purposes, just with credits to authors.

## 2 Methodology and Implementation

Distributed computing system implementation implies the coexistence of different computing and communication tasks on heterogeneous modules. These modules interface each other through well-known standard protocols, with known performances. But the overall performance can not be obtained just by these data. In real systems, unexpected interactions and synchronization requirements, both software and hardware, often arise, leading to real performance degradation.

Firmware and hardware changes needed to satisfy system requirements in a late design phase are a problem, too, as they can invalidate previous choices (memory footprint, clock frequencies, battery sizes, etc.).

Extensive simulations through the so-called *digital-twins* are then sometimes used, but they need extra-time, and require the availability of suitable simulation models and input data.

To overcome estimated performance uncertainties, IoT nodes are often oversized, leading to undue cost increments. An alternative approach, where feasible, is to adapt overall system specifications to the obtained ones in the deployed system.

The proposed solution is to develop a flexible and scalable system, able to output real performance data for IoT nodes and gateways. It must be used in the lab, at the beginning, but also on-the-field, for a more precise assessment.

The basic idea is that IoT nodes are typically built up by the following building blocks:

- Mixed-signal microcontrollers
- Specialized hardware accelerators
- Sensors and actuators
- Communication interfaces, wired or wireless
- Power sources.

Also gateways are typically coincident with nodes, or they include a subset of the preceding building blocks.

Performance data collection must be extracted from the running hardware with minimal interaction, to not change significantly the system behaviour. This implies the usage of as fast as possible hardware methodologies, linking the measurement appliances to the application. Minimal required measurements can be summarized as such:

– Time difference measurement
– Event count
– Power consumption.

Hardware and software events can be monitored just toggling a simple digital I/O line, as this is the technique which guarantees less overhead at all, up to a single clock cycle, if carefully implemented. Power consumption must be externally measured through sensing of the power supply rails.

The VirtLAB board [6] is a low cost and small size open system, developed by the authors to overcome teaching restrictions during the CoVid pandemic. It contains two sections: a user and a master ones. The user section contains a mixed-signal microcontroller and a low power FPGA, connected through a 32 bit I/O bus. Moreover, an Arduino compatible shield socket is connected to the bus, allowing to seamlessly connect existent shields. The master section is built-up by an MCU, an FPGA, volatile (DRAM) and not volatile (QSPI flash) storage. This section has access to the 32 bit user I/O bus, too, and is then able to perform time measurements on it. Current drawn by the user MCU and user FPGA can be measured in real time by the master section, too.

The hardware features of the VirtLAB board can then be easily mapped to what is required according to the preceding lists:

– Hardware modules—These blocks are built up in the user section of the Virt-LAB board.
  • Mixed-signal microcontrollers—The user MCU perfectly accomplishes this task. Moreover, it is a low power device, which can be put in several sleep modes, mimicking the desired behaviour of a real IoT MCU.
  • Specialized hardware accelerators—The user FPGA can be used to implement these kinds of functions (DSP, AI, etc.).
  • Sensors and actuators—Several Arduino shields are available, with a variety of sensor and actuator types (environmental, inertial, motor control, etc.). Even if a specified device wasn't available, it would be relatively easy to design a Specialized shield to be interfaced to the VirtLAB board.
  • Communication interfaces, wired or wireless. Arduino shields are available for the more common communication standards, ranging from BLE and similar short range protocols, up to LoRa and GSM/LTE, for long range communications. Wired communications shield are available, too, for the more common industrial standards (RS232/485, CAN, EtherCAT, etc.)
  • Power sources—The VirtLAB board has an independent power supply, which derives local voltages from a simple 5V USB input.
– Measurements—Data collection and analysis is implemented in the master section of the VirtLAB board.
  • Time difference measurement—Two different solutions are available on the board. The first one, less complex but with limited performances, is implemented using the master MCU hardware timer. It is an integrated peripheral able to record the exact time at which a selected input state is changed. Up to 4 different inputs can be sampled at the same time. Time

resolution is 12.5 ns, due to the internal 80 MHz clock speed, but there is a limit to the maximum frequency of events, as a DMA memory access is needed to record current time in response to an input event, leading to a bandwidth in the order of 10 MB/s. The second solution, a high performance one, is based on dedicated circuitry implemented inside the master FPGA. In this case, all the 32 bit of the I/O bus can be monitored and recorded in a high speed HyperRAM. For sake of simplicity, resolution is held at 12.5 ns, but there is no limit to the recording speed, as data transfer toward memory will be on a high speed bus with burst transfers, guaranteeing up to 160 MB/s bandwidth.

- Event count—The two solutions implemented for time difference measurements are used for event counting, as well.
- Power consumption—Power supply lines to the user MCU and to the user FPGA have independent sense resistors and amplifiers for each rail. The outputs are sampled by dedicated ADCs channels in the master MCU, and the converted results are saved in memory through DMA accesses. Bandwidth is 350 kHz, limited not only by the current sense amplifiers, but by power decoupling capacitors, too. These are unavoidable, due to the strict requirements on power supply voltage ripple of the FPGA.

## 3    Usage Examples

To assess the feasibility of the methodology described in the preceding section, an IIoT node has been implemented on the designed system, and the methodology has been used to extract real performances.

The first case of study has been the development of an EtherCAT slave node. The EtherCAT standard uses Ethernet technology in real time industrial applications, using specialized hardware to overcome the timing uncertainties present in the Ethernet protocol [7]. A dedicated EtherCAT Slave Controller integrated circuit (ESC), connected to the network, is interfaced to a common microcontroller. A critical factor, impacting overall system performances, is the time lag between hardware events and the reaction of the firmware running on the MCU. But datasheets are not accurate enough to analyse the impact of hardware and firmware latencies on real behaviour.

In our test implementation, a LAN9252 Microchip ESC is hosted on an Arduino shield, and the application microcontroller has been replaced by the VirtLAB user MCU. The implemented system is shown in Fig. 2.

The ESC uses an SPI bus to communicate with the MCU, through a register based interface. The application runs on the MCU, but a FreeRTOS kernel is used, to mitigate system timing complexity. EtherCAT packets are read from the network, are decoded by the ESC, and then generate hardware events to the MCU (interrupts). These events trigger the activation of firmware real-time tasks. The developed system has been used to extract accurate timing information about application behaviour. Data are sent to a PC in textual form, through a USB based virtual serial.

**Fig. 2.** VirtLAB board (green) with EtherCAT shield (red) installed

Latencies due to hardware and firmware interactions has been extracted, and some of the data are visible in Tables 1 and 2.

**Table 1.** Latency of FreeRTOS interrupt to task, with different strategies

| Metrics | Value |
|---|---|
| Interrupt to task delay with binary semaphores | 10 µs |
| Interrupt to task delay with direct to task notifications | 8.7 µs |
| Interrupt to task delay with stream buffers | 19.2 µs |
| Interrupt to task delay with polling on static shared variable | 320 ns |

**Table 2.** Timings of ESC registers read, with direct to task notifications

| Metrics | Total time | SPI bus time |
|---|---|---|
| SPI direct register read | 10 µs | 15 µs |
| SPI indirect register read | 62 µs | 15 µs |
| SPI indirect register read with auto-increment | 40 µs | 15 µs |

The obtained results have been verified through the usage of standard laboratory equipment, such as oscilloscope, and multimeter.

## 4   Conclusions and Future Work

In this paper, a methodology to assess real world IoT system performances has been demonstrated. A first usage example, with quantitative data relative to

a timing analysis, has been shown. In the next future, other examples will be developed, including power consumption data, too, to further validate and refine the designed firmware and hardware. A GUI running on the PC will be a near future development, too.

# References

1. Sethi P, Sarangi SR (2017) Internet of Things: architectures, protocols, and applications. J Electr Comput Eng (gen. 2017), 1–25. https://doi.org/10.1155/2017/9324035
2. Vashi S, Ram J, Modi J, Verma S, Prakash C (2017) Internet of Things (IoT): A vision, architectural elements, and security issues. In: 2017 international conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC), pp 492–496. https://doi.org/10.1109/I-SMAC.2017.8058399
3. Lee B, Cooper R, Hands D, Coulton P (2019) Value creation for IoT: challenges and opportunities within the design and development process. In: Living in the Internet of Things (IoT 2019). London, UK, pp 1–8. https://doi.org/10.1049/cp.2019.0127
4. Krčo S, Pokrić B, Carrez F (2014) Designing IoT architecture(s): A European perspective. In: IEEE World Forum on Internet of Things (WF-IoT). Seoul, Korea (South), pp 79–84. https://doi.org/10.1109/WF-IoT.2014.6803124
5. European Open Science Cloud (EOSC). www.eosc-portal.eu/
6. Ruo Roch M, Martina M (2022) VirtLAB: a low-cost platform for electronics lab experiments. Sensors 22:4840. https://doi.org/10.3390/s22134840
7. Rostan M, Stubbs JE, Dzilno D (2010) EtherCAT enabled advanced control architecture. In: IEEE/SEMI advanced semiconductor manufacturing conference (ASMC). San Francisco, CA, USA, pp 39–44. https://doi.org/10.1109/ASMC.2010.5551414

# A Compact Continuous Analyzer of Particulate Matter Radioactivity

Christian Riboldi[1,2], Daniele M. Crafa[1,2], Carlo Fiorini[1,2], and Marco Carminati[1,2(✉)]

[1] Dipartimento di Elettronica, Informazione e Bioingegneria, Politecnico di Milano, Milan, Italy
marco1.carminati@polimi.it
[2] Istituto Nazionale di Fisica Nucleare, Milano, Italy

**Abstract.** We present a microcontroller-based in-flow airborne spectrometer to quantify the presence and activity of radioactive isotopes in particulate matter (PM). It is based on a hollow CsI(Tl) scintillator inside which air flows, thanks to a fan, and PM is captured by a micrometric filter placed in the middle, thus allowing full solid angle capture of gamma rays. The instrument is very compact, fitting in a 10 cm cube and thus being compatible with several unmanned flying vehicles, including an experimental rocket inside which it will be launched as scientific payload to evaluate the robustness of its design during flight. Here we report the mechanical and electronics (for SiPM readout) design, as well as a preliminary characterization. Simulations have allowed to find the optimal compromise in the number and position of SiPMs with an expected energy resolution better than 7% at 662 keV.

**Keywords:** Gamma Spectroscopy · Silicon PhotoMultiplier · Microcontroller

## 1 Introduction

The continuous measurement of the radioactivity of airborne particulate matter (PM) is routinely performed by national authorities for environmental protection, often managing networks of fixed ground stations [1]. This allows to monitor the effects of release in the atmosphere of radioactive material due, for instance, to nuclear accidents or weapons. The traditional technique is based on sampling large volumes of air pumped through filters that collect PM. Periodically, the Teflon filter is then placed in a shielded box (the so-called Marinelli beaker) coupled to a gamma-ray spectrometer to measure the presence and activity of the captured radionuclides. Automatic samplers exist, but their cost and bulkiness limit their diffusion.

In this paper we present the design and characterization of a portable gamma-ray spectrometer for continuous monitoring of radioactivity of PM and fast deployment. This work is characterized by two key elements of novelty: (i) the hollow scintillator design, allowing high sensitivity and continuous-flow sensing and (ii) the compactness thanks to a volume of 10 cm × 10 cm × 10 cm that enables its transportation in the atmosphere as payload in drones, balloons and even experimental rockets (Fig. 1). In fact, this project was selected as scientific payload of the experimental rocket Gemini developed by the

Skyward association of students of Politecnico di Milano that will compete in the EuRoC international competition in October 2023. Figure 1c shows Pyxis, the previous version of the rocket, highlighting the location of the payload that will descend after reaching the apogee with an independent parachute).



**Fig. 1.** Fast deployment modes for the proposed instrument: **a** on a balloon for long (days) missions, on a **b** drone or **c** rocket for short ones.

Several radionuclides of interest can be monitored; we chose to focus on the most relevant ones for environmental monitoring: $^{131}$I (364 keV), $^{106}$Ru (622 keV) and $^{137}$Ce (662 keV), limiting the energy range to fit the $^{137}$Cs photopeak with suitable energy resolution better than 10%, typically provided by low-cost handheld monitors. Anyway, a gain switch allows to double the energy range, extended to 2 MeV to include also $^{134}$Cs (1.6 MeV).

## 2 Instrument Design

The design of the instrument, dictated by the CubeSat volume constraints, is shown in Fig. 2: the hollow scintillator (60 mm in diameter and 40 mm in thickness, custom designed by us and manufactured by the company Proteus, USA) is coupled to SiPMs along four flattened windows. The other surfaces are covered with a diffusive coating to contain the scintillation light in the crystal. This design was inspired by a continuous crystal for preclinical PET [2]. CsI(Tl) was chosen as compromise between energy resolution and sensitivity. The PM filter is blocked with a plastic holder in the middle of the hole thanks to a step in the internal diameter (from 14 to 10 mm). A centrifugal fan (with a flow rate of 0.4 m$^3$/min) is placed behind the scintillator creating the aspiration of air from the front inlet and where a solid-state flow rate sensor (FSR3000 by Renesas) is placed to compute the total amount of scanned air. The selection of this fan represents again a compromise between achievable flow rate (the higher to better, to shorten the

air sampling time or maximize the acquired volume during a fixed scan time) and the power dissipation. A 6.6 V two-cell LiFe battery compliant with the weight (85 g) and capacity (1450 mAh) requirements of the rocket was selected for this prototype.

The architecture of the electronics (Fig. 3) is an evolution of the one we developed for a compact USB-powered gamma-ray spectrometer for scrap metal screening [3]. Here the mainboard contains the control microcontroller (STM32), an SD card to store data acquired during flight, a MEMS inertial measurement unit with triaxial accelerometer and gyroscope (used to estimate the module attitude, in particular during the rocket flight, and thus be able to correlate the detected events with the different orientations and phases of the mission) and the DC/DC converters to bias the SiPM (~30 V) and the fan (15 V). Four identical SiPM boards are coupled to the flattened windows in the reflective coating, hosting SiPM, input transimpedance amplifier (TIA), filter and peak stretcher, along with a temperature sensor (to compensate the SiPM gain drifts) and an LED (emitting at 465 nm) to inject light pulses to calibrate and periodically test the acquisition chain.



**Fig. 2.** Module design: hollow scintillator coupled to a centrifugal fan for air aspiration and in-flow spectroscopy of PM captured in the inner filter (Pb shielding not shown).

In order to size the readout chain and select the number of SiPMs to be coupled to the scintillator as in [3] (we selected arrays of 2 SIPM model AFBR-S4N66P024M by Broadcom with area of 13.5 mm × 6.54 mm), we performed numerical simulations by means of the software package ANTS2 [4], specifically developed for Anger cameras, to compare different configurations. The results of simulations are reported in Fig. 4: the full coverage (a) of the flattened windows (24 SiPM) would give an ideal energy resolution of 5.1% at 662 keV. The represents the best possible result. On the other hand,

**Fig. 3.** Architecture of the instrument electronics consisting of a motherboard with MCU and four identical SiPM front-end boards.

a single SiPM per side would give 8.8% (c), while 3 SiPM per side would give 6.8% (b). The latter was chosen as the optimal compromise, in line with the performance of other portable scintillator-based spectrometers with resolution around 10% [3]. Within geometry (b), the exact location of the 3 SiPMs in the window (in the middle, centered or edges) was also simulated. Since no significant difference in energy resolution emerged (see bottom spectra), the edges were chosen in order to grant better mechanical stability. As shown in the spectra at the bottom right of Fig. 4, each TIA connected to 3 SiPMs will ideally collect an average of 6000 photoelectrons for $^{137}$Cs (the centroid of the Gaussian peak in the histogram of collected events) and, correspondingly, the chain is sized to read a maximum of 8000 photons. An optical glue matching the refraction index of the scintillator is used to attach the SiPM to the flat face of the crystal.

Figure 5 shows the realized boards: the circular motherboard is directly plugged to the four SiPM boards attached to the scintillator. An auxiliary board hosting the SD memory and the USB connector for PC interfacing is located on the side of the cube for easier accessibility. The event processing time (filtering and acquisition with a 12-bit ADC, rebinning the spectrum histogram to the standard of 10 bits) is expected to be in the order of tens of microseconds, leading to an expected max. Count rate in the ~kcps range, definitely compatible with environmental monitoring applications. Considering the battery capacitance and the power consumption, we expect an operating time of about 3 h. For longer missions a bigger battery could be employed.

The external frame is made of water-cut carbon fiber plates, while the internal mechanical supports, holding the scintillator, boards and fan, are 3D printed. Figure 6 shows a mechanical dummy in the assembly phase.

## 3 Conclusions

We have presented the novel design of very compact (fitting in a single CubeSat unit) spectrometer based on a hollow scintillator for continuous-flow capture and analysis of airborne dust in terms of radioactivity. The instrument is conceived, in particular, for

**Fig. 4.** Simulations of the energy resolution at 662 keV achievable for different configurations of SiPM: the optimal trade-off (6.8%) is with 12 SiPM placed at the edges (**b**).



(a) Motherboard                (b) SiPM Board                (c) Interfaces Board

**Fig. 5.** System PCBs: **a** motherboard hosting the MCU and DC/DC converters, **b** SiPM board hosting the analog front-end, **c** interface PCB with SD flash memory and USB ports.

fast deployment in case of accidents entailing the release in the atmosphere of radioactive particulate matter. The mechanical and electronics design were custom tailored on the specifications, in particular challenging volume ($10 \times 10 \times 10$ cm$^3$) and weight constraints. In fact, for the rocket competition, the maximum payload weight is 1.3 kg met by this design with the exception of the lead shielding required to reject external effects. Simulations have allowed to choose the optimal trade-off between the number of SiPM (12 arranged into 4 tile boards coupled to the scintillator) and corresponding energy resolution (7% at 662 keV, aligned with the state of the art [3]).

In conclusion, we envision the embedment of machine learning in the instrument MCU [5] for real-time processing of the acquired spectra for automatic identification of the isotopes [6].

**Fig. 6.** Mechanical dummy of the CubeSat prototype to check assembly of all components.

# References

1. Bem H et al (2002) Determination of radioactivity in air filters by alpha and gamma spectrometry. Nukleonika 47(2):87–91
2. Freire M et al (2021) Experimental validation of a rodent PET scanner prototype based on a single LYSO crystal tube. IEEE Trans Radiat Plasma Med Sci 6(6):697–706
3. Carminati M et al (2022) Handheld magnetic-compliant gamma-ray spectrometer for environmental monitoring and scrap metal screening. Sensors 22(4):1412
4. Morozov A et al (2012) ANTS—a simulation package for secondary scintillation Anger-camera type detector in thermal neutron imaging. J Instrum 7(8):P08010
5. Buonanno L et al (2020) A directional gamma-ray spectrometer with microcontroller-embedded machine learning. IEEE J Emerg Sel Topics Circ Syst 10(4):433–443
6. Kamuda M, Stinnett J, Sullivan CJ (2017) Automated isotope identification algorithm using artificial neural networks. IEEE Trans Nucl Sci 64(7):1858–1864

# Data Acquisition System for a 28 nm Flash-ADC Based Programmable Front End Channel for HEP Experiments

Andrea Galliani[1,2(✉)] , Luigi Gaioni[1,2] , and Gianluca Traversi[1,2]

[1] Universitá degli Studi di Bergamo, Dipartimento di Ingegneria e Scienze Applicate, I-24044 Dalmine (BG), Italy
andrea.galliani@unibg.it
[2] INFN, Sezione di Pavia, Via Bassi 6, I-27100 Pavia, Italy

**Abstract.** This work discusses the flash ADC-based front-end designed for future, high-rate pixel detector applications.The paper includes a description of the front-end blocks within the illustration of the submitted for fabrication prototype chip configuration blocks and the readout structure.

## 1 Introduction

Future High Energy Physics (HEP) experiments, where Total Ionizing Doses (TIDs) beyond 1 Grad ($SiO_2$) are expected and where detectors will experience unprecedented levels of particle rates, demand for advanced pixel readout processors. The 65 nm CMOS based RD53 pixel chips [1] currently represent the cutting edge of HEP's pixel front-end electronics. The 28 nm CMOS is the next key industrial node and it is currently the main focus of the HEP microelectronics community for future developments. With the aid of such technologies, designers may make pixel cells smaller and increase the speed of readout logic and I/O circuits. It also brings along a significant improvement in terms of radiation hardness [2]. The INFN Falaphel project, whose ultimate goal is the integration of silicon photonics modulators with fast, rad-hard circuits in a 28 nm process, is carrying out the front-end design activity. The project aims to develop front-end circuits suitable for use in possible future upgrades of the HL-LHC experiments, as well as to meet the requirements of the Future Circular Collider experiments [3]. The development of 28 nm CMOS front-end circuits for future, high-rate pixel detector applications is achieved through the design and comparison of two different architectures: the first one is under development and implements the time-over-threshold method in order to convert the analog signal from the preamplifier, while the latter integrates a zero dead-time front-end channel. It is composed of a Charge Sensitive Amplifier (CSA) stage that converts the signal

from the detector in a voltage signal [4], which is then digitized with a 2-bit, in-pixel flash-ADC. A prototype chip has been submitted for fabrication at the end of 2022, and it is currently under an early stage characterisation. This paper briefly describes the analog channel and then investigates in depth the blocks included in the developed Data Acquisition System describing the configuration structures integrated in the pixel and the methods to manage data to and from the chip.

## 2   Analog Channel

Figure 1 shows the schematic of the analog front-end channel. The Charge Sensitive Amplifier stage consists of a forward gain stage surrounded by a fast feedback including a capacitance $C_F$ and an NMOS transistor $M_F$, which is biased to regulate $C_F$'s discharge current, and a slow feedback loop capable of compensating for the detector leakage current. When a signal is delivered by the detector, the CSA output is approximated with a step signal with amplitude $Q_{in}/C_F$ and a linear return-to-baseline with a slope of $\sim -I_F/C_F$, with $I_F$ being the current generated by $M_F$. The leakage compensating feedback loop around the gain stage, locks the CSA output voltage to the gate-to-source voltage of the input device of the gain stage. If a detector current is present, the CSA output tends to move. Nonetheless, the gate of the $M_L$ transistor is driven in such a way to carry the whole leakage current, avoiding the saturation of the CSA. The low-pass $R_L - C_L$ filter guarantees low-frequency operation for the $M_L$ transistor. The CSA output is AC-coupled to a trio of clocked, auto-zeroed comparators used for building the 2-bit, in-pixel flash ADC. Such a comparator is based on the cascade of two common-source amplifiers and two inverters. Every comparator requires a negative step threshold signal. The CSA output and the threshold signal are injected at the gate of the comparator input transistor $M_1$, synchronous with the falling edge of the 40 MHz clock signal. The clock controls the opening and closing of the switches $S_1$ and $S_2$. When the two switches are open, the node P behaves like a summing node. If the CSA signal amplitude is larger than the threshold, then a positive going signal is present at the $M_1$ gate, amplified by two common source stages and fed to the output inverters that consolidate the output logic levels.

This front-end channel targets an environment where a very high rate of events is expected and where the probability of having hits in adjacent bunch crossings is not negligible. In order to test this zero dead-time front-end, an in-pixel read-out block has been integrated with the capability to store events from two subsequent bunch crossing periods of 25 ns. Thus, in cascade to the comparator, a block that controls the double hit detection mechanism has been integrated. It receives a signal from the comparator as input, (IN), it reads a global reset signal (RESET_B) and provides two logical output signals, OUT_1, which indicates the detection of a first hit, and OUT_2, that signals the presence of a second hit.

**Fig. 1.** Charge sensitive preamplifier with leakage compensation feedback and a comparator stage

## 3   The Prototype

The prototype chip is composed of a $4 \times 8$ matrix of front-end channels (Fig. 2). Because every channel shares an input line and six output lines, the matrix will be tested enabling one pixel at a time. Thus a programmable structure integrated in every channel is required. It allows to manage custom configurations through a set of switches activated with digital control signals. The main setting is the enable or disable of the input line. In case of disable, the input line is grounded so that the channel is not sensitive to any disturbances on the injection bus. It is also possible to connect the input node of the CSA with a desired detector-emulating capacitor (25, 50 or 75 fF). Moreover a circuit that emulates the presence of a sensor leakage current (LKG mosfet) is also integrated. Finally these settings allow the management of the output lines. As a matter of fact every double-hit detection stage exposes two output signals. Each of them is connected to the proper bus line. In order to decouple every pixel from the others, the ability to disconnect the channels output from the output bus is required. Therefore every output signal is driven by a column driver built with a tristate inverter. Considering that the 2-bit flash ADC is composed of a trio of comparator, and their output is processed by a double-hit detection stage, every channel has six output signals. Thus six digital control signals are needed to connect or disconnect their column driver from the output bus. Therefore each channel integrates a 10 bit shift register to manage these configurations (Fig. 3). In Table 1 is shown the relationship between the n-th output of the shift register and the controllable setting of the channel. The shift register output of any pixel is connected to the input of the next one, then a 320 bit serial signal is needed to configure the matrix.

**Fig. 2.** Diagram of 32 pixel matrix. In the foreground is shown the channel of a pixel within its building blocks and configuration signals. It is highlighted the presence of shared buses for both the input and output signals



**Fig. 3.** Diagram for the 10 bit shift register integration with the analog frontend channel

**Table 1.** Shift register code and control signal pairing

| SR bit | Enabled signal |
|--------|----------------|
| 0 | OUT2_1_EN |
| 1 | OUT2_2_EN |
| 2 | OUT1_1_EN |
| 3 | OUT1_2_EN |
| 4 | OUT0_1_EN |
| 5 | INJ_EN |
| 6 | C1_EN |
| 7 | C0_EN |
| 8 | LKG_EN |
| 9 | OUT0_2_EN |

# 4 Data Acquisition System

The digital output of the comparator is the only element that can be used to evaluate the main analog performance parameters of the prototype chip. In order to get these information, is necessary to properly configure the matrix using the Data Acquisition System (DAQ) that has been developed. It is mainly composed of a Tektronix 5334 Data Timing Generator (DTG) connected through a GPIB-USB-HS interface to a PC and a PCB motherboard designed to integrate the matrix carrier-board and a Arduino micro-controller board, which is needed to manage the configuration bits required by every pixel and to read the outputs of the prototype chip. A GUI Python-based software has been developed specifically for managing the data readout, the configuration of the prototype chip and sending the required settings to the Data Timing Generator. The exchange of the readout and configuration data between the motherboard and the PC is made through the serial port of the Arduino board. Considering a typical testing flow, it is necessary to automate these operations:

1. configure the 10 bit shift register of every pixel with the desired settings generating the clock signal to slide the proper digital serial configuration signal. Figure 4 shows the example of the signals involved in the configuration of a pixel.
2. set the injection signal amplitude and the 40 MHz clock form the DTG. Considering that it is necessary to iterate the injections several times, the first is realized through a square wave which amplitude can span in the range of (3–300 mV). The latter is a 900 mV amplitude square wave.
3. read any output signal that may be present and reset the double-hit detection stage (that is working as a latch) in order to be ready for the next injection event.

In order to keep everything consistent, the frequency of the reading and reset operations is set by the injection signal. The rising edge of this signal triggers the reset operations. The falling edge triggers the reading operations. Actually the DTG only offers two output channels that are used for the 40 MHz clock and the low voltage injection signal. The Arduino board needs 3.3 V digital signals in order to be sensitive to external interrupts. For this reason it has been designed the stage shown in Fig. 5. The circuit is composed of a non inverting amplifier that increase the amplitude of the injection signal from the DTG and than it compares it to a threshold set by the trimmer connected to the positive input of the LM311N comparator. The threshold should be put to a value higher than (2.5 V) in order to clone the injection signal and provide a signal that can be interpreted by the Arduino as digital. At this stage, the system can be configured to execute a chip scan to calculate the Equivalent Noise Charge and the threshold dispersion that distinguishes this front-end circuit. These parameters are obtained interpolating the comparator response to an increasing range of input signals with the so-called s-curve [5] for every pixel.

**Fig. 4.** An example of the configuration signal pattern for a pixel. In this picture only the six outputs and the injection settings are enabled



**Fig. 5.** Schematic of the signal adapter from the DTG to the microcontroller

## 5   Conclusions

The design and the main technological choices for the Data Acquisition System that will be used in the characterization of the flash-FE architecture, developed in the framework of the INFN Falaphel project in a 28 nm CMOS technology, has been described. The DAQ system is under calibration and the characterization of the prototype chip has started. The system provides significant benefits in terms of testing speed and workflow convenience, significantly accelerating the characterization of the prototype chip and increasing efficiency in data analysis.

# References

1. Gaioni L (2019) Test results and prospects for RD53A, a large scale 65 nm CMOS chip for pixel readout at the HL–LHC. Nucl Instrum Methods 936:282–285
2. Bonaldo S et al (2020) Ionizing-radiation response and low-frequency noise of 28-nm MOSFETs at ultrahigh doses. IEEE Trans Nucl Sci 67(7):1302–1311
3. Abada A et al (2019) FCC-hh: the hadron collider: future circular collider conceptual design report volume 3. Eur Phys J: Spec Top
4. Gaioni L et al (2023) 28 nm CMOS analog front-end channels for future pixel detectors. In: Nuclear instruments and methods in physics research section a: accelerators, spectrometers, detectors and associated equipment, vol 1045, pp 167609. ISSN 0168–9002
5. Galliani A et al (2023) DAQ system for the readout of a flash-ADC based front-end channel matrix. In: 2023 12th international conference on modern circuits and systems technologies (MOCAST). Athens, Greece, pp 1–4

# Preliminary Development of a Full-Digital Smart System for Chest Auscultation and Further Internet of Medical Things Framework

Matteo Zauli[1](✉), Lorenzo Mistral Peppi[1], Valerio Antonio Arcobelli[2],
Luca Di Bonaventura[3], Valerio Coppola[1], Sabato Mellone[2],
and Luca De Marchi[2]

[1] ARCES, University of Bologna, 40136 Bologna, Italy
`matteo.zauli7@unibo.it`
[2] DEI, University of Bologna, 40136 Bologna, Italy
`luca.di-bonaventura@st.com`
[3] STMicroelectronics, ARCES, Bologna, Italy

**Abstract.** In this article, we report the design and initial testing of a device which paves the path to the realization of a low-cost wearable device for detecting heart and lung tones for telemedical and long-term monitoring applications. In the current practise, the capability to assess chest and lung sound characteristics heavily depends on the experience of the clinician and a device which minimizes the human factor is highly in demand. In particular, this work addresses the development of a smart device based on the novel approach of digital Micro Electro-Mechanical Systems (MEMS) microphones capable of operating in a wide frequency range and low noise ratio. The aim is to propose in future developments Artificial intelligence-based algorithms for patient monitoring, a medical diagnostic decision support system, belonging to the Internet of Medical Things framework.

**Keywords:** Auscultation · Stethoscope · Wearable · Digital acquisition system · Telemedicine · IoMT

## 1 Introduction

Smart devices for real-time remote monitoring of the patient can make a big contribution to clinical practice, in particular, the COVID-19 pandemic highlighted the need for more pervasive and efficient approaches in telemedicine [1].

Auscultation is the practice devoted to listening to internal body sounds, usually using an acoustic stethoscope. In particular, chest auscultation is an important part of the patient's physical examination and an essential part of

clinical diagnosis used by doctors to investigate the condition of the heart [2] and lungs [3]. The main sounds in the chest can be divided into three categories: breath sounds (airflow through the tracheobronchial tree), heart sounds (turbulence due to valves snapping shut), and murmurs (rapid and choppy blood flow). Those noises can be divided into health or disease indicators. This is why is fundamental to determine the characteristics of the sounds to design a digital acquisition system acting suitable measurements, moving from an acoustic stethoscope to a digital one. As an example of chest audio signal features, the work of heart valves along with the blood flow into arteries and veins produces acoustic signals in the range of 20–1000 Hz [4], while respiratory sounds occupy a frequency band between 100 and 5000 Hz [5].

In order to design a suitable acquisition system, the careful selection of the frequency range is crucial, as can be observed in the wearable patches presented in literature [6,7] and commercial digital stethoscopes, i.e., eKuore and Thinklabs One. Those devices acquire audio signals using a single audio transducer in a frequency band around 500–600 Hz for eKuore and [7], and up to 2 kHz for Thinklabs and [6]. While, as far as Signal to Noise Ratio is concerned, except for [7] for which it is not indicated, the other devices settle between 12 and 15 dB.

In this work, we want to propose a preliminary study about a digital acquisition system for remote chest tones monitoring, with improvements over the devices presented above in terms of bandwidth and Signal-to-Noise Ratio (SNR). Moreover, the presented device features multichannel acquisition capabilities. The designed device acts as a USB PC peripheral and supports up to 4 microphone sources with continuous data streaming to the PC for real-time auscultation and data storage. Featuring a mic array is important to perform adaptive noise suppression, which, in turn, allows to extend the acquisition band up to the maximum frequencies occurring in the human body chest, those related to respiratory sounds. Otherwise, unwanted high-frequency noise, due to speech sounds and motions, may interfere with auscultation. In addition, in the proposed prototype an on-board audio filtering has been evaluated, as well as post-processing on Personal Computer side powered by MATLAB software and Python scripts. Overall the designed device features of a frequency band between 10 and 3000Hz with an SNR of 16 dB.

## 2   Material and Methods

The auscultation device exhibits the capacity to acquire up to 4 distinct channels concurrently, connected via USB to a laptop, and is under the control of a MATLAB interface script. For post-processing purposes, the device employs the same software and Python scripts to apply signal processing, facilitating the visualization and evaluation of the acquired measurements' quality. In addition, to improve the acquisitions, the sensors have been enclosed in 3D printed bell structures and closed by a Littmann® stethoscope diaphragm. The diaphragm enhances the acquisition sound quality thanks to a better sound pressure transmission.

## 2.1   Prototype

The realized prototype is based on Printed Circuit Board stacking techniques, to have a common computational layer and two different signal condition layers. The computational layer, named "MCU layer", hosts an STMicroelectronics STM32L552RE Microcontroller Unit (MCU). Conversely, the upper layer focuses on the sound sensor choice performing different signal conditioning tasks up to the main component configuration, a dual operational amplifier by Texas Instruments, OPA2197. Four different types of microphones were evaluated: (i) piezoelectric, i.e. the "Murata Piezoelectric Diaphragms 7BB-35-3L0", (ii) electret condenser, the "Primo EM272Z1", (iii) analog MEMS, the "STMicroelectronics MP23ABS1", (iv) digital MEMS, the "STMicroelectronics IMP34DT05". These transducers were selected for hearth and pulmonary activity monitoring because of their small form factor, non-invasiveness, and signal quality. Two signal condition layers have been designed to process the output of the sensors: a "Piezoelectric Layer" and an "Analog Microphones Layer", while the digital mic can be connected directly to the "MCU layer".

When analog sensors are used, the acquired signals are sent from the signal conditioning layers to the computational one. Each signal is provided as input to the Analog-to-Digital Converter (ADC) embedded in the MCU, performing a digital conversion with a sampling frequency of 44 ksample/second/channel at a 12 bits resolution. To match the ADC dynamic range, two different preamplifier layers were developed, obtaining the ideal conditions for each type of transducer, making it possible to customise the gain, currently around 22 dB in both cases, and the input impedance, thus maximizing the SNR.

The digital microphone has a sample rate of 44 kHz too and unlike the analog solution, it makes use of the DFSDM peripheral (Digital Filter for Sigma-Delta Modulator). This peripheral benefits from 4 input channels and features also a filtering function such as 4 Sinc filters (order from 1 to 5) and 4 integrators. Morever, the IMP34DT05 microphone has a Pulse Density Modulated data output whose format enables a direct connection to DFSDM which can process the audio signal up to 24-bit resolution from the mentioned digital mic. Lastly, DFSDM decreases the computational load, because it can perform processing (filtering) and data transfer (taking advantage of Direct Memory Access) autonomously. Such solution is depicted in Fig. 1. That Figure shows the "MCU layer" with the connection toward the digital mic embedded in an ad-hoc designed 3D printed PLA (Polylactic Acid) plastic bell. The decision to house the microphone in that bell has been taken to mimic the performances of conventional acoustic stethoscopes. In Fig. 1, it can be observed also the Littmann® stethoscope diaphragm. In the same Figure, the "Connectors to upper layers" are shown, i.e. 2 female headers (20 pins in 2 rows) used to stack board such as the ones designed for the analog microphones interfacing (the above-mentioned "Piezoelectric Layer" and "Analog Microphones Layer"), and to plug directly the digital solution.

The final microphone choice has been followed thanks to a specific test setup based on a loudspeaker LG 6400GSMC01A placed behind a water-filled polymer balloon. On that balloon the selected microphones, each one encapsulated in a

**Fig. 1.** Prototype parts

suitable 3D printed bell, one after the other were placed on the polymer surface and two kinds of audio recordings were done: playing a white noise track from the speaker and the other one in a complete silence mode. In addition, two commercial digital stethoscopes have been evaluated: Thinklabs One and eKuore (see Sect. 3).

## 2.2   Software Architecture and Signal Processing

In the current study, the designed software (mainly implemented in MATLAB) provides the capability to (i) acquire data from the external device, (ii) conduct real-time audio listening, (iii) locally store acquired data, and (iv) process the collected data. Such operations are enabled through serial communication established between MATLAB and the connected hardware prototype via a USB cable. The main operations flow of the MATLAB routine is given by a set of messages types (see Fig. 2) : (i) "START", is the command to start data acquisition, (ii) "DATA", data streaming for monitoring and listening, (iii) "STOP", is the command to terminate the auscultation recording. The digital processing of the acquired signals was implemented in Python using the SciPy library. A sequence of two pass-band Butterworth filters was utilized, aimed at isolating the cardiac heartbeat signal, which typically resides within the range of 20 to 150 Hz [8]. The selected bandwidth of the filters was intentionally wider to ensure that no components of the signal of interest were excluded (ranging from 10 Hz



**Fig. 2.** Acquisition system

to 400 Hz, with a stopband attenuation of greater than or equal to 40 dB). Additionally, a 14 dB gain was applied to each signal and, for visualization purposes, the signals were normalized by dividing each time series by its maximum value. SNR has been used for quantifying the results.

## 3   Results

This preliminary work reports the design and development of a full-digital device test bench for multi-channel digital auscultation with wearable devices that can support both analog and digital sensors. Moreover, after the initial tests, an innovative oriented digital MEMS mic design was developed, and a dedicated hardware filtering stage was implemented basing on DFSDM peripheral programming. Tests were conducted to verify the SNR values and frequency range (Table 1) of the different solutions, including the two mentioned commercial devices. Those values and the features exhibited by the sensor's datasheet were instrumental to select the transducer most suited for the auscultation task, which is the digital MEMS microphone IMP34DT05. The presented device is characterized by 4 channels, 24-bit resolution, sample rate at 44 kHz, SNR of 16 dB, connectivity via USB 2.0 Full-Speed, main board size of $40 \times 48$ mm, and a chest-piece of diameter 43 mm, height 11 mm.

**Table 1.** SNR comparison ($dB$)

|  | Primo EM272Z1 | MP23ABS | MP34DT05 | Piezo | One | eKuore |
|---|---|---|---|---|---|---|
| SNR | 11 | 14 | 16 | 7 | 14 | 12 |
| Upper cutoff frequency | 1200 | 1100 | 3000 | 2900 | 700 | 600 |

The performances of "Murata Piezoelectric Diaphragms 7BB-35-3L0" and "STMicroelectronics IMP34DT05", can be qualitatively observed in Fig. 3 which represents: the time-series acquired with MEMS and with piezoelectric sensor, subplots (a) and (b), respectively. The implemented signal processing steps are shown in Fig. 3c: (i) raw signals acquisition (shown in blue), (ii) pass band filter application, and (iii) gain adjustment implementation (see Sect. 2.2 for more details on the processing), with the ultimate processed outcome indicated by black traces. As anticipated, the MEMS solution demonstrated superior performances compared to other solutions, specifically in terms of signal amplitude and noise level. This can be observed in Fig. 3a, b, since the MEMS channel shows more prominent first heart sounds (S1) resulting in an SNR of 16 dB, in contrast to the piezoelectric channel which only achieves an SNR of 7 dB.

**Fig. 3.** Processed signal

## 4   Conclusion

The prototype described in this work is designed to record and listen to chest tones in real-time and store them for future use such as auscultations by an experienced clinician or for developing AI-based algorithms for long-term monitoring of the patient. Thus, this device is a crucial first step towards offering new telemonitoring and services to support continuity of care approaches.

## References

1. Taha A, Shehadeh M, Alshehhi A, Altamimi T, Housser E, Simsekler M, Alfalasi B, Al Memari S, Al Hosani F, Al Zaabi Y and Others (2022) The integration of mHealth technologies in telemedicine during the COVID-19 era: a cross-sectional study. Plos One 17:e0264436
2. Montinari M, Minelli S (2019) The first 200 years of cardiac auscultation and future perspectives. J Multidiscip Healthc 183–189
3. Cottin V, Richeldi L (2014) Neglected evidence in idiopathic pulmonary fibrosis and the importance of early diagnosis and treatment. Eur Respir Rev 23:106–110
4. Gupta P, Moghimi M, Jeong Y, Gupta D, Inan O, Ayazi F (2020) Precision wearable accelerometer contact microphones for longitudinal monitoring of mechano-acoustic cardiopulmonary signals. NPJ Digit Med 3:19
5. Bohadana A, Izbicki G, Kraman S (2014) Fundamentals of lung auscultation. N Engl J Med 370:744–751
6. Lee S, Kim Y, Yeo M, Mahmood M, Zavanelli N, Chung C, Heo J, Kim Y, Jung S, Yeo W (2022) Fully portable continuous real-time auscultation with a soft wearable stethoscope designed for automated disease diagnosis. Sci Adv 8:eabo5867
7. Klum M, Leib F, Oberschelp C, Martens D, Pielmus A, Tigges T, Penzel T, Orglmeister R (2019) Wearable multimodal stethoscope patch for wireless biosignal acquisition and long-term auscultation. In: 2019 41st annual international conference of the IEEE engineering in medicine and biology society (EMBC), pp 5781–5785
8. Phua K, Chen J, Dat T, Shue L (2008) Heart sound as a biometric. Pattern Recognit 41:906–919

# Reducing Energy Consumption in NB-IoT by Compressing Data and Aggregating Transmission

Mohamad Kheir El Dine[1(✉)], Hussein Al Haj Hassan[2], Abbass Nasser[3], Chamseddine Zaki[4], Azza Moawad[5], and Ali Mansour[1]

[1] ENSTA Bretagne, LabSTICC UMR CNRS 8265, 29806 Brest Cedex 9, France
mohamad.kheir@ensta-bretagne.org
[2] CCE Department, American University of Science and Technology, Beirut, Lebanon
[3] Computer Science Departement, American University of Culture and Education, ICCS-Lab, Beirut, Lebanon
[4] College of Engineering and Technology, American University of the Middle East, Egaila 54200, Kuwait
[5] Department of Electronics and Communication Engineering, Arab Academy for Science, Technology and Maritime Transport (AASTMT), Cairo, Egypt

**Abstract.** The rapid expansion of Internet of Things (IoT) applications necessitates the utilization of efficient Low Power Wide Area Network (LPWAN) technologies. Among these technologies, Narrowband IoT (NB-IoT) has emerged as a highly promising solution due to its energy efficiency, long-range communications, and anticipated scalability. In this study, we conducted simulations to examine the impact of various factors on energy usage in NB-IoT devices. Taking into consideration the significance of energy efficiency in IoT devices, particularly in hard-to-reach locations, we study the effects of data compression on energy consumption. The obtained results revealed that implementing data compression techniques can significantly reduce energy consumption, specifically in areas with limited transmission coverage. Additionally, we study the benefits of expanding data transmission intervals to conserve energy. The findings of this paper emphasize the importance of adopting energy-efficient practices in the development of IoT technology.

**Keywords:** IoT · LPWAN · NB-IoT · Energy consumption

## 1 Introduction

The term "Internet of Things (IoT)" refers to electronic systems that include sensors and actuators controlled by software to enable massive devices connectivity through the Internet. Needless to say, the number of IoT devices is growing quickly each year were it is predicted that there will be 30.9 billions by 2025 [1]. In the realm of IoT technology, there exist two distinct categories: the first

one encompasses short-range applications with high energy consumption, exemplified by technologies such as Bluetooth Low Energy (BLE), ZigBee, Z-Wave, and Wi-Fi, which are designed for specific uses and limited proximity. The second category comprises Low Power Wide Area Network (LPWAN) technologies including LoRa, and Sigfox were developed to fulfill the demand for long-distance communication. LPWAN is a generic term for any network designed to communicate wirelessly with lower power than other networks such as cellular, satellite. Sigfox and LoRAWAN are unlicensed technologies that work in the spectrum give a ton of device connectivity [2]. The Third Generation Partnership Project (3GPP) has developed two IoT technologies: Long-Term Evolution (LTE)-M and NB-IoT. Both technologies operate on the licensed spectrum. Unlicensed IoT demands building new infrastructures; however, NB-IoT reuses LTE infrastructures. That means using NB-IoT is simpler and less expensive than using unlicensed technologies. The 3GPP Release 13 presents the specifications of NB-IoT. The key specifications of NB-IoT technology include low uplink delay, a high Maximum Coupling Loss (MCL) of 164 dB, extended UE battery life, and support for a large number of devices per square kilometer [3].

NB-IoT uses a bandwidth of 180 kHz, and can be deployed in one Physical Resource Block (PRB) of the available LTE bandwidth , or in the Guard-band, or occupy the GSM system's released spectrum. The NB-IoT system was previously planned to operate in a Frequency-Division-Duplexing (FDD) based mode in the 3GPP Release 13 standard. Thanks to 3GPP Release 15, it is now possible to use NB-IoT in Time-Division Duplexing (TDD) mode. The deployment of Narrowband Internet of Things (NB-IoT) holds great potential for various applications, like smart metering, smart cities, localization purposes, agriculture and smart building [4]. NB-IoT has two primary states: connected state and idle state, in the Connected State the UE maintains a control link with the network and compose from the Synchronization (Sync), Transmission and Reception (TX/RX), Connected Mode Discontinuous Reception (CDRX) and Release. The idle state in NB-IoT refers to the state where the device is not actively connected to the network but is still able to receive paging messages [5], it is composed from the Extended Discontinuous Reception (EDRX) and Power Saving Mode (PSM).

Recent research has focused on IoT technologies, including studies on NB-IoT power consumption and device lifespan. An analytical model proposed in [6] using Markov chains revealed that NB-IoT devices can extend battery life up to 10 years with a minimal 10-second delay for various scenarios. In [7], NB-IoT device measurements were conducted to assess their energy consumption, with slight parameter adjustments leading to notable energy savings. In studies [8,9], the same devices were used in different networks, with emphasizing energy efficiency in the Telekom Malaysia network and highlighting the importance of firmware setup for achieving promised device lifetime. [6] conducted measurements in a Spanish commercial network, showing energy savings using the Release Assistance Indication (RAI) flag during Extended Discontinuous Reception (eDRX).

In this paper, we present energy efficiency techniques in NB-IoT with focus on the impact of compressing data and expanding data transmission interval on the energy consumption in NB-IoT module. Specifically, our study aims to prove experimentally that expanding the period of collecting messages considerably affects the energy efficiency, same to data compression in bad coverage area.

The rest of the paper is structured as follows: In Sect. 2 we present a scenario work of an active sensor in NB-IoT and in Sect. 3 we present our proposed energy model and in Sect. 4 we simulate the performance and highlight the essential impact of compressing and expanding Data transmission interval on energy consumption and finally the conclusion is given in Sect. 5.

## 2   Scenario of an Active Sensor in NB-IoT

In this section, we present a scenario work of an active sensor in NB-IoT:

The UE sends a request message to eNB inquiring connection to the network. After receiving an acknowledgment from the eNB, the UE sends a request message to synchronize the SIM card. When the UE receives an acknowledgment, RRC setup makes a transition from RRC Idle mode to RRC connected mode. When the data is ready to be transmitted or received, the UE enters in the Connected state. In the Connected state, the UE can send and receive the data. Regarding the uplink, the UE has 3 options:

– Uplink with RAI 000: The UE sends the packet data and waits for the acknowledgment from the eNB. When the time of TX/RX phase ends, the UE enters in CDRX phase. If the UE receives a download packet, it reenters in the TX/RX phase. Otherwise, the UE releases the Connected state then enters in the Idle state.
– Uplink with RAI 200: The UE sends the packet data and directly releases the Connected state, then enters in the Idle state.
– Uplink with RAI 400: The UE sends the message and enters in the CDRX phase. If the UE receives an acknowledgement message from the base station, it leaves the Connected state and enters in the Idle state.

Regarding the downlink, the UE receives the data and enters the CDRX phase. After that, the UE releases the Connected state and enters in the Idle state. In the Idle state, the UE enters EDRX. If the UE receives a download packet from the eNB, it enters again the Connected state. Otherwise, and after ending the EDRX, the UE enters the PSM (deep sleep) (Fig. 1 shows an illustrative flowchart).

## 3   Proposed Saving Energy Model

In our research, we performed an investigation on the nRF9160 rev2 sensor of battery voltage 3.7 V, by utilizing the online power profile for LTE available

**Fig. 1.** Scenario work of UE in NB-IoT

at Devzone online power profile LTE.[1] This platform provided us with valuable information regarding the electrical current and duration spent in different phases. It should be noted that the values of electric current and time vary according to multiple factors, such as the specific module, the Extended Coverage Level (ECL), the data size, and the repetition factor. By identifying and addressing specific aspects that contribute to higher energy usage, we can effectively work towards optimizing power efficiency. Through our investigation, we successfully established correlations between each phase and the factors influencing the electric current and time within it.

The NB-IoT technology encompasses various phases of UE operation, each influenced by specific parameters that impact energy consumption. These phases can be categorized as follows:

$E_{sync}$: The current I and time T are constant in synchronization. Therefore, the Energy E of network synchronization is constant (E = Pt).

$E_{SIMsync}$: current I is constant, but the time T depends on the module (Table 1).

$E_{RRC}$: The current in the RRC setup changes when the ECL (0–1–2) is changed, and the time changes when the module is changed (Table 2).

$E_{TAU}$: The current in TAU depends on the ECL (0–1–2), but the time is constant (Table 3).

$Erelease$: The current in TAU depends on the ECL (0–1–2) whereas the time is constant (Table 4).

---

[1] https://devzone.nordicsemi.com/power/w/opp/3/.

$E_{TX/RX}$: Current and time of data transmission or data receiving depend on the data size (Table 5).

$E_{CDRX}$ and $E_{EDRX}$: $I_{sleep}$ is constant while $I_{listening}$ depends on the repetition of the Physical Downlink Control Channel (PDCCH). On the other hand, $T_{sleep}$ and $T_{listening}$ depend on the interval time of CDRX.
EPSM: In PSM phase, the current depends on the module. The time is between 2 s and 413 days (Table 6).

**Table 1.** Energy consumption of SIM sync phase with different modules, I is the current in mA, T is the time in ms and E is the energy in Joule

| Module | I (mA) | T (ms) | E (J) |
|---|---|---|---|
| nRF9160 rev1 | 6.093 | 734.600 | 0.0165 |
| nRF9160 rev2 | 6.093 | 483 | 0.0108 |

**Table 2.** Energy consumption of RRC phase with different ECL, 0 is good coverage to 2 bad coverage

| ECL | I (mA) | T (s) | E (J) |
|---|---|---|---|
| 0 | 34.425 | $149.743 \times 10^{-3}$ | 0.0190 |
| 1 | 36.294 | $149.743 \times 10^{-3}$ | 0.0201 |
| 2 | 61.808 | $149.743 \times 10^{-3}$ | 0.0342 |

**Table 3.** Energy consumption of TAU phase with different ECL

| ECL | I (mA) | T (s) | E (J) |
|---|---|---|---|
| 0 | 34.425 | $931.115 \times 10^{-3}$ | 0.1185 |
| 1 | 36.294 | $931.115 \times 10^{-3}$ | 0.1250 |
| 2 | 61.808 | $931.115 \times 10^{-3}$ | 0.2129 |

Overall, the parameters that affect the energy consumption are the operator and module, packet size, Extended Coverage Level (ECL), Release Assistance Indicator (RAI), and repetition of packet.

One can say that the total energy consumption of a sensor in NB IoT is:

$$\begin{aligned}
E_{Total} &= E_{sync} + E_{SIMsync} + E_{RRCsetup} + E_{TAU} + E_{release} + E_{TX/RX} + E_{CDRX} \\
&+ E_{EDRX} + E_{PSM} = U.(I_{sync}T_{sync} + I_{simsync}T_{simsync} + I_{RRC}T_{RRC} + I_{TAU}T_{TAU} \\
&+ I_{release}T_{release} + I_{TX/RX}T_{TX/RX} + N.CDRX.(I_{sleep}T_{sleep} + I_{listening}T_{listening}) \\
&+ N.EDRX.(I_{sleep}T_{sleepEDRX} + I_{listening}T_{listening}) + I_{PSM}T_{PSM}).
\end{aligned} \quad (1)$$

where $N_{CDRX} = T_{CDRX}/CDRX$ and $N_{EDRX} = T_{EDRX}/EDRX$ intervals

**Table 4.** Energy consumption of Release phase with different ECL

| ECL | I (mA) | T (s) | E (J) |
|-----|--------|-------|-------|
| 0 | 17.683 | $527.939 \times 10^{-3}$ | 0.0345 |
| 1 | 17.523 | $527.939 \times 10^{-3}$ | 0.0342 |
| 2 | 15.345 | $527.939 \times 10^{-3}$ | 0.0299 |

**Table 5.** Energy consumption of TX/RX phase with different data size

| packet size | I (mA) | T (s) | E (J) |
|-------------|--------|-------|-------|
| 20 | 35.388 | 962.224 | 0.1259 |
| 100 | 35.741 | 1011 | 0.1336 |
| 200 | 36.145 | 1072 | 0.1433 |

## 4   Simulations

We conducted a Python-based simulation to evaluate the energy consumption of an nRF9160 rev 2 sensor. We modified the sensor's configuration to investigate the effects of expanding the data transmission interval and data compression on reducing energy consumption. To conduct this study, we designed three scenarios:

1. In the first scenario, a UE sends a message of 100 bytes every hour. The energy consumption by the UE per hour is 497 J. Thus, the total energy consumption after 10 h is 4970 J.
2. In the second scenario, the UE collects information over a span of 10 h and subsequently sends a single message of 1000 B. The energy consumption of the UE in this scenario is measured to be 875 J.
3. In the third scenario, the Release Assistance Indicator (RAI) is changed from 000 to RAI 200. Following the transmission of the message, the UE enters an idle state. Notably, the energy consumption of the UE decreases to 753 J in this particular case.

Figures 2 and 3 represent the difference of energy consumption between each phase of UE in the first and second scenario to study the impact of Expanding Data Transmission Interval on the energy consumption in NB-IoT. In Fig. 3, the Power-Saving Mode (PSM) phase has the longest time interval and also the

**Table 6.** Energy consumption of PSM phase with different modules

| Module | I (mA) | T |
|--------|--------|---|
| nRF9160 rev1 | $4 \times 10^{-3}$ | [2 s–413 days] |
| nRF9160 rev2 | $2.7 \times 10^{-3}$ | [2 s–413 days] |

highest energy consumption. However, it's important to note that this phase effectively conserves energy because it operates in a deep sleep mode characterized by significantly reduced energy consumption.



**Fig. 2.** Distribution of energy consumption in different phase of UE that send 100 B.



**Fig. 3.** Distribution of energy consumption in different phase of UE that send 1000 B.

We noticed that expanding the data transmission intervals demonstrates a notable reduction in energy consumption, and thus yielding favorable results on the lifespan of battery. Furthermore, optimizing the configuration contributes to decreasing the energy consumption of the UE.

On the other hand, our study focused on examining the impact of message size on energy consumption in various ECLs. We conducted experiments by sending messages of 100 and 1000 B on two different ECLs (ECL 0 and ECL 2), and then compressing these messages to 50 and 500 B respectively in order to investigate the effects of message compression. The energy consumption for each message in each ECL is presented in Table 7.

Our findings revealed that the energy consumption remains unaffected when the message sent by the UE is small, regardless of the ECL. However, in the case of large message sizes, energy consumption is significantly reduced when the messages are compressed, particularly as the ECL increases.

## 5   Conclusion

In this paper, we present an overview of NB-IoT technology with a focus on rationalising energy consumption to enjoy greater efficiency in NB-IoT. We

**Table 7.** Energy consumption of packet size

| Size | ECL = 0 | ECL = 2 |
|---|---|---|
| 100 B | 497.8 J | 2226.1 J |
| 50 B | 493.8 J | 2222 J |
| 1000 B | 569.6 J | 2298.1 J |
| 500 B | 529.7 J | 2258.3 J |

have develop an application model to implement a simulation that measures the amount of energy consumption at each stage in NB-IoT under different conditions using a calculator developed in python. In the end, we conclude that the amount of energy consumption is not affected by varying the packet size except for one case which is ECL 2 (bad coverage). We can also save energy by reducing the number of times in which the data is sent. To clarify, messages can be collect and sent at once between larger time interval taking into consideration the efficiency of work. In addition to what is mentioned, the most appropriate configuration can be chosen to decrease the energy consumption in NB-IoT technology. We can also develop effective ways to reduce the average latency in NB-IoT.

# References

1. Vailshery LS (2022) Internet of things (iot) and non-iot active device connections worldwide from 2010 to 2025
2. Pérez M et al (2022) Coverage and energy-efficiency experimental test performance for a comparative evaluation of unlicensed lpwan: Lorawan and sigfox. IEEE Access 10:97183–97196
3. García-martín JB, Torralba A (2021) Model of a device-level combined wireless network based on nb-iot and IEEE 802.15.4 standards for low-power applications in a diverse iot framework. Sensors 21:6
4. Routray S, Akanskha E, Sharmila KP, Sharma L, Ghosh AD, Pappa M (2021) Narrowband iot based support functions in smart cities, pp 1459–1464
5. Yassine F, El Helou M, Bazzi O, Lahoud S (2021) Investigation on narrowband iot link adaptation with rate and energy objectives, pp 412–417
6. Andres-Maldonado P et al (2019) Analytical modeling and experimental validation of nb-iot device energy consumption. IEEE Internet Things J 6:6
7. Dissecting energy consumption of nb-iot devices empirically. IEEE Internet Things J 8:1224–1242 (2021)
8. Yeoh C, Man A, Ashraf Q, Samingan A (2018) Experimental assessment of battery lifetime for commercial off-the-shelf nb-iot module
9. Martinez B, Adelantado F, Bartoli A, Vilajosana X (2019) Exploring the performance boundaries of nb-iot. IEEE Internet Things J 6:5702–5712

# Design and Development of New Wearable and Protective Equipment for Human Spaceflights

Maurizio Pellegrini[1], Giuseppe Coviello[1], Giuseppe Brunetti[1], Francesco Angelini[2], Ilario Lagravinese[2], Giorgia Manca[2], Flavio Augusto Gentile[2], Roberto Vittori[1], and Caterina Ciminelli[1(✉)]

[1] Optoelectronics Laboratory, Politecnico di Bari, Via Orabona 4, 70125 Bari, Italy
`caterina.ciminelli@poliba.it`
[2] REASPACE, Via Demarinis, 16, 70021 Acquaviva Delle Fonti (BA), Italy
`f.gentile@reaspazio.com`

**Abstract.** Exposure to microgravity poses challenges to the astronaut's musculoskeletal system, causing temporary issues that can be addressed through exercise routines. For long-duration space missions, wearable technology with sensors to measure and stimulate muscular activity is crucial. The proposed active suit system is a custom-made intra-vehicular wearable designed to counteract muscle and bone mass loss in microgravity. It consists of an anti-allergic, antibacterial suit with differential compression capabilities, an electronic system to detect body joint movements, an electrostimulation system, and an electronic control unit (ECU) for motion data analysis and stimulation regulation. The suit aims to reduce muscle atrophy by providing intelligent electrostimulation, simulating Earth-like gravity conditions. System calibration on Earth and adaptive neuromuscular stimulation in space are essential steps to maintain astronauts' muscle tone during extended space missions.

**Keywords:** Microgravity · Suit · Electrostimulation · Electromyography

## 1 Introduction

During permanence in extra-atmospheric environments, human physiology indeed undergoes significant adaptations to cope with the challenges posed by the absence of Earth's atmosphere and gravity. These adaptations are crucial for the survival and well-being of astronauts during their missions in space [1]. Physiological changes in astronauts during their stay in space can be approached from both microscopic and macroscopic standpoints. At the former level, mainly focused on the cellular environment, microgravity involves the cells' change in size, shape, volume, and adherence properties. This is caused by the influence of mechanical unloading on the cytoskeletal structures which alters the transcription, translation, and organization of its proteins [2]. At the macroscopic level, the health effects involve cardiovascular (cardiac atrophy, redistribution

of blood volume to the head and torso, increased renal excretion of salt and water and reduction of plasma levels), neuro-vestibular (disorientation, cold sweat, vomiting, facial pallor, nausea and stomach awareness all defined by the term "space motion sickness") and musculoskeletal systems. The normal function of bone and muscle tissues is the most affected. Microgravity causes muscle atrophy, especially in the higher load-bearing muscles, due to fibre loss, diminished protein synthesis and increasing protein degradation. The skeletal system suffers from the demineralization of the bone in microgravity, due to changes in the blood-forming cells and decoupling of reabsorption and formation [2].

To mitigate some of these effects and ensure the well-being of astronauts, space agencies carefully monitor their health during space missions. They provide specialized exercise routines, nutrition plans, and other countermeasures to help minimize the impact of these physiological changes. However, they are time-consuming and limit the astronauts' routine operations. Therefore, studying and overcoming these adaptations is vital for planning future long-term space missions, such as those to Mars or beyond, to ensure the health and safety of astronauts during extended stays in space.

To monitor vital parameters in extra-atmospheric conditions with the presence of microgravity and ionizing radiation, a wearable multisensory system has been designed and developed. The importance of this interdisciplinary research lies in ensuring continuous monitoring of astronauts' health during space missions, measuring, and counteracting possible shifts from the physiological condition, also given increasingly limited interactions of the crew with the ground staff, laying the foundations for missions of increasingly longer duration. As the first step of the development, the main goal was to reduce the phenomenon of muscle atrophy by providing intelligent electrostimulation which compensates for the differences between the muscular effort measured in normal gravity conditions and that in microgravity conditions.

## 2  Medical Research During Spaceflights

Periodic flights to orbital and sub-orbital heights take place, organized by national and international agencies or private companies, or with public-private joint collaborations. In any scenario, whether for commercial or scientific research purposes, these missions achieve the objective of conducting experiments encompassing a wide range of scientific fields, including medicine and engineering. One of the most recent examples of this is the Virgin Galactic sub-orbital flight that took place in June 2023. The goal of the medical research was observing the change of heart rate of the passengers during acceleration and an attempt to measure cosmic radiation in the upper atmosphere [3].

Remarkable results have been achieved in the framework of a joint mission between NASA and Space X in October 2021, that aimed at investigating, among countless research topics, several nutrition plans and efficient exercise routines, the measurement of physiological and psychological parameters of astronauts at the International Space Station (ISS) [4]. The latter aimed at increasing the efficiency of exercise to prevent muscle atrophy and resulting bone loss using electrical muscle stimulation (EMS) using the EasyMotion non-invasive wearable skin suit [4]. This suit initiates involuntary contractions of global musculature using dry electrodes while being controlled by the crew's iPad aboard the ISS [5]. Therefore, the development proposed here appears to be perfectly timely given the growing interest in studying and counteracting the problems

caused by microgravity conditions and the potentially increasing number of flights in which the system under development could find use.

## 3  Wearable Sensing and Stimulation System in Microgravity Conditions

The active suit system proposed here (see Fig. 1) is a custom-made intra-vehicular wearable designed to counteract the loss of muscle and bone mass experienced by astronauts in microgravity conditions. It consists of a suit made with anti-allergic, antibacterial materials; an electronic system for detecting the movement of the main body joints (sensors, flexometers); an electrostimulation system which acts based on the movements detected; an electronic control unit (ECU) that analyses motion data and regulates the stimulations. The ECU is placed within a carbon fibre case. The main objective of the suit is to reduce the phenomenon of muscle atrophy by providing intelligent electrostimulation which compensates for the differences between the muscular effort measured in normal gravity conditions and that developed in microgravity conditions. The idea is that, in the first realization step, a system calibration is performed on the Earth and electromyographic data collection of the astronaut's joint angles is conducted. Compared to systems like the EasyMotion, the proposed active suit is not limited to a rigid set of exercises controlled by an external source simulating Earth-like gravity in activities of daily living thanks to the above-mentioned ECU.



**Fig. 1.** The suit with detailed flexometers (in red) and electrodes (in blue), for a limited number of sensing and stimulation points.

### 3.1  The Suit Fabric

The Polypropylene fabric offers a unique combination of advanced performance and comfort, ensuring the astronaut's maximum protection. Distinctive features of the fabric are low hysteresis, for reduced deformation and optimal maintenance of its properties even after long periods of use; antibacterial properties, which help maintain a hygienic

and safe environment for astronauts in space, reducing the risk of contamination and infection; differential compression, for proper redistribution of body fluids while in space, which helps to mitigate the negative effects of microgravity on the circulatory system; intelligent thermal regulation, capable of adapting to the temperature variations in the space and, so, maintaining an optimal body temperature; lightness and flexibility, for complete freedom of movement and a lightweight feel while in use.

### 3.2   The Sensing and Stimulation System

All sensors are directly applied to the fabric, ensuring anatomic placement that precisely follows the user's movements. These sensors interact in real-time with electrical stimulation that varies based on the speed or specific weight of onboard equipment, recognized through vision sensors or codes. A control unit can easily deactivate each sensor and/or electrode. Moreover, the system can manage various stimulation mappings, such as muscle tone maintenance, strengthening and relaxation, increased load during workouts, capillarization and post-trauma rehabilitation.

The flexometers are 2-D measuring tapes, based on capacitive effect. They are positioned on the suit and connected to the ECU, through an I2C bus. The measuring tape and the connection point to the cable are enclosed in a special silicone case in order to protect electrical connections and avoid any accidental out-of-the-plane movement. All electrode pads will be placed on the targeted muscles. The majority of antigravity or tonic-postural muscles, as well as some phasic muscles, will be connected and controlled by the ECU. Given the application, the electrodes need to be dry and at the same time galvanically isolated from the user [6]. Dry electrodes help to impede the dermatologic effects of electrolytic gel, indirectly caused by using the more common wet electrodes. For the purposes of the first experiment, electrodes placed on the biceps (one pair per arm) and on the quadriceps femoris (two pairs per muscle) have been used, for a total number of 12. All electrodes are connected directly to the ECU, which provides a current of low to moderate intensity (max 30 mA), with a frequency appropriate to the expected stimulation effect, to the type of muscle and with a biphasic behaviour. The electrodes are made of a specific silicone that does not require gel, with the addition of a conductive material. The central unit in this system performs continuous monitoring of the flexometers at a frequency of 100 Hz. Muscle activity is determined by comparing the measured degrees with individual muscle-specific thresholds; for instance, the biceps may have a threshold of 30 degrees. When the threshold is exceeded, the corresponding muscle stimulation is activated, and the current and the stimulation frequency depend on the preselected program. This customizable approach allows for the implementation of various stimulus mappings, catering to different muscle stimulation requirements.

The suit in this final version will monitor biometric, electromyography muscle parameters and relative limb movement in microgravity. To these aims, current development is focused on the integration of a heath rate monitor, a portable blood pressure monitor, a respiratory rate monitor, a thermometer IC, a $SpO_2$ measurement IC, electromyography sensors, acceleration and movement 6-dofs IMU, a stabilizer system and a bioelectrical impedance analysis machine.

### 3.3 The Electronic Control Unit

For each couple of electrodes, the driving section consists of an improved Howland current pump with a buffer in the feedback path, using an instrumentation amplifier (INA), that can be used in an inverting or noninverting configuration, depending on the input signals, and an OPAMP in a voltage follower configuration that acts as INA reference. The circuit load voltage is the input of the buffer. This design enables the current source to drive a wide range of load resistances, making it versatile for various applications that require a current source capable of bipolar (source or sink) operation, such as electrostimulation applications. In this context, the primary objective is to provide a controlled and safe current for stimulating biological tissues. The Improved Howland current pump with buffer in the feedback loop [7] presents two operational amplifiers (op-amps), with the op-amp in the feedback path configured as a buffer to eliminate the feedback current, improving the thermal noise performance and simplifying the overall design.

Figure 2 shows the block diagram of the custom-designed and developed board. It consists of a microcontroller from STMicroelectronics that drives the Howland circuit's control signals through square waves generated by the PWM protocol and appropriately filtered. For each muscle group, only one Howland circuit and two electrodes connected to the board's reference and the circuit's output are required. A dual power supply of $\pm 18$ V is used to provide both sourcing and sinking capabilities for the current. Additionally, a 2-D flexometer is connected to the same IIC bus for each muscle group. In this specific implementation, six muscle groups are controlled with a maximum programmable current of $\pm 50$ mA per muscle group through software. The waveform employed is a balanced biphasic waveform with a pulse duration of approximately $250\,\mu s$ and its programmable repetition frequency ranges from 5 Hz to 50 Hz. Lastly, the board is powered by three 18650 lithium batteries, space-approved, arranged in a series configuration. Each battery has a nominal voltage of 3.60 V, with a maximum current rate of 4000 mA, ensuring functionality for more than 30 h.



**Fig. 2.** Block diagram of the ECU.

### 3.4 Experimental Results

Figure 3 illustrates the generated signals for driving a single muscle group, specifically the biceps. The green and blue signals represent the input signals already filtered at the input of the Howland circuit, while the yellow signal represents the signal at the positive electrode. Figure 3a, b show the waveforms for 20 mA and 30 mA stimulation, respectively. The electrode/muscle group has an impedance of approximately 270 Ω. As known, the latter varies due to various factors such as sweat, adhesion, body temperature, etc., and the proposed circuit is designed to be independent of such variations. To test this, in vitro laboratory tests were conducted by replacing the electrodes with electronically variable resistance, and the actual current delivered by the stimulation circuit was measured. It was experimentally verified that the system effectively responds to changes in impedance, ensuring a constant stimulation current over time, even in the presence of sudden variations.



**Fig. 3.** Generated signals for driving a muscle group (biceps): green and blue are the filtered signals at the input of the Howland circuit, yellow is the signal at the positive electrode; **a** 20 mA and **b** 30 mA stimulation.

## 4   Conclusions

The paper proposes an active suit system to counteract muscle and bone loss in astronauts during microgravity. The suit uses sensors, flexometers, and electrostimulation to simulate Earth-like gravity conditions. The suit's sensors interact in real-time with electrical stimulation, performed by Howland current pump, tailored to specific needs. Lab tests confirm the suitability of the Electrical Control Unit for controlled and safe current stimulation. The suit shows promise for astronaut health during space exploration and long-duration missions.

## References

1. Nguyen N, Kim G, Kim K-S (2020) Effects of microgravity on human physiology. Korean J Aerosp Environ Med 30(1):25–29
2. Bradbury P et al (2020) Modeling the impact of microgravity at the cellular level: implications for human disease. Front Cell Dev Biol 8:00096
3. Virgin Galactic has sold 800 tickets to the edge of space. The first customers just took fight. https://edition.cnn.com/2023/06/29/travel/virgin-galactic-launch-italian-air-force-scn/index.html. Accessed 01 July 2023

4. Crew-3 Astronauts Launch to Space Station Alongside Microgravity Research. https://www.nasa.gov/mission_pages/station/research/news/Crew-3-Launch-to-ISS-Alongside-Microgravity-Research/. Accessed 30 June 2023
5. ElectroMyoStimulation (EMS body suit) to Improve Exercise Outcome in Space Missions. https://www.nasa.gov/mission_pages/station/research/experiments/explorer/Investigation.html?#id=8526. Accessed 02 July 2023
6. Morun C, Lake S (2018) Systems, articles, and methods for capacitive electromyography sensors. U.S. Patent No 10,042,422
7. Texas Instruments (2008) An-1515 a comprehensive study of the Howland current pump. Application report SNOA474A

# Batteries

# Electrochemical and Thermal Modelling of a Li-Ion NMC Pouch Cell

Aljon Kociu[✉], Luca Pugi, Lorenzo Berzi, Edoardo Zacchini, Massimo Delogu, and Niccolò Baldanzini

Department of Industrial Engineering, University of Florence, Florence, Italy
aljon.kociu@unifi.it

**Abstract.** Lithium-ion battery cells are considered the best solution for energy storage systems in Battery Electric Vehicles (BEVs), due to their high energy density, high lifecycle, and low aging impact. Battery performances are strictly connected to operating conditions, especially temperature. Modern Battery Thermal Management Systems (BTMSs) optimize battery usage in terms of temperature to achieve high performance and high safety of the battery. In this paper an electrochemical-thermal coupled 3D model will be developed and discussed. The electrochemical model used is a pseudo-bidimensional model (P2D) which will output the heat source for the thermal model in order to get from this one the temperature of the battery when charging and discharging. The model data will be validated through experimental data coming from a 30 Ah NMC pouch cell used in automotive applications.

**Keywords:** Electrochemical · Thermal · P2D · NMC

## 1 Introduction

Lithium-ion batteries are the most promising technology for vehicle energy storage systems (ESSs). Their complex electrochemical nature and, above all, the strong temperature dependency during operation are subject of this study. In this paper a 3D electrochemical model coupled with a 3D thermal model is proposed. Both electrochemical and thermal model are used to study a high capacity NMC pouch cell. The electrochemical model with its temperature-dependent parameters was validated through charge-discharge cycles to accurately reproduce the operating conditions of the battery cell. The heat sources are distinguished and characterized through an advanced electrochemical approach. The coupled model provides information on the current and voltage profiles that cause the least thermal stress thus effectively reproducing the distribution and temperature gradients in the battery. Identifying these quantities helps in the development of adequate cooling and/or heating systems in order to use the battery cells at optimal temperature ranges. Once the behaviour of a single cell is established, it is then possible to develop models that represent the entire battery pack, so an even more accurate and reliable result is reached. This process optimizes battery thermal management system (BTMS) algorithms, improves lifespan, performance, and safety of future energy storage systems in electric vehicles.

## 2   Electrochemical and Thermal Model

A battery model lies on of a group of equations representing battery performance and characteristic curves with a certain degree of approximation. Battery models presented in literature can be classified in three main categories: the physics-based electrochemical models, the electrical equivalent circuit models (ECM) and the data driven models [1]. In this paper an accurate 3D electrochemical model coupled with a 3D thermal model will be proposed and discussed. The coupled electrochemical-thermal battery model is a high-efficient and powerful tool to study lithium-ion battery cell characteristics. Electrochemical models typically consider a single unit of the entire battery cell (a single sheet) as study object. The battery sheet includes current collectors, positive electrode (cathode), negative electrode (anode) and separator, as shown in Fig. 1 [2]. The main electrochemical models presented in literature are the pseudo-bidimensional (P2D) model, the single-particle (SP) model and the extended SP (ESP) model [3, 4]. Although the pseudo-bidimensional (P2D) model is time consuming, it represents the state of the art in terms of battery electrochemical modelling, with its main advantage to be found in high accuracy. Not only, the P2D model can dialogue with aging models to estimate internal battery cell parameters [5]. Thus, it will be implemented in this study.

The electrochemical model provides the heat generation volumetric rate as output, that represents the input for the coupled thermal model in terms of heat source, allowing the study of thermal characteristics and improve simulation accuracy. The thermal model can calculate the temperature of the cell which itself will become the new input for the electrochemical model to give an accurate estimate of the heat generation rate (see Fig. 2). While the electrochemical 3D model considers a single sheet of the cell, the thermal 3D model considers the whole battery as a single material which properties are averaged by all the sheets that make up the whole cell. Finally, simulations are executed, and finite element model is used to solve the equations.

## 3   Definition of the Model

As mentioned previously, the electrochemical model adopted in this study is a pseudo bidimensional model (P2D). The main difficulty associated with the model is the parameterization, however this allows to characterize the primary dynamics of a lithium-ion battery cell thanks to its implicit versatility which allows to incorporate the thermal part. When defining the electrochemical model a few assumptions are made, such as the porous electrode theory [6], mass transportation and charge transportation, and the concentrated solution theory [7]. As regards thermal modelling, heat generation of collectors, heat dissipation in battery tabs, and irradiation of the battery are neglected. Moreover, it is assumed that the current flows uniformly in each cell and the distribution of heat generation is equivalent for each of them.

The aim of the work is to build a model that is more accurate than a 1D electrochemical model but at the same time limit the increase of computational cost typical for a 3D model. Both electrochemical and thermal model have been validated through the experimental data gained at Moving Lab, University of Florence. The validation consisted in a charge-discharge 1C cycle of the battery under investigation.

**Fig. 1.** Lithium-ion battery cell schematization [2]



**Fig. 2.** Electrochemical and thermal model interaction overview [6].

## 4   Model Implementation

### 4.1   Electrochemical Model Implementation

As mentioned in the previous sections, a single sheet of the entire battery cell has been considered. This consists of half positive current collector (since collectors are shared between sheets), one positive electrode, one electrolyte separator, one negative electrode and half negative current collector, as depicted in Fig. 3.

The input parameters are provided manually to the solver so that they fitted the battery cell under examination. The battery cell itself was designed maintaining the same length and width but modifying the depth so that it reproduced a single sheet. Material used for every layer of the cell sheet are derived from an extended literature research and extended simulation, results are shown in Table 1. Once the model parameters are defined the simulation could be performed.

### 4.2   Thermal Model Implementation

In the thermal model three different scenarios have been analysed: single sheet cell, one-to-one scale battery cell with and without heat sinks in both surfaces of the pouch. The

**Fig. 3.** Composition of a single battery cell sheet.

**Table 1.** Materials selected for each component.

| Component | Material |
|---|---|
| Negative current collector | Copper—Cu |
| Anode (negative electrode) | Graphite—LixC$_6$ MCMB |
| Electrolyte separator | LiPF$_6$ in 3:7 EC:EMC |
| Cathode (positive electrode) | LiNi$_{0.6}$Mn$_{0.2}$Co$_{0.2}$O$_2$—NMC622 |
| Positive current collector | Aluminium—Al |

different scenarios depicted have been assessed with different convective coefficients to evaluate battery performance in different working condition in terms of temperature and its gradient. All the heat sources are assumed as volumetric heat sources in the 3D geometry. Moreover, conduction is considered the heat transfer mechanism inside the battery cell, while convection is the heat transfer mechanism between the battery cell and the environment. As for the materials, in the first single-sheet scenario materials associated to the various components are identical to those of the electrochemical model. In further one-to-one scale scenarios, ad-hoc material was characterized. The new material considers the three main properties concerning the developing of the thermal analysis (density, heat capacity and thermal conductivity) which have been averaged based on real materials found in battery cells, Table 2.

**Table 2.** Ad-hoc material designed.

| Property | Value |
|---|---|
| Heat Capacity [J/(kg·K)] | 839.4 |
| Density [kg/m$^3$] | 2010.9 |
| Thermal Conductivity [W/(m·K)] | 1.28 |

Once the thermal model is defined, a time-dependent study was performed on the different scenarios. The simulation receives heat generation values as an input from

the electrochemical model and it outputs the battery temperature as an outcome. The results obtained are validated through experimental data of the pouch cell afterwards. The battery cell temperature is controlled and kept constant in a climatic chamber (25 °C) and undergoes a convective flux in its whole surface. To mimic this flux, the model implements a convective coefficient h = 15 W/m$^2$*K which was considered appropriate since the climatic chamber does not represent a natural convection environment due to the periodical activation of a fan used to maintain the temperature in the chamber constant. Air cooling is simulated in the battery surface.

## 5   Results and Model Validation

Model tests were performed simulating charge and discharge cycles with 1C current rate. Charge and discharge are performed between ranges of 3.4–4.2 V for the previous one and 4.2–3 V for the latter. A coupled electrochemical-thermal 3D simulation is time-consuming, so to simplify the simulation and make it more efficient the time step was 20 s. In the next figures, voltage versus SoC is plotted in different situations.

Figure 4 compares the electrochemical model and experimental data in a charge and discharge cycle. Despite the electrochemical model considering a single sheet of the cell, it can accurately describe battery charge-discharge curves within a modest error.



**Fig. 4.** Experimental data versus model.

The simulated average surface temperature of the battery cell while charging and discharging with different convective heat transfer coefficients (5–10-15 W/m$^2$*K) is shown in Fig. 5. The simulation highlights the effectiveness of the heat flux coefficient on battery thermal performance. Considering the working condition in which experimental data was gained (climatic chamber with static temperature set at 25 °C) the higher coefficient appeared to be the most suitable to properly represent battery temperature changes during cycling.



**Fig. 5.** Heat fluxes.

Finally, the results of validation process are shown in Fig. 6. Despite the simplification introduced (a single sheet in the electrochemical model and an averaged single material for the thermal model), the system developed in this study appears to accurately represent real temperature curves through charge and discharge cycles.



**Fig. 6.** Validation through experimental data.

## 6   Conclusions and Future Developments

This work shows that an accurate coupled electrochemical-thermal model with reduced computational cost can be developed to be implemented in real-time applications. To achieve this, a single sheet for the electrochemical model and an averaged material for the thermal model have been considered, which impact is under investigation. The model can be an important tool to examine different batteries in multiple scenarios of interest with a high degree of confidence. Further understanding and research on battery material properties may result in better accuracy of the model. The implementation of State of Charge and State of Health estimation filters is the next step of this work [8].

## References

1. Wang Y et al (2020). A comprehensive review of battery modeling and state estimation approaches for advanced battery management systems. https://doi.org/10.1016/j.rser.2020.110015
2. He CX, Yue QL, Wu MC, Chen Q, Zhao TS (2021) A 3D electrochemical-thermal coupled model for electrochemical and thermal analysis of pouch-type lithium-ion batteries. Int J Heat Mass Transf 181. https://doi.org/10.1016/j.ijheatmasstransfer.2021.121855
3. Doyle M, Fuller TF, Newman J (1992) Modeling of galvanostatic charge and discharge of the lithium/polymer/insertion cell. J Electrochem Soc 140:1526
4. Romero-Becerril A, Alvarez-Icaza L (2011) Comparison of discretization methods applied to the single-particle model of lithium-ion batteries. J Power Sources 196:10267–10279. https://doi.org/10.1016/j.jpowsour.2011.06.091
5. An F, Zhou W, Li P (2021) A comparison of model prediction from P2D and particle packing with experiment. Electrochim Acta 370. https://doi.org/10.1016/j.electacta.2021.137775
6. Schweiss R (2020) Validation of 1D porous electrode theory using steady-State measurements of flooded electrodes at variable electrolyte compositions. Chem Eng Sci 226. https://doi.org/10.1016/j.ces.2020.115841
7. Vynnycky M, Borg KI (2010) On the application of concentrated solution theory to the forced convective flow of excess supporting electrolyte. Electrochim Acta 55:7109–7117. https://doi.org/10.1016/j.electacta.2010.06.070

8. Locorotondo E, Lutzemberger G, Pugi L (2021) State-of-charge estimation based on model-adaptive Kalman filters. Proc Inst Mech Eng. Part I: J Syst Contr Eng 235:1272–1286. https://doi.org/10.1177/0959651820965406

# Low-Cost Configurable Electronic Load for Lithium Ion Batteries Testing

Niccolò Nicodemo[(✉)], Roberto Di Rienzo, Alessandro Verani,
Federico Baronti, Roberto Roncella, and Roberto Saletti

Dipartimento di Ingegneria dell'Informazione, Università di Pisa, Via Caruso 16,
56122 Pisa, Italy
niccolo.nicodemo@phd.unipi.it

**Abstract.** Characterization data are essential to improve state estimation algorithms for lithium-ion batteries. Unfortunately, only bigger companies can afford extensive test campaigns, as they require expensive and specific equipment. A Low-Cost, Open-Design, and highly configurable electronic load for battery testing is proposed in this paper. The developed equipment is able to sink a desired current profile configured by using a host computer Python interface. A preliminary experimental validation of the designed electronic load is performed obtaining promising results. This equipment could help smaller companies and laboratories to perform battery characterization tests and increase the amount of available data for state estimation algorithms improvement. For this reason, it is released as Open Hardware/Software system and is freely available to the research community.

**Keywords:** Lithium-ion batteries · Battery characterization ·
Low-cost · Open-design

## 1 Introduction

Advancements in Lithium-Ion Batteries (LIBs) technology considerably boosted vehicle electrification in the last years, paving the way to reduce greenhouse gas pollution. However, LIB technology still presents significant open challenges [1], such as the estimation of the charge left in the battery and its performance degradation due to aging. Improving the estimation accuracy of such quantities would address some of the main factors that prevent a wider diffusion of electric vehicles, e.g. the range anxiety [2]. Unfortunately, expensive test campaigns and large amount of data are required to improve battery state estimation algorithms [1,3]. Those limitations often make research on LIBs a prerogative of big companies. On the other hand, small laboratories rely on test setup based on general purpose equipment and specific setup configuration to reduce the power requirements and thus equipment cost [4–6]. Unfortunately, component shortage caused by the breakout of Covid pandemic is still affecting the research community [7,8] reducing the availability of laboratory equipment and further

hampering the research. The Open Hardware approach emerges as an important tool in this scenario for carrying out scientific researches in adverse conditions and under-resourced areas [9,10]. This work aims to provide Low-Cost Open Hardware testing equipment for battery characterization. The developed equipment draws a controlled current from a device under test, so it can be used to discharge a battery with controlled current profiles. In addition, the proposed equipment can be used to charge a battery if an auxiliary battery setup is used as suggested in [6]. Similar works have been already presented in the literature [11], but they are focused on reaching higher performance and accuracies rather than low cost, and configurability. Instead, the proposed equipment features a modular structure which allows it to cover different current ranges and different accuracies. In addition, the design is provided to be freely modified according to the desired budget and required performance. The equipment comes with an open-source software interface with which it can be controlled. The developed software allows the users the execution of laboratory tests and custom current profiles to emulate any real-world applications. Both Hardware and Software design can be downloaded from [12] with the hope of encouraging the contribution of smaller companies and laboratories to LIB related research.

## 2   Hardware and Software Platform Description

The architecture of the developed equipment is shown in Fig. 1 and consists of two parts: a control board, and a current sink circuit. The control board drives the current sink circuit, measures the main system quantities, and manages the communication with the control interface. The current sink circuit is based on a well-known topology [11,13,14]. It is composed of an operational amplifier OA, a shunt resistor $R_{sh}$ and a n-channel MOSFET $M_{Out}$. A desired voltage $V_{sh}$ is imposed across $R_{sh}$ by the operational amplifier and the MOSFET $M_{Out}$. A current $I_s$ is consequently forced in $R_{sh}$ as a function of $V_{sh}$ according to the Ohm's Law. The range and accuracy of $I_s$ mainly depend on the MOSFET and the shunt resistor characteristics. Low values of $R_{sh}$ keep low its power dissipation but enhance the effects of offset errors resulting in inaccurate $I_s$ values. Instead, the device under test voltage $V_{dev}$ has no effects on the controlled current values as long as it allows a correct polarization of $M_{Out}$.

$R_{sh}$ and $M_{Out}$ can arbitrarily be changed according to the application requirements. For this reason, $R_{sh}$ and $M_{Out}$ are placed in a module sepa-



**Fig. 1.** General architecture

rated from the control board. This module is hereinafter referred to as "output module", as shown in Fig. 1. An IXTN660N04T4 N-Mos from IXYS and a WSBS5216L1000JK14 100 $\mu\Omega, 5\%$ shunt resistor from Vishay have been considered for the first implementation of the equipment. The N-Mos is capable of managing a DC current of 100 A with a maximum $V_{ds}$ of 5 V. It is mounted on a 10 cm $\times$ 7 cm $\times$ 4 cm heatsink.

The control board uses a 16-bit Digital-to-Analog Converter (DAC) from Texas Instrument, the DAC8501, to generate the voltage $V_{sh}$. In particular, $V_{sh}$ is generated as a partition of the DAC output voltage in order to match the small voltage drop across $R_{sh}$ and still exploit the full DAC output dynamic. The default partition ratio $A_{DAC}$ is set to 1/101 but it can easily be changed by modifying the resistance values of the DAC output voltage divider.

A 16-bit Delta-Sigma Analog-to-Digital Converter (ADC) from Microchip, the MCP3462R, is used to measure $V_{dev}$ and $V_{sh}$ whose dynamic ranges are properly conditioned by Precision Operational Amplifiers OPA2156 from Texas Instruments. Default gain values $A_{ADC_{sh}}$ and $A_{ADC_{dev}}$ are set to 101 and 1/101, respectively, but can easily be changed by modifying the resistance values of the conditioning stage. Both DAC and ADC use a Serial Peripheral Interface (SPI) as communication protocol. A MCP2210 USB-to-SPI converter from Microchip Electronics handles the communication between the control board and the Software application on a computer. The communication is electrically insulated by means of a MAX14850 from Maxim Integrated.

A coarse estimated cost of the developed equipment, including only the components and the Printed Circuit Board costs, is around 110 \$, 50 \$ of which are due to the specific output module described above. This cost is much lower than the price of commercial electronic loads with a similar maximum current value. For example, the cost of LD400 from Aim TTi (60 A), EL 3080-60 B from Electro-Automatik (60 A), and BK8540 from B&K Precision (30 A) are around 1400 \$, 1300 \$, and 700 \$, respectively. These prices are obtained by consulting Farnell Italia Srl distributor in September 2023.

## 2.1   Software Platform

An open-source software application to control the equipment has been developed using the Python3 language and the Graphical User Interface (GUI) design module PyQt5. The application allows the user to control the device and log the quantities measured by the equipment on a host computer. The application offers two different control modalities: Manual and Automatic. The user can set a specific value of output current in Manual mode. Instead, a specific output current profile can be executed in Automatic mode. In addition, the application allows the user to create and edit test profiles that are saved as Tab-delimited text files. The test file is composed of two columns: the step requested current, and the step time length. Instead, the adjustable equipment parameters, i.e, sampling time, $R_{sh}$, $A_{DAC}$, and $A_{ADC_x}$, can be configured by editing the configuration file "Configuration.py". The minimum step time length and sampling time are set to 100 ms.

**Fig. 2.** Software architecture

## 3    Experimental Setup and Results

A preliminary evaluation of the developed equipment was carried out using the experimental setup shown in Fig. 3a. It is composed of a QPX1200SP DC power supply from AimTTi used to emulate a lithium-ion cell under test. A NI cDAQ-9178 equipped with one NI 9228, one NI 9219, and a 1 mΩ shunt resistor was used as a reference measurement system. The NI 9228 features eight 24-bit/±60 V differential voltage inputs that are used to measure $V_{dev}$, $V_{sh}$, and the voltage across the 1 mΩ reference shunt resistor. The NI9219 analog input module is used to monitor the MOSFET $M_{Out}$, the heatsink, and the ambient temperatures by means of thermocouples. Data acquisition from NI modules was performed at 2 kHz sampling rate.

First, a thermal characterization test was performed. It consists in sinking a constant current of 25 A with a $V_{dev}$ of 1.5 V for 15 min and using the subsequent thermal relaxation to identify the parameters of a 2RC Foster thermal model for the N-Mos assembled on the heatsink. The obtained model allowed us to estimate the Safe Operating Area (SOA) for the DC operations of the equipment. The SOA is reported in Fig. 3b, considering a maximum N-Mos junction temperature of 150 °C and an ambient temperature of 25 °C. It is worth pointing out that this limitation is only due to the components chosen for the output module and the characteristics of the used heatsink.



(a)          (b)

**Fig. 3.** Test Setup (**a**) and estimated SOA for DC operations (**b**)

A constant current discharge phase of a lithium-ion cell is considered as the first case study for the equipment validation. The voltage $V_{dev}$ is set in the working ranges of various lithium-ion cell chemistries (1.5–4.5 V) while the developed equipment was set to sink a current sweep profile. The sweep is composed of 10 s constant current steps whose value is increased of 1 A in each step up to the maximum value allowable by the estimated SOA. Figure 4a shows the relative errors of the set and measured currents. The set error is obtained by comparing the theoretical current value with the reference measure system one. Instead, the measure error is evaluated as the difference between and the measured current by the reference system current with the developed equipped one. The results show set and read errors lower than 0.5% for all the tested $V_{dev}$, with current values higher than 7 A. As expected, higher errors are obtained for low current values, where offset errors have more impact. The maximum absolute set and read errors are 80 and 56 mA, respectively.

The Urban Dymamometer Driving Schedule (UDDS) was considered as second case study. It aims to evaluate the equipment capability to accurately generate a standard automotive application profile. Only the discharge phases of the UDDS cycle were considered, i.e., the effects of regenerative braking were ignored. In particular, the profile was scaled to have a maximum peak current of 50 A. $V_{dev}$ was set to 3.7 V to emulate the typical scenario for an Nickel-Manganese-Cobalt LIB in Electric Vehicles. The UDDS ideal current profile $I_{UDDS}$ was programmed on the equipment. Figure 4b shows $I_{UDDS}$, the current measured by the reference measure system $I_{ref}$, and by the developed equipment one $I_{meas}$. The accuracy of reproducing the ideal current profile was quantified by comparing the time integral of $I_{UDDS}$, $I_{ref}$, and $I_{meas}$. The difference between the $I_{UDDS}$ and $I_{ref}$ charge values, and the difference between $I_{ref}$ and $I_{meas}$ ones resulted to be around 53 and 58 C, i.e. less than 0.8% of the total charge moved by the ideal profile.



(a)                    (b)

**Fig. 4.** Set and measure relative errors (**a**) and results for UDDS profile (**b**)

# 4    Conclusion

A low-cost and highly configurable electronic load for lithium-ion battery testing is proposed in this paper. The equipment is able to sink a controlled current value with a set and measurement accuracies lower than 0.5% for almost all the test conditions. Moreover, the developed equipment is able to reproduce standard driving cycles obtaining a very good accuracy. A simple thermal model of the equipment is developed to identify its Safe Operating Area, whose limitations appear to be related only to the thermal management of the chosen output module power components, rather than to the equipment architecture. However, the modular architecture of the developed equipment allows to easily adapt the current and power test ranges by changing the output module components. Finaly, the developed device and its software platform are provided Open Hardware and Open Software to allow smaller companies and laboratories to contribute to LIB related research.

# References

1. Wang Z, Feng G, Zhen D, Gu F, Ball A (2021) A review on online state of charge and state of health estimation for lithium-ion batteries in electric vehicles. Energy Rep 7:5141–5161
2. Liu X, Zhao F, Geng J, Hao H, Liu Z (2023) Comprehensive assessment for different ranges of battery electric vehicles: Is it necessary to develop an ultra-long range battery electric vehicle? iScience 26:106,654
3. dos Reis G, Strange C, Yadav M, Li S (2021) Lithium-ion battery data and where to find it. Energy AI 5:100,081
4. Carloni A, Baronti F, Di Rienzo R, Roncella R, Saletti R (2018) Open and flexible li-ion battery tester based on python language and raspberry pi. Electronics 7:454
5. Vergori E, Mocera F, Somà A (2018) Battery modelling and simulation using a programmable testing equipment. Computers 7:20
6. Nicodemo N, Di Rienzo R, Verani A, Baronti F, Roncella R, Saletti R (2023) Low-cost lithium-ion battery characterization setup based on auxiliary batteries. In: Applications in electronics pervading industry, environment and society. Springer Nature, Switzerland, pp 157–162
7. Greenstein S (2021) Shortages of integrated circuits. IEEE Micro 41:86–88

8. Pennisi S (2022) Pandemic, shortages, and electronic engineering. IEEE Circuits Syst Mag 22:41–49
9. Stirling J, Bowman R (2021) The covid-19 pandemic highlights the need for open design not just open hardware. Des J 24:299–314
10. Webb H, Bezuidenhout L, Nurse JR, Jirotka M (2019) Lab hackathons to overcome laboratory equipment shortages in Africa: Opportunities and challenges. In: Conference on human factors in computing systems—proceedings
11. Weßkamp P, Haußmann P, Melbert J (2016) 600-A test system for aging analysis of automotive li-ion cells with high resolution and wide bandwidth. IEEE Trans Instrum Meas 65:1651–1660
12. BatteryLabUnipi: https://github.com/batterylabunipi/Low-Cost_Electronic_Load
13. Zamora M (2020) Precision current sources and sinks using voltage references application report precision current sources and sinks using voltage references. Application Report—SNOAA46
14. TI: Tips and tricks for designing with voltage references. AN—SLYC147A (2021)

# In-Home Measurement and Analysis of Health-Related Physiological Parameters

# Remote Healthcare System Based on AIoT

Alberto Cabri[1,3]($\boxtimes$), Stefano Rovetta[2,3], Francesco Masulli[2,3], Akshi Sharma[2], Pier Giuseppe Meo[4], and Mario Magliulo[4,5]

[1] University of Milan, Computer Science Department, Via Celoria 18, Milan, Italy
`alberto.cabri@unimi.it`
[2] University of Genoa, DIBRIS, Via Opera Pia 13, Genoa, Italy
[3] Vega Research Laboratories srl, Via Ippolito D'Aste 7, Genoa, Italy
[4] Hassisto srl, Via Mercato 16e, Mugnano di Napoli, Naples, Italy
[5] CNR - Istituto di Biostrutture e Bioimmagini, Via De Amicis 95, Naples, Italy

**Abstract.** Life expectancy in recent years has sensibly increased and age related problems in elderly people have followed a similar trend. Being able to find innovative solutions to enable senior population to maintain their quality of life despite the presence of chronic illnesses has become crucial for high quality ageing. The opportunities offered by the technological advancement with remote assistance applications, wearable devices, and *Artificial-Intelligence-Of-Things (AIoT)* architectures are of paramount importance to improve the services in healthcare facilities by adding the power of Artificial Intelligence to Internet-Of-Things devices. An experimental framework has been deployed to two residential homes in collaboration with two italian companies to collect and analyze data in order to actively monitor the vital signs of their guests, predict critical situations and identify significant clusters or communities.

**Keywords:** Remote assistance · Healthcare · Artificial-intelligence · Internet-of-things · Smart devices

## 1 Introduction

In developed countries, average life expectancy has sensibly increased in the last few years thanks to improved healthcare services, active prevention of diseases and pathologies and the availability of new drugs. As a result, we are witnessing an increase in geriatric population and consequently in a group of diseases that are directly related to ageing, such as degenerative diseases. As reported in [1], the exponential growth of elderly people requiring treatment for chronic life-long diseases will dramatically impact the cost of healthcare and the ability to provide the necessary treatments that, at the time being, can be effectively supported by innovative smart devices and technologies. In fact, new portable or wearable devices which integrate disparate monitoring sensors are nowadays available for everyone at low cost, making it possible to constantly assess patients' health

status through a Health Monitoring System (HMS) that not only limits hospitalization and medical staff intervention but also cuts the waiting lists improving the consultation effectiveness while reducing the overall healthcare expenditure [2]. When the HMS makes use of smart devices to record and track patients' status it is termed Smart Health Monitoring System (SHMS) and can be general purpose (GHMS) when multiple generic vital signs are recorded or Remote (RHMS) if data are collected at a remote location and transferred to healthcare facilities for medical follow-up analysis. Another non-negligible benefit comes from the availability of mobile devices that can provide reliable network connections and computational power to perform data collection and timely onsite preliminary analysis via direct interaction with the medical staff. The advent of IoT devices is fostering the availability of new assets for the development of a brand new personalized care that can improve the quality of life for elderly people.

## 2    State of the Art

As presented in [3], various bluetooth wireless devices, such as blood pressure monitor or pulse oximeter for oxygen saturation, were used in combination with an IrDA (Infrared Data Association) connected blood glucose meter and an electronic thermometer to gather patient's data into a set-top-box, responsible for data transmission over a secure network connection. An integrated Wearable and Mobile HMS was setup to record vital signs that were manually checked by registered personnel and evaluated according to standard ward procedures, whereas the researchers were responsible for chasing inaccuracies and time delays. A fuzzy logic system where the significant physiological parameters are properly weighted for a particular condition severity, was used to interpret heart rate, blood pressure, pulse rate, temperature and oxygen saturation (SpO2) and detect cases of bradycardia. The authors of [4] propose the adoption of tabled-based applications to make their patients' data timely accessible to caregivers and provide the basic routine functions to clinicians; the application, which is required to be simple, is organized into five sections for patient related information, biomedical parameters, patient's historical record, medical notes and contact information. Similarly, wireless sensors are used by [5] to relieve chronic patients from the burden of intrusive instruments. The authors propose a three stage system based on an Arduino board for sensing, filtering and transmitting the signals over internet and later display the relevant results with a computer-based or mobile-based application. The authors of [6] assert that a decrease in health related quality of life (HRQoL) is correlated to frailty or pre-frailty status, which is why they recommend an HMS capable of achieving remote continuous monitoring without active patient's intervention by integrating Artificial Intelligence with *Internet of Things* (AIoT). IoT features in medical systems should be designed in a human-centric perspective, considering human beings as critical *components* of the system and focusing on two use cases (remote elderly monitoring and smart ambulance) that combine emerging technologies with best healthcare

practices. Patients and caregivers become *actors* in the new cyber-physical system, with specific context-driven tasks and critical issues. One recurrent goal of most papers in literature, such as [3,8], is the adoption of devices that do not interfere with patient's regular life, which is becoming more and more possible with smart bracelets and other integrated sensors.

## 3 System Architecture

Our project involved multiple *actors*, playing different roles in the operational interaction with the selected hardware devices for data sensing and collection. The patients are housed in two residential home-like accommodation (RSA Minoretti and RSA Valpolcevera) based in Genoa (Italy) that provide specialized medical care for people who cannot be cared at home. The total number of patients involved in the present study is 25, with 10 males and 15 females. The technological infrastructure was provided by Hassisto srl, a spin-off of the National Research Council (CNR), that developed an innovative software platform for e-Health, connecting multiple bluetooth devices to a hub that transmits the anonymized patients' data to a central database for continuous monitoring through a dedicated mobile application. This platform integrates a wide variety of professional medical devices, such as ECG or EEG, but also commercial wearable sensors that can measure vital signs, other biomedical data and physical activity in real-time. The primary sources of data were the IoT devices worn by each patient who was identified by a unique anonymous code to comply with the European General Data Protection Regulation (GDPR). Despite the multiplicity of smart devices supported by the hardware platform, in this initial project we adopted only two sensors and the transmission hub, not to be too invasive in the elderly patients' lives.

### 3.1 Spovan H03 Wristband

Unlike other smartwatches, the primary aim of this waterproof wristband is the collection of health data: it is and advanced fitness and health tracker which provides on-wrist ECG, blood oxygen metering, blood pressure, sleep and heart rate variability (HRV) monitoring [9].

### 3.2 The Sleeping Band

The sleeping band is a non-invasive device, which is positioned between the sheets and the mattress, at the level of user's rib cage, to monitor the quality of sleep, breathing and heart beats when the patient is in bed. It can also identify some specific movements or situations that might require the caregiver intervention, such as in case of seizure, prolonged sleep apnea or sudden chaotic movements.

### 3.3    The Communication Hub

The communication hub, namely the *Hassisto Gateway*, is based on a customized Android TV box that provides the necessary bluetooth interface towards monitoring devices and hosts an HTTP Web server, which allows for the remote connection to the collected data [6].

## 4    Data Analytics

All data collected by the selected devices are stored on a remote software platform that provides a dashboard to manage patients, devices and the relevant measurements which can be supplemented by the residential home personnel with additional medical information about each patient. Moreover, measurement data can be imported, via a REST API and a Python Request module [10], from a second portal and analyzed in a Pandas DataFrame [11] to perform various machine learning tasks, such as statistical analysis. The medical staff requested to aggregate patient's data at weekly, hourly and daily resolutions therefore the collected measurements were averaged over the requested time span, also to take into account the different measurement time intervals of each sensor device. Our focus was set on six sensors: breath, heart rate, maximum and minimum blood pressure, saturation and steps. In order to find relevant information and possibly detect abnormal situations we used the radar plots to visualize the overall status of a patient and compare it with other patients having similar health conditions. To have an informative insight of patients' health status from the collected data, cluster analysis was performed as it allows to group subjects by similar symptoms and assess the presence of possible deviation from the community profile: this was achieved by means of the graded possibilistic clustering algorithm [12], which can iteratively track outliers and adapt to concept shifts and drifts in non-stationary data.

### 4.1    Clustering Algorithms

Clustering is a family of unsupervised learning algorithms whose aim is automatic discovery of similarities among data. There are hard and soft clustering algorithms which differ by the degree of membership to the clusters that each data point may have. Specifically, in hard clustering a data point can be assigned to one cluster only with binary membership $\mu_i \in \{0, 1\}$, whereas in the soft one each data point can have a $\mu_i \in [0, 1]$ with the optional additional probabilistic constraint that $\sum_i \mu_i = 1$. The elbow method [13] was initially used to estimate the optimal number of clusters $k$ but it did not yield a significant result, therefore an empirically determined value $k = 3$ was set.

Four different algorithms, two hard and two soft methods, were considered and assessed against our dataset: k-means [14], k-medoids [15], Fuzzy C-means [16] and Graded Possibilistic Clustering [17]. Data wrangling and appropriate scaling of our dataset was necessary to compensate the different ranges of the measured

quantities, moreover, it was important that the cluster center be a true patient, which in principle makes the *k-medoids* clustering method favorable against the other presented algorithms. Radar plots were implemented via D3—Data Driven Documents - to provide single patient's hourly, daily and weekly profiles and clustering for a population of patients over a defined time span. Examples of the generated radar plots can be seen in Fig. 1.

## 4.2 Results and Discussion

Data from a four week period was selected and further divided into two blocks of two weeks each to be able to appreciate any possible change in the centroids. Standard scaler normalization was performed before running the tests and two clustering techinques, commonly used in similar problems, were used: Fuzzy C-means for soft clustering and K-means for hard clustering.   In Fig. 1, it is evident



**Fig. 1.** Comparison of clustering models built for Fuzzy C-means (left) and k-means (right)

that both algorithms produce nearly identical results over the two considered periods respectively. The limited number of patients in this study represents a serious limitation on the reliability of the results but we could demonstrate the usefulness of the adopted radar plots in conveying significant information about patient's communities to the medical staff who will in future be able to deliver more targeted therapies to their elderly guests.

## 5 Conclusions

This paper presented a preliminary study of an AIoT system aimed at active monitoring of vital signs for frail and elderly people hosted in residential homes. We used wearable sensors and sensored beds to continuously collect streams of data for a long period and perform clustering analysis over several time frames. The medical staff and caregivers can visualize a set of radar plots enabling a synoptic view of a single patient or groups of similar patients.

# References

1. Esposito M, Minutolo A et al (2018) A smart mobile, self-configuring, context-aware architecture for personal health monitoring. Eng Appl Artif Intell 67:136–156. https://doi.org/10.1016/j.engappai.2017.09.019

2. Baig MM, Gholamhosseini H (2013) Smart health monitoring systems: an overview of design and modeling. J Med Syst 37(2). https://doi.org/10.1007/s10916-012-9898-z

3. Baig MM, GholamHosseini H, Connolly MJ, Kashfi G (2014) Real-time vital signs monitoring and interpretation system for early detection of multiple physical signs in older adults. In: IEEE-EMBS international conference on biomedical and health informatics (BHI), Valencia, Spain, pp 355–358. https://doi.org/10.1109/BHI.2014.6864376

4. Baig MM, GholamHosseini H, Salete Sandini Linden M (2015) Tablet-based patient monitoring and decision support systems in hospital care. In: International conference of the IEEE engineering in medicine and biology society. https://doi.org/10.1109/embc.2015.7318585

5. Ali NS, Alyasseri ZA, Abdulmohson A (2018) Real-time heart pulse monitoring technique using wireless sensor network and mobile application. Int J Electr Comput Eng (IJECE) 8(6):5118–5126. https://doi.org/10.11591/ijece.v8i6.pp5118-5126

6. Tramontano A, Scala M, Magliulo M (2019) Wearable devices for health-related quality of life evaluation. Soft Comput 23(19):9315–9326. https://doi.org/10.1007/s00500-019-04123-y

7. Kotronis C et al (2019) Evaluating internet of medical things (IoMT)-based systems from a human-centric perspective. Internet Things 8:100125. https://doi.org/10.1016/j.iot.2019.100125

8. Esposito M et al (2018) A smart mobile, self-configuring, context-aware architecture for personal health monitoring. Eng Appl Artif Intell 67:136–156. https://doi.org/10.1016/j.engappai.2017.09.019

9. SPOVAN-Shenzhen Tianpengyu Technology Co., Ltd. High Quality ECG Smart bracelet Smart Band Spovan H03 Wholesale—Shenzhen Tianpengyu Technology Co., Ltd. 2019. https://www.spovan.com/products-detail-150445. Accessed 19 Feb 2023

10. Python Requests module. https://github.com/psf/requests

11. The pandas development team: pandas-dev/pandas: Pandas, Zenodo (2020). https://doi.org/10.5281/zenodo.3509134

12. Abdullatif A et al (2016) Graded possibilistic clustering of non-stationary data streams. In: Lecture notes in computer science. https://doi.org/10.1007/978-3-319-52962-2_12

13. Thorndike RL (1953) Who belongs in the family?. In: Psychometrika, vol 18, pp 267–276. https://doi.org/10.1007/BF02289263

14. Lloyd SP (1982) Least squares quantization in PCM. IEEE Trans Inf Theory 28(2):129–137. https://doi.org/10.1109/tit.1982.1056489

15. Kaufman L, Rousseeuw PJ (1990) Finding groups in data. In: An introduction to cluster analysis. Wiley-Interscience

16. Dunn JC (1973) A fuzzy relative of the ISODATA process and its use in detecting compact well-separated clusters. J Cybern 3(3):32–57. https://doi.org/10.1080/01969727308546046
17. Masulli F, Rovetta S. Soft transition from probabilistic to possibilistic fuzzy clustering. IEEE Trans Fuzzy Syst 14(4):516–527. https://doi.org/10.1109/TFUZZ.2006.876740

# Oral Health Phenotype of Postmenopausal Women Using AI

Rodion Kraft[1](✉), Juan A. Ortega[2] , Áurea Simón-Soro[2] , Natividad Martínez[1] ,
and Luis González-Abril[2]

[1] Reutlingen University, 72762 Reutlingen, Germany
rodion.kraft@student.reutlingen-university.de
[2] University of Seville, 41012 Sevilla, Spain

**Abstract.** Menopause is the permanent cessation of menstruation occurring naturally in women's aging. The most frequent symptoms associated with menopausal phases are mucosal dryness, increased weight and body fat, and changes in sleep patterns. Oral symptoms in menopause derived from saliva flow reduction can lead to dry mouth, ulcers, and alterations of taste and swallowing patterns. However, the oral health phenotype of postmenopausal women has not been characterized. The aim of the study was to determine postmenopausal women's oral phenotype, including medical history, lifestyle, and oral assessment through artificial intelligence algorithms. We enrolled 100 postmenopausal women attending the Dental School of the University of Seville were included in the study. We collected an extensive questionnaire, including lifestyle, medication, and medical history. We used an unsupervised k-means algorithm to cluster the data following standard features for data analysis. Our results showed the main oral symptoms in our postmenopausal cohort were reduced salivary flow and periodontal disease. Relying on the classical assessment of the collected data, we might have a biased evaluation of postmenopausal women. Then, we used artificial intelligence analysis to evaluate our data obtaining the main features and providing a reduced feature defining the oral health phenotype. We found 6 clusters with similar features, including medication affecting salivation or smoking as essential features to obtain different phenotypes. Thus, we could obtain main features considering differential oral health phenotypes of postmenopausal women with an integrative approach providing new tools to assess the women in the dental clinic.

**Keywords:** Menopause · Artificial intelligence · Phenotype

## 1 Introduction

Menopause is the permanent cessation of menstruation that occurs naturally in women's aging, generally between the ages of 45 and 55. This biological transition is accompanied by various symptoms, such as mucosal dryness, increased weight and body fat, and alterations in sleep patterns. Oral symptoms associated with menopause, such as reduced saliva flow, can lead to dry mouth, ulcers, and changes in taste and swallowing patterns.

Several soft tissue lesions, including benign mucosal pemphigoid, lichen planus, candidiasis, and pemphigus vulgaris, have been identified as menopause-associated disorders. Despite this knowledge, the oral health phenotype of postmenopausal women remains largely unexplored.

The aim of this study is to determine the oral phenotype of postmenopausal women by examining factors such as medical history, lifestyle, and oral assessments, utilizing artificial intelligence algorithms for analysis. To achieve this goal, 100 postmenopausal women attending the Dental School of the University of Seville were included in the study.

Participants completed a comprehensive questionnaire that covered various aspects of their lifestyle, medication use, and medical history. In addition, an oral examination was conducted to assess dental and periodontal status, salivary flow, and oral pH levels. The collected data was then analyzed using artificial intelligence algorithms, including k-nearest-neighbor and decision trees, to classify the data and identify which features had the most significant impact on the oral health phenotypes observed in the study participants.

By employing machine learning algorithms, the study aims to provide a better understanding of the specific oral health challenges faced by postmenopausal women. This understanding may help dental practitioners and healthcare professionals tailor their approach to better address the unique needs of this patient population. Furthermore, the results may contribute to the development of targeted interventions or preventive measures aimed at improving oral health and overall well-being in postmenopausal women.

## 2   Related Work

Marne et al. [1] use a breast cancer dataset for applying k-means clustering to group data based on similarities in the features. After clustering, the decision tree algorithm was applied to each cluster to create classification models for predicting breast cancer. The decision tree algorithm was chosen for its simplicity and ability to handle both continuous and categorical data. The performance of the proposed method was evaluated using accuracy, sensitivity, specificity, and precision. The results demonstrated that the combined use of decision tree and k-means clustering techniques significantly improved the accuracy of breast cancer prediction, outperforming other existing methods. The study shows that combining decision tree and k-means clustering techniques can effectively classify and predict breast cancer, potentially aiding healthcare professionals in early detection and treatment of the disease.

Rajaguru et al. [2] use a breast cancer dataset to compare the performance of k-nearest-neighbor and decision trees using accuracy, sensitivity, specificity, and precision. The results showed that the decision tree algorithm had higher accuracy, sensitivity, and specificity compared to the KNN algorithm. However, the KNN algorithm demonstrated higher precision than the decision tree algorithm. This study provides insights into the effectiveness of the decision tree and k-nearest neighbor algorithms in classifying breast cancer. While both algorithms have their advantages and limitations, the decision tree

algorithm outperforms the KNN algorithm in terms of accuracy, sensitivity, and specificity. These findings can be valuable for healthcare professionals and researchers seeking to improve early detection and treatment of breast cancer using machine learning techniques.

Shin et al. [3] compared two machine learning algorithms by using three medical datasets from the UCI Machine Learning Repository, which included the Pima Indians Diabetes dataset, the Breast Cancer Wisconsin dataset, and the Hepatitis dataset. These datasets contained various features related to each disease, such as patient demographics, medical history, and test results. The decision tree algorithm and the bagging algorithm were applied to the datasets, and their performances were evaluated using accuracy, precision, recall, and F1-score. The bagging algorithm was chosen for comparison due to its ability to improve the performance of base classifiers, like decision trees, by combining multiple models and reducing overfitting.

The results indicated that the bagging algorithm outperformed the decision tree algorithm in terms of accuracy for all three datasets. Additionally, the bagging algorithm showed better performance in terms of precision, recall, and F1-score for the Pima Indians Diabetes and Breast Cancer Wisconsin datasets. However, the decision tree algorithm had slightly better results in terms of precision and recall for the Hepatitis dataset.

Zeng et al. [4] propose a medical and health data classification method based on machine learning to improve the accuracy and efficiency of medical diagnosis and treatment. The authors aim to address the challenges associated with the rapidly increasing volume of medical data and the need for effective tools to analyze and classify such data. Here the combination of machine learning algorithms, such as decision trees, support vector machines, and k-nearest neighbors were used to classify medical and health data. They also utilized a feature selection algorithm to select the most relevant features, reducing the dimensionality of the data and improving the classification performance. The algorithms were tested on three different datasets and evaluated by accuracy, precision, recall and F1-score. The results demonstrated that the proposed method outperformed traditional classification methods and achieved high accuracy and efficiency in classifying medical and health data. The authors suggest that the method can be applied to various medical fields, including disease diagnosis, treatment planning, and prognosis prediction. In conclusion, this study presents a machine learning-based method for medical and health data classification, which has the potential to improve the accuracy and efficiency of medical diagnosis and treatment processes. By employing a combination of machine learning algorithms and feature selection techniques, the proposed method offers a promising approach to handling the challenges associated with the increasing volume of medical data.

## 3  Methods and Results

*Design of the study and sample selection*

The population of this study consisted of 100 patients from the University of Seville. The inclusion criteria were women in any phase of menopause aged 40 years or older.

Exclusion criteria were incapacitating disease for obtaining clinical data or biological samples independently, women under 40 years of age, men, and women who, meeting the inclusion criteria, did not sign the informed consent for the study. To ensure anonymity, each participant was labeled with an identification code to preserve the affiliation data, separate from the clinical data.Prior to the oral and dental clinical data collection, a basic health questionnaire was carried out in which aspects of lifestyle-exercise, diet, toxic habits (tobacco and alcohol), coffee-allergies, medical and family history, autoimmune diseases and medication were collected. In addition, a menopause questionnaire was included to determine the stage of menopause and the use of hormone therapy.

*Clinical data collection*

For the collection of oral and dental clinical data, the clinical history facilitated the systematic and complete evaluation of the stomatognathic apparatus, through codes associated with the most frequent pathologies. The function of the temporomandibular joint, state of the dentition, caries, fillings, tooth loss, fissures, fractures, hypoplasia, teeth acting as bridge or crown abutments and unerupted teeth were recorded. Finally, the community periodontal index was performed to determine the periodontal status.

*Data analysis*

For the experiments we used k-nearest-neighbor (knn), decision tree and support vector machine (svm). At first, we preprocessed the data by filtering features like specific diseases, which occur rarely. We combined them to one binary feature which says that there is a disease. After this we selected the most important features for the menopause and used them as the main feature for the prediction of the classification with the three different algorithms. These features are the menstrual status (MS), the SA-Flow and Periodontal DX. The prediction was made by using the python package scikit-learn which contains pre-implemented methods for knn, decision tree and svm. For prediction it is important for defining the feature which should be predicted. Starting with the KNN we defined two different variants, one with three neighbors and one with seven by using the elbow method. Then we predict the features for the menstrual status. Here we get an accuracy of 66% and we get all the patients which are correct or wrong predicted.



**Fig. 1.** KNN 3 Correct and wrong prediction of MS values

After the prediction we looked over the common features of the correct predicted patients. In all figures the red dots are the wrong predicted ones and the blue dots are the correct predicted once. In Fig. 1 it was Liver disease and hyposalivation medication. For the same procedure we do for the prediction for the SA Flow there we get an accuracy of 56.6% (Fig. 2; Table 1).

**Fig. 2.** KNN 3 Correct and wrong predicions of th SA Flow

**Table 1.** KNN3 SA Flow

| SA Flow values | Amount of patients | Common features |
|---|---|---|
| 0 | 9 | ST Flow, Cardiovascular, liver disease, hyposalivation medication, probiotic |
| 1 | 8 | Diet, Hormo10L (DM), Gastrointesti10L, Brain disorder, Liver disease, RE10L Disease/Kidney, probiotic |

And for the periodontal dx feature we get an accuracy of 30%. There is no correct prediction for the values 0, 1, 4. As the most correct predictions are with the value 2 we look just at these common features and get Liver disease and probiotic (Fig. 3).



**Fig. 3.** KKN 3 Correct and wrong prediction of the periodontal dx

We extend the KNN to 7 neighbors. For the prediction of the MS feature we do not get any common features, but the SA-Flow prediction gives us an accuracy of 63% and 19 common features which are listed in Table 2.

**Table 2.** KNN7 SA Flow

| SA Flow values | Amount of patients | Common features |
|---|---|---|
| 0 | 9 | Previous cigarettes/day, Drain disorder, liver disease, RE10L Disease/Kidney, Probiotic |
| 1 | 10 | Diet, Brain Disorder, RE10L Disease/Kidney, Probiotic |

There is not just a different amount of common features but there are also different common features. For the decision tree prediction of the menstrual status the accuracy is 73%. The most predicted value is status 2 but it has no common features (Fig. 4).



**Fig. 4.** Decision Tree correct and wrong menstrual status



**Fig. 5.** Decision Tree correct and wrong prediction of SA Flow

The prediction of the SA Flow with decision trees has an accuracy of 73% and gets common features for the value 0 which are Diet, Brain Disorder and RE10L Disease/Kidney and for the value 1 Liver disease. The prediction of the periodontal dx is just 30% and again almost just the value 2. The common features here are diet and liver disease. At last, we predict the classification with a support vector machine. Here the Menstrual status got an accuracy of 60% but without common features. The SA-Flow gets the best results with an accuracy of 73%.



**Fig. 6.** SVM Correct and wrong predictions of SA Flow

The common features for value 0 are liver disease and hyposalivation medication and for value 1 Diet, Hormo10l (DM), gastrointessti10L, Brain disorder, Liver disease and RE10L Disease/Kidney.

## 4 Conclusion

As seen in Figs. 5 and 6 the best results are gotten from the prediction of the SA-Flow by the decision tree and the SVM. The features which show the most importance for this prediction seem to be the Liver disease and RE10L Disease/Kidney or diet feature.

The accuracy of the periodontal prediction is so low that the only thing which should be remembered is that diet and liver disease are the common features of the correct predicted value 2. That also can be a result that the periodontal dx cannot be predicted with the given features as it has no similarities. With the gained knowledge it may be a possibility to rework the data and try similar predictions just with the common features of the experiments, so that the results are equals we can say that these are the most important features to predict the correct values. Also, it may be an option to clean the data and remove features where most of the patients don't have any variances like "Active Infection". The goal also would be to improve the accuracy of the prediction. As some of the common features like diet are important for different features another possibility may also be to remove them and look on the influence of the change.

# References

1. Marne S, Churi S, Marne M (2020) Predicting breast cancer using effective classification with decision tree and K means clustering technique, In: 2020 international conference on emerging smart computing and informatics (ESCI), Pune, India, pp 39–42. https://doi.org/10.1109/ESCI48226.2020.9167544
2. Rajaguru H, S R SC (2019) Analysis of decision tree and k-nearest neighbor algorithm in the classification of breast cancer. Asian Pac J Cancer Prev 20(12):3777–3781. https://doi.org/10.31557/APJCP.2019.20.12.3777. PMID: 31870121; PMCID: PMC7173366
3. Tu MC, Shin D, Shin D (2009) A comparative study of medical data classification methods based on decision tree and bagging algorithms. In: 2009 Eighth IEEE international conference on dependable, autonomic and secure computing, Chengdu, China, pp 183–187. https://doi.org/10.1109/DASC.2009.40
4. Yu Z, Fuchao C (2021) Medical and health data classification method based on machine learning. J Healthc Eng 2021:5. Article ID 2722854. https://doi.org/10.1155/2021/2722854
5. Sorrell S (2009) The rebound effect: definition and estimation. In: Evans J, Hunt L (eds) International handbook on the economics of energy, Cheltenham, Edward Elgar

# Prototyping a Compact Form Factor Module for Physiological Measurement with Multiple Applications During the Daily Routine

Erik Stahl, Mostafa Haghi$^{(\boxtimes)}$, Wilhelm Daniel Scherz, and Ralf Seepold

HTWG Konstanz - University Applied Sciences, 78462 Konstanz, Germany
mostafa.haghi@htwg-konstanz.de

**Abstract.** With the advancement in sensor technology and the trend shift of health measurement from treatment after diagnosis to abnormalities detection long before the occurrence, the approach of turning private spaces into diagnostic spaces has gained much attention. In this work, we designed and implemented a low-cost and compact form factor module that can be deployed on the steering wheel of cars as well as most frequently touch objects at home in order to measure physiological signals from the fingertip of the subject as well as environmental parameters. We estimated the heart rate and SpO2 with the error of 2.83 bpm and 3.52%, respectively. The signal evaluation of skin temperature shows a promising output with respect to environmental recalibration. In addition, the electrodermal activity sensor followed the reference signal, appropriately which indicates the potential for further development and application in stress measurement.

**Keywords:** Continuous health monitoring · Dynamic point of perception · Ubiquitous computing

## 1 Introduction

Health-related quality of life (HRQoL) encompasses various dimensions, including physical, mental, emotional, and social functioning. Research has shown that HRQoL, as a broad measure, can predict both morbidity and mortality [1].

According to World Health Organization (WHO) statistics, three of the top ten causes of death related to diseases are cardiac conditions, three are respiratory conditions, and one is associated with physical-related activities. All three categories involve parameters that can be measured using medical and non-medical sensors and devices [2].

Having said that behavioral, environmental, physiological, and psychological domains have been addressed in some other works as the four influencing domains of health. The intensive correlation and interaction of the four domains are discussed in earlier investigations. Whether the aim of monitoring is long-term health prognostics and short-term emergency detection and assistance, the processing and decision-making are subject to data acquisition from multi-sensing systems in multi-domains [3].

Health-related signals must be measured in the places where people spend most of their time, i.e., home, work and the means of transport used to commute between home and work. Wearable devices, such as smartwatches are frequently used to measure biosignals. However, to some extent obtrusiveness have introduced some other alternative approaches in private places (e.g., car and home), using non-contact, remote, or touch-based devices where wearing is not mandatory [4].

In this work, we aim to develop and implement a compact form factor, low-cost, and multi-signal measuring prototype with multiple applications for measuring HRQoL. The proposed module is operating standalone and measures several physiological and non-physiological parameters.

## 2   Materials and Methods

### 2.1   Requirement analysis

We designed a low-cost, compact, integrative, and multi-parameters measuring module from physiological and non-physiological influencing health domains. We integrated photoplethysmography (PPG), skin temperature (ST), electrodermal activity (EDA) as well as air temperature, humidity, pressure, and Volatile organic compounds (VoCs) sensors from physiological and environmental domains. We extracted heart rate (HR) and oxygen blood (SpO2) from the PPG sensor which addresses the two of vital signs. We addressed the emotional state recognition using EDA and assessed the environmental conditions using the environmental sensors [5,6].

### 2.2   Sensor Deployment

Although the environmental measurements are not dependent on the touch the physiological measurements require adequate location of sensor deployment. Therefore, we identified the fingertip as the optimal location of touch for data acquisition. This is reflected in PPG, ST, and EDA sensors deployment [5–7]. It is not optimal to measure EDA only on one finger, but the limited size of the module and setup restricts us to one finger. To improve the signal quality, the electrodes of the EDA sensor are as far apart as possible and deployed on the fingertip.

### 2.3   Electronics Design and Layout

We used Easy-EDA in order to design the electronics circuit, simulate, and build the PCB layout. The module is battery-powered using the coin cell CR2032 appropriately positioned by a battery holder. We marked the region of fingertip touch on the module which is $18 \times 11$ mm. In the contact region, we deployed the electrodes of EDA, PPG, and ST sensors. The EDA electrodes are at the largest distance within the contact region. The circles marked with 4 are the surfaces on which the electrodes are subsequently mounted. PPG and ST sensors

are mounted in the middle of the contact region side by side to ensure appropriate contact ( see Fig. 1). The environmental sensor is deployed outside of the contact region as it does not require a direct touch. We utilize inter-integrated communication (IIC) and analog-to-digital converter (ADC) to support the ST, PPG, environmental as well as EDA sensor communication with the embedded system, respectively.



**Fig. 1.** The 3D model of the module in two layers.

### 2.4   Data Collection

We utilized the board-to-board connection to attach the module to Arduino Nano RP 2040 Connect. The combination of the module and the embedded system structured a two-layer device with a dimension of 44 × 22 mm. The embedded system supports ADC and I2C wire data transmission protocols to read out all sensors. It also has a 16 MB storage and a Bluetooth and Wi-Fi modules [4,5].

### 2.5   Experimental design

We recruited three subjects in order to validate the module and sensors. The subjects were healthy including 2 males and 1 female and received the consent form. We measured data from all sensors at the same time. We used Plux sensors as the golden standard. The data from the module was compared to the reference sensors after the synchronization for further processing. We asked subjects to watch a movie including 5-min neutral, 5-min sad, and 5-min joyful sections. We fixed the module to the left hand, the second finger of the subject, and asked the subject to wear the Plux sensors on the right hand. Prior to the experiment, the subject was asked to relax for 5 min.

## 3   Results

In total, four sensors of PPG, ST, EDA, and environmental were deployed which delivered eight parameters HR and SpO2, ST, skin conductance response (SCR),

and air temperature, pressure, humidity, and VoC, respectively. We collected the data for a total duration of 45 min and segmented the physiological signals into portions of 30 s for processing and analysis. The sensors' data acquisition was performed at the sampling rates of 100, 1, and 0.1 Hz, respectively for PPG and EDA, ST, and environmental sensors.

We processed 135 segments of PPG and EDA signals. ST and environmental sensors were not subject to signal processing. The environmental data of subject 2 during 15-minute experiments shows boundary changes of 26–29 °C, 971.8–971.65 hPa, and 74–63 RH%. At the same time the ST swings between 28 and 31 °C (see Fig. 2).

ST output differs from a regular body temperature. However, considering both changes in the ST sensor and calibration with respect to the environmental sensor might re-correct the output. The close distance between the two sensors of ST and environmental indicates the influence of fingertip heat on the environmental sensor.



**Fig. 2.** Example of skin temperature and all environmental parameters.

**Table 1.** The MAE of HR with respect to the reference sensor indicates the max and min errors 3,9 and 1,61, respectively.

| Subject   | 1   | 2    | 3    |
|-----------|-----|------|------|
| MAE (bpm) | 3,9 | 3,08 | 1,61 |

The PPG signal was analyzed using the Python library HeartPy. We analyzed the HR of the signal in 30 s intervals and calculated the difference to the signal of the reference sensor for every interval. The mean average error (MAE) of subjects is 2,83 bpm (see Table 1).

We applied the same analysis to the SpO2 measurements, calculated from the PPG Signal. The average MAE of all subjects is 3,52% (see Table 2).

**Table 2.** The estimated SpO2 from the module shows the reliability and consistency of the measurement.

| Subject | 1 | 2 | 3 |
|---------|------|------|------|
| MAE (%) | 3,59 | 3,17 | 3,79 |

**Table 3.** Add caption.

| | $EDR_{lat}$ | $EDR_{amp}$ | $EDR_{rist}$ | $t_{max}$ | $RT_{(50\%)}$ | $RT_{(63\%)}$ |
|---|---|---|---|---|---|---|
| EDA (module) | 30,5 | 1,2 | 30,5 | 60,9 | 25,7 | 38,1 |
| EDA (reference) | 21 | 28,7 | 27,9 | 49 | 10,5 | 12,9 |

The EDA signal from the designed module and the reference sensor is depicted in Fig. 3. Even though the signal requires further analysis, however, the pattern of the reference signal is followed by the module, appropriately. We analyzed some shorter intervals of the signal depicted in Fig. 3 with the Python library from biosignals, in significant terms (see Fig. 3).



**Fig. 3.** EDA signal visualization for the sad video.

## 4   Discussion

The proposed module addresses the measurement of four domains to facilitate the mutual interaction measurement in the later stage for more complex diseases. Currently, the prototype measures HR, SpO2, and ST (i.e., physiological domain), EDA (i.e., emotional state recognition: psychological domain), air temperature, pressure, humidity, and VoC (i.e., environment domain), and raw data of inertial measurement unit (i.e., behavioral domain). In recent years, with increasing the attention to turning private spaces into diagnostic spaces and

implementing the dynamic point of perception as a fundamental step in contributing to a smart healthy city, the module can contribute to both approaches. Due to the compact form factor and low cost, it can be within the interest of users to attach it to frequently touched objects such as on the steering wheel of a car or bicycle with which one travels distances and needs to hold the steering. This could be considered an unobtrusive method that does not impose any additional workload on the users during their daily routine. Besides, the application could be extended to objects at home such as a door handle or light switches (Table 3).

This work presented is a validation of the limited number of healthy subjects. In future studies, it is required to recruit more subjects with particular applications.

## 5 Conclusions

We designed and implemented a compact form factor, low-cost, and multi-signal acquisition prototype with applications in HRQoL measurement. It measures and acquires the physiological signals from the fingertip. We could measure HR with as small an error of 2,83. The stability and low error of SpO2 indicated its reliability and applicability. In future applications, the respiration rate also could be modeled and estimated based on the PPG signal. The prototype meets the requirements of an efficient module to be deployed on the steering wheel of cars to measure stress. As an IoT module it can also play its role in developing a dynamic point of perception in smart healthy city.

## References

1. WHOQOL—Measuring Quality of Life— The World Health Organization. https://www.who.int/tools/whoqol. Accessed 07 Mar 2023
2. WHO statistics. https://www.who.int/news-room/fact-sheets/detail/the-top-10-causes-of-death. Accessed 07 Mar 2023
3. Haghi M, Spicher N, Wang J, Deserno TM (2022) Integrated sensing devices for disease prevention and health alerts in smart homes. In: Accident and emergency informatics. IOS Press, pp 39–61
4. Haghi M, Danyali S, Ayasseh S, Wang J, Aazami R, Deserno TM (2021) Wearable devices in health monitoring from the environmental towards multiple domains: a survey. Sensors 21(6):2130 Mar 18
5. Kyriacou PA (2022) Introduction to photoplethysmography. In: Photoplethysmography. Academic Press, pp 1–16
6. Boucsein W, Boucsein W (2012) Principles of electrodermal phenomena. Electrod Activity 1–86
7. Jarchi D, Salvi D, Tarassenko L, Clifton DA (2018) Validation of instantaneous respiratory rate using reflectance PPG from different body positions. Sensors 18(11):3705

# Simultaneous Detection of pH, Antioxidant Capacity and Conductivity Through a Low-Cost Wireless Sensing Platform

Riccardo Goldoni[1,2(✉)], Andrea Ria[3], Daniela Galimberti[4,5], Paola Dongiovanni[6], Lucanos Strambini[2], and Gianluca Tartaglia[5,7]

[1] Department of Electronics, Information and Bioengineering (DEIB), Politecnico Di Milano, Milan, Italy
riccardo.goldoni@polimi.it

[2] CNR-Istituto Di Elettronica E Di Ingegneria Dell'Informazione E Delle Telecomunicazioni, Torino, Italy

[3] Department of Information Engineering, University of Pisa, Pisa, Italy

[4] Department of Biomedical, Surgical and Dental Sciences, University of Milan, Milan, Italy

[5] Neurology-Neurodegenerative Diseases, Fondazione IRCCS Cà Granda, Ospedale Maggiore Policlinico, Milan, Italy

[6] Medicine and Metabolic Diseases, Fondazione IRCCS Cà Granda, Ospedale Maggiore Policlinico, Milan, Italy

[7] UOC Maxillo-Facial Surgery and Dentistry, Fondazione IRCCS Cà Granda, Ospedale Maggiore Policlinico, Milan, Italy

**Abstract.** Portable diagnostic tools enable the next generation of in-home sensing and biosensing devices for screening of pathologies and monitoring of physiological parameters. In this context, the reduction in cost of off-the-shelf electronics components is rapidly advancing the development of affordable and easy-to-use portable platforms. In this work we describe the development of a multisensing platform capable of simultaneously or concurrently performing three electrochemical measurements that can provide information in real time about the levels of important physiological markers, namely pH, antioxidant capacity and conductivity of biofluids. The developed device can be wirelessly powered via an inductive charging module and an onboard battery and transmits data via Bluetooth to a mobile device for rapid and simple data sharing and processing. By integrating multiple types of electrochemical sensors into an inexpensive, pocket-sized diagnostic tool, we demonstrate a low-cost approach to point-of-care health monitoring using open-source hardware that can be readily deployed in resource-limited settings. These types of inexpensive sensing modalities have the potential to provide actionable health information as initial screening markers for the onset of several systemic diseases.

**Keywords:** Point-of-care devices · Electrochemical sensors · pH · Antioxidant capacity · Conductivity · Salivary diagnostics

# 1   Introduction

The convergence of inexpensive printed circuit boards, miniaturized electronics, and electrochemical sensing technologies has enabled the rapid development of miniaturized and portable diagnostic devices [1]. Electrochemical sensors provide a fast, affordable, and accurate way to detect several biomarkers in various body fluids, including saliva [2–6], opening up opportunities for point-of-care diagnostics and in-home physiological monitoring. However, for these devices to benefit in-home healthcare, they must be simple to operate, affordable, and able to be deployed in resource-limited settings. One way to achieve that would be to have an onboard battery for field use that could be wirelessly recharged for a better user experience and the reduction of unnecessary wiring. Additionally, wireless data transmission to mobile devices should allow data sharing and remote processing without dedicated cables or equipment. In this work, we demonstrate the feasibility of rapidly prototyping diagnostic devices using off-the-shelf (OTS) components. Our customized platform is based on an ARM microcontroller with an integrated analog front end optimized for electrochemical sensing. This enables real-time monitoring of three sensors, two of which are custom electrochemical sensors: an iridium oxide-based pH sensor and a cerium oxide nanoparticle-based antioxidant capacity sensor. These sensors measure key physiological parameters as a first step toward developing inexpensive and scalable tools for assessing oral health at a population level. Our device consists of a print-ed circuit board populated with commercial electronic components as well as our custom-built electrochemical sensors. Data from the three sensors are transmitted digitally via Bluetooth to a mobile device for collection, visualization, and sharing. By leveraging such an "open hardware" development paradigm based on commercial off-the-shelf components, we aim to demonstrate a rapid and low-cost pathway for prototyping future diagnostic platforms.

# 2   Proposed System

The structure of the proposed system is reported in Fig. 1. The building blocks of the developed platform are the following:

- Sensor elements, including a pH sensor and an antioxidant capacity sensor
- Electronic platform, featuring the ADuCM355 as the integrated IC which includes an ARM M3 chip and the AD5940 dedicated analog front-end, optimized for chemical sensing. The embedded ADC is a 16-bit SAR ADC, which can work in either low-power (800 kSPS) or high-power mode (1.6MSPS). The electronics board features Bluetooth communication and wireless powering via inductive charging.
- Mobile device, which is used to operate the platform as well as to read and process the output data.

## 2.1   Electrochemical Sensors

The reagents used in the study were analytical-grade chemicals purchased from Sigma-Aldrich (St. Louis, MO, USA). The chemical reagents used included potassium hexacyanoferrate (III) ($K_3[Fe(CN)_6]$), potassium hexacyanoferrate (II) ($K_4[Fe(CN)_6]$),

**Fig. 1.** Building blocks of the wireless sensing system for at-home point-of-care diagnostics

sodium chloride (NaCl), potassium chloride (KCl), sodium phosphate dibasic ($Na_2HPO_4$), potassium phosphate monobasic ($KH_2PO_4$). Cerium (IV) oxide NP (20 wt% colloidal dispersion in acetic acid 2.5 wt%, d = 30–60 nm), which was acquired from Merck (Darmstadt, Germany). Phosphate-buffered saline (PBS) solutions were prepared by dissolving 8 g of NaCl, 0.2 g of KCl, 1.44 g of $Na_2HPO_4$, and 0.245 g of $KH_2PO_4$ in 1 L of ultrapure deionized water (DI). The iridium oxide solutions for electrodeposition were obtained from 4.5 mM iridium tetrachloride, 0.3% hydrogen peroxide and 40 mM oxalic acid dehydrate. The solutions are then left to stabilize at room temperature for one day. A set of screen-printed carbon electrodes (SPCE) of model DRP-11L were purchased from Dropsens and consisted of a working carbon electrode, counter carbon electrode and Ag/AgCl reference electrode. A second set of screen-printed electrodes (SPE) of model 250AT were purchased from Dropsens and consisted of a gold working electrode, a platinum counter electrode and an Ag/AgCl reference electrode.

## 2.2 Electronic Platform

The readout electronic module was assembled from available OTS components. The main board that was employed is the EVAL-ADuCM355QPZ, purchased from Analog Devices. A Bluetooth module HC-06 (Guangzhou HC Information Technology Co., Ltd.) and a 3.7 V 1800 mAh LiPo battery (EEMB LP953450) were purchased from DigiKey. An inductive power transfer charging kit 3.3V 500 mA and a USB-C MicroLipo 3.7 V/4.2 V 100 mA were purchased from Adafruit Industries LLC.

## 3 Experimental Results

The board was assembled with discrete components as listed in the description of the proposed system. The assembly is shown in Fig. 2, which demonstrates the modularity of the setup, where different sensors can be added and removed depending on the desired application. The form factor, while it could be further reduced integrating the different components into a custom portable circuit board (PCB) is still limited, which allows for portability and field deployment. Three buttons on the board allow for the possibility to run each test—pH, antioxidant capacity and conductivity—independently, if needed, for maximum flexibility of the system. The electrochemical sensors constitute the disposable element of the sensing system, as they are only intended to be used for a single measurement.

**Fig. 2.** Board experimental setup that includes the ADuCM355 IC, LiPo battery with microLiPo charging module for wireless power transfer via inductive charging kit, HC-06 Bluetooth module and multiple sensors.

The electrochemical sensors used in the study were developed following different approaches. The nanoceria based sensor followed a drop-casting modification, which has been reported in literature [7–9] where a diluted dispersion of cerium oxide nanoparticles was deposited on the surface of an activated SPCE. The activation procedure consisted of the drop casting of 50 µL of a PBS solution followed by the application of a constant potential of 1.4 V for 300 s.

The iridium oxide-based sensor instead followed an electrodeposition protocol. An iridium oxide solution was prepared following the protocol first described by Yamanaka [10]. Different voltage windows and scan rates for the electrodeposition were evaluated, following a Design of Experiments (DoE) approach, in order to select the optimal parameters for the final electrodeposition protocol. The iridium oxide solution exhibited a storage time of at least one month, which was confirmed by UV-Vis spectroscopy.

Since the conductivity sensor has not been developed yet, a dummy cell with known resistance has been used instead to test the developed system in a simulated condition of known conductivity of the system.

Each developed sensor was tested in replicates (n = 3) using both a laboratory grade potentiostat (PalmSens 4, PalmSens) and the newly developed electronic readout device, allowing us to compare the performance of the two instruments for analytical applications. A solution of a redox probe, potassium ferro/ferricyanide, was used to characterize the SPCE electrode modified with nanoceria. The pH sensor, based on iridium oxide modified screen printed gold electrode, was in turn characterized using standard commercial pH solutions of pH 4, 7 and 10. The conductivity dummy cell has been characterized by running electrochemical impedance spectroscopy and extracting the resistance value from the obtained data. The results are reported in Fig. 3.

The results from the characterization tests show how the performance exhibited by the developed portable platform almost matches the performance of a laboratory grade potentiostat in all of the reported tests.

**Fig. 3.** Analytical performance of developed electrochemical sensors and conductivity dummy cell. **A** Characterization of nanoceria-modified SPCE via cyclic voltammetry (CV) over a range of different scan rates, from 25 mV/s to 200 mV/s. **B** Characterization of dummy cell used for conductivity tests simulation. **C** Characterization of iridium oxide-modified gold SPE via open circuit potentiometry (OCP) using standard commercial solutions of different pH values, 4, 7 and 10 respectively.

# 4   Conclusion

This study represents an initial step toward the development of a miniaturized, wireless diagnostic device for point-of-care testing. The present work has resulted in a first prototype of a portable system capable of simultaneous potentiometric, voltammetric and amperometric measurements, powered wirelessly for in-home use. Devices like the one reported here have the potential to enable continuous, in-home monitoring of key physiological parameters. The ability to non-invasively track biomarkers over time could allow earlier disease detection, improved treatment efficacy monitoring, and reduced hospital visits. Miniaturized, user-friendly devices are key to making this type of frequent, in-home testing feasible. The diagnostic tool from this study represents an important step toward that goal.

# References

1. Ria A, Cicalini M, Manfredini G, Catania A, Piotto M, Bruschi P (2022) The SENSIPLUS: a single-chip fully programmable sensor interface. In: Saponara S, De Gloria A (eds) Applications in electronics pervading industry, environment and society. Springer International Publishing, Cham, pp 256–261
2. Goldoni R, Farronato M, Connelly ST, Tartaglia GM, Yeo W-H (2021) Recent advances in graphene-based nanobiosensors for salivary biomarker detection. Biosens Bioelectron 171:112723
3. Goldoni R et al (2022) Salivary biomarkers of neurodegenerative and demyelinating diseases and biosensors for their detection. Ageing Res Rev 76:101587
4. Goldoni R et al (2021) Malignancies and biosensors: a focus on oral cancer detection through salivary biomarkers. Biosensors 11(10):396
5. Poli PP et al (2023) Detection and sensing of oral xenobiotics in edentulous patients rehabilitated with titanium dental implants: Insights from a scoping review. J Prosthetic Dentistry
6. Dongiovanni P et al (2023) Salivary biomarkers: novel noninvasive tools to diagnose chronic inflammation. Int J Oral Sci 15(1):27
7. Tortolini C, Bollella P, Zumpano R, Favero G, Mazzei F, Antiochia R (2018) Metal oxide nanoparticle based electrochemical sensor for total antioxidant capacity (TAC) detection in wine samples. Biosensors 8(4):108
8. Andrei V, Sharpe E, Vasilescu A, Andreescu S (2016) A single use electrochemical sensor based on biomimetic nanoceria for the detection of wine antioxidants. Talanta 156:112–118
9. Goldoni R et al (2023) Quality-by-Design R&D of a novel nanozyme-based sensor for Saliva antioxidant capacity evaluation. Antioxidants 12(5):1120
10. Yamanaka K (1989) Anodically electrodeposited iridium oxide films (AEIROF) from alkaline solutions for electrochromic display devices. Jpn J Appl Phys 28(4R):632

# On the Assessment of Gray Code Kernels for Motion Characterization in People with Multiple Sclerosis: A Preliminary Study

Matteo Moro[1,2,3(✉)], Maria Cellerino[4,5], Matilde Inglese[4,5], Maura Casadio[1,3], Francesca Odone[1,2,3], and Nicoletta Noceti[1,2,3]

[1] DIBRIS, University of Genova, Genova - IT, Italy
matteo.moro@unige.it
[2] MaLGa-Machine Learning Genoa Center, Genova - IT, Italy
[3] RAISE-Robotics and AI for Socio-economic Empowerment Ecosystem, Genova - IT, Italy
[4] DINOGMI, University of Genova, Genova - IT, Italy
[5] IRCCS Ospedale Policlinico San Martino, Genova - IT, Italy

**Abstract.** Motor deficits in the lower limbs are common in people with multiple sclerosis (MS), impacting mobility and quality of life. Objective and quantitative metrics are crucial for effective identification and monitoring of motor deficits. Recent advancements in computer vision and human pose estimators allow for automatic extraction of movement information from video data, offering potential insights into human motion patterns. This exploratory study investigates the use of Gray-Code Kernels (GCKs) in characterizing gait patterns in individuals with advanced-stage MS compared to age- and sex-matched unimpaired controls. The preliminary results obtained demonstrate the promising potential of combining GCKs and pose estimators in characterizing gait patterns, warranting further investigation in this area.

**Keywords:** Human motion analysis · Gray-code kernels · Multiple sclerosis · Human pose estimation

## 1 Introduction

Motor deficits in the lower limbs are common symptoms experienced by individuals with multiple sclerosis (MS), impacting their ability to walk and reducing overall mobility, balance, and quality of life [6]. MS is a potentially disabling disease of the central nervous system, whose symptoms vary greatly among people and depend on the location and severity of nerve fiber damage [7].

To monitor motor deficits effectively, objective and quantitative metrics are crucial [15]. Such metrics can enable early and accurate assessments, personalized medical decisions, and tailored rehabilitation treatments. In recent decades, significant progress has been made in this research field, with state-of-the-art technology offering useful and precise quantitative measures based on markers or wearable systems [4]. However, these technologies may interfere with natural movements, cause discomfort, and require specialized personnel [3].

Thanks to the advancements in the computer vision domain, RGB cameras recently emerged as a suitable alternative to marker-based approaches, also thanks to the recent development of efficient and accurate algorithm for human pose estimation. The goal of human pose estimation is to determine the spatial positions and relationships of body keypoints in images or videos. This process involves inferring the articulated human body's pose, comprising joints and rigid parts, through image-based observations. Pose estimation can be performed in 3D or 2D, and the literature offers various approaches to address this problem [5,10,16]. Human pose estimators have been successfully applied to the problem of gait analysis in motor control and rehabilitation [8,9], and combined with data-driven techniques they provide a valuable tool for quantitative human motion analysis in natural environments, such as home-like settings [14]. Pose estimators provide a compact representation of pose, that can be used for a coarse analysis of human movements in a video analysis setting. In this respect, given the sparsity of the representation, it is beneficial to combine body keypoints with denser representations, able to capture the dynamic evolution of pixels over time and provide an intuitive visual representation where salient parts in the image are highlighted. This can be seen as a sort of attention mechanism, that can guide the following steps in the analysis. Unfortunately, also RGB cameras present drawbacks, e.g., privacy, legal and ethical considerations and data storage. However, the effect of these drawbacks can be reduced with appropriate regulations.

Inspired by these considerations, in this work, we present an exploratory study to assess the potential of Grey-Code Kernels (GCKs) [11,12] in characterizing gait patterns in people with MS. GCKs, proposed by [1], are a family of filters that can be used to efficiently project the content of images and videos. They provide a way to represent motion information in a form that is easily interpretable and computationally efficient. GCKs have been already adopted in the field of human action recognition [13] and they seem particularly suitable to characterize simple tasks.

The combination of GCKs and pose estimation allows for the automatic extraction of movement information from video data, offering richer insights about human motion patterns without interfering with natural movements or causing discomfort to the subjects. A global dynamic representation for each frame is computed by pooling the features obtained by the projection with the whole family of filters. In the pooling maps relevant spatio-temporal motion cues are highlighted. For our experiments, we acquired and analysed videos of five people with advanced-stage MS and five age- and sex-matched unimpaired

controls. We considered people with an Expanded Disability Status Scale (EDSS) equal to or greater than 4.

Our preliminary analysis shows promising results in easily characterizing motion patterns. In particular, we observe that statistics on the pooling maps reflect impaired mobility and restricted movements in the leg of people with MS. In contrast, unimpaired controls demonstrated a relatively more consistent motion distribution with higher values in the maximum pooling, indicating more balanced and regular gait patterns. Such results set the basis for further investigations.

## 2   Materials and Methods

### 2.1   Dataset

We acquired videos (resolution $1280 \times 720$, at 30 fps) of the lateral gait of 5 MS people (5 females, mean age  standard deviation: 47  6, mean disease duration standard deviation: 12.5  10 years, EDSS mean [minimum, maximum]: 5.3 [4, 6]). The EDSS reflects the level of disability in MS, ranging from 0 (normal neurological exam) to 10. We used only one computer webcam to minimize the invasiveness of the setup and to explore the possibility of extracting quantitative information with cheap tools widely available. In the same way, we acquired also videos of the gait of 5 unimpaired controls (age and sex-matched with MS people).

The task consisted of three different walking conditions: (i) normal walking, (ii) toe walking, and (iii) tandem walking. For each walking condition, participants were asked to walk following straight paths, from one side of a room to the opposite, for 6 trials. The path was 4 m long.

The study conforms to the standard of the declaration of Helsinki and was approved by the local ethical committee (CER Liguria, No. 222REG2017). All participants provided written informed consent prior to participation in the study.

### 2.2   Our Pipeline

In this subsection we highlight the main steps we perform to analyze the videos (summarized in Fig. 1) and assess the potential of GCKs representations to characterise gait patterns.

- **_Pose estimation_**. We employ Openpose [2], a well known architecture adopted in other studies [14] in the analysis and characterization of human motion. Openpose provides accurate and real-time human pose estimation by detecting and localizing multiple keypoints. Among them, we focus on the knee and on the foot (those more meaningful to describe the gait) of the leg visible from the lateral acquisition view-point.

– **Gray-code kernels computation**. To efficiently represent motion features in the video frames, we generated GCKs following the implementation described in [12] to compute the video projection and build the maximum and average pooling maps.

– **Metrics extraction**. In the maps representing the maximum and the average pooling of the GCKs, we focus on a neighborhood of the knee and of the foot locations detected with Openpose. In such neighborhoods, we computed statistics (i.e., mean, maximum, minimum, and median) of the values of the maps. The choice of these metrics allow us to capture different aspects of the motion behavior, such as average motion intensity, extreme motion events, and typical motion characteristics. Furthermore, it allows us to compare the gait behaviour of the two population (people with MS and unimpaired controls).



**Fig. 1.** Summary of the main steps of the video analysis

## 3    Experimental Results

Our experimental results reveal that the maximum of the maximum pooling of Gray-code kernels (GCKs) plays a crucial role in effectively highlighting gait patterns in individuals with MS compared to unimpaired controls. To define the neighborhood of the GCKs maps to focus on, we select a number of pixels that allow us to consider the whole anatomical points: with the resolution of the camera as reported in the methodological section, we empirically select a square of 40 pixels side. By utilizing this specific statistical metrics, we are able to discern differences in the motion patterns of the knee and foot between the two groups and among all the three different gait tasks. In particular, we report the most meaningful result: Fig. 2 reports the histograms of the maximum values extracted from the maximum pooling of each map for each person in a neighborhood of the foot.

**Fig. 2.** Histograms of the maximum values over time of the maximum pooling of the GCKs in a neighborhood of the foot. *Left*: behaviour of people with MS. *Right*: behaviour of unimpaired controls. The three rows represent the three different gait tasks: normal, toe, and tandem walk. The x axis represent the values extracted from the GCKs maps (normalized between 0 and 1). The y axis represent the probability associated with each value

In the case of maximum pooling, the maximum values within the neighborhood around each keypoint indicate the presence of prominent and intense motion events. We observe that individuals with MS exhibit less repetitive and more distributed values (see Fig. 2). On the other hand, unimpaired controls demonstrate more consistent and smoother motion patterns. Regarding the differences across walking condition, both in people with MS and unimpaired controls, it is possible to notice that lower maximum values are associated with tandem walk.

## 4    Discussion and Conclusion

The observed differences in gait patterns between individuals with MS and unimpaired controls may be relevant in the clinical context. Motor deficits in the lower limbs are a common symptom experienced by individuals with MS, impacting their ability to walk and reducing overall mobility and balance. Our approach, based on the combination of human pose estimation and GCKs and subsequent computation of meaningful metrics, may be useful to provide a quantitative and objective assessment of these motor deficits.

The behaviour highlighted in Fig. 2 may reflect impaired mobility and restricted movements in the leg of people with MS. In contrast, unimpaired controls demonstrated a more consistent motion distribution with higher values in the maximum pooling, indicating more balanced and regular gait patterns.

The analysis we performed leave some open problems that we will address in future research. Firstly, the choice of the neighborhood size for computing GCKs metrics can significantly impact the results. While a larger neighborhood may capture more contextual information, it might introduce noise and reduce the

sensitivity to subtle motion changes. Conversely, a smaller neighborhood may miss relevant motion patterns. In this direction, consideration and further investigation on the pose estimator adopted are needed to reduce the effect of possible mispredicitons. Then, a deeper analysis on the maps obtained with different kernels is needed. Furthermore, our approach may be affected when applied to different datasets with varying characteristics, such as different camera angles, lighting conditions, or diverse populations. Ensuring robustness and generalizability across various scenarios remains a challenge and requires the acquisition of additional data. Lastly, a comprehensive comparison with other state-of-the-art methodologies (e.g., shoes-based gait analysis) would be necessary to assess the actual potential of the pipeline.

# References

1. Ben-Artzi G, Hel-Or H, Hel-Or Y (2007) The gray-code filter kernels. IEEE Trans PAMI 29(3):382–393
2. Cao Z, Simon T, Wei SE, Sheikh Y (2017) Realtime multi-person 2d pose estimation using part affinity fields. In: IEEE CVPR, pp 7291–7299
3. Carse B, Meadows B, Bowers R, Rowe P (2013) Affordable clinical gait analysis: an assessment of the marker tracking accuracy of a new low-cost optical 3d motion analysis system. Physiotherapy 99(4):347–351
4. Colyer SL, Evans M, Cosker DP, Salo AI (2018) A review of the evolution of vision-based motion analysis and the integration of advanced computer vision methods towards developing a markerless system. Sport Med-Open 4(1):1–15
5. Dang Q, Yin J, Wang B, Zheng W (2019) Deep learning based 2d human pose estimation: a survey. Tsinghua Sci Technol 24(6):663–676
6. Fritz NE, Marasigan RER, Calabresi PA, Newsome SD, Zackowski KM (2015) The impact of dynamic balance measures on walking performance in multiple sclerosis. Neurorehabilitation Neural Repair 29(1):62–69
7. Lublin FD, Reingold SC, Cohen JA et al (2014) Defining the clinical course of multiple sclerosis: the 2013 revisions. Neurology 83(3):278–286
8. Moro M, Marchesi G, Hesse F, Odone F, Casadio M (2022) Markerless versus marker-based gait analysis: a proof of concept study. Sensors 22(5) (2022)
9. Moro M, Marchesi G, Odone F, Casadio M (2020) Markerless gait analysis in stroke survivors based on computer vision and deep learning: a pilot study. In: ACM symposium on applied computing. pp. 2097–2104 (2020)
10. Munea, T.L., Jembre, Y.Z., Weldegebriel, H.T., Chen, L., Huang, C., Yang, C.: The progress of human pose estimation: A survey and taxonomy of models applied in 2d human pose estimation. IEEE Access 8 (2020)

11. Nicora E, Noceti N (2022) Exploring the use of efficient projection kernels for motion saliency estimation. In: ICIAP, pp 158–169
12. Nicora E, Noceti N (2022) On the use of efficient projection kernels for motion-based visual saliency estimation. Front Comput Sci 4
13. Nicora E, Pastore VP, Noceti N (2023) Gck-maps: A scene unbiased representation for efficient human action recognition. In: International conference on image analysis and processing. Springer, Berlin, pp 62–73
14. Setti F, Avogaro A, Cunico F, Rosenhahn B, Markerless human pose estimation for biomedical applications: a survey. Front Comput Sci 5
15. Williams KL, Brauer SG (2022) Walking impairment in patients with multiple sclerosis: The impact of complex motor and non-motor symptoms across the disability spectrum. Aust J Gen Pract 51(4):215–219
16. Xu Y, Zhang J, Zhang Q, Tao D (2022) Vitpose: Simple vision transformer baselines for human pose estimation. NeurIPS 35:38571–38584

# Short Contributions

# Crack Auscultation in Asphalt Pavements Using Computer Vision

Emanuel A. Cortez Médici[(✉)], Ricardo Petrino, and Ramón Cortez

Universidad Nacional de San Luis, LEIS, Departamento de electrónica, Av. Ejército de los Andes 950, D5700 San Luis, Argentina
{emanuelacortezm,rpetrinodlv}@gmail.com, arcortez@unsl.edu.ar

**Abstract.** Auscultation of cracks in asphalt pavements plays a fundamental role for ensuring transportation infrastructure maintenance and longevity. This paper presents a system that combines neural network algorithms to detect, classify, and segment pavement cracks in order to produce reports of the concentration of cracks in asphalt pavements. The proposed approach generates a choropleth map that allows road inspectors to quickly determine the state of the pavement and get more information on the type of cracks present on the pavement. By implementing this system, continuous auscultation of asphalt pavements becomes feasible, contributing to effective infrastructure management and maintenance practices.

**Keywords:** Pavement crack detection · Artificial vision · Convolutional neural networks

## 1 Introduction

Pavement crack auscultation is a critical task in ensuring the maintenance of transportation infrastructure, since cracks let the water penetrate from the surface to underlying layers reducing pavement life and worsening ridability [1]. Traditional methods often require specialized equipment and manual inspection, making them time-consuming and costly [2].

This paper presents a practical system that exploits the proven capabilities of convolutional neural networks to detect, classify and segment asphalt cracks [3–6]. By combining a *You Only Look Once* (YOLO) and a U-Net neural network into a single program, the system is capable of processing a recording of a pavement surface and a position logfile generated by a smartphone's GPS, and produce a report that can be used by road inspectors to determine the condition of the pavement.

Figure 1 shows the workflow of the obtained system. By mounting a smartphone on a vehicle, pavement conditions can be recorded with ease in real-time. The smartphone's camera captures the pavement surface, while the GPS sensor logs the vehicle's position at regular intervals.

Afterwards, the data is transferred to a computer for processing where the recorded video is analyzed frame by frame using the YOLO algorithm. Combining YOLO V5 detection and Norfair tracking [7], each detected crack is identified with a unique number, so the same crack is not detected twice. The bounding boxes are then used to extract the cracks from the whole image and, with a U-Net neural network [8], the crack is segmented, extracting as a result the shape of the crack to automatically calculate the thickness and longitude of the crack.

Finally, a database is created using the ID of each crack, storing the processed data, and a choropleth map is generated from it in order to graphically show the information of the pavement condition regarding cracks, and allowing the user to inspect each section of the pavement. The system was tested on two roads of the province of San Luis, Argentina.



**Fig. 1.** Workflow of the designed system

## 2   Pavement Recording

The mounting was designed to go on top of a truck bed at a fixed height that allows recording the entire lane width. For the auscultation of the pavement, two resolutions were used, one being of $1920 \times 1080$ and a higher resolution of $3840 \times 2160$, both at 60 frames per second. Approximately 125 min of video was recorded on a first run in order to generate the database needed to train the neural networks.

# 3   Crack Detection and Classification

In a first pass, the video is analyzed using the YOLO neural network, as this is more efficient than using the U-Net network to segment each frame. In this way, parts of the pavement where there are no cracks are not processed by the U-Net network. Consequently, YOLO detects and classifies three different types of cracks—longitudinal cracks, transverse cracks and alligator cracks. The U-Net network is only responsible for the segmentation of each detected crack.

For this purpose, a custom dataset was made using images obtained from recordings of different roads. The videos were decomposed into images, reaching a total of 1329 images. With the help of an annotating tool, *ImLabel*, approximately 1400 transversal cracks, 800 longitudinal cracks, and more than 500 alligator cracks were labeled. The training was conducted over 100 epochs and the neural network achieved a recognition rate of approximately 84%. The computer used for the training consisted on a AMD 5900X, a GPU Nvidia RTX 3080, 12GB VRAM and 32GB of DDR4 RAM memory. On average, the training took around 10 min per epoch using batches of 16 images.

In addition, it is necessary to identify each detected crack with a unique number and to retain this number until the crack is no longer present in the video. A SORT algorithm was used for this purpose, specifically the Norfair algorithm, as it offered better performance than other options such as KCF or CSRT. With this identification process, it is possible to analyze each crack individually and create a database to store all the data processed. Figure 2 shows how the system is able to assign an ID number and keep it as the crack moves through the recording.

**Fig. 2.** Tracking and identification of a detected crack.

# 4   Crack Segmentation

The thickness and the longitude of a crack are key aspects to take into account when evaluating the condition of the pavement. Therefore, it is necessary to extract from the bounding box the crack and calculate its dimensions. Figure 3 shows the result of processing an asphalt crack with a U-Net in order to completely extract the crack from the image and then calculate its proportions. For that purpose, the U-Net network was trained for 250 epochs using de Crack500 dataset, where it achieved a efficiency rate of 79%. The platform used for the training was the same as the one used for the training of the YOLO network.

The resulting image is just a matrix where a values equall to 1 represent the crack, while 0 represent the rest of the pavement. By using morphological operations and calculating the proportion of a image pixel in real life, it is possible to compute the thickness and longitude of each crack.



**Fig. 3.** Detected crack segmented by the U-Net network.

## 5   Map Generation

The last step is to display all the collected and processed data in a visual and understandable way. A choropleth map allows a quick overview of the concentration of cracks in the pavement. Figure 4 shows the resulting map and the Popup windows that opens when a section of the map is clicked. The window shows information of each crack in the zone as well as the segmented image.



**Fig. 4.** Resulting choropleth map [left] and Popup window [right]

## 6   Experimental Results

The system was used on two different roads located in the province of San Luis, Argentina, in order to test how the system performed against varying lightning conditions, speed of the vehicle and pavement textures.

### 6.1  *Ruta Provincial N.° 15*

This road was chosen due to its deteriorated state. The system was then evaluated by testing it on five sections of the road, varying the driving speed of the vehicle. Figure 5 shows the different maps obtained for each section of the road.



**Fig. 5.** Results obtained for the *Ruta Provincial N.° 15*.

In total, the pavement auscultation for the 4 sections took approximately 90 min. The processing time of the data was greatly influenced by the resolution of the analised video. With FullHD videos, the YOLO network was able to analise 24 frames per second, while with 4 K videos, that number decreased to 10 frames per second. Thus, taking into account videos recorded at 60 fps, the detection and classification process took 2.5 times longer than the video's duration for FullHD videos and 6 times longer for 4 K videos. Regarding the segmentation process, it was influenced by the image resolution but also by the size of the extracted bounding box to analise. However, it can be said that on average the inference time for each image was around 1 s. The average vehicle speed that allowed the pavement to be scanned with sufficient quality to properly detect, and segment cracks was 12 km/h. Although speeds higher than 15 km/h allowed cracks to be properly detected, the motion blur caused by the vibrations of the vehicle hampered the segmentation performed by the U-Net.

In terms of resolutions, the one that yielded the most favorable results was 1920 × 1080. This was primarily due to the considerably faster processing speed and the significant impact of motion blur on high-resolution images."

### 6.2  *Ruta Nacional N.° 146*

Figure 6 shows the obtained map of the road auscultation for the *Ruta Nacional N.° 146*. The map shows a low concentration of cracks compared to the maps obtained for the previous test. This is consistent with reality, as this pavement is in better conditions.

Although the neural network responsible for detection was able to correctly detect most of the cracks present in the pavement, the U-net was not able to successfully segment the detected cracks. Since the detection and classification

did not fail, the system was still able to generate a choropleth map showing the condition of the pavement and the concentration of cracks along the road.



**Fig. 6.** Map obtained for *Ruta Nacional N.° 146*.

Therefore, by using two different neural networks for different specific tasks, it is possible to obtain a more flexible system that is easier to improve. One implemented solution that drastically improved the performance of the system was the addition of a better motion dampener to the support on which the camera is placed. This reduced the motion blur in the video and allowed for better segmentation. In order to solve the segmentation problem, another possibility is to include more training data containing thinner cracks, images with motion blur or with shadows from nearby objects.

## 7   Conclusion

In conclusion, this paper presents an innovative pavement crack auscultation system that harnesses the power of smartphones and neural networks. By using a smartphone for data collection and processing the recorded information on a computer, the system provides a cost-effective and efficient solution for pavement maintenance. Subsequent analysis using neural networks ensures accurate crack detection and segmentation, enabling informed decision-making for infrastructure management.

# References

1. Rahbar-Rastegar R (2017) Cracking in asphalt pavements: impact of component properties and aging on fatigue and thermal cracking
2. Batac G, Ray M (1988) French strategy for preventive road maintenance: why and how?
3. Hu G, Hu B, Yang Z, Huang L, Li P (2021) Pavement crack detection method based on deep learning models. Wirel Commun Mob Comput 2021:1–13
4. Lau S, Chong E, Yang X, Wang X (2020) Automated pavement crack segmentation using U-net-based convolutional neural network. IEEE Access 8:114892–114899
5. Di Benedetto A, Fiani M, Gujski L (2023) U-Net-based CNN architecture for road crack segmentation. Infrastructures 8:90
6. Kheradmandi N, Mehranfar V (2022) A critical review and comparative study on image segmentation-based techniques for pavement crack detection. Constr Build Mater 321:126162
7. Alori J, Descoins A, Javier LF, KotaYuhara FD, Moises tryolabs/norfair: v2.2.0. (Zenodo, 2023, 1). https://doi.org/10.5281/zenodo.7504727
8. Ronneberger O, Fischer P, Brox T (2015) U-NET: convolutional networks for biomedical image segmentation

# The SensiTag: An Innovative BAP RFID TAG for Environmental Multi-sensing

Andrea Ria$^{(\boxtimes)}$, Andrea Motroni, Francesco Gagliardi, Massimo Piotto, and Paolo Bruschi

Department of Information Engineering, University of Pisa, Pisa, Italy
`andrea.ria@ing.unipi.it`

**Abstract.** This work presents the SensiTag, a novel Battery Assisted RFID Tag with multi-sensing capabilities enabled by a recently introduced versatile integrated sensor interface (the Sensiplus). The suggested tag enables a variety of sensing capabilities by using both the built-in sensors within the Sensiplus and external commercial sensors that can be connected through the interface. This study demonstrates, for the first time, the Sensiplus platform ability to conduct temperature measurements and communication using commercial RFID systems at the UHF band. The SensiTag performance was assessed through experimental testing on the prototype.

**Keywords:** RFID · Sensing · CMOS · Sensor interface

## 1 Introduction

In recent years, the Internet of Things (IoT) has experienced significant advancements. One of the emerging goals is to establish a network of interconnected sensors capable of collecting and exchanging environmental information [1–3]. By replacing wired connections with wireless interfaces, installation and maintenance costs are substantially reduced, and network flexibility is enhanced. RFID (Radio Frequency IDentification) technology, specifically passive UHF (Ultra High Frequency) RFID, has played a crucial role in developing passive labels that identify objects and store extensive information [4, 5] also for industrial environments [6]. These labels can be remotely read by an external reader connected to an antenna. However, the passive nature of these devices poses limitations when they are required to perform environmental monitoring tasks in addition to their regular functions [7, 8]. In fact, the addition of sensors increases the tag power consumption, often requiring the use of batteries. Although incorporating batteries introduces drawbacks, it enables the creation of wireless sensor nodes with advanced intelligence and a wide range of functionalities applicable in various domains. Modern RFID chips are designed to incorporate internal sensors and offer interfaces that facilitate the integration of external sensors and microcontrollers [9]. The RFID reader, operating under the constraints of the EPC UHF Gen2 Air Interface Protocol [10], can be programmed to interact with the tag memory of the RFID chip. This enables

the reader to send requests to read sensor data and receive the corresponding responses. Some commercial readers even allow the RFID chip to operate as a bridge, acting as a remote master node for SPI communications. However, in some instances, a separate microcontroller connected to the tag is necessary to serve as the master node.

In this work, we propose a Battery Assisted Passive UHF Tag based on an innovative low-power/low-voltage general purpose integrated sensor interface called Sensiplus (SP-IC) [11]. In this way, the Sensiplus ability to interface with a wide range of commercial sensors and to take advantage of sensors integrated within the chip itself, is advantageously combined with the UHF RFID method for communication. The SP-IC has been developed by SENSICHIPS s.r.l. [12] in collaboration with the University of Pisa. The SP-IC has been already employed in different sensing applications [13, 14] but it is the first time where it is experimentally tested in a RFID scenario.

## 2   Proposed System

A simplified block diagram of the proposed SensiTag is shown in Fig. 1, enclosed in the dashed line. It consists of three main blocks: a BAP RFID TAG which embeds the RFID EM4325 chip [15], the Arduino Nano V3 platform and the SP-IC. The RF communication is ensured by means of an Impinj RFID R420 reader connected to a A4030L Times-7 linearly polarized antenna radiating at a frequency of 865.7 MHz, placed at a distance of 50 cm from the SensiTag. The communication between the BAP RFID TAG and the SP-IC is managed by Arduino via SPI protocol with dedicated chip-select pins: $CS_{RF}$ and $CS_{SP}$. A 3.7 V Lithium-Ion cell is connected to the BAP which includes a 3.3 V LDO in order to deliver the supply voltage to all the circuits. A commercial Negative Temperature Coefficient (NTC) resistor (Panasonic ERT-JZEV104F) is connected to the SP-IC. At the same time, in this work, the SP-IC internal temperature sensor is used.



**Fig. 1.** Simplified block diagram of the proposed system (SensiTag).

### 2.1   BAP RFID Tag

The BAP RFID TAG is a custom FR4 PCB based on the EM4325 RFID chip, a popular integrated circuit for general purpose communication with passive or battery assisted UHF RFID tag. A dipole antenna has been designed considering the standard FR4 PCB fabrication parameters (i.e., the substrate thickness and the relative dielectric constant).

An L-shaped LC matching network was necessary to obtain the maximum power transfer from the antenna to the chip. The EM4325 is programmed for BAP mode and SPI-slave configuration.

## 2.2 The Sensiplus (SP-IC)

The SP-IC, whose complete functionalities and characteristics are described in [11], is a general purpose, single chip, mixed-signal integrated sensor interface, designed and fabricated with the standard UMC 180 nm 3.3 V CMOS process. It is capable of interfacing with up to 16 different sensors, by stimulating them with programmable ac and/or dc current and voltage stimulus. The corresponding response (i.e., a voltage or a current signals, respectively) is processed by a versatile readout channel. In particular, the SP-IC can perform Electrochemical Impedance Spectroscopy (EIS) with high accuracy, enabling the reading of resistive and/or capacitive sensors. All the SP-IC functionalities are handled by the integrated digital finite state machine, which also manages digital communication with external devices through standard SPI or I2C protocols. A proprietary single-wire digital communication protocol is also available. As far as the integrated SP-IC sensors are concerned, they consist in a bandgap-based temperature sensor, a light and UV sensor and a humidity sensor.

## 2.3 The Arduino Nano Platform

In this application, the Arduino Nano V3 platform is configured as SPI-master in order to communicate with the BAP RFID TAG and the SP-IC. In particular, the microcontroller (MCU) embedded in the Arduino platform accesses the EM4325 user memory locations where the RFID reader writes the op-code of the corresponding requested measurement. In this paper, as simple demonstration examples, we implemented two different measurements: (i) a temperature measurement by means of the SP-IC internal temperature sensor, (ii) a temperature measured with a commercial NTC resistor connected to the SP-IC port. Clearly, different types of measurements can also be performed, using different kinds of sensors connected to the SP-IC. Once the desired measurement has been selected, the MCU initiates the sequence of instructions to be transmitted to the SP-IC. The MCU receives the measurement results from the SP-IC and stores them in a dedicated EM4325 user memory location, remotely readable by the RFID reader. Figure 2 shows the flowchart of the software implemented in Arduino to manage of the measuring process.



**Fig. 2.** Flowchart of the software implemented onto Arduino Nano platform.

## 3   Experimental Results

A photograph of the proposed system is shown in Fig. 3, where the main components are indicated. The system size is 15 cm × 5 cm. It is worth noting that the SP-IC is encapsulated within a 44-pins JLCC package allowing the user to access to all the pins. A reduction of the SensiTag size can be accomplished by designing an ad-hoc PCB with both the SP-IC and the MCU directly soldered on it. However, possible size reduction is limited by the dipole antenna length, designed to work at a frequency of 865.7 MHz on a standard FR4 PCB substrate.



**Fig. 3.**  Photograph of the proposed system with the main components indicated.

In the following, we present the results of experimental tests concerning temperature measurements by means of both the SP-IC internal temperature sensor and the mentioned commercial NTC resistor connected to the SP-IC. For both the experimental tests, the temperature was imposed and controlled by a Peltier cryostat. The temperature was swept in the range 10–85 °C. For each temperature value set by the cryostat, the RFID reader required one measurement with the internal sensor of the SP-IC and one with the NTC resistor. For each applied temperature, 10 samples were recorded.

Figure 4 shows the value of the measured resistance as a function of the applied temperature, compared with the characteristic reported in the datasheet of the component. A good agreement is clearly visible, showing the SP-IC capability of interfacing resistive sensors while being embedded in a RFID platform.

Figure 5 shows the temperature measured with the SP-IC internal sensor (red squared dots) and the temperature derived from measurements reported in Fig. 4 (blue triangled dots), compared with the ideal case (black line). For the SP-IC internal sensor, a linear fit (red dashed lines) gives the following relationship:

$$T_{meas} = 1.0127 \cdot T_{applied} - 0.466 \tag{1}$$

The small slope and offset errors affecting the measured temperature, which can be noticed in (1), suggest that the temperature can be measured by the SP-IC with sufficient accuracy for many applications. As far as the temperature measured with the NTC resistor is concerned, the linear fitting (blue dashed line) highlighted a temperature offset value close to 5 °C. This error may be due to the fact that the two employed sensors are in two physically different positions and the reheating process is not uniform on the PCB.



**Fig. 4.** SP-IC measured NTC resistance as function of the cryostat applied temperature.



**Fig. 5.** Measured temperature with the SP-IC temperature sensor (red) and with the external NTC resistor (blue)

## 4   Conclusions

In this work, the innovative BAP RFID SensiTag has been described. For the first time, the Sensiplus interface has been used in a battery-assisted UHF RFID tag to perform temperature measurements, demonstrating its versatility. The employment of the Sensiplus combined with a RFID interface paves the way for a promising development of lightweight wireless sensor nodes to be deployed in a variety of environments.

# References

1. Al-Fuqaha A, Guizani M, Mohammadi M, Aledhari M, Ayyash M (2015) Internet of Things: a survey on enabling technologies, protocols, and applications. IEEE Commun Surv Tutor 17(4):2347–2376. https://doi.org/10.1109/COMST.2015.2444095
2. Alioto M (2017) Enabling the Internet of Things. Springer International Publishing. https://doi.org/10.1007/978-3-319-51482-6
3. Landaluce H, Arjona L, Perallos A, Falcone F, Angulo I, Muralter F (2020) A review of IoT sensing applications and challenges using RFID and wireless sensor networks. Sensors 20:2495. https://doi.org/10.3390/s20092495
4. Ria A, Michel A, Singh RK, Franchina V, Bruschi P, Nepa P (2020) Performance analysis of a compact UHF RFID ceramic tag in high-temperature environments. IEEE J Radio Freq Identif 4(4):461–467. https://doi.org/10.1109/JRFID.2020.2998008
5. Motroni A, Bernardini F, Buffi A, Nepa P, Tellini B (2022) A UHF-RFID multi-antenna sensor fusion enables item and robot localization. IEEE J Radio Frequ Identif 6:456–466. https://doi.org/10.1109/JRFID.2022.3166354
6. Franchina V, Ria A, Michel A, Bruschi P, Nepa P, Salvatore A (2019)A Compact UHF RFID ceramic tag for high-temperature applications. In: 2019 IEEE international conference on RFID technology and applications (RFID-TA), Pisa, Italy, pp 480–483. https://doi.org/10.1109/RFID-TA.2019.8892217
7. Motroni A, Ria A, Strambini L, Nepa P (2022) Experimental assessment of passive UHF-RFID sensor tags for environment and kinematic data. In: 2022 IEEE 12th international conference on RFID technology and applications (RFID-TA), Cagliari, Italy, pp 201–204. https://doi.org/10.1109/RFID-TA54958.2022.9924095
8. Costa F, Genovesi S, Borgese M, Michel A, Dicandia FA, Manara G (2021) A review of RFID sensors, the new frontier of internet of things. Sensors 21:3138. https://doi.org/10.3390/s21093138
9. Colella R, Sabina S, Mincarone P, Catarinucci L (2023) Semi-passive RFID electronic devices with on-chip sensor fusion capabilities for motion capture and biomechanical analysis. IEEE Sensors J 23(11):11672–11681. https://doi.org/10.1109/JSEN.2023.3267540
10. EPC UHF Gen2 Air Interface Protocol. https://www.gs1.org/standards/rfid/uhf-air-interface-protocol. Last Accessed 30 June 2023
11. Ria A, Cicalini M, Manfredini G, Catania A, Piotto M, Bruschi P (2022) The SENSIPLUS: a single-chip fully programmable sensor interface. In: Saponara S, De Gloria, A (eds) Applications in electronics pervading industry, environment and society. ApplePies 2021. Lecture Notes in Electrical Engineering, vol 866. Springer, Cham. https://doi.org/10.1007/978-3-030-95498-7_36
12. www.sensichips.com. Last Accessed 30 June 2023
13. Motroni A, Ria A, Cecchi G, Nepa P (2023, June) Robot-based UHF-RFID joint SAR localization and tag sensing. In: 2023 8th international conference on smart and sustainable technologies (SpliTech). IEEE, pp 1–4
14. Ria A, Manfredini G, Gagliardi F, Vitelli M, Bruschi P, Piotto M (2023) Online high-resolution EIS of lithium-ion batteries by means of compact and low power ASIC. Batteries 9:239. https://doi.org/10.3390/batteries9050239
15. www.emmicroelectronic.com/product/epc-and-uhf-ics/em4325. Last Accessed 30 June 2023

# Digital Twins for Remote ECG Monitoring

Stelios Ninidakis, Miltos D. Grammatikakis[✉], and George Kornaros

Department ECE, Hellenic Mediterranean University, 71410 Heraklion, Greece
{sninidakis,mdgramma,kornaros}@cs.hmu.gr

**Abstract.** Digital twins provide virtual replicas of medical systems, reshaping the E-Health industry towards smart, intelligent, and proactive services. With most existing biosensor solutions, patient data is transmitted over Bluetooth to a gateway that in turn connects via Ethernet to a file server (or cloud) for storage and analysis. In this work, a performance twin of a biosensor is developed, generating self-similar ECG data at increased rates compared to the actual device. In addition, we design a biosensor digital twin representation in the open-source Eclipse Ditto twin management system. We evaluate the scalability of the traditional file server and the digital twin approach. Our results indicate that the file server scales slightly better than the digital twin since Ditto operates in soft real-time only for rates below 4096 pulses/s. The digital twin's energy costs from the client/server viewpoints are also much higher.

**Keywords:** Biosensor · Digital Twin · ECG Monitoring · Energy · Soft Real-Time

## 1 Introduction

The concept of digital twins (DT) has evolved since its inception in the early 2000s [1]. Initially focused on improving product lifecycle management, DT gained momentum with advancements in IoT, cloud computing, and data analytics. The aerospace industry, particularly NASA [2], played a significant role in early adoption, by applying digital twins to simulate spacecraft operations. Industry 4.0 and E-Health have also embraced DT for optimizing manufacturing and enhancing patient care. As technology continues to mature, DT expands into new domains, offering real-time monitoring, data processing, predictive analytics, and data-driven decision-making.

DT can be categorized into different types based on their application and complexity. These types provide a semantic framework for understanding different manifestations and functionalities. A *performance twin* focuses on real-time monitoring and analysis of operational performance [3]. Performance twins continuously gather data from sensors and devices, comparing them to expected or desired metrics. They enable organizations to detect anomalies, identify inefficiencies, and make data-driven decisions to optimize operations and improve outcomes. They are common in E-Health, industrial manufacturing, and infrastructure management.

Eclipse Ditto is the most popular open-source DT platform [4]. It is a robust framework that integrates the Mongo-DB database for managing DT. It supports a wide range of operations, such as DT creation, deletion, access, search, or update via Web Sockets or HTTP, e.g., using POST/PUT/GET/DELETE calls. MQTT, AMQP, and CoAP protocols are supported via Eclipse Hono. Commercial DT platforms like Microsoft Azure DT, GE Predix, and Siemens Mindsphere cater to specific use cases.

In this paper, we develop a performance twin (called device twin) operating at higher ECG rates (>256 pulses/s) than the actual pulse sensor STMicro BodyGateway (BGW). This enables evaluating the scalability of modern E-Health facilities equipped with higher-rate sensors. We also examine real-time ECG monitoring when the sensor's digital twin is represented in Eclipse Ditto, i.e., when Ditto gathers ECG data from the actual BGW sensor [5], or its device twin. Comparing Ditto with a traditional ad-hoc file server approach, Ditto operates in soft real-time only for rates below 4096 pulses/s. The DT energy measured at the client/server is also higher. For the file server, we have considered a lock-free concurrent queue (called CQ) and lock-based circular FIFOs (called list) at the Bluetooth-to-Ethernet gateway. Our results indicate that CQ is slightly faster for large biosensor rates, above 1024 pulses/s.

We next examine how our work relates to state-of-the-art. Experimental work related to DT in E-Health is practically non-existent. Costa et al. evaluate software aging by deploying a vehicle DT using a Kubernetes platform on PCs [6]; for the DT, they use ad-hoc mathematical models without Ditto. Ala-Laurinaho et al. designed a DT of an industrial crane on a Raspberry Pi 4 using commercial frameworks, including Siemens Mindsphere [7]. Value pairs are sent from connected subsystems to an API gateway using HTTP GET/PUT requests, and packet latency is measured using Wireshark. GETs and PUTs appear to have similar latency. Differences from our study are (a) our API service runs directly on the server, (b) our use case focuses on E-health, and (c) we examine soft real-time rather than packet delays. Lee et al. examine the dynamic management of massive IoT devices using Kubernetes and DT services based on Eclipse Ditto with Hono protocol adapters [8]. In their use case, a 16-core server runs Ditto, while 3 similar Hono nodes load balance requests among multiple adapters; each adapter converts HTTP requests into AMQP 1.0 for communicating with Ditto. Their system saturates at 700 requests/s when dozens of IoT devices (emulating mobile robots) generate 10 to 100 events/s each. Although this study examines Ditto and Hono's limitations in massive DT creation/access, it uses Kubernetes for dynamic management and PCs for emulating actual sensors. Finally, Kamath et al. explore representation, real-time acquisition, data analytics, and visualization of a smart factory using Eclipse Ditto and Hono, Apache Kafka, Influx DB, and Grafana on a Linux server with Intel Xeon E5606 and 64 GB RAM [9]. With 100 Things created by a laptop, Ditto is the bottleneck; it takes more than twice the time (0.45s per request) compared to Influx DB and Grafana.

Section 2 discusses the embedded platform and software architecture. Section 3 focuses on the experimental framework comparing the file server approach with the digital twin. Finally, Sect. 4 provides a summary and discusses future work.

## 2   Biosensor, BT-Ethernet Application, and Device/Digital Twins

The STMicro BGW biosensor is a single-lead patch device operating at 128 or 256 pulses/s, supporting Bluetooth. It is like the Preventice BodyGuardian mini patch used for recording ECG data on the cloud for almost 1 million patients in the USA.

Based on the BGW sensor, we have developed a performance twin (called *device twin*) that extends the BGW operating range to 512, 1024, 2048, or 4096 pulses/s). The device twin abstracts the BT interface, and sends self-similar ECG data at higher rates, making it amenable for performance/scalability studies. The accuracy of the device twin is very good since the real BGW device rarely misses its real-time deadline; in this case, it transmits data using special delayed packets (not modeled).



**Fig. 1.**   The BT-to-Ethernet sw stack with the biosensor and device/digital twin.

We have also developed a BT-to-Ethernet gateway on top of an ARMv7 Odroid XU4 single-board computer. The gateway manages communication among the BGW biosensor (or its device twin) and the file server or Eclipse Ditto (managing Things). As shown in Fig. 1, the BT-to-Ethernet gateway uses 5 software components.

**Initialization functions** include (a) the *BT Connect Function* which uses `bluez` tools to pair with the real BGW device, or (b) the *Things Init (Boot) Function* which registers (using HTTP POST) the DT to Eclipse Ditto (stored in Mongo-DB). The JSON file that partially describes the BGW DT prototype is shown later. It is like a Web of Things mapping of a physical IoT device [10].

The **BT Writer Thread** transmits commands using BT packets to configure the BGW device to transmit periodic ECG data at a specified rate, i.e., 128 or 256 samples of ECG data per sec (with 12-bit precision); notice that the BGW device can also support the transmission of bioimpedance and accelerometer data.

The **BGW Reader Thread** uses the rfcomm library to connect and periodically receive raw BT data from the real BGW device at 128, or 256 pulses/s [11]. Then, it periodically extracts ECG values from raw BT data using bit operations and inserts them either into a CQ or list structure. Alternatively, the BGW Device Twin Reader Thread emulates the BGW device for higher rates, inserting self-similar ECG values to a CQ, since this structure is faster.

The **BT-to-Ethernet Sharing Thread** periodically dequeues biometric data from the CQ (or resp. The list) and transmits them either to (a) a file server (an Odroid XU3 board) or (b) transfers them via HTTP PUT to the Eclipse Ditto Server (a Dell PowerEdge T330 with 32GB RAM and 4 HDD drives, running Ubuntu 22.04.2 LTS and Eclipse Ditto 3.3.0 from [4]) via a 2.1 Gbit/s TP-Link Archer C5400 router. Transmission occurs when ECG data exceeds the current BGW rate (or 512, for rates above 1024 pulses/s); this optimization improves network performance.

In the BGW twin (JSON), acl is the authentication section. Attributes/features describe static/dynamic (name, value) pairs with defaults (or latest values).

```
{"thingId":                      "attributes": {
"org.eclipse.ditto.example:odroid",   "manufacturer": "STMicro",
  "acl":{                             "gateway": "ecg-ditto"},
    "ditto": {"read": true,      "features": {
             "write": true,          "BGW": {
             "admin": true             "properties": {
           },                            "ECG Array": [ 90.9, ...],
    "odroid": {"read": true,            "date":"2023-07-16",
              "write": true,            "time":"12:00:00",
              "admin": false            "samplingRate": 4096
            }                      } } } }
  },
```

## 3   Scalability of File Server and Device Twin Use Cases



**Fig. 2.** The cumulative rate of ECG samples at the file server for CQ and list.

In our first use case, we consider traditional file server performance for both the concurrent queue and lock-based queue implementation of the BT-to-Ethernet gateway. As discussed, gateway input is either from a real BGW device (for ECG rates of 128 or 256 pulses/s) or from its device twin (for higher rates) [12].

Large fluctuations are experienced in the instant server rate for small BGW rates of 128 and 256 pulses/s, since BGW transmits pulses in bursts. However, as shown in Fig. 2, we concentrate on the cumulative average rate at the file server, which is more decisive from a soft real-time viewpoint. This rate is computed using the number of ECG samples processed and the inter-arrival latency. During the period of measurements (60 s), both implementations (CQ and list) sustain even the higher BGW rates in real-time; graphs for smaller rates are omitted.

In addition, we have profiled power consumption at the gateway during file server access. For both queues, power consumption is small (~0.2 W on average, with rare peaks at 0.6 W). Maximum power at the file server is much higher (~50W).



**Fig. 3.** The rate of PUT Operations on the Eclipse Ditto. PUT is slower than GET.

In our second use case, we focus on the rate of PUT operations at the Eclipse Ditto server, when the BGW (or its device twin) transmits at 128–4096 pulses/s. As shown in Fig. 3, the server can update the ECG array of the BGW DT stored in Eclipse Ditto in real-time, for all sampling rates below 4096 pulses/s. Note that for 1024 or 2048 pulses/s, there is a larger variation, however, the cumulative rate eventually matches the BGW rate.

Power consumption at the gateway is 3.3–3.7 W (with peaks up to 5.3 W). Power consumption at the Ditto server is again much higher (in the region of ~90 W).

Further experiments with more simultaneously transmitting BGW devices show that soft real-time is possible when up to two BGW devices operate at 128 pulses/s. The same is true, even if simultaneous GET operations occur since GET is much faster than PUT: 8 times faster for a single transmitting device, or twice faster for two devices; graphs for this experiment are omitted due to lack of space.

## 4 Conclusions and Future Work

With Healthcare 4.0, the widespread availability of smart, wearable medical sensors and related gateway, server, and digital twin platforms can improve patient health by continuously collecting and analyzing patient data. Therefore, modern E-Health architectures must support the reliability, power efficiency, and real-time requirements of medical monitoring and decision-support applications. In this context, this study provides a first glimpse at the performance, scalability, and power consumption of traditional ad-hoc file servers and digital twins.

In the future, it would be interesting to examine the system efficiency for different data processing and propagation methods and analytics, including compression and filtering for heart-beat detection and classification.

Moreover, the symbiosis of performance with security and data privacy in digital twin systems for E-Health requires careful integration of existing application workloads with efficient security patterns, mechanisms, protocols, and policies [13]. Hence, in this direction, it is interesting to develop reliable authentication protocols, that use standalone, tamper-proof cryptographical circuits, and provide configurable high-level security policies among sensor, gateway, and server entities [14, 15].

## References

1. Grieves MW (2023) Digital twins: past, present, and future. In: Crespi N, Drobot AT, Minerva R (eds) The digital twin. Springer
2. Li L, Aslam S, Wileman A et al. Digital twin in the aerospace industry: a gentle introduction. IEEE Access 10:9543–9562
3. Khan LU, Han Z, Saad W et al (2022) Digital twin of wireless systems. IEEE Commun Surv Tutor 2230–2254
4. Eclipse Ditto. 18 July 2023. https://projects.eclipse.org/projects/iot.ditto
5. BodyGateway, ST Micro. 18 July 2023. https://usermanual.wiki/ST-Microelectronics-S-R-L/MHBGW1.Product-literature
6. Costa J, Matos R, Araujo J et al. Software aging effects on Kubernetes in container orchestration systems for digital twin cloud infrastructures of urban air mobility Drones 7(1), 35, 1–22
7. Ala-Laurinaho R, Autiosalo J, Nikander A et al. Data link for the creation of digital twins. IEEE Access 8:228675–228684
8. Lee J, Kang S, Chun I-G (2021)mIoTwins: design and evaluation of mIoT framework for private edge networks. In: International conference information communication technology convergence (ICTC), Jeju Island, Korea, pp 1882–1884
9. Kamath V, Morgan J, Ali MI (2020) Industrial IoT and digital twins for a smart factory. In: Global Internet of Things Summit (GIoTS), Dublin, Ireland, pp 1–6
10. Grammatikakis MD, Ninidakis S, Kornaros G et al (2022) Towards efficient gateways and servers for biosensors In: International conference applications in electronics pervading industry, environment and society, pp 346–352

11. Web of Things (WoT), 18 July 2023. https://www.w3.org/WoT/
12. Corral-Acero J, Margara F, Marciniak M et al (2020) The "Digital Twin enables the vision of precision cardiology." Eur Heart J 41(48):4556–4564
13. Hathaliya JJ, Tanwar S (2020) An exhaustive survey on security and privacy issues in healthcare 4.0. Comput Commun 153:311–335
14. Grammatikakis MD, Petrakis P, Papagrigoriou A et al (2015) High-level security services based on a hardware NoC Firewall module. In: International workshop on intelligent solutions in embedded systems (WISES), Ancona, Italy, pp 73–78
15. Chukwu E, Garg L (2020) A systematic review of blockchain in healthcare: frameworks, prototypes, and implementations. IEEE Access 8:21196–21214

# A Deeply Quantized Classifier for Very Low Resolution ToF Imaging

Danilo Pietro Pau[1(✉)], Welid Ben Yahmed[2], and Jeffrey M. Raynor[3]

[1] FIEEE, STMicroelectronics, System Research and Applications, Agrate Brianza, Italy
Danilo.pau@st.com
[2] Universita degli Studi di Trento, Trento, Italy
[3] STMicroelectronics, Imaging, Edinburgh, UK

**Abstract.** This paper presents the design of tiny deeply quantized neural networks suitable for super integration within the Time-of-Flight, low-resolution, image sensor. They were aimed to achieve ultra low complexity with adequate classification accuracy. First floating-point models were studied to process 8-bits images $2 \times 2$ pixel resolution, downsampled from $8 \times 8$ images of a public dataset. Data were acquired with an off-the-shelf Time-of-Flight sensor. Next, many deeply quantized networks were designed from scratch by using QKeras quantization training-aware schema. 8 and 6 bits pixel depth were generated by the sensor. Ternary, 6, and 8 bits quantizers were used along with convolutions and dense layers. Experimental results have shown that the proposed deeply quantized models achieved in maximum an accuracy of 77.79% compared to 88.48% for the floating-point models. The model size of the latters ranged from 92,320 bytes for the largest floating-point model to 151 bytes for the tiniest 6-bits deeply quantized model.

**Keywords:** Deeply quantized neural networks · Time-of-flight sensor · Low-resolution image classification · QKeras

## 1  Introduction

Since 2019, the TinyML[1] community demonstrated practical applications of Machine Learning (ML) based on resource-constrained artificial neural networks (NNs). MLCommons/Tiny[2] industry benchmarks have proven that NN tiny inferences can consume energy as low as a few tens of $\mu$J. Recently Deeply Quantized Neural Networks (DQNNs) [1] delivers un-compromised accuracy (w.r.t. floating point 32bits (fp32) NNs) with the lowest hardware cost. This paper investigated the applicability of DQNNs to tiny Time-of-Flight (ToF) sensors. A possible application is detecting hand-hovering above the sensor, enabling actions such as initiating system power-up or transition to sleep mode. For such usage, a high resolution sensor is unnecessary and against a compact solution.

---

[1] https://www.tinyml.org/.
[2] https://mlcommons.org/en/inference-tiny-10/.

## 2   Problem Definition and Requirements

The challenge was to design DQNN models tiny enough to be super-integrated within the ToF sensor. The DQNN models shall require very limited on-chip memory. High classification accuracy shall be achieved w.r.t. fp32 NNs. The computational complexity, the number of multiplications and additions, shall be the very low and use few bits. Any fp32 NN shall be deployable on a low-power off-the-shelf micro-controller (e.g. such as the low-power STM32U5[3]) Microcontroller Unit(MCU). The DQNNs shall be deemed for further hardware implementation.

## 3   Related Works

Reference [2] presented a self-powered module that combined energy harvesting and gesture sensing using small low-cost photo-diodes. The recognition module included a MCU running a robust and lightweight algorithm that could accurately recognize finger gestures even in the presence of ambient light with different intensities, directions, and fluctuations. The results demonstrated that the system achieved high recognition accuracy (>96%) in all tested settings.

Tiny Deep image recognition was performed in [3] by using visual attention condensers. The AttendNets was presented as a low-precision, compact deep NN tailored for on-device image recognition. Experimental results on the ImageNet50 dataset showed that AttendNets lowered 3 times the number of multiplication and addition operations and improved the accuracy by 7.2%. It required 4.17 times less parameters and minimized the memory to store the weights by 16.7 times if compared to MobileNet-V1.

Reference [4] introduced a very tiny re-configurable NN hardware able to compensate for the sensor drift measures. The circuit required $0.0373 \text{ mm}^2$ in 130 nm CMOS technology and a dynamic power of $1.07 \text{ µW}$. These results enabled its integration in the same sensor package for pressure measures compensation.

## 4   Dataset and preprocessing

The dataset shapes a case study where data are captured using a low resolution ToF.

The original public dataset is composed of ToF[4] depth images sampled by a miniature mobile robot navigating in a sand-like terrain [5]. Each image has a resolution of 8x8 pixels. The task involves classifying the images into two classes: Background and Robot. It comprises a total of 4,150 samples, with 2,062 samples belonging to the Background class and 2,088 samples to the Robot class. The bit-depth is 8 bits per pixel, where each pixel value represents the depth information. The dataset isn't normalized.

---

[3] https://www.st.com/en/microcontrollers-microprocessors/stm32u535-545.html.
[4] https://www.st.com/en/imaging-and-photonics-solutions/vl53l5cx.html.

Pre-processing was applied to the dataset. First, the image resolution was re-sampled to $2 \times 2$ pixels. The number of samples per class was balanced resulting in a total of 4,124 samples for the two classes (Background and Robot).

For the 6-bits and 8-bits DQNNs, the pixel values were normalized, respectively, between –32 and 31 and between –128 and 127.

## 5   Proposed Neural Networks

The NN architectures were built on top of convolution blocks. In total, three sets of models were investigated:

1. floating point 32-bits (fp32).
2. fixed point 8-bits (named 8q).
3. fixed point 6-bits (named 6q).

### 5.1   Floating-Point 32-Bit Models

The main goal was to achieve maximum accuracy. The model footprint, and the number of parameters, were tuned w.r.t. the amount of training samples. Models varied the number of convolution blocks, the number, and size of the filters. Some NNs have been designed with an activation layer and/or a batch normalization step. Non-linearity was the Rectified Linear Units (ReLU). Also fully connected layers were used and the Softmax layer computed the probabilities for the two output classes. The general structure of the fp32 NNs is reported in Table 1.

**Table 1.** fp32 NNs configurations

| Series block | Convolution 1D | Fully connected |
|---|---|---|
| Series details | Number of filters filters size: 2,3 or 4 batch normalization ReLU | Output size ReLU, Softmax |

### 5.2   6 and 8-bits DQNNs

Using quantization-aware training with Qkeras, the DQNN models were designed trying to keep the same levels of accuracy w.r.t. the fp32 models.

Different combinations of numbers of convolution blocks, fully connected layers, numbers of kernels, and sizes of kernels were explored. Moreover, different QKeras quantization functions were used covering fixed number of bits (e.g. 3, 6, 8), ternary (e.g. –1, 0, 1), and powers of two (–8 to 8). Table 2 sums up the configurations and the different parameters for the proposed DQNNs. The DQNNs differ by the fixed number of bits of the quantization function being either 6 or 8 bits for 6-bits DQNNs and 8-bits DQNNs respectively.

**Table 2.** DQNNs configurations

| Series block | Convolution 1D | Fully connected | Quantizer func-tions |
|---|---|---|---|
| Series details | Kernel quantizer, bias quantizer quantized activation ReLU | Kernel quantizer | Quantized_bits, ternary, quantized_po2, quantized_relu |

## 6    Accuracy Studies

Performance evaluation was performed by 5-fold validation with one fold being used for testing. Across the test fold for the different data splits, the mean accuracy and the variance were computed. The training was performed using a cross-entropy loss function. Figure 1 shows that the largest mean accuracy value achieved across fp32 models was 88.48% corresponding to the model f10 with a variance of 2%. Figure 2 shows that the greatest mean accuracy value across DQNN models was 77.79% and was achieved by the model f10 with a variance of 2%. The fp32 models achieved a significant higher accuracy than the DQNNs models.



**Fig. 1.** Mean accuracy and variance with respect to the number of parameters for fp32 models

**Fig. 2.** Mean accuracy and variance with respect to the number of parameters for DQNN models

## 7 Complexity Analysis

The complexity of the fp32 models was reported in Table 3 in terms of the number of multiply-accumulate (MACC) operations, random access memory (RAM), read-only memory (ROM) and inference time measured on the Arm Cortex-M33 MCU running at 160 MHz using an off-the-shelf cloud-based deployment (on the MCU) service.[5] The MCU embedded RAM was 768 KiB and the embedded flash was 2 MB.

## 8 Conclusions

Binary classifiers of $2 \times 2$ ToF images of a miniature robot were studied. Several 6 and 8-bits DQNNs were compared to fp32 models. When comparing these models in the same range of number of parameters as reported in Table 4, the gap of the maximum accuracy between DQNNs and fp32 models was 5.58% whereas the quantized models were significantly lower in ROM size. Future work will use knowledge distillation in order to take benefit from the higher accuracy achieved by the fp32 models.

---

[5] https://stm32ai-cs.st.com/home.

**Table 3.** Complexity of fp32 models

| Series model | Series MACC | Series ROM size [Bytes] | Series RAM size [Bytes] | Series Inference time [ms] |
|---|---|---|---|---|
| f1 | 46928 | 77440 | 4456 | 2.005 |
| f4 | 46096 | 88904 | 4456 | 2.113 |
| f5 | 38968 | 63896 | 4856 | 1.857 |
| f6 | 43064 | 67992 | 4856 | 2.019 |
| f7 | 17784 | 42712 | 4600 | 1.094 |
| f8 | 24680 | 55696 | 5624 | 1.532 |
| f9 | 47440 | 92320 | 6040 | 2.185 |
| f10 | 44504 | 72204 | 6968 | 2.087 |
| f11 | 31128 | 58832 | 6840 | 1.603 |
| f12 | 39864 | 67208 | 6232 | 1.883 |
| f16 | 6224 | 28728 | 5804 | 0.551 |
| f17 | 10096 | 29420 | 4784 | 0.633 |
| f19 | 10760 | 29200 | 5200 | 0.794 |
| f22 | 3080 | 20304 | 4128 | 0.353 |
| f24 | 4040 | 22724 | 5052 | 0.400 |

**Table 4.** DQNNs and fp32 models in the same range of parameters

| Series model | Series parameters | Series ROM size [Bytes] | Series mean accuracy | Series maximum accuracy |
|---|---|---|---|---|
| 6q_1 | 2816 | 2028 | 77.79% | 78.42% |
| 8q_1 | 2816 | 2704 | 73.38% | 77.21% |
| f19 | 2818 | 29200 | 81.26% | 82.79% |

# References

1. Coelho Jr CN, Kuusela A, Zhuang H, Aarrestad T, Loncar V, Ngadiuba J, Pierini M, Summers S, Ultra low-latency, low-area inference accelerators using heterogeneous deep quantization with QKeras and hls4ml. http://arxiv.org/abs/2006.10159v1
2. Li Y, Li T, Patel RA, Yang X-D, Zhou X (2018) Self-powered gesture recognition with ambient light. In: Proceedings of the 31st annual ACM symposium on user interface software and technology (UIST '18). Association for Computing Machinery, New York, NY, USA, pp 595–608. https://doi.org/10.1145/3242587.3242635
3. Wong A, Famouri M, Shafiee MJ (2020) AttendNets: tiny deep image recognition neural networks for the edge via visual attention condensers. arXiv:2009.14385
4. Licciardo GD, Vitolo P, Bosco S, Pennino S, Pau D, Pesaturo M, Di Benedetto L, Liguori R (2023) Ultra-tiny neural network for compensation of post-soldering thermal drift in MEMS pressure sensors. In: 56th edition of the IEEE international symposium on circuits and systems (IEEE ISCAS 2023)
5. Jan P, škulj G, Esnault C, Puc J, Vrabič R, Podržaj P (2023) Miniature mobile robot detection using an ultra-low resolution time-of-flight sensor dataset. IEEE Dataport

# On Enhancing the Throughput of the Latched Ring Oscillator TRNG on FPGA

Riccardo Della Sala[✉] and Giuseppe Scotti

Università La Sapienza, Via Eudossiana 18, Rome 00184, Italy
`riccardo.dellasala@uniroma1.it`

**Abstract.** In this summary, we introduce a new design of the Latched Ring Oscillator (LRO), true-random-number-generator (TRNG) on a 7-Series FPGA, demonstrating the portability of the architecture. The novel design, combined with a novel sampling strategy, allows to enhance performance of the previous work, speeding up the throughput of about 265 times. The proposed LRO-TRNG effectively utilizes both meta-stability, manufacturing variations and accumulated jitter as sources of entropy, resulting in excellent levels of unpredictability and randomness. The TRNG's performance has been validated with considering NIST and AIS-31 tests. Measurement results indicate that the LRO-TRNG achieves an estimated entropy of approximately 7.99984 per byte (as per the T8 test of the AIS-31) and a throughput of 201 Mbits/s using a 450MHz clock. When compared to existing TRNGs, the LRO-TRNG surpasses most of previously published designs in terms of the throughput attaining a good trade-off among speed and FPGA resources usage.

**Keywords:** TRNGs · FPGA · Ring oscillators · Metastability · Jitter

## 1 Introduction

Nowadays many essential services used in daily life, such as data storage, booking and banking, rely on authentication protocols, which should meet hard requirements in terms of device security. Several authentication protocols have been proposed in recent literature [1], in order to effectively protect sensitive data on devices through the usage of Physical Unclonable Functions (PUFs) [2–5] for the key generation and True Random Number Generators (TRNGs) [6]. For example, the Privacy-Preserving mutual authentication protocol requires the simultaneous presence of both PUF and TRNG primitives [7–9]. TRNGs are among the most hard-to-design building blocks requiring the adoption of ad-hoc strategies to assure good entropy performance. In addition, deploying TRNGs on portable devices presents significant challenges due to the multitude of requirements that must be met simultaneously.

The use of Field-Programmable Gate Arrays (FPGAs) has gained significant popularity in prototyping and deploying embedded systems. As a result, more

and more devices are leveraging the capabilities offered by FPGAs. TRNGs typically rely on dynamic entropy sources, such as jitter accumulation [10], thermal noise [11], and metastability [6]. However, the exploitable entropy sources on FPGA devices are somewhat limited. One of the commonly used entropy sources in TRNG architectures is jitter in ring oscillators (ROs) [6,12–15]. The fundamental principle behind RO-TRNGs is to simultaneously stimulate ring oscillators with different nominal frequencies and to sample their outputs at a different sampling frequency. This process generates a random output sequence that depends on jitter accumulation. The different ring-oscillators outputs can be XOR-ed between each other in order to enhance the entropy of the generated bitstream [6], at the cost of reduced Throughput. It is important to note that these designs often require post-processing encoding schemes (e.g., Von Neumann encoding [16]) to pass statistical tests, and therefore introduce additional FPGA resources and power consumption.

In a previous study [12], the metastability of latched ring oscillator was utilized as an entropy source combined with the accumulated jitter. The architecture was implemented on a Xilinx Spartan-6 FPGA with a reference clock of 50 MHz, which significantly limited the throughput of the TRNG which resulted about 0.76 $MBits/s$.

In this work, we present the porting of the LRO TRNG on a Xilinx Artix-7 FPGA. As in the previous work [12], the proposed TRNG implementation harnesses both the jitter in ring oscillators and the metastability of D-latches as entropy sources. However, differently than [12], the design proposed in this work instantiates 128 identical LRO cells which are sampled in parallel to better exploit manufacturing variations. Extensive investigations were conducted to explore the interconnections among the Artix-7 FPGA elements and to understand the characteristics and physical behaviors of the built-in blocks. This research provided valuable insights into the jitter accumulation mechanism and allowed us to fully exploit metastability, optimizing the TRNG's performance in terms of entropy and randomness while minimizing hardware resources utilization.

## 2    The Latched Ring Oscillator

The Latched Ring Oscillator (LRO) architecture is depicted in Fig. 1a. The single LRO is composed by an inverter (NOT) gate and a D-type latch with an asynchronous reset pin $R$. The output of the latch Q is fed back through an inverter to the its D terminal. with this arrangement, the circuit is in a ring oscillator (RO) configuration when the gate pin $G$ of the latch (connected to the $START$ signal) goes high and the reset pin $R$ of the latch connected to the $RESET$ signal is low. The startup sequence of the proposed cell is described below:

– Reset-phase: during this phase, the $RESET$ signal is high and the output of the latch is forced to 0, thus its input is forced to 1, due to the inverter behavior.

– RO-phase: during this phase the $RESET$ signal goes low and at the same time the $START$ signal goes high, initializing the oscillation.
– Sampling-phase: after $N_{clk}$ clock cycles in which the ring oscillator was allowed to oscillate, collecting Jitter, the $START$ signal goes low forcing the latch to be opaque (in its holed state) and thus locking the oscillation. During this phase the LRO takes advantage of the metastability of the latch and of the Jitter accumulated during the oscillation.



**Fig. 1.** The Latched Ring Oscillator topology (**a**) and proposed scheme to enhance the throughput of the latched ring oscillator (**b**).

Due to its simplicity, the single LRO is very compact, and is able to fit in only 1/4 of a FPGA Slice. Since the focus of this paper is on achieving high throughput from this simple but effective cell, we changed our strategy, and focused on the parallel sampling of 128 identical LRO cells, according to the scheme in Fig. 1b. The throughput of the architecture in Fig. 1b can be expressed as:

$$Throughput = \frac{f_{clk} \cdot N_{cells}}{N_{clk} + 1} \tag{1}$$

where it has been denoted with $N_{cells}$ the number of instantiated cells and with $N_{clk}$ the number of clock cycles after which the output is sampled. For the proposed design $N_{cells}$, and $N_{clk}$ have been set to 128 and 10 respectively, whereas a clock frequency of 450 MHz has been chosen. We have taken into account also the reset phase by considering an additive clock cycle. With these values the throughput in Eq. 1 results to be 5.24 Gbit/s.

However, as demonstrated in [12], the performance of a single LRO are not good enough to fulfill statistical requirements, and the outputs of several nominally identical LROs have to be combined to collect enough entropy in order to pass NIST and AIS-31 tests. In fact, due to the effect of mismatch and process variations on routing delays, each LRO oscillates with a frequency that follows a normal distribution $\mathcal{N}(\mu_{ring}, \sigma_{ring}^2)$. If the outputs of the 128 LRO cells are combined, after $N_{clk}$ clock cycles the output bitstream will depend on the Jitter accumulation of the different LROs, and also on the oscillation frequency of the different LRO instances (which is related to mismatch and process variations). By combining these effects a very high throughput bitstream, with good statistical performances can be extracted from the scheme of Fig. 1b. The main

post-processing approaches to combine the outputs of the different LRO cells in a single bitstream are the Von Neumann encoding and the XOR-combining. Both these approaches allow to improve the statistical performances, but result in an output bitstream with lower throughput.

## 3   FPGA Implementation

The proposed architecture has been integrated on an Artix-7 FPGA with an ad-hoc strategy to implement the macro cell of the latched ring oscillator, reported in Fig. 2. Indeed, in order to make the most of mismatch and process variations, the design should guarantee that all the cells are nominally identical. In this way the variation of the clock frequency will be given just by the given site in which the cell is instantiated. The 8 bit macro and the 16 identical 8 bits macros are reported in Fig. 2a and b respectively. As it can be observed, the adopted routing strategy allowed to obtain very similar routing connections and thus also very similar clock frequencies for the different LROs. The average oscillation frequency of the different LROs has been evaluated to be approximately equal to 370 MHz. For this computation the routing delay from the output of the inverter to the input of the latch, the routing delay from the output of the latch to the input of the inverter and finally the propagation delay of the inverter and of the latch have been evaluated in the Xilinx Vivado environment. Being the inverter implemented through a look up table (LUT), its delay is strongly related to the input gate on which the operation is performed. However, this information is not given by the producer and we considered the delay from the generic input to the output of the LUT. Anyway, it has also to be remarked that these delays have been extracted with considering the SLOW_MAX corner and thus extracted oscillation frequencies are a worst-case scenario.



**Fig. 2.** Macro of 8 bits of the LRO (**a**) 16 identical 8 bits macros instantiated on FPGA (**b**).

## 4   Statistical Performance and Evaluation

Main performances of TRNGs include the throughput, the resources usage in terms of Slice/bit, the biasing and complexity of the bitstream evaluated through ad-hoc statistical tests such as NIST 800-90B, and entropy evaluation, for example through the AIS-31 T8 test or through the NIST 800-22A.

**Table 1.** 800-90B NIST tests

| Test name | Result | Pass | Test name | Result | Pass |
|---|---|---|---|---|---|
| Frequency | 0.5982 | ✓ | Block Frequency | 0.3724 | ✓ |
| Cumulative Sums | 0.3481/0.1573 | ✓ | Runs Test | 0.7345 | ✓ |
| Longest Run | 0.1102 | ✓ | Rank | 0.7754 | ✓ |
| FFT | 0.4398 | ✓ | Non Overlapping | 0.4651 | ✓ |
| Overlapping | 0.1735 | ✓ | Universal | 0.3379 | ✓ |
| Approximate Entropy | 0.3135 | ✓ | Random Excursion | 0.0349/0.2674 | ✓ |
| Random Excursion Variant | 0.01735/0.4673 | ✓ | Serial | 0.0137/0.6681 | ✓ |
| Linear Complexity | 0.5982 | ✓ | | | |

As a first characterization step of the proposed TRNG we performed the NIST 800-90B tests on the bistream generated without any post-processing technique, and from the results it was clear that some kind of post-processing was needed to improve the statistical properties and to satisfy NISTs' requirements. Then, we considered two Von Neumann rounds on the extracted bitstream, which have been found sufficient to correct the biasing of the bitstream and to increase the bitstream complexity. On the post-processed bitstream NIST tests were conducted and results of 800-90B tests have been reported in Table 1. In addition, AIS-31 tests were performed on the post-processed bitstream and the resulted byte entropy (T8 test result) has been found to be about 7.99984 .

Finally, the average compression ratio due to the two Von Neumann rounds has been evaluated on 20 bitstreams, and it has been found to be about 26, resulting in an overall throughput of about 201 [MBit/s] ($\frac{450 \cdot 128}{(10+1) \cdot (26)}$ [Mbit/s]).

## 5    Comparison with the Literature

The proposed TRNG design has been compared with others from literature, and results are summarized in Table 2. Table 2 shows that the proposed design is able to guarantee the highest throughput (T.P.) with respect to the recent literature. However, if we look at the FOM reported in [12] defined as:

$$FOM = \frac{TP}{O.F.} \cdot (Bits/Slice) \tag{2}$$

the proposed design achieves a good value of the FOM, but the highest value is achieved is achieved by [17] which obtains the best trade-off among resource consumption and throughput.

**Table 2.** Comparison table

|  | Bits/Slice | TP [MBit/s] | O.F. [MHz] | FOM [bits/Slice] | Device |
|---|---|---|---|---|---|
| This work | 0.1538 | **201** | **450** | 0.0687 | Artix-7 |
| [12] | **1** | 0.76 | 50 | 0.0152 | Spartan-6 |
| [18] | 0.0065 | 100 | 400 | 0.0016 | Virtex-5 |
| [17] | 0.25 | 100 | 100 | **0.25** | Artix-7 |
| [19] | 0.0714 | 100 | 400 | 0.0179 | Virtex-6 |

TP: Throughput, O.F.: Operating Frequency, FOM: Figure of Merit.

## 6    Conclusion

In this work we ported the LRO TRNG on a 7-Series FPGA, demonstrating the portability of the architecture. In addition a novel design has been carried out to enhance throughput performance of the previous work, speeding up the throughput of about 265 times. Despite its small size, the presented LRO-TRNG effectively utilizes both meta-stability and accumulated jitter as sources of entropy and make the most of process and mismatch variations by collecting 128 bits in parallel from 128 different sites. Performances have been validated with considering NIST and AIS-31 tests. Measurement results indicate that the proposed LRO-TRNG design achieves an estimated entropy of approximately **7.99834** per byte (as per the T8 test of the AIS-31) and a throughput of about 201 Mbits/s using a 450MHz clock.

## References

1. Wang Y, Liang H, Wang Y, Yao L, Yi M, Huang Z, Lu Y (2022) A reconfigurable PUF structure with dual working modes based on entropy separation model. Microelectron J 124:105445
2. Della Sala R, Bellizia D, Scotti G (2021) A novel ultra-compact FPGA PUF: the DD-PUF. Cryptography 5:23
3. R. Della Sala, Bellizia D, Scotti G (2022) A lightweight FPGA compatible weak-PUF primitive based on XOR gates. IEEE Trans Circuits Syst II: Express Briefs 69:2972–2976
4. Della Sala R, Bellizia D, Centurelli F, Scotti G (2023) A monostable physically unclonable function based on improved RCCMs with 0–1.56% native bit instability at 0.6–1.2 V and 0–75° C. Electronics 12:755
5. Della Sala R, Scotti G (2023) A novel FPGA implementation of the NAND-PUF with minimal resource usage and high reliability. Cryptography 7:18
6. Della Sala R, Bellizia D, Scotti G (2022) High-throughput FPGA-compatible TRNG Architecture exploiting multistimuli metastable cells. IEEE Trans Circuits Syst I: Regul Pap 69:4886–4897
7. Japa A, Majumder MK, Sahoo SK, Vaddi R (2021) Tunnel FET-based ultra-lightweight reconfigurable TRNG and PUF design for resource-constrained internet of things. Int J Circuit Theory Appl 49:2299–2311

8. Della Sala R, Scotti G (2023) Exploiting the DD-cell as an ultra-compact entropy source for an FPGA-based Re-configurable PUF-TRNG architecture. IEEE Access 11:86178–86195
9. Della Sala R, Scotti G (2022) The DD-cell: a double side entropic source exploitable as PUF and TRNG. In: 2022 17th conference on Ph.D research in microelectronics and electronics (PRIME). IEEE, pp 353–356
10. Prada-Delgado MA, Martínez-Gómez C, Baturone I (2020) Auto-calibrated ring oscillator TRNG based on jitter accumulation. In: 2020 IEEE international symposium on circuits and systems (ISCAS). IEEE, pp 1–4
11. Matsuoka S, Ichikawa S, Fujieda N (2021) A true random number generator that utilizes thermal noise in a programmable system-on-chip (PSoC). Int J Circuit Theory Appl 49:3354–3367
12. Della Sala R, Bellizia D, Scotti G (2021) A novel ultra-compact FPGA-compatible TRNG architecture exploiting latched ring oscillators. IEEE Trans Circuits Syst II: Express Briefs 69:1672–1676
13. Li X, Stanwicks P, Provelengios G, Tessier R, Holcomb D (2023) Jitter-based adaptive true random number generation circuits for FPGAs in the cloud. ACM Trans Reconfigurable Technol Syst 16:1–20
14. Li X, Stanwicks P, Provelengios G, Tessier R, Holcomb D (2020) Jitter-based adaptive true random number generation for FPGAs in the cloud. In: 2020 international conference on field-programmable technology (ICFPT). IEEE, pp 112–119
15. Fischer V, Bernard F, Bochard N, Varchola M (2008) Enhancing security of ring oscillator-based TRNG implemented in FPGA. In: 2008 international conference on field programmable logic and applications. IEEE, pp 245–250
16. Avaroğlu E, Tuncer T, Özer AB, Ergen B, Türk M (2015) A novel chaos-based post-processing for TRNG. Nonlinear Dyn 81:189–199
17. Addabbo T, Fort A, Moretti R, Mugnaini M, Takaloo H, Vignoli V (2020) A new class of digital circuits for the design of entropy sources in programmable logic. IEEE Trans Circuits Syst I: Regul Pap 67:2419–2430
18. Cherkaoui A, Fischer V, Fesquet L, Aubert A (2013) A very high speed true random number generator with entropy assessment. In: Cryptographic hardware and embedded systems—CHES 2013. Springer, Berlin, pp 179–196
19. Wang X, Liang H, Wang Y, Yao L, Guo Y, Yi M, Huang Z, Qi H, Lu Y (2020) High-throughput portable true random number generator based on jitter-latch structure. IEEE Trans Circuits Syst I: Regul Pap 68:741–750

# An Advanced Customizable Circuit Simulator to Investigate Memristor Dynamics

Riccardo Moretti[✉], Tommaso Addabbo, Ada Fort, and Valerio Vignoli

Department of Information Engineering and Mathematics, University of Siena, Siena, Italy
riccardo.moretti@unisi.it

**Abstract.** In this paper we present a circuit simulator based on MAT-LAB, developed in a dedicated way for the analysis of simple circuit topologies given by the series between a memristor described by the Stanford model and a current limiter. The simulator is built taking into account critical aspects regarding the numerical integration of the Stanford model. The performed simulations demonstrate the ability of the simulator to manage the considered critical aspects, confirming the reliability of the results it provides.

**Keywords:** Stanford Memristor · Circuit Simulator · MATLAB

## 1 Introduction

A non-volatile memristor can be defined as a nonlinear device which possesses an internal state that determines its resistance. The dynamics of this internal state is influenced by the applied excitation and encompasses an infinite number of equilibrium points within a specific range [1]. Exploiting this dynamical property, the memristor resistance can be varied, enabling the utilization of this device as a programmable resistor or as a component for memory storage [2]. To investigate memristor dynamics, accurate device numerical models are mandatory. In particular, from a circuit design point of view, models should allow the retrievement of information regarding the relationship between the characteristics of applied stimuli and the resulting memristor state dynamics, given an initial state condition. Among the different solutions presented in literature, a large number of models are deterministic [3–6], based on a relation and a differential equation of the form:

$$\begin{cases} f_m(i, v, \mathbf{M}) = 0, \\ \frac{d\mathbf{M}}{dt} = \mathbf{G}_m(i, v, \mathbf{M}), \end{cases} \tag{1}$$

where $\mathbf{M} \in \mathbb{R}^n$ is the internal state variable and $i, v \in \mathbb{R}$ are the memristor current and voltage, respectively. One of the most relevant proposals among this

class of models is represented by the Stanford model [7]. In our research activity we investigate the dynamical aspects of non-volatile memristors programming, referring to the Stanford model representation [8]. We refer to the generic circuit topology shown in Fig. 1, given by the series between a memristor and a generic two-terminal nonlinear resistor acting as a current limiter. In detail, we relate the static characteristics of the limiting circuit to the dynamical characteristics of the memristor. In the investigation process, a fundamental tool for the verification of hypotheses is their validation through simulations based on the numerical integration of the circuit. In this sense, a simple simulator is not enough, as it does not allow the implementation of the preliminary static analyses of the circuit and the comparison of the simulations with the fixed premises. For this purpose, despite the fact that a Stanford model official description in Verilog-A compatible with SPICE is available [9], we developed a dedicated work environment, capable of configuring the simulation parameters on the basis of the preliminary theoretical analyses, accurately simulating the circuit, and interfacing the simulation results with the initial hypotheses.



**Fig. 1.** Generic representation of a current-limited memristor circuit topology.

In this work, we present our dedicated simulation environment based on MATLAB, purposely designed to provide low-level control of simple circuit topologies referable to the scheme shown in Fig. 1, in which the current limiter can be implemented through sub-circuits with different complexities. The paper is organized as follows. In Sect. 2 the critical aspects of the Stanford model with respect to its simulation over time are presented. In Sect. 3 the developed simulation environment is described. Experiments, conclusion and references conclude the paper.

## 2   Stanford Model

Stanford model concerns filamentary devices characterized by a two-dimensional state $\mathbf{M} \in \Omega = \{(g, T) \in [g_{\min}, g_{\max}] \times \mathbb{R}^+\}$, where $g$ is the distance from the tip of the conductive filament to the top electrode of the memristor (i.e., the gap

distance) and $T$ is the device temperature. Without going into detail about the model equations, we focus for our purposes on the differential equation describing the gap distance dynamics:

$$\frac{dg}{dt} = -v_0 \left( \exp\left( \frac{A(v,g)}{T} \right) - \exp\left( \frac{B(v,g)}{T} \right) \right) - n(g), \qquad (2)$$

where $A(v,g) \propto -vg^3$ and $B(v,g) \propto vg^3$. $n(g) \in N_{[g_{\min},g_{\max}]}(g)$, i.e. it belongs to the normal cone to the gap domain $[g_{\min}, g_{\max}]$, according to which when $g = g_{\min}$ the gap time derivative can only be zero or positive, and when $g = g_{\max}$ the gap time derivative can only be zero or negative [10]. This equation is problematic from the viewpoint of its numerical integration for two main aspects. The first problem is its stiffness. The trend of (2) can be approximated by a hyperbolic sine. This entails that the derivative of the gap passes from almost negligible values for low voltages to very high values as soon as a threshold condition is exceeded. This behavior is compatible with the condition of stiffness, for which explicit numerical integration methods (such as the explicit Runge–Kutta methods) or implicit methods of higher order than the second become unstable [11]. The second aspect to deal with is the presence of the normal cone. When the gap distance reaches one of the extremes of its domain, the simulator must intervene by verifying the value of the derivative as a function of the current state of the circuit, canceling it in case the gap distance would exceed its domain.

## 3    Simulation Environment

The proposed simulator has an object-oriented structure, in which the circuit and its components are defined as classes providing properties and methods necessary for their simulation. In the proposed simulation environment objects are instantiated within scripts, which allow for the proper dimensioning of the parameters, the configuration and launch of the simulations, the collection and plotting of the results. Compared to a commercial simulator, the developed work environment allows the user to intervene on the sizing of the circuit, the parameterization of the simulation and the analysis of the results with dedicated functions, oriented towards a low-level control. Specifically, the simulator: implements algorithms for sizing circuit components based on a static analysis; allows to carry out transient simulations with full control over all the parameters of the numerical integration method, also including the possibility of performing parametric simulations, which can be extended to each element of the circuit; has an event state machine for the correct simulation of the memristor dynamics; performs intermediate and post-simulation analysis and plotting of the simulated signals. The structure of the simulator is designed in a modular perspective, so as to allow the user to add or remove features according to the needs. The result is a versatile tool, which can be adapted to simulate different circuit topologies with respect to those presented in this paper, with a limited effort by the user. Due to lack of space, in the following we limit the presentation to an in-depth analysis of the solutions developed to manage the chief critical aspects related

to the Stanford memristor model, introduced in Sect. 2. For a further insight on the simulator, the reader can refer to the source code available in [12].

### 3.1  Stanford Memristor Class

The `stanfordmemristor` class contains the information required to simulate a Stanford model-described memristor. The model parameters are defined as class properties, while the equations defining the device behavior have been implemented as class methods. The gap distance time derivative is computed by the `gaptimederivative` method. The method takes as inputs the memristor voltage, the gap distance value, the memristor temperature value and a special variable `gregion`. This variable defines the behavior of the derivative at the extremes of the domain of the gap. The management of this variable is entrusted to the circuit class containing the memristor, according to an approach described in Sect. 3.2.

### 3.2  Circuit Class

Each circuit is defined by a class, according to a template which guarantees their interchangeability within the scripts. The components of the circuit are defined through their respective classes as properties, which must be instantiated at the time of construction of the object. The template imposes that the methods of each class are defined with the same name, the same input variables and the same output variables. What distinguishes one circuit from the other are only the operations carried out within the methods. This structure allows to analyze different circuits using scripts by modifying only the name of the constructor of the object, keeping all function calls unchanged. The transient simulation of the circuit is carried out through the `transient` method. To handle the stiffness of the gap distance time derivative, the method performs the numerical integration of the circuit using the MATLAB function `ode15s`. `ode15s` is a multistep solver which can use the backword differentiation formulas (corresponding to the SPICE Gear method), and belongs to the set of stable solvers with respect to stiff integration problems [13]. The `transient` method also handles the variable `gregion` described in Sect. 3.1. The method defines for the `ode15s` function four events that interrupt its execution before the configured stop time is reached. Events are crossings of the gap distance at the following values:

1. $g_{\min}$ with negative gradient;
2. $g_{\min} + \varepsilon$ with positive gradient;
3. $g_{\max} - \varepsilon$ with negative gradient;
4. $g_{\max}$ with positive gradient.

$\varepsilon$ is a multiple of the machine epsilon corresponding to the minimum variation of the gap distance detectable by the function `ode15s` starting from $g_{\min}$ and $g_{\max}$. The method initializes the variable `gregion` according to the initial value of the gap distance (`"gmin"` if $g \leq g_{\min} + \varepsilon$, `"gmax"` if $g \geq g_{\max} - \varepsilon$, or an empty

string), after which it calls `ode15s` . When exiting the function, the method checks if an event has occurred, updates `gregion` and restarts the simulation from the moment of the event. With reference to the list defined above, events 1 and 4 correspond to the achievement of one of the extremes of the domain of the gap distance. As a result, `gregion` is updated to the values `"gmin"` or `"gmax"`, respectively, forcing the vanishing of the derivative in the corresponding direction out of the domain. Conversely, events 2 and 3 indicate that the gap distance has moved into the interior of the domain, thus allowing positive and negative values of its derivative.

## 4    Experiments and Results

To check the previously described simulator properties, the case of a memristor in series with a linear resistor has been considered. Figure 2 shows the time simulation of the sequence of two pairs of excitation pulses (the first negative and the second positive), assuming an initial value of the gap distance within its domain. The resistor is sized to ensure that a pulse carries the gap distance from one domain end to the other. From the red line in Fig. 2b it is possible to appreciate the stiffness of the equation which describes the time derivative of the gap distance. In fact, the derivative reaches extremely high peaks in correspondence with the excitation voltage level transition. Together with this aspect, we also appreciate the ability of the simulator to manage the dynamics of the gap distance in correspondence with the extremes of its domain. Notwithstanding the high values of its derivative, upon reaching the upper or lower extreme, the gap distance enters a saturation condition, from which it exits only in correspondence with the next inversely biased pulse. These results confirm the correct behavior of the simulator, proposed as an alternative to others based on Verilog-A models
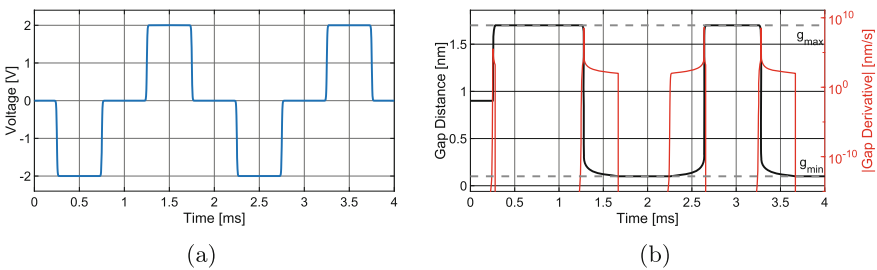


(a)

(b)

**Fig. 2.** Transient simulation of the series between the memristor and a linear resistor. **a** Excitation voltage applied to the series. **b** Black line: gap distance; red line: gap distance time derivative absolute value; grey dashed lines: gap distance domain boundaries.

## 5    Conclusion

In this paper we presented a circuit simulator based on MATLAB, developed in a dedicated way for the analysis of simple circuit topologies given by the series between a memristor described by the Stanford model and a current limiter. The simulator was developed in the context of a research activity oriented towards the study of memristor programming techniques. In this activity the analysis of transient simulations is combined with numerical investigations based on nonlinear equations solvers. This entails the need for a simulation environment that goes beyond the classic SPICE simulators, also requiring the implementation of dedicated pre- and post-simulation processing features. The simulator was built taking into account critical aspects regarding the numerical integration of the Stanford model. Experiments demonstrated the ability of the simulator to properly manage the considered critical aspects, confirming its reliability. In its present form the simulator is suitable for the simulation of low-complexity programming circuits based on current limiting. Future activities envisage the extension of the simulator functionalities, aimed at its use for the simulation of more complex memristor-based circuits.

## References

1. Chua L (2018) Five non-volatile memristor enigmas solved. Appl Phys A 124:1–43
2. Soni K, Sahoo S (2022) A review on different memristor modeling and applications. In: 2022 international mobile and embedded technology conference (MECON). IEEE, pp 688–695
3. Kvatinsky S, Friedman EG, Kolodny A, Weiser UC (2012) TEAM: threshold adaptive memristor model. IEEE Trans Circuits Syst I: Regul Pap 60(1):211–221
4. Kvatinsky S, Ramadan M, Friedman EG, Kolodny A (2015) VTEAM: a general model for voltage-controlled memristors. IEEE Trans Circuits Syst II: Express Briefs 62(8):786–790
5. Prodromakis T, Peh BP, Papavassiliou C, Toumazou C (2011) A versatile memristor model with nonlinear dopant kinetics. IEEE Trans Electron Devices 58(9):3099–3105
6. Jiang Z, Wu Y, Yu S, Yang L, Song K, Karim Z, Wong HSP (2016) A compact model for metal-oxide resistive random access memory with experiment verification. IEEE Trans Electron Devices 63(5):1884–1892
7. Guan X, Yu S, Wong HSP (2012) A SPICE compact model of metal oxide resistive switching memory with variations. IEEE Electron Device Lett 33(10):1405–1407
8. Addabbo T, Moretti R (2023) Static analysis of current limiting techniques for accurate memristor programming. In: 2023 IEEE international symposium on circuits and systems (ISCAS). IEEE, pp 1–5
9. Jiang Z, Wong HSP (2014) Stanford university resistive-switching random access memory (RRAM) verilog-a model. https://nanohub.org/publications/19/1
10. Di Marco M, Forti M, Moretti R, Pancioni L, Innocenti G, Tesi A (2022) Convergence of a class of delayed neural networks with real memristor devices. Mathematics 10(14):2439
11. Dahlquist GG (1963) A special stability problem for linear multistep methods. BIT Numer Math 3(1):27–43

12. Moretti R (2023) Current limited memristor. https://github.com/RiccMore/CurrentLimitedMemristor
13. Shampine LF, Reichelt MW (1997) The MATLAB ODE suite. SIAM J Sci Comput 18(1):1–22

# Monitoring Hardware True Random Number Generators with Artificial Neural Networks: Problem Modeling and Training Dataset Generation

Tommaso Addabbo[1]([✉]), Gian Domenico Licciardo[2], Riccardo Moretti[1], Alfredo Rubino[2], Filippo Spinelli[1], Valerio Vignoli[1], and Paola Vitolo[2]

[1] Department of Information Engineering and Mathematics, University of Siena, 53100 Siena, Italy
tommaso.addabbo@unisi.it, riccardo.moretti@unisi.it, filippo.spinelli@unisi.it, valerio.vignoli@unisi.it

[2] Department of Industrial Engineering, University of Salerno, 84084 Fisciano, Italy
gdlicciardo@unisa.it, arubino@unisa.it, pvitolo@unisa.it

**Abstract.** We investigate the possible application of Artificial Neural Networks (ANNs) to monitor and test hardware True Random Number Generators (TRNGs). The preliminary results have been obtained developing original investigation and design tools, including the characterization of optimized statistical estimators, considering a class of TRNGs based on binary Markov Chains. Due to the large amount of data required during the design of the ANN we developed an original flexible Database Management Systems to manage and generate training and testing data samples organized in a Object Oriented Database.

**Keywords:** Hardware true random number generators · Entropy testing · Artificial neural networks

## 1 Introduction

Machine learning (ML) is a continuing evolving field that has achieved significant prominence in recent years. Nowadays, the application of Artificial Neural Networks (ANNs) to hardware security [1] and cryptography [2–4] is an emerging research topic [5,6]. In this brief, we investigate the possible application of ANNs to the assessment and monitoring of hardware True Random Bit Generators (TRBGs), targeting low/medium-complexity solutions suitable for hardware-based lightweight cryptography.

Previous research introduces machine learning classifiers focusing on specific statistical features in data sequences, in particular on highly sensitive statistical tests designed to validate the output of cryptographic True Random Number Generators (TRNGs) [7,8]. In contrast to the previous approaches, we focus on the direct estimation of the TRNG information source entropy, which involves

a more flexible technical framework that is well-suited for TRNGs monitoring, development and validation [9,10]. In detail, we target a test system suitable to detect static or dynamic degradation of the entropy with respect to a given threshold level.



**Fig. 1.** Flowchart describing the reference design flow followed in our research.

The reference design flow followed in our research is shown in Fig. 1. Once given the class of TRBGs under investigation, the statistical feature (e.g., in our case, the source entropy and the classifier specification and constraints (e.g., entropy threshold level, classification accuracy, speed, computational complexity, hardware resources and target device), an iterative design flow is followed, targeting the system specifications. The design flow shown in Fig. 1 involves several sub-problems to be addressed, including the design and characterization of specific statistical estimators, considering the class of TRBGs under investigation.

In this work we present in detail the research activity carried out to investigate the latter mentioned sub-problem, and the original design tools developed to perform the activity 6 in the flowchart (i.e., "6. Training and Test Dataset Creation"), considering a specific class of TRNGs based on Markov Chains.

## 2    Design and Characterization of Specific Statistical Estimators

The TRNGs considered in this work were based on Markov Chains, that are low-complexity stochastic models that can offer a variety of dynamical behaviors suitable to emulate different statistical defects affecting actual entropy sources [10]. In detail, a binary Markov Chain has been considered, defined by a 2 $\times$

2 parametric state transition probability matrix dependent on two parameters only [10]. On the other side, the statistical feature considered in this work is the Average Shannon Entropy (ASE) [10] that results particularly effective and based on low-complexity calculations [9].

Approximated estimations of the ASE can be obtained with parametric estimators, with accuracy that depends on the length of the observed random sequences, the length of the binary words setting the source alphabet and the source entropy [9]. As a result, the first problem that we investigated was to identify an adequate sizing of the ASE estimator, considering the need to label hundreds of thousands of sources for the creation of the training and testing dataset.

Accordingly, we investigated the impact on the ASE estimation relative error of the observed random sequence length, as a function of the nominal source entropy. As shown in Fig. 2, for random sequence lengths larger than $3 \cdot 10^6$ bits, the relative error is lower than 0.4% for binary words of 12bits (analysis performed evaluating 100 observations for each parametric point). For the same random sequence lengths and ASE levels greater than 0.3 bit/symb the relative error resulted lower than 0.25%. This means, in absolute terms, a worst case estimation error of 0.0025 bit/symb for the highest source entropy levels ($\approx 1$ bit/symb).



**Fig. 2.** The impact on the ASE estimation relative error of the observed random sequence length, as a function of the nominal source entropy. Error bars report the maximum and minimum occurrences (100 observations for each parametric point).

Once the random sequence length has been set to $3 \cdot 10^6$ bits, a similar analysis has been carried out considering the impact on the ASE estimation relative error of the binary words size ranging between 4 and 12 bits. Results have been reported in Fig. 3 (analysis performed evaluating 100 observations for each parametric point). For entropy levels greater than 0.3 bit/symb the

**Fig. 3.** Flowchart describing the reference design flow followed in our research. Error bars report the maximum and minimum occurrences (100 observations for each parametric point).

relative estimation error turned out to be lower than 0.5% considering sizes of the alphabet binary words greater than or equal to 10 bits.

Using the designed estimator, the authors could perform a complete characterization of the parametric Markov Chain model, obtaining a numerical relation between the source entropy and the two parameters of the binary Markov Chain.

## 3 Development of an Object-Oriented TRNGs Database Management System for ANN Training and Testing

Nowadays, Database (DB) Management Systems (DBMSs) are crucial for efficient data storage, recovery and management, enabling scalability and integration in distributed systems. Object Oriented Programming (OOP) supports modular, reusable and maintainable code, enhancing software design principles such as abstraction, inheritance and polymorphism. Joining DBs data persistency and OOP coding results in robust and scalable applications called Object Oriented Databases (OODs).

Due to the large amount of data required during the design of the TRBG monitoring ANN, we developed an original flexible DBMS to manage and generate training and testing data samples organized in a OOD. In the considered application an ANN has been designed to classify different kinds of images obtained by pre-processing random data (operation 4 in Fig. 1) [10]. Accordingly, the designed DB has been organized in three tables associated to the TRBG models, the pre-processed images and the estimated statistical features (the ASE, in this work) used to label the random sources, respectively. The latter two tables are related to the first table by a *1 to many* relationship, as shown in Fig. 4.

Once the database architecture has been defined, adopting the OOP approach the class objects have been developed in *Matlab* environment. Table records are managed by specific methods (i.e., *getters* and *setters*) providing the functionality of standard queries.



**Fig. 4.** Entity-relationship model describing the DB architecture.

## 4    Results & Conclusion

Using the developed design tools and following the flowchart shown in Fig. 1 the authors could first design an optimized ANN built and trained in Python language, using TensorFlow and Qkeras frameworks. We also performed a preliminary investigation of the hardware complexity involved by the proposed architecture using high level synthesis tools in the Xilinx and Matlab environments targeting a Xilinx Zynq Ultrascale Multiprocessor System on Chip (See Fig. 5, also reporting information about the hardware consumption and timing results). Referring to an entropy threshold level of 0.98 bit/symb, the relative classification error of the classifier resulted not greater than 0.3%.



**Fig. 5.** Preliminary investigation of the hardware complexity involved by the proposed architecture using high level synthesis tools in the Xilinx and matlab environments targeting a Xilinx Zynq ultrascale multiprocessor system on chip.

In conclusion, we have investigated the possible application of Artificial Neural Networks to monitor and test entropy sources. The preliminary results have been obtained developing original investigation and design tools, including the characterization of optimized statistical estimators, considering the class of TRBGs based on binary Markov Chains. Also, due to the large amount of data required during the design of the TRBG monitoring ANN (i.e., in the activity "6. Training and Test Dataset Creation" reported in the flowchard in Fig. 1), we developed an original flexible DBMS to manage and generate training and testing data samples organized in a OOD.

# References

1. Duan S, Li Z, Luo Y, Sun M, Wang W, Lin XS, Xu X (2021) Machine learning in hardware security. Springer International Publishing, Cham, pp 111–146
2. Brunetta C, Picazo-Sanchez P (2022) Modelling cryptographic distinguishers using machine learning. J Cryptogr Eng 12(2):123–135
3. Truong ND, Haw JY, Assad SM, Lam PK, Kavehei O (2019) Machine learning cryptanalysis of a quantum random number generator. IEEE Trans Inf Forens Secur 14(2):403–414
4. Meraouche I, Dutta S, Tan H, Sakurai K (2021) Neural networks-based cryptography: a survey. IEEE Access 9:124727–124740
5. Feng Y, Hao L (2020) Testing randomness using artificial neural network. IEEE Access 8:163685–163693
6. Zhu S, Ma Y, Li X, Yang J, Lin J, Jing J (2020) On the analysis and improvement of min-entropy estimation on time-varying data. IEEE Trans Inf Forens Secur 15:1696–1708
7. Nagy I, Suciu A (2021) Randomness testing with neural networks. In: 2021 IEEE 17th international conference on intelligent computer communication and processing (ICCP), pp 431–436
8. Li H, Zhang J, Li Z, Liu J, Wang Y (2023) Improvement of min-entropy evaluation based on pruning and quantized deep neural network. IEEE Trans Inf Forens Secur 18:1410–1420
9. Addabbo T, Fort A, Moretti R, Mugnaini M, Papini D, Vignoli V (2022) A stochastic algorithm to design min-entropy tuning controllers for true random number generators. IEEE Trans Circuits Syst I: Regul Pap 69(5):2084–2094
10. Spinelli F, Moretti R, Addabbo T, Vitolo P, Licciardo GD (2023) Low-complexity machine learning architecture for hardware-aware true random number generators assessment and continuous monitoring. In: 2023 18th conference on Ph.D. research in microelectronics and electronics (PRIME), pp 221–224

# Highly-Efficient Galois Counter Mode Symmetric Encryption Core for the Space Data Link Security Protocol

Luca Crocetti[1(✉)], Francesco Falaschi[2], Sergio Saponara[1], and Luca Fanucci[1]

[1] Department of Information Engineering, University of Pisa, Via G. Caruso 16, 56122 Pisa, Italy
{luca.crocetti,sergio.saponara,luca.fanucci}@unipi.it
[2] IngeniArs S.r.l, Via Ponte a Piglieri 8, 56121 Pisa, Italy
francesco.falaschi@ingeniars.com

**Abstract.** In the last decades, the space sector has been the subject of significant technological improvements and investments from both government agencies and private companies, generating an increase in data rates and volumes of exchanged data. Accordingly, the security threats and the number of documented cyberattacks have grown. In order to meet the requirements of space applications, the Consultative Committee for Space Data Systems (CCSDS) has issued and maintained a series of reports and recommendations over the years, including a set of standards aimed at efficiently exploiting the communication channels. In this work, we present the implementation of an Advanced Encryption Standard – Galois/Counter Mode (AES-GCM) core on space-grade FPGAs, that is compliant with the latest CCSDS security standards and outperforms the state-of-the-art in terms of resource efficiency.

**Keywords:** Space security · Cybersecurity · Hardware security · CCSDS · 352.0-B-2 · 355.0-B-2 · SDLS · Protocol · AES · GCM · CTR · Confidentiality · Integrity · Authentication · Authenticated encryption · Space-grade · FPGA

## 1 Introduction

The advent of the *New Space Economy* [1] is favoring the ambitious intervention of private investors in the space sector, which for a long time enjoyed predominant institutional support, having the effect of pushing to the limits the requirements related to both data rate speed and security of the communication links [2]. In 2020, research by Morgan Stanley [3] estimated that the global space industry could generate revenues of more than 1 trillion by 2040 from the current 350 billion. In this scenario of rapid and significant growth, cybersecurity has become a central issue in the debate on the development of space applications [4,5] indeed, the number of cyberattacks has significantly grown [6].

The Consultative Committee for Space Data Systems (CCSDS) has issued and maintained a series of reports and standards to meet the requirements of

space applications over the years. These requirements are applicable to both the satellite-to-ground (*downlink*) and ground-to-satellite (*uplink*), actually showing an average data rate per satellite from 2.04 Gbps to 22.74 Gbps [7]. Remote controls sent through the *uplink* as *telecommand* (TC) packets need protection against unauthorized tampering; scientific payload data transmitted through the *downlink* as *telemetry* (TM) packets require confidentiality. The CCSDS reports 350.0-G [8] and 350.1-G [9] provide an overview of possible threats and introduce multiple options for implementing security mechanisms. In particular, [8] presents the space link reference model and compares it against the OSI model (Fig. 1).



**Fig. 1.** CCSDS space link reference model with security options [8]

According to [10], we focused our work on the development of a hardware accelerator to address the space security requirements at the Data Link layer, with the goal to implement a solution that is both high-speed and efficient in terms of supported data rate per used logic resources.

## 2 Space Data Link Security

The corresponding layer in the CCSDS space link reference model is the Space Data Link Security (SDLS) layer, which is regulated by the CCSDS standards 352.0-B [11] and 355.0-B [12]. The SDLS introduces a frame format integrating a Security Header and an optional Security Trailer, respectively, before and after the payload (i.e. Frame Data in Fig. 2), and to be set accordingly to the employed cryptographic algorithm.

| Transfer Frame without SDLS | Space Link Protocol Frame Headers | Frame Data | | | Space Link Protocol Frame Trailers |
|---|---|---|---|---|---|
| Transfer Frame with SDLS | Space Link Protocol Frame Headers | Security Header | Frame Data (secured by authentication and/or encryption) | Security Trailer | Space Link Protocol Frame Trailers |

**Fig. 2.** Transfer Frame structure with/without SDLS fields [12]

Between the algorithms specified by [11,12], this work presents the implementation of the Advanced Encryption Standard – Galois/Counter Mode (AES-GCM) [13] because it offers a comprehensive solution for addressing all the security requirements for space communications, providing at the same time confidentiality, data integrity, and data authentication (i.e. authenticated encryption). The implementation of the other algorithms will presented in other works.

## 3   Implementation of a Hardware AES-GCM Core

The AES-GCM algorithm is a mode of operation of the AES cipher [14] that exploits the AES—Counter Mode (AES-CTR) algorithm [15] for data encryption and decryption, and Galois multipliers for the MAC generation and verification. We designed the core in SystemVerilog, by exploiting the AES engine proposed in [16–18] for high-performance computing applications. In particular, we modified and optimized the engine for the implementation on space-grade FPGAs. The main modification concerned the S-box, the substitution circuit of the substitution-permutation network on which the AES is based [14]. Since the release of the AES standard [14], the optimization of this circuit gained particular interest from researchers, because it is the most expensive in terms of area and critical delay among all the operations of the AES algorithm. All the main works typically converge on solutions based on field arithmetic [19,20] instead of solutions based on Look-up Tables (LUTs), because they lead to more compact and more efficient implementations on standard-cell technologies. Anyway, in recent years the trend inverted for FPGA implementations thanks to technological advances on these devices [21], and we opted for this last.
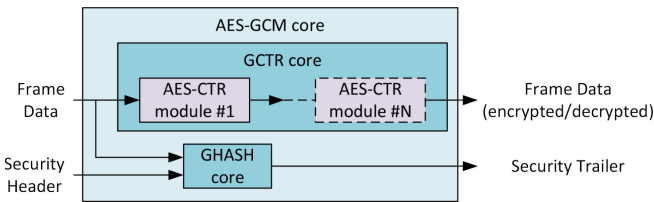


**Fig. 3.** Outline of the implemented AES-GCM core

The resulting hardware accelerator is shown in Fig. 3, and it is mainly composed of a GCTR core for encryption/decryption of data and a GHASH unit

dedicated to the generation of the MAC. The GCTR core is built on top of an underlying AES-CTR module, that can be replicated in cascade in order to form a pipeline for maximizing the throughput without affecting the maximum frequency. The GHASH module is based on a Galois multiplier and works in parallel to the GCTR core. In addition, the core can be customized through other synthesis parameters to select the supported key size (i.e. only 128 bits, only 256 bits, or both), varying the utilization of logic resources accordingly.

## 4   Results and Comparison with the State-of-art

We implemented the AES-GCM core on different space-grade FPGAs, in particular on a Kintex Ultrascale KU060 and a Virtex 4QV by Xilinx/AMD, and an RTG4 by Microsemi. The implementation results shown in Table 1 refer to the proposed AES-GCM core configured to support both 128-bit and 256-bit keys and instantiate a single AES-CTR module.

**Table 1.** Implementation results on space-qualified FPGAs

| Space-qualified FPGA | Frequency | Resource utilization | Throughput | Power consumption |
|---|---|---|---|---|
| KU060 | 224.37 MHz | 665 CLBs | 2.871/2.051 Gbps | 239 mW |
| RTG4 | 71.39 MHz | 5967 LEs | 913.8/652.7 Mbps | 161 mW |
| Virtex 4QV | 139.7 MHz | 2832 Slices | 1.788/1.277 Gbps | – |

CLBs (Configurable Logic Blocks), LEs (Logic Elements), and Slices in Table 1 are the programmable logic resources of the corresponding FPGA device; the column Frequency indicates the maximum frequency ($f_{max}$), and the throughput is indicated, respectively, for the encryption (or decryption) process with 128-bit/256-bit keys. Being 128 bits the length of data blocks, the throughput ($THR$) can be calculated as:

$$THR = \frac{128 \cdot f_{max}}{CPD} \tag{1}$$

$CPD$ in Eq. 1 stands for *Clock per Data* and indicates the number of clock cycles required to generate an output data block. Such value depends on the number $N$ of instantiated AES-CTR modules inside the GCTR core according to the formula:

$$CPD = \left\lceil \frac{L}{N} \right\rceil \tag{2}$$

$L$ in Eq. 2 is the latency expressed in clock cycles (cc) and corresponds to 10/14 cc, respectively, for 128-bit/256-bit keys. The maximum number of cascadable internal AES-CTR modules is $N = 10$ (only 128-bit keys), or $N = 14$ (only 256-bit keys or both 128-bit/256-bit keys), maximizing the throughput in each case ($CPD = 1$). For instance, the maximum achievable throughput for the implementation on the FPGA KU060 is 28.71 Gbps, which gives the possibility

to support the maximum average data rate per satellite (22.07 Gbps).

The state-of-the-art related to the implementation on space-grade FPGAs of AES-GCM cores for the CCSDS standards [11,12] offers only one (recent) work [22]. The authors of [22] use a Virtex 4QV FPGA and propose a pipelined architecture with 7 stages, similar to ours when configured with $N = 7$ and only 256-bit keys. Therefore, we configured our AES-GCM core accordingly to make a fair comparison. The results are reported in Table 2 highlighting the better resource efficiency of our proposal. In addition, our solution is able to double the throughput without doubling the resources. Indeed, if we configure our core with $N = 14$, only the number of the AES-CTR modules shown in Fig. 3 is doubled (i.e. only the GCTR core is doubled), while the GHASH core that works in parallel to the GCTR core is not modified. As the cascading of the AES-CTR modules does not affect the maximum frequency, the resulting (doubled) throughput of 17.88 Gbps can be calculated from Eqs. 1 and 2 by using $L = 14$ (due to the usage of 256-bit keys), $N = 14$ and the maximum frequency of 139.7 MHz reported in Table 2.

**Table 2.** Comparison with the state-of-the-art

|  | Fmax | Resource Utilization | Throughput | Resources Efficiency |
|---|---|---|---|---|
| [22] | 139.725 MHz | 16393 Slices | 8.942 Gbps | 0.546 Mbps/Slice |
| This work | 139.7 MHz | 15488 Slices | 8.941 Gbps | 0.577 Mbps/Slice |

## 5    Conclusions

In this work, we present the implementation of a configurable AES-GCM core in SystemVerilog to protect space communications by satisfying the security requirements specified by the CCSDS Space Data Link Security protocol. The proposed core can be configured to find the most efficient solution in terms of throughput per resource utilization, and it outperforms the state-of-the-art in terms of efficiency, being able also to support the average data rates of modern space applications with throughputs up to 28.71 Gbps on space-grade FPGAs.

Future works will include the evaluation of the proposed core for Automotive Ethernet security applications [17] and the analysis of Side-Channel vulnerabilities [23].

# References

1. Economic Cooperation and Development (OECD): The Space Economy at a Glance (2014). https://doi.org/10.1787/9789264217294-en. Accessed June 2023
2. Orlova A, Nogueira R, Chimenti P (2020) The present and future of the space sector: a business ecosystem approach. Space policy, vol 52, art. 101374
3. Stanley M (2020) Space: investing in the final frontier. Morgan Stanley (2020)
4. Mulder LCCP, Siegel JT (2021) The future of security in space: a thirty-year US strategy
5. Naja G, Mathieu C (2015) Space and security in Europe. Handbook of space security. Springer, pp 371-3-83
6. Manulis M, Bridges CP, Harrison R, Sekar V, Davis A (2021) Cyber security in new space: analysis of threats, key enabling technologies and challenges. Int J Inf Secur 20(3):287–311. Springer
7. Del Portillo I, Cameron BG, Crawley EF (2019) A technical comparison of three low earth orbit satellite constellation systems to provide global broadband. Acta Astronaut 159:123–135
8. CCSDS: The application of security to CCSDS protocols. Informational report CCSDS 350.0-G-3 (2019)
9. CCSDS: security threats against space missions. Informational report CCSDS 350.1-G-3 (2022)
10. Baldanzi L, Crocetti L, Di Matteo S, Fanucci L, Saponara S, Hameau P (2019) Crypto accelerators for power-efficient and real-time on-chip implementation of secure algorithms. In: 26th IEEE international conference on electronics, circuits and systems (ICECS). IEEE, pp 775–778
11. CCSDS: CCSDS cryptographic algorithms. Recommended standard CCSDS 352.0-B-2 (2019)
12. CCSDS: space data link security protocol. Recommended standard CCSDS 355.0-B-2 (2022)
13. NIST: recommendation for block cipher modes of operation: Galois/Counter Mode (GCM) and GMAC. Special Publication (SP) 800-38D (2007)
14. NIST: advanced encryption standard (AES). Federal information processing standards (FIPS) publication 197 (2001)
15. NIST: recommendation for block cipher modes of operation: methods and techniques. Special Publication (SP) 800-38A (2001)
16. Nannipieri P, Di Matteo S, Baldanzi L, Crocetti L, Zulberti L, Saponara S, Fanucci L (2022) VLSI design of advanced-features AES cryptoprocessor in the framework of the European processor initiative. IEEE Trans Very Large Scale Integr (VLSI) Syst 30(2):177–186. https://ieeexplore.ieee.org/document/9631958
17. Carnevale B, Falaschi F, Crocetti L, Hunjan H, Bisase S, Fanucci L (2015) An implementation of the 802.1AE MAC security standard for in-car networks. In: IEEE 2nd world forum on internet of things (WF-IoT). IEEE, pp 24–28
18. Nannipieri P, Crocetti L, Di Matteo S, Fanucci L, Saponara S (2023) Hardware design of an advanced-feature cryptographic tile within the European processor initiative. IEEE Trans Comput
19. Singha TB, Palathinkal RP, Ahamed SR (2023) Securing AES designs against power analysis attacks: a survey. IEEE Internet Things J
20. Mozaffari-Kermani M, Reyhani-Masoleh A (2011) Efficient and high-performance parallel hardware architectures for the AES-GCM. IEEE Trans Comput 61(8):1165–1178

21. Khairallah M, Chattopadhyay A, Peyrin T (2017) Looting the LUTs: FPGA optimization of AES and AES-like ciphers for authenticated encryption. In: 18th international conference on cryptology in India (INDOCRYPT). Springer, pp 282–301
22. Muraleedharan D, Daniel SK (2020) An efficient IP core of consultative committee for space data systems (CCSDS) recommended authenticated cryptographic algorithm. In 2020 8th international symposium on digital forensics and security (ISDFS). IEEE, pp 1–6
23. Crocetti L, Baldanzi L, Bertolucci M, Sarti L, Carnevale B, Fanucci L (2019) A simulated approach to evaluate side-channel attack countermeasures for the advanced encryption standard. Integration 68:80–86

# A Pedestrian Detection Method Based on YOLOv7 Model

Binglin Li[1], Qinghe Zheng[1(✉)], Xinyu Tian[1], Abdussalam Elhanashi[2], and Sergio Saponara[2]

[1] School of Intelligence Engineering, Shandong Management University, Jinan 250357, China
15005414319@163.com

[2] Department of Information Engineering, University of Pisa, 56122 Pisa, Italy

**Abstract.** Pedestrian detection has extensive applications in computer vision, such as intelligent transportation and security surveillance. YOLOv7 is an object detection model that achieves real-time object detection by dividing the image into grids and predicting bounding boxes and classes for each grid. This study explores and optimizes pedestrian detection based on the YOLOv7 model. Firstly, a large-scale pedestrian dataset is used to train and fine-tune the model to improve detection accuracy and robustness. The YOLOv7 model demonstrates powerful pedestrian detection performance, achieving a high average precision (AP) of 0.92 and a recall rate of 0.85. Experimental results indicate that the pedestrian detection method based on YOLOv7 achieves good results in terms of accuracy and speed, and it has high practical value and application prospects.

**Keywords:** Object detection · Pedestrian detection · YOLOv7

## 1 Introduction

Object detection is a task in the field of computer vision that aims to enable computers to automatically identify and locate objects in images, providing a foundation for many applications such as intelligent surveillance, autonomous driving, and medical image analysis [1, 2]. The significance of object detection lies in enhancing the computer's understanding and classification of complex scenes, thereby improving human-computer interaction and reducing human workload.

Pedestrian detection, as an important research direction in computer vision, has received wide attention. With the rapid development of intelligent transportation and security surveillance, accurate and efficient pedestrian detection algorithms are crucial for real-time object recognition and tracking [3, 4]. However, pedestrian detection still faces a series of challenges due to the complexity of factors such as pedestrian poses, scales, and occlusions. Pedestrians may appear in various poses, leading to changes in their appearance and shape in images, making them difficult to detect and recognize. Pedestrians can have different scales and proportions, requiring the use of multi-scale detection algorithms to detect pedestrians of different sizes. Additionally, consideration needs to be given to different image resolutions for pedestrian detection in different

environments. Pedestrians may be occluded by other pedestrians or objects, making them difficult to detect. In such cases [5, 6], occlusion detection algorithms [7, 8] and corresponding processing techniques [9, 10] are needed to detect pedestrians. In recent years, object detection algorithms based on deep learning have made significant breakthroughs, and the YOLO (You Only Look Once) model has attracted attention due to its real-time performance and accuracy. The YOLO model transforms the object detection task into a regression problem by dividing the image into grids and predicting each grid's targets, significantly improving detection speed. YOLOv7, as an improved version of the YOLO series, further optimizes accuracy and efficiency, making it a powerful tool for pedestrian detection.

This paper aims to conduct research on pedestrian detection based on the YOLOv7 model. We utilize a large-scale pedestrian dataset for model training and fine-tuning to improve the accuracy and robustness of pedestrian detection. Additionally, we propose some improvement methods to enhance detection performance in challenging situations encountered in pedestrian detection.

## 2    YOLOv7 for Pedestrian Detection

YOLOv7 is an extended version of the YOLO (You Only Look Once) model, aimed at achieving real-time object detection without sacrificing accuracy. Similar to YOLOv4, YOLOv7 is based on the Darknet neural network framework, which includes the backbone network, intermediate network, and head network. In terms of pedestrian detection, YOLOv7 demonstrates excellent results due to its real-time performance and high detection accuracy. By utilizing pre-trained weights and transfer learning, the YOLOv7 model can be fine-tuned to detect pedestrians in different scenarios.

### 2.1    YOLO V7 Model Structure

YOLO v7 network model as shown in Fig. 1. It is mainly includes input (Input), backbone network (Backbone), head (Head) three parts, input layer to input picture resize of 640x640 size, input into the backbone network, Backbone module consists of several CBS convolutional layers, ELAN convolutional layers and MPConv convolutional layers, of which SBC consists of convolutional layers, batch normalization layers (Batch Normalization (BN), SiLU activation function, used to extract image features at different scales. The CBS convolutional layer is an improvement on Batch Normalization. Batch normalization refers to the standardization of the data in each batch, while the CBS convolutional layer is to standardize and scale between different feature channels, so that the features between different channels are more independent and balanced, thereby improving the generalization ability and robustness of the model. ELAN convolutional layer is a convolutional layer that introduces attention mechanism, which weights the features of different channels by calculating the similarity between different channels in the feature map, so that important features can be better captured and utilized, thereby improving the accuracy and efficiency of the model. The MPConv convolutional layer is a mixed-precision convolutional layer that accelerates the training and inference process of the model by using half-precision floating-point numbers to reduce the amount of computation and storage space for convolution computation.

**Fig. 1.** Specific structure of YOLOv7.

## 2.2 Optimization Method

Adaptive Moment Estimation (Adam) is a commonly used optimization algorithm for training deep learning models. The specific principle of Adam algorithm is as follows:

1. For the current parameter value $\theta_t$, calculate its corresponding gradient using batch stochastic gradient descent (mini-batch SGD), $g_t$.
2. Calculate first-order moment estimation: maintain an exponentially weighted moving average variable $m_t$, which is used to estimate the first-order moment, i.e. mean, of the current gradient. Specifically, $m_t$ is updated as shown by

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_t \tag{1}$$

3. Calculation of the second-order moment estimate: Similarly, a exponentially weighted moving average variable $v_t$ is maintained to estimate the second-order moment (variance) of the current gradient. The update equation for $v_t$ is given by

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)g_t^2 \tag{2}$$

where $\beta_2$ is a hyperparameter ranging between 0 and 1, commonly set to 0.999.

4. Calculation of the adaptive learning rate: Based on the estimates of the first and second-order moments, the adaptive learning rate for each parameter is computed. Specifically, first, the bias-corrected values of $\widehat{m}_t$ and $\hat{v}_t$ are calculated as

$$\widehat{m}_t = \frac{m_t}{1 - \beta_1^t} \tag{3}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t} \tag{4}$$

Then, the adaptive learning rate $\alpha_t$ for each parameter is computed as

$$\alpha_t = \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \tag{5}$$

where $\eta$ represents the initial learning rate, and $\epsilon$ is a very small constant used to avoid division by zero errors caused by extremely small denominators.

5. Parameter Update: Finally, the parameters are updated according to

$$\theta_{t+1} = \theta_t - \alpha_t \cdot \widehat{m}_t \tag{6}$$

The Adam combines momentum and adaptive learning rate to adapt to update frequencies and variations. Its basic idea is to adjust the learning rate for each parameter at each time step by estimating the first and second-order moments of the gradient.

## 3 Experiments and Evaluation

### 3.1 Experimental Setup

To validate the effectiveness of the YOLOv7 model used in this paper, we conducted experiments using a part of COCO dataset for training (2000 images for the training set, 200 for the validation set, and 200 for the test set). We selected pedestrians in different backgrounds as the research objects. The model was implemented using the PyTorch framework and trained and inferred efficiently using GPU acceleration with CUDA and cuDNN libraries. The experiments were conducted on a Windows 11 system with a batch size of 4 and 60 epochs. The experimental platform is consisted of AMD 6800HS CPU and 16GB of memory. The model was trained and predicted on an RTX3060 LAPTOP GPU with 6GB of VRAM. In the testing phase, we evaluated the trained model on a separate test set and recorded the detection results.

### 3.2 Evaluation Metrics

To evaluate the performance of the pedestrian detection method, we used commonly used evaluation metrics. These metrics include accuracy, recall, and average precision. Accuracy measures the overall correctness of detected pedestrian instances, while recall quantifies the method's ability to identify true positives. Average precision (AP) provides a comprehensive measure of accuracy at different recall levels, thus evaluating the model's performance comprehensively.

To evaluate the performance of our pedestrian detection method, we adopted a standard evaluation protocol. We computed the precision-recall curve, as shown in Fig. 2. Then we reported the mean average precision (mAP) at different IoU (intersection over union) thresholds in Fig. 3.
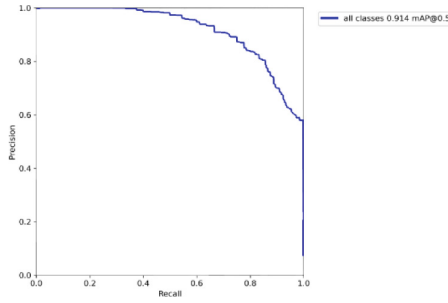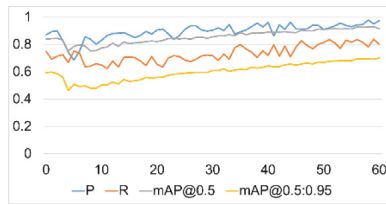
**Fig. 2.** Precision-recall curve.



**Fig. 3.** The mAP at different IoU in each epoch.

### 3.3 Experimental Results

We present the experimental results of the YOLOv7 model on the pedestrian detection task. The model's performance is evaluated based on metrics such as accuracy, recall, and average precision (AP).

During the training phase, with only a partial COCO dataset and training for 60 iterations, the YOLOv7 model achieved an AUC of 0.914, indicating that it successfully learned to detect pedestrians from the training dataset. In the testing phase, the YOLOv7 model demonstrated strong performance in pedestrian detection. It achieved a high average precision (AP) of 0.92, indicating its ability to accurately locate pedestrians in various real-world scenarios. The model also showed a recall rate of 0.85, indicating its ability to detect most true positive pedestrian instances. The loss values of the training set (top three plots) and the test set (bottom three plots) are shown in Fig. 4. In the experiments, we evaluated the impact of different image sizes on the speed and quality of object detection. As shown in Fig. 5, larger image sizes mean the ability to detect smaller objects with higher precision, but it also increases the processing time.

The experiments are shown in Fig. 6, were conducted using images taken from four different focal lengths of the Galaxy S21 Ultra camera. Figures (a–e) show the inference results with an input size of 680 × 480, while Figures (f–j) show the results with an input size of 2560 × 1920. Figure (a) and Figure (f) were captured with a 13 mm wide-angle lens, and the model successfully detected pedestrians in the distorted regions caused by the wide-angle lens. Figures (b–c), as well as Figures (d–e), compare the results of different focal lengths of the same scene. Figure (b) has a focal length of 72 mm, Figure (c) has a focal length of 230 mm, Figure (d) has a focal length of 72 mm, and Figure (e) has a focal length of 24 mm. The model performs well under various focal lengths, even

in situations with poor image quality or average lighting conditions. Comparing Figures (a–e) with Figures (f–j), it can be observed that increasing the input image size allows for better detection of pedestrians that are farther away (meaning smaller in the image), but it may result in missed detections of pedestrians closer to the camera. Overall, the method performs well under different focal lengths and various scenes.
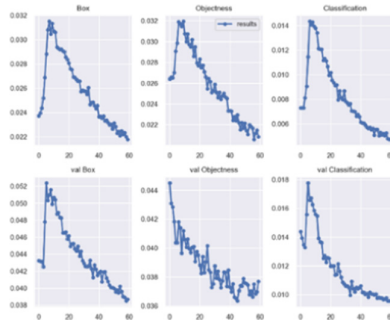


**Fig. 4.** The loss values of the training set and the test set.
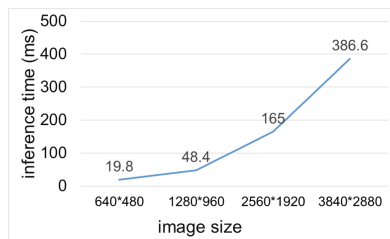


**Fig. 5.** The inference time of different image size.



**Fig. 6.** Experimental results. **a** scene 1, **b** scene 2, **c** scene 3, **d** scene 4, **e** scene 5, **f** scene 6, **g** scene 7, **h** scene 8, **i** scene 9, **j** scene 10

## 4  Conclusion

Pedestrian detection is an important task in the field of computer vision. Its significance lies in helping computers identify pedestrians in images or videos and perform operations such as tracking and counting. It can be applied to various practical scenarios. For example, in pedestrian safety, it can improve pedestrian safety in urban traffic. In video surveillance, pedestrian detection can help monitor densely populated areas, identify abnormal behaviors, and more. In this paper, we deployed the YOLOv7 object detection model on a computer to achieve pedestrian detection and analyzed the entire process from training to inference. In summary, our experimental results demonstrate the effectiveness and potential of the YOLOv7 model in pedestrian detection. The method achieves high accuracy and real-time performance, making it valuable in various applications such as surveillance systems and autonomous driving.

## References

1. Jiao L, Zhang F, Liu F et al (2019) A survey of deep learning-based object detection. IEEE access 7:128837–128868
2. Ouyang W, Wang X (2013) Joint deep learning for pedestrian detection. In: Proceedings of the IEEE international conference on computer vision, pp 2056–2063
3. Zheng Q, Yang M, Tian X et al (2020) A full stage data augmentation method in deep convolutional neural network for natural image classification. Discrete Dyn Nat Soc
4. Zhang Q, Yang M, Zheng Q, Zhang X (2017) Segmentation of hand gesture based on dark channel prior in projector-camera system. In: IEEE/CIC ICCC, Qingdao, China, pp 1–6
5. Li J et al (2019) Dynamic hand gesture recognition using multi-direction 3D convolutional neural networks. Eng Lett 27(3):490–500
6. Zhao L, Li S (2020) Object detection algorithm based on improved YOLOv3. Electronics 9(3):537
7. Zheng Q, Zhao P, Wang H, Elhanashi A, Saponara S (2022) Fine-grained modulation classification using multi-scale radio transformer with dual-channel representation. IEEE Commun Lett 26(6):1298–1302
8. Jiang N, Fu F, Zuo H, Zheng X, Zheng Q (2020) A municipal PM2.5 forecasting method based on random forest and WRF model. Eng Lett 28(2):312–321
9. Pathak AR, Pandey M, Rautaray S (2018) Application of deep learning for object detection. Procedia Comput Sci 132:1706–1717
10. Zheng Q, Yang M, Tian X, Wang X, Wang D (2020) Rethinking the role of activation functions in deep convolutional neural networks for image classification. Eng Lett 28(1):80–92

# Dynamic Capture Algorithm Based on Visual Background Extractor (Vibe)

Dali Qiao[1], Qinghe Zheng[1(✉)], Xinyu Tian[1], Abdussalam Elhanashi[2], and Sergio Saponara[2]

[1] School of Intelligence Engineering, Shandong Management University, Jinan 250357, China
15005414319@163.com
[2] Department of Information Engineering, University of Pisa, 56122 Pisa, Italy

**Abstract.** Computer vision and motion capture have gradually developed, and the detection of moving objects has always been very important. Vibe is a simple and efficient algorithm with low computation, good real-time performance, fast speed. There will be ghosting when detecting images in the foreground, which allows people to observe the trajectory of objects, but it will also affect the image display at the next moment to a certain extent. Vibe is divided into two steps: initializing the background model and updating the background model. By adjusting the secondary sampling factor, very few sample values can cover all background samples, store a set of values for each pixel that used to be at the same location and its neighbors. And then compare the pixel values in this collection with the current pixel values, to determine if the pixel belongs to the background and adapt the model by randomly selecting which values to replace from the background model. Finally, when a pixel is found to be part of the background, its value is propagated to the background model of neighboring pixels. The results of the vibe test are not affected by the speed at which the object moves, but by light.

**Keywords:** Moving object detection · Vibe algorithm · Background subtraction technique · Quadratic sampling factor

## 1 Introduction

The field of computer vision has been gradually developing, among which moving object detection has always been a very hot research field in machine vision, optical flow method, frame difference method, background difference method, etc. are more commonly used. The optical flow method [1, 2] requires a lot of computing and special hardware support to determine the appropriate amount of optical flow of pixels in the image to achieve the detection of moving objects. The frame difference method [3] uses interval or consecutive frames of images to detect objects, although the principle is simple and convenient to calculate, but the detected image contour is intermittent and there are holes inside. The background difference method [4] differs the current video frame image from the background image to obtain a moving object image, but the requirements for the background image are higher. Motion targets are not allowed in the background image

and need to be updated in real time to adapt to changes that you can make. Among them, the VIBE algorithm has been studied and applied by people because of its small amount of calculation, fast operation speed and better detection effect, and has made certain contributions to the field of computer vision. It is the process of identifying foreground (moving objects) and background (relatively static) through video sequences.

It lays the foundation for subsequent tasks such as target tracking, and is widely used in intelligent surveillance [5], unmanned [6], aerospace [7], and defense military [8]. In recent years, multimedia communication technology and computer networks have developed rapidly, and the information transmission rate has also increased. People are no longer satisfied with the existing methods of communicating information such as pictures, text, and voice. Motion capture is widely used in military, education, public security, medical and other fields. Due to the inherent characteristics of video surveillance, it is necessary to collect video data around the clock [9, 10]. This has led to a sharp increase in video data, which has brought great challenges to the storage and use of information.

The performance of the existing image keyframe feature point matching method cannot meet the social demand for video information, so the image key feature point matching method based on vibe algorithm is proposed. Vibe algorithm, also known as visual background extraction algorithm, has the advantages of low computation, good real-time, fast speed, simple and efficient. The Vibe algorithm can quickly detect image keyframes, provide accurate basis for image key feature point matching, improve the matching accuracy of feature point matching methods, and enhance the rationality of matching speed adjustment parameters.

## 2 Vibe Algorithm

### 2.1 Background Model Initialization

The background model of the Vibe algorithm only needs the first frame image in the video to be initialized, which mainly includes three parts: pixel sample set initialization, foreground pixel segmentation and pixel update. The Vibe algorithm adopts the update strategy of randomly replacing background pixel samples, which can cope with the uncertainty of pixel transformation to a certain extent.

First, $N$ matrices of the same size as the image are created, assuming a certain pixel A in the first image, and randomly select $n$ pixel values in the pixel and its 8 neighborhoods to fill in the corresponding position of pixel $a$ of $n$ matrices. That is, the initialization of one pixel is completed, and the initialization of the background modulus is completed after the same operation is performed on all pixels. In subsequent detection, each pixel has $N$ sample values corresponding to the background model, as shown in Eq. (1).

$$M(a) = \{n_1, n_2, ..., n_{N-1}, n_N\} \tag{1}$$

### 2.2 Foreground Pixels Extraction

The process of segmenting the foreground image by the Vibe algorithm is the process of classifying all pixels in the current frame image, and completing the foreground image

segmentation at the same time as the classification is completed. The principle of pixel classification is to compare a pixel with a sample in the corresponding sample set to determine whether it is a foreground attraction.

The $F_1$ and $F_2$ are the color distance, $d(x)$ is the pixel value of pixel $b$, and $r$ is the color distance radius. $d_i$ is the pixel value in the sample set, $i = 1, 2, …, N$. Color distance refers to the gap between two colors, usually the larger the distance, the greater the difference between the two colors; Conversely, the more similar the two colors are. When the total number of samples $Vr(d(x))$ at a distance from $d(x)$ less than $r$ is greater than the set threshold $m$, that is, if the matching value of pixel $b$ and the background is large enough, it is judged that pixel $b$ is the background point; The opposite is the former attraction. The foreground segmentation principle is shown in Eqs. (2) and (3).

$$c = V_r(d(x)) \cap \{d_1, d_2, ..., dn\} \tag{2}$$

$$b = \begin{Bmatrix} g_b, c \geq m \\ g_f, c \leq m \end{Bmatrix} \tag{3}$$

where $c$ is the number of samples less than R away from $v(x)$; $b$ is the pixel currently judged; $g_b$ is the background point; $g_f$ is the foreground point.

Foreground pixel segmentation: In two-dimensional Euclidean space, calculates the Euclidean distance between the current pixel and the pixel in the corresponding sample set. If greater than the radius value $R$, the matching number is incremented by 1. When the number of matches is greater than the match threshold, the pixel is the background, otherwise it is the foreground.

## 2.3  Updates to the Background Model

Vibe algorithm updates have randomness in space-time. When updating to a new image, if the pixel corresponding to a certain location of the image is judged to be a background pixel, the background model needs to be updated. Background pixels are used to randomly replace any pixel value in the same location of the background model and its eight neighborhoods. Repeat the previous steps to complete the update of the image background model.

In the background model, the probability of a sample being retained is correlated with the historical value, and the number of matches $p$ is incremented by one for each sample detected as foreground gf. After sample $b$ is detected as a foreground 50 times in a row, the secondary sampling factor $\Psi$ is updated to the background $g_b$ with a probability of $1/\Psi$, and the judgment principle is shown by Eq. (4).

$$b = \begin{Bmatrix} g_b, p > 50 \\ g_f, p \leq 50 \end{Bmatrix} \tag{4}$$

As shown in Eqs. (5) and (6), $\Psi$ is the time sampling factor, and when a pixel is classified as a background point, a probability of $1/\Psi$ is used to update a random sample

in the corresponding sample set. At the same time, the background has a probability of $1/\Psi$ to change one of its eight neighborhoods.

$$b_\Psi = g_b, \Psi \in (0, b) \tag{5}$$

$$b_{\Psi+r} = g_b, r \in [-1, 1] \tag{6}$$

## 3   Experiment

### 3.1   Dataset Introduction

720p, 30fps video shot with iQOO9 was tested in different brightness, location, and angle. In the video, there are long-distance object movement, close-range object movement, rapid movement of objects, moving objects blocking each other, and sudden movement of light.

### 3.2   Experimental Environment

We use Dell G15 5520 PC, 12th Gen Intel® Core(TM) i7-12700H 2.30 GHz processor, system type 64-bit operating system, x64-based processor. The number of background samples $N = 10$; the judgment threshold $m = 2$; the initial value of color distance radius $R = 20$; and the time sampling factor $\Psi = 16$.

### 3.3   Experimental Results and Analysis

**Experiments at different brightness levels**. The experimental results of lens shake in Fig. 1, the experimental results of lens fixation (background initialization not completed) in Fig. 2, and experimental results in lens fixation (background initialization completed) in Fig. 3 are shown. When the lens is not fixed, the background is misjudged as the foreground due to the shaking of the lens. Even if the lens is fixed, if the background is not initialized, it will be misjudged as the foreground. As the light decreases, less and less parts of foreground objects are detected; As the light increases, the more details of the foreground object are detected. But there will be backgrounds around the foreground that are misjudged as foregrounds.



a                                   b                                   c

**Fig. 1.** Experimental results of lens shake. **a** brightness (low), **b** brightness (medium), **c** brightness (high).

**Fig. 2.** Experimental results with lens fixation (background not initialized). **a** brightness (low), **b** brightness (medium), **c** brightness (high).



**Fig. 3.** Experimental results with lens fixation (background initialization complete). **a** brightness (low), **b** brightness (medium), **c** brightness (high).



**Fig. 4.** Experimental results when suddenly brightening. **a** before the change, **b** change, **c** after the change.



**Fig. 5.** Experimental results of sudden darkening. **a** before change, **b** change, **c** after change.



**Fig. 6.** Experimental results of relatively large hours of objects in the lens frame. **a** small, **b** medium, **c** large.

**Fig. 7.** Experimental results at different object speeds. **a** fast, **b** medium, **c** slow.

**Experiments when brightness changes suddenly**. Figure 4 shows the sudden brightness and Fig. 5 shows the sudden dark. When the brightness of light changes suddenly, the part of the brightness that changes is misidentified as a moving object. However, in the end, the miscalculated picture will still be corrected, and the smaller the brightness change, the faster the correction.



**Fig. 8.** Background clutter. **a** very messy, **b** messy, **c** concise.



**Fig. 9.** Experimental results when shooting angles are different. **a** looking down, **b** looking up, **c** looking up.



**Fig. 10.** Experimental results when moving an object. **a** less, **b** medium, **c** more.

**Experiments in which objects in the lens are relatively large hours**. This is shown in Fig. 6: The experimental results of the relatively large hours of objects in the lens screen. When the relative size of the object in the lens frame is smaller, the less influence on other objects, and its own contour is more perfect; Conversely, the larger the relative

size, the greater the impact on the surrounding picture, the more incomplete its own contour, and the white phenomenon will occur when the moving area is too large.

**Experiments at the speed of an object**. Figure 7 shows the experimental results at the speed of the object. The speed of movement of objects ranges from fast afterimage to slow movement, and the detected picture is very similar. However, the slower the object moves, the less misjudgment about the object's shadow, which can be detected regardless of speed.

**Experiments with varying levels of background clutter**. Figure 8 shows the experimental results with the degree of background clutter. The more cluttered the background, the greater the effect of the ghosting phenomenon on the foreground object, destroying the structure and contour of the foreground object. Conversely, the simpler the background, the less affected the object is affected by the background.

**Experiments at different angles**. Figure 9 shows the results of experiments with different shooting angles. According to previous experiments, the less complex the background is when looking down, the less impact it has on the image, and the more complete the image is. The light source is mostly above, the brighter the image, the more complete it is, the background is mostly more complex when looking up, the light is dark when looking up, the outline of the object is incomplete and the background is messy.

**Experimental results under different number of target conditions**. Figure 10 shows the experimental results of how many objects are moved. No matter how many moving objects are moved, moving objects can be detected. When there is a small amount of movement, the image can be clearly detected, but because of the ghosting phenomenon, the more objects moved, especially when overlapping occurs, the detected image is very important.

## 4   Conclusion

The field of computer vision has been gradually developing, among which moving object detection has always been a very hot research field in machine vision. Vibe is divided into two steps: initializing the background model and updating the background model. By adjusting the quadratic sampling factor, very few sample values cover all background samples. Stores a set of values for each pixel that used to be at the same location and its neighbors, and then compares the pixel values in this collection to the current pixel values to determine if the pixel belongs to the background and adapt the model by randomly selecting which values to replace from the background model.

The vibe algorithm is affected by light, the relative size of objects within the lens, objects blocking each other, background, viewing angle, and camera movement, and is more susceptible to noise. However, the speed at which the object moves does not affect the detection of the system. When the light is dark or changes suddenly, it will cause the system to misjudge; Ghosting when the background is complex can cause images to overlap. How to completely detect the foreground target and reduce the error caused by light, background, and ghost shadow will be the direction of subsequent research.

# References

1. RoshanA, Zhang Y (2019) Moving object detection using spatial correlation in lab colour space. ISPRS—Int Archiv Photogramm, Remote Sens Spatial Inf Sci. ISSN 2194-9034
2. Zhao N, O'Connor D, Basarab A, Ruan D, Sheng K (2019) Motion compensated dynamic MRI reconstruction with local affine optical flow estimation. IEEE Trans Biomed Eng 1558–2531
3. Zheng Q, Yang M, Tian X et al (2020) A full stage data augmentation method in deep convolutional neural network for natural image classification. Discr Dyn Nat Soc
4. Zhang Q, Yang M, Zheng Q, Zhang X (2017) Segmentation of hand gesture based on dark channel prior in projector-camera system. In: IEEE/CIC ICCC, Qingdao, China, pp 1–6
5. Li J et al (2019) Dynamic hand gesture recognition using multi-direction 3D convolutional neural networks. Eng Lett 27(3):490–500
6. Xin X, Huiping L, Fengping H (2015) The flocs target detection algorithm based on the three frame difference and enhanced method of the Otsu. Int J Comput Theory Eng. ISSN 1793-8201
7. Zheng Q, Zhao P, Wang H, Elhanashi A, Saponara S (2022) Fine-grained modulation classification using multi-scale radio transformer with dual-channel representation. IEEE Commun Lett 26(6):1298–1302
8. Jiang N, Fu F, Zuo H, Zheng X, Zheng Q (2020) A municipal PM2.5 forecasting method based on random forest and WRF model. Eng Lett 28(2):312–321
9. Noor W, Vincent S, Anggi PR (2015) Background subtraction berbasis self organizing map untuk deteksi objek bergerak. Jurnal Systemic: Inf Syst Inf J 1(1):42–51
10. Zheng Q, Yang M, Tian X, Wang X, Wang D (2020) Rethinking the role of activation functions in deep convolutional neural networks for image classification. Eng Lett 28(1):80–92

# Car Recognition Based on HOG Feature and SVM Classifier

Xuanming Zhu[1], Qinghe Zheng[1](✉), Xinyu Tian[1], Abdussalam Elhanashi[2], Sergio Saponara[2], and Pierpaolo Dini[2]

[1] School of Intelligence Engineering, Shandong Management University, Jinan 250357, China
15005414319@163.com

[2] Department of Information Engineering, University of Pisa, 56122 Pisa, Italy

**Abstract.** Public places such as parks, pedestrian lanes, and dangerous construction sites are often prohibited for vehicles to enter, but there are still some cars that cause hidden safety hazards due to driver negligence or deliberate entry. At the same time, there are problems of high cost and low efficiency in car detection and warning in public places completely relying on manpower. Therefore, in view of the security risks caused by blind car driving in some public places, this paper builds a car recognition model based on directional gradient histogram (HOG) and support vector machine (SVM) to realize the car recognition function. This model is used to judge whether there is a car entering the recognition range, and is applied to the car recognition in a specific public environment, so as to warn cars entering the forbidden area and feed back to the staff for corresponding treatment. Through the experiment and the analysis of sample data confusion matrix, the accuracy and specificity of the model are up to 92.7% and 92.6%, respectively, which proves that the model has a high recognition rate and can be well applied to the automobile recognition task.

**Keywords:** Public Place · Driving Safety · SVM + HOG · Car Recognition

## 1 Introduction

In recent years, due to the surge of car ownership per capita, driving safety problems in public places occur frequently, and a considerable part of accidents occur because drivers ignore or ignore the warning of traffic control [1, 2]. If real-time warning and feedback related information can be given to the staff in critical areas, and penalties can be taken to warn the driver, the occurrence of such accidents can be reduced [3].

To realize the recognition of cars, it is necessary to rely on many methods of image processing in computer vision system, such as image transformation [4, 5], image enhancement [6, 7], image segmentation [8–10]. The image background is complex, the light is strong or weak, and the shooting angle is poor, which will lead to the image processing results falling short of expectations and affect the recognition effect.

Based on OpenCV, this paper uses HOG and SVM to construct and verify the relevant model. The training set with multiple views, different lighting environments and multiple

scenes is used to construct the model, so that it still has a high recognition rate in the complex recognition environment. The model can be used to identify cars in public environment to reduce the occurrence of related accidents.

## 2   Some Work Based on HOG Feature Extraction and SVM

### 2.1   HOG Features Extraction

The HOG feature is a feature descriptor used for object detection in computer vision and image processing. It constructs features by calculating and counting the gradient direction histogram of the local area of the image.

The input image is grayed, normalized and Gamma corrected to improve the robustness of the image to light changes and reduce the interference of noise. Gamma is calculated as follows:

$$I(x, y) = I(x, y)^{\gamma} \tag{1}$$

*In this equation: I represents the image, (x, y) represents the pixels of the image, and $\gamma$ is usually 1/2.* Since the gradient is a vector, its amplitude and direction should be calculated. The formula for calculating the x direction gradient of each pixel is as follows:

$$G_x(x, y) = H(x + 1, y) - H(x - 1, y) \tag{2}$$

The formula for calculating the gradient in y direction is as follows:

$$G_y(x, y) = H(x, y + 1) - H(x, y - 1) \tag{3}$$

*In this equation: H (x, y) is the pixel value.* The formula for calculating the gradient value is as follows:

$$G(x, y) = \sqrt{G_x(x, y)^2 + G_y(x, y)^2} \tag{4}$$

The formula for calculating the direction is:

$$G(x, y) = \tan^{-1}\left(\frac{G_y(x, y)}{G_x(x, y)}\right) \tag{5}$$

The magnitude of the gradient value represents the change trend of the object boundary related features in both horizontal and vertical directions. The Fig. 1 shows the procedure for extracting HOG features.

Main parameters of HOG feature extraction: detection window size (winSize), block size (blockSize), block sliding step (blockStride), cell unit size (cellSize), and the number of histogram bin (nBins). Cell is the smallest unit of HOG feature extraction, and the feature descriptor is obtained by combining the gradient or edge direction histogram of each pixel in the cell. Relevant parameters are shown in Fig. 2.

In Fig. 2, the moving step is a cell size of $8 \times 8$, The number of bin in the histogram is 9, Then each row of blocks needs to be moved 6 times, a total of 7 blocks, Each column of blocks needs to be moved 14 times, a total of 15 blocks, There are 105 blocks, and each block is composed of 4 cells, so there are 420 cells, so the actual HOG feature dimension is 3780 dimensions.

**Fig. 1.** The HOG feature extraction procedure



**Fig. 2.** The HOG feature extraction related parameters

## 2.2 SVM Model

The SVM are a kind of generalized linear classifier which can classify data according to supervised learning.

The basic idea of SVM learning is to solve for separated hyperplanes that correctly partition the training data set and have the greatest geometric spacing. A set of training data and labels such as $\{(x_1, y_1), (x_2, y_2) \dots (x_N, y_N)\}$ is assumed. $x_i \in R^n$, $y \in \{+1, -1\}$, $i = 1, 2, \dots N$, $x_i$ is the ith eigenvector, and $y_i$ is the class label. A label equal to $+1$ is considered as a positive class and a label equal to $-1$ is considered as a negative class. Then the Hyperplane can be represented by $\omega^T x + b = 0$.

Some understanding of classification problem:

(1) The training set is linearly separable:

For $\{(x_i, y_i)\}$, $i \in 1 \sim N$, $\exists(\omega, b)$ make $\forall i = 1 \sim N$. We have: if $y_i = +1$, then $\omega^T x_i + b \geq 0$, if $y_i = -1$, then $\omega^T x_i + b < 0$, finally $y_i[\omega^T x_i + b] \geq 0$.

(2) Maximum geometric Margin:

Vector $x_0$ to the hyperplane $\omega^T x + b = 0$ distance of $d = \frac{|\omega^T x_0 + b|}{\omega}$. The $(a\omega, ab)$ by $a$ zoom $(\omega, b)$, make on the receiving support vector $x_0$ have $|\omega^T x_0 + b| = 1$, the support vector and plane distance $d = \frac{1}{\omega}$. Then maximizing $d$ means minimizing $\omega$. The SVM margin is shown in Fig. 3.

**Fig. 3.** The SVM classification margin

The problem to be optimized by the support vector machine for linear problems is obtained as follows:

$$\text{minimize: } \frac{1}{2}\omega^2 \tag{6}$$

$$\text{subject to: } [y_i\omega^\mathrm{T}x_i + b] \geq 1 \ (i = 1 \sim N) \tag{7}$$

In the case of non-linearly separable data, SVM uses kernel function to map input variables into a high-dimensional feature space, and makes them linearly separable in the high-dimensional space. Then, the optimal classification hyperplane is constructed in this high-dimensional space.

The problem to be optimized when dealing with nonlinear problems:

$$\text{minimize: } \frac{1}{2}\omega^2 + C\sum_{i=1}^{N}\xi_i \tag{8}$$

$$\text{subject to: } \begin{cases} y_i[\omega^\mathrm{T}\varphi(x_i) + b] \geq 1-\xi_i \ i = 1 \sim N \\ \xi_i \geq 0 \end{cases} \tag{9}$$

*In this equation:* $\xi_i$ *as the slack variables, C to predetermined parameters,* $C\sum_{i=1}^{N}\xi_i$ *as regular,* $\varphi(x)$ *for high-dimensional mapping x.*

## 3   Experimental Procedure and Evaluation

### 3.1   Experimental Environment

The experimental environment on which this paper is based is shown in Table 1.

### 3.2   Dataset Introduction

The training samples of this classifier come from a variety of actual road images, including several colors of objects, such as sedan, buses, trucks, roadside trees, road surfaces, guardbars, etc. The training samples are gray images of $64 \times 128$ pixels, of which 400

**Table 1.** Experimental environment

| Experiment environment | Configuration/Version |
|---|---|
| Operating system | 64-bit Windows11 |
| Processor | Intel(R) Core(TM) i7-8750H @ 2.20GHz 2.21 GHz |
| Compiling software | Visual Studio 2022 |
| OpenCV | OpenCV4.5.5 |

positive samples and 800 negative samples are included. The 500 positive and negative samples are used in the test set. Positive samples are real car images taken in different scenes, while negative samples are images that do not contain cars associated with the positive sample environment. Some processed positive and negative sample data are shown in Fig. 4.



(a) Positive sample



(b) Negative sample

**Fig. 4.** Some processed positive and negative sample data

### 3.3 Experimental Results and Analysis

In the experiment, 1000 Image in any scene were randomly selected, including 500 images containing cars and 500 images without cars for mixed detection. The test results are shown in Table 2.

The accuracy rate, error rate, precision rate, recall rate, specificity and F-measure in the general indicators of model evaluation are analyzed as shown in Table 3.

**Table 2.** Data confusion matrix

| Truth | Predicted | |
|---|---|---|
| | Positive | Negative |
| Positive (500) | TP (473) | FN (27) |
| Negative (500) | FP (37) | TN (463) |

In the table: TP: true example, FN: false negative example, FP: false positive example, TN: true negative example

**Table 3.** Evaluation indicators

| Evaluation indicators | Computational equation | Value (%) |
|---|---|---|
| Accuracy | $\frac{TP+TN}{TP+FN+FP+TN}$ | 93.6 |
| Error rate | $\frac{FP+FN}{TP+FN+FP+TN}$ | 6.4 |
| Precision | $\frac{TP}{TP+FP}$ | 92.7 |
| Recall | $\frac{TP}{TP+FN}$ | 94.6 |
| Specificity | $\frac{TN}{TN+FP}$ | 92.6 |
| F-Measure | $\frac{2\,Recall*Precision}{Recall+Precision}$ | 93.6 |

According to the data of each evaluation index in Table 3, it is easy to know that this model has high accuracy. It is expected that this model can perform the task of car recognition in relevant areas well.

In order to further verify the generalization ability of the model, 140 images of different types of car images in various scenarios were selected for separate detection, and a total of 1400 images were selected for mixed detection. The recall rate of common car type is shown in Table 4.

The recall rate of racing car and taxi can reach 90.0% and 95.0% after retraining with hard example feature extraction. It can be seen that the model can improve its generalization ability and accuracy after corresponding training, and has the ability of future learning enhancement, which can be better applied to car recognition tasks.

**Table 4.** The recall rate of common car type

| Car type | Recall (%) |
|---|---|
| Bus | 92.9 |
| Family sedan | 96.4 |
| Fire engine | 92.1 |
| Heavy truck | 92.1 |
| Jeep | 95.7 |
| Minibus | 94.3 |
| Racing car | 86.4 |
| SUV | 99.3 |
| Taxi | 82.9 |
| Truck | 90.7 |
| All(mixture) | 92.3 |

## 4   Conclusion

Based on the phenomenon of high incidence of driving safety problems in some public places, this paper trains the car recognition model based on HOG and SVM to judge whether there are cars passing through the recognition area, and then make corresponding treatment to reduce the occurrence of accidents. Compared with the traditional warning measures such as manual inspection and placing ordinary warning signs, the application of this model can reduce labor cost, improve efficiency and enhance the warning effect. According to the analysis of the experimental results, all the indexes of the model have good performance, which proves that the model has the corresponding ability of car recognition.

However, the model still lacks the corresponding terminal practice experience, and the future research direction will gradually shift to the realization of terminal equipment and the processing of recognition results. For example, the speed is judged according to the driving distance and driving time of the car entering and leaving the recognition range, and the pure visual speed detection scheme is realized to help the driving safety in public places.

## References

1. Guo X, Hao Q, Yang F (2022) Apple multi-target detection method based on improved HOG and SVM. Foreign Electron Meas Technol 41(11):154–159. https://doi.org/10.19652/j.cnki. femt.2204286
2. Zhao J, Wang Y (2020) Unsafe behavior recognition based on image recognition technology. Saf Environ Eng 27(1):158–165. https://doi.org/10.13578/j.cnki.issn.1671-1556.2020.01.024
3. Zheng Q, Yang M, Tian X et al (2020) A full stage data augmentation method in deep convolutional neural network for natural image classification. Discrete Dyn Nat Society

4. Zhang Q, Yang M, Zheng Q, Zhang X (2017) Segmentation of hand gesture based on dark channel prior in projector-camera system. In: IEEE/CIC ICCC. Qingdao, China, pp 1–6
5. Li J et al (2019) Dynamic hand gesture recognition using multi-direction 3D convolutional neural networks. Eng Lett 27(3):490–500
6. Geng Q, Zhao H, Fanhua Y, Wang Y, Zhao H (2018) Vehicle recognition algorithm based on improved HOG feature extraction. Chin Optics 11(02):174–181
7. Zheng Q, Zhao P, Wang H, Elhanashi A, Saponara S (2022) Fine-grained modulation classification using multi-scale radio transformer with dual-channel representation. IEEE Commun Lett 26(6):1298–1302
8. Jiang N, Fu F, Zuo H, Zheng X, Zheng Q (2020) A municipal PM2.5 forecasting method based on random forest and WRF model. Eng Lett 28(2):312–321
9. Pathak AR, Pandey M, Rautaray S (2018) Application of deep learning for object detection. Procedia Comput Sci 132:1706–1717
10. Li Z, Liu W (2017) A moving vehicle detection method based on local HOG fature. J Guangxi Normal Univ (Natural Science Edition) 35(3):1–13. https://doi.org/10.16088/j.issn.1001-6600.2017.03.001

# A Microcontroller-Based System
# for Human-Emotion Recognition with Edge-AI
# and Infrared Thermography

Maria Gragnaniello ⓘ, Alessandro Borghese⁽✉⁾ ⓘ, Vincenzo Romano Marrazzo ⓘ,
Giovanni Breglio ⓘ, Andrea Irace ⓘ, and Michele Riccio ⓘ

Department of Electrical Engineering and Information Technology (DIETI), University of
Naples Federico II, 80138 Naples, Italy
`alessandro.borghese@unina.it`

**Abstract.** Infrared thermography has shown great promise as a diagnostic method
for health care, providing useful information on a person's physiological, and
pathological state. Recently, the use of artificial intelligence combined with
infrared technology has boosted the adoption of thermal imaging in various appli-
cations and has been proposed to recognize human emotion by measuring facial
skin temperature. However, its application has been limited to laboratory settings
due to demanding computational and hardware resources. In this scenario, this
work presents the design and development of a portable system based on a low
power microcontroller implementing an optimized Edge-AI solution for binary
emotional state classification using minimal hardware resources. The recognition
of happiness and sadness emotional states induced by audiovisual stimuli serves as
a case-study for feasibility assessment. Thermal images, produced by an uncooled
and low-cost thermal sensor, along with electrocardiogram, are acquired and pro-
cessed with an Arm® Cortex®-M4 microcontroller. A simple, yet effective neural
network has been developed, optimized, and deployed to run the emotion detec-
tion algorithm in real time. The complete system has been experimentally verified
and results in terms of accuracy and hardware constraints are discussed. Specifi-
cally, by employing a dataset consisting of 60 infrared videos, an accuracy of 80%
was achieved with a resource occupation of 3.4 kB of RAM and 76.4 kB of flash
memory.

**Keywords:** Edge-AI · Embedded systems · Infrared thermography

## 1 Introduction

Infrared thermography (IRT) has become a reliable technology for a wide range of med-
ical applications, allowing for contactless acquisition of temperature [1]. The possibility
for such a measurement technique to serve as viable diagnostic method descends from
the fact that several information about a person's physiological, pathological, and emo-
tional state can be derived from the body temperature distribution [2, 3]. Studies have
also demonstrated the effectiveness of thermography in recognizing human emotions

[4, 5]. Detecting emotional reactions in humans can be difficult due to people hiding their true emotions in public, especially negative ones [6]. There are several application examples where emotion recognition is based on facial expressions, voice, and posture. Nevertheless, the manifestation of emotion through these channels can be masked. Instead, physiological signals, such as the activity of the autonomous nervous system, reflect the intrinsic emotional state of the person, and are immune to social masking [7].

To address this issue, various approaches have been proposed in the literature, including functional Infrared Thermal Imaging (fITI). This method is recognized as a promising technique for measuring emotional unconscious responses by analyzing thermal variations on the facial skin [8].

State-of-the-art applications are based on high-end IR cameras and Convolutional Neural Networks (CNNs), which are incompatible with low-power portable devices because of greater memory and computation-time requirements than Neural Networks (NNs).

This paper proposes a new approach built upon the integration of Edge-AI technique and a microcontroller-based system. A classifier able to recognize the happy and sad emotional states is developed as a case-study. The use of a low-power and cost-effective microcontroller enables the use of this solution in large-scale portable applications, including hospitals, clinics [9], workplaces, and home-based monitoring. Despite the absence of high-performance hardware, the system achieves a promising initial accuracy thanks to the ECG signals sampled synchronously alongside the IR stream.

The proposed embedded system is based on an Arm® Cortex®-M4 microcontroller that simultaneously handles (i) the acquisition of thermal images and ECG signal; (ii) signal processing and features extraction; (iii) NN inference. Significant optimizations were applied to the hardware resources to meet the time constraints necessary for real time operation. Additionally, an analysis on the impact of data quantization on the network performance was conducted.

## 2   Hardware Specifications

The system architecture, designed and developed as depicted in the block diagram of Fig. 1, consists of three fundamental parts: a Microcontroller Unit (MCU), an IR thermal camera module (IR module) and an analog front-end (AFE) to sense the ECG (ECG module).

The adopted MCU (STM32F411CEU6) is equipped with an Arm® Cortex®-M4 CORE processor with 84 MHz clock frequency; it embeds high-speed memories offering 512 kB of flash memory and 128 kB of SRAM [10].

The AFE is based on the AD8232 chip [11], a fully integrated single-lead ECG front-end with a common-mode rejection ratio of 80 dB and a gain of 100.

The IR thermal stream is captured through the FLIR Lepton 2.5 sensor, a radiometric-capable long-wave infrared (LWIR) camera. It features a focal plane array of 80 × 60 active pixels with a 15 µm pitch and 9 Hz frame rate. The IR module detects infrared radiation in the wavelength range of 8–14 µm exploiting uncooled microbolometer technology.

**Fig. 1. a** Block diagram of the prototype; **b** time management diagram depicting the time constraints for data acquisition and elaboration.

The data transfer between the MCU and the IR module is implemented through the Video over Serial Peripheral Interface (VoSPI) Mode 3 protocol [12]. The time interval between two consecutive frames, about 27 ms, sets a time constraint for the MCU, which must complete the data processing (i.e., signal processing and inference) to maintain synchronization and avoid frame loss. To this purpose, the video stream from the sensor is managed via Direct Memory Access (DMA), as shown in Fig. 1b, thus allowing the NN inference to be performed every 360 frames (40 s-long time window) while keeping the devices synchronized.

## 3 Neural Network Design and Training: Methodology and Results

The dataset used to train and validate the designed NN was produced at the Electrothermal Characterization Laboratory of the University of Naples "Federico II". It consists of 60 five-minutes-long IR videos featuring volunteers who were subjected to audio-visual stimuli. Such videos were recorded by a high-performance IR camera which can produce videos with higher frame rate and image resolution than the FLIR Lepton 2.5 (9 Hz and $60 \times 80$ pixels). Therefore, the quality of the videos making the dataset was preliminarily lowered to match the quality of the videos recorded by the FLIR Lepton 2.5. Two pictures of the same subject are reported in Fig. 2.

In particular, the picture of Fig. 2a results from the resolution degradation of a frame produced by the high-performance camera, while the frame of Fig. 2b was captured by the FLIR Lepton 2.5.



**Fig. 2.** Example of (**a**) a resolution-reduced frame from the training dataset and (**b**) a frame acquired by the FLIR Lepton 2.5 (b).

The computational complexity was mitigated by the dimensionality reduction of thermal data, enabling successful emotional state classification through synchronous ECG signal sampling. 1-D time series were extracted from the image stream, originating from four distinct Regions of Interest (ROIs) shown in Fig. 2: forehead (1), cheeks (2,4), and philtrum (3). For each ROI, the time evolution of spatial average temperatures and the standard deviation were calculated.

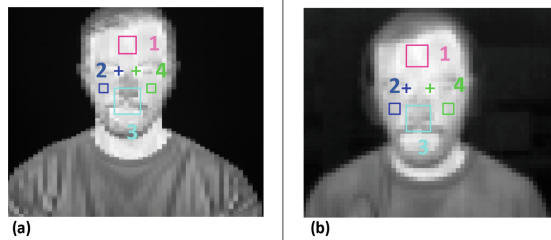These signals, together with the ECG signal, operated as inputs for the NN designed and trained on the Edge Impulse platform [13]. The impulse (i.e., a preprocessing stage whose output is fed to the NN) consists of a spectral analysis block that operates on 375 training windows, each spanning 40 s, to analyze the data. Specifically, the Fast Fourier Transform (FFT) with a length of 512 points was applied to each window. A classifier was employed to process the results of the spectral analysis block and make predictions on the two emotional labels: happy and sad. The NN adopted is a Fully Connected Network consisting of 2 dense layers, with 30 and 20 neurons, respectively, and using the activation function LeakyReLU ($\alpha = 0.2$). Additionally, the dropout technique was employed with a rate of 0.2 to mitigate overfitting. In Fig. 3a, the training data is visualized, with correctly classified instances shown in green and misclassified instances in red. Figure 3b provides a summary of the estimated performance of the network on the target MCU. The expected inference time is 58 ms, with a hardware usage of 9.5 kB of RAM and 256.3 kB of flash memory. The achieved accuracy is 99.1%. A further NN training was conducted by changing the data quantization from 32 to 8 bit. The results of the so-trained NN are summarized in Fig. 3c, which shows that significantly lower hardware and time constraints can be achieved at the cost of a slightly lower accuracy (98.2%). Therefore, the selected NN was the one operating on 8-bit data and its performance on the validation dataset are summarized in Table 1. In this case, the systems exhibited an accuracy of 80%, with an F1 score of 0.73 for the "happy" output and 0.85 for the "sad" output.



**Fig. 3.** **a** Classified training data mapped on a reduced order 2-D space. Summary of the estimated on-device performance of the network on the target MCU for **b** 32-bit and **c** 8-bit quantized data.

## 4   System Validation

A demonstrative prototype was developed based on the hardware of Fig. 1a. The software packages were exported from Edge Impulse and implemented on the microcontroller. Figure 4a showcases the prototype including the clamp electrodes for ECG. The IR

**Table 1.** Confusion matrix referring to the 8-bit quantized data.

|          | Happy | Sad   | Uncertain |
|----------|-------|-------|-----------|
| Happy    | 74.7% | 21.3% | 4%        |
| Sad      | 17%   | 83%   | 0%        |
| F1 Score | 0.73  | 0.85  |           |

module and the MCU are housed in a case for ease of mounting on the desktop monitor providing the audiovisual stimuli. Figure 4b describes the experimental protocol for testing the application where the subject faces the IR module and seats at 1 m from it.



**Fig. 4.** **a** Final prototype showing the IR inference system mounted on a desktop monitor along with the clamp electrodes for ECG and **b** subject measurement protocol.

Sections 2 and 3 emphasized the time constraints imposed by the FLIR Lepton 2.5 sensor and the network inference time predictions on the chosen device. Figure 5 demonstrates the temporal performance of the prototype, illustrating its effectiveness in maintaining synchronization. Specifically, in Fig. 5a, channel 4 provides information about the time required by the MCU to acquire and process a sequence of 360 frames, which is approximately 41 s. On the other hand, channel 7, as detailed in Fig. 5b, illustrates the time necessary for the inference itself, which is about 9 ms. The latter fits within the time window between two frames acquisition (27 ms) thus allowing the application to run without interruption of streaming of IR frames.



**Fig. 5.** Time diagram profiling the microcontroller activity; **a** time required for acquiring and processing 360 frames, **b** time required for the NN inference.

## 5   Conclusion

In this paper, we have introduced a first release of a low cost and small-size Edge-AI system prototype designed for emotions recognition using infrared thermography and synchronized ECG sampling. The system utilizes a single microcontroller to handle both the IR video stream from a low-power thermal sensor and perform inference. The neural network, developed and trained on the Edge Impulse platform, is a fully connected network achieving an 80% accuracy with an inference time of 21 ms, consuming 3.4 kB of RAM and 76.4 kB of flash memory. This proves that the hardware resources and timing constraints required for network inference are compatible with microcontrollers, eliminating the need for more powerful hardware such as dedicated FPGAs and SoCs. These promising results set the stage for further improvement in the NN's performance and advance towards the realization of a fully developed contactless emotion recognition system. This highlights the system as a versatile monitoring device, applicable not only in clinical and automotive environments but also as a safety instrument in emotionally charged workplaces, effectively monitoring high-stress environments.
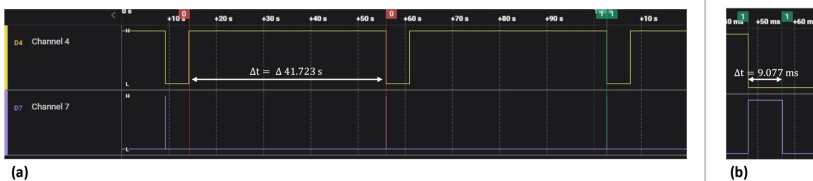
## References

1. Lahiri BB et al. Medical applications of infrared thermography: a review. Infrared Phys Technol 55(4):221–235
2. Kim Y et al. Remote heart rate monitoring method using infrared thermal camera. Int J Eng Res Technol
3. Anaya-Isaza A, Zequera-Diaz M. Detection of diabetes mellitus with deep learning and data augmentation techniques on foot thermography. IEEE Access
4. Jamal KMS, Kamioka E. Emotions detection scheme using facial skin temperature and heart rate variability. In: MATEC web conferences, vol 277, p 020377
5. Cruz-Albarran IA et al. Human emotions detection based on a smart-thermal system of thermographic images. Infrared Phys Technol 81:250–261
6. Hassani S et al. Physiological signal-based emotion recognition system. In: 2017 4th IEEE ICETAS, pp 1–5
7. Selvaraj J et al. Classification of emotional states from electrocardiogram signals: a non-linear approach based on hurst. Biomed Eng Online 12(1):44
8. Kosonogov V et al. Facial thermal variations: a new marker of emotional arousal. PloS One 12(9):e0183592
9. Gasparro R et al. Thermography as a method to detect dental anxiety in oral surgery. Appl Sci 11:12
10. 'Arm Cortex-M4—Microcontrollers—STMicroelectronics. https://www.st.com/content/st_com/en/arm-32-bit-microcontrollers/arm-cortex-m4.html. Accessed 10 May 2023
11. 'ad8232.pdf'. https://www.analog.com/media/en/technical-documentation/data-sheets/ad8232.pdf. Accessed 28 Jun 2023
12. Lepton_Engineering_Datasheet_Rev200.pdf. https://cdn.sparkfun.com/assets/f/6/3/4/c/Lepton_Engineering_Datasheet_Rev200.pdf. Accessed 10 May 2023
13. Edge Impulse. https://edgeimpulse.com/. Accessed 28 Jun 2023

Andrei Boiko[1(✉)] , Maksym Gaiduk[1] , Natividad Martínez Madrid[2] ,
and Ralf Seepold[1]

[1] Ubiquitous Computing Lab, Department of Computer Science, HTWG Konstanz—University
of Applied Sciences, 78462 Konstanz, Germany
{Andrei.Boiko,Maksym.Gaiduk,Ralf.Seepold}@htwg-konstanz.de
[2] School of Informatics, Reutlingen University, 72762 Reutlingen, Germany
Natividad.Martinez@reutlingen-university.de

**Abstract.** The monitoring of a patient's heart rate (HR) is critical in the diagnosis of diseases. In the detection of sleep disorders, it also plays an important role. Several techniques have been proposed, including using sensors to record physiological signals that are automatically examined and analysed. This work aims to evaluate using a contactless HR monitoring system based on an accelerometer sensor during sleep. For this purpose, the oscillations caused by chest movements during heart contractions are recorded by an installation mounted under the bed mattress. The processing algorithm presented in this paper filters the signals and determines the HR. As a result, an average error of about 5 bpm has been documented, i.e., the system can be considered to be used for the forecasted domain.

**Keywords:** Accelerometer · Ballistocardiography · Contactless Measurement · Heart Rate

## 1 Introduction

Continuous monitoring of a patient's vital signs is essential for many chronic diseases. This can be done at home or in the hospital. The intrusion into the patient's life associated with such monitoring often becomes an obstacle to everyday life. A solution to this problem can be found in a non-contact monitoring system's development and practical application. This involves installing a system in the patient's bed in a hospital or even in a regular bed at home so that vital signs can be monitored anywhere [1]. Heart rate (HR) is one of the vital signs of interest, as changes in it can indicate severe disorders such as arrhythmia or bleeding [1, 2]. This signal can also be used as part of sleep monitoring, including vital signs monitoring, which is essential given the increasing risk of cardiovascular disease [3], or for automatic sleep stage scoring [4].

Polysomnography (PSG) is the most accurate and leading method for assessing sleep patterns [5]. However, this method is expensive, cumbersome, and time-consuming, so

it can't be done at home. In this case, developing a low-cost system for ambulatory or home use is particularly interesting [6].

Many measurement techniques can be used to obtain information about cardiac activity [7]. So, a review of non-contact cardiorespiratory monitoring systems can be found in [8]. Several methods have been proposed to measure ballistocardiography (BCG) signals during sleep. These include using different sensors, their number, and their position for measurement. Some have used load cells [9], piezoelectric or pressure sensors [10, 11], and others. The accelerometers usage in various measurement fields is increasing rapidly with the development of highly sensitive sensors based on micro-electromechanical systems (MEMS) [12]. Unfortunately, using accelerometers to measure BCG signals remains underexplored [13].

This paper aims to investigate a contactless accelerometer system for non-invasive HR monitoring. The main focus is to evaluate the system's performance concerning different sleeping positions. This is done by comparing the HR values recorded by the subjects with the reference measurement system of PSG.

## 2 Proposed System

The ability to work autonomously is a prerequisite for developing our system. This is a requirement for the system's hardware and software. This ensures that the processes of collecting and processing data can be carried out without using external equipment. This capability makes the system affordable for home use, with high accuracy, unobtrusiveness, and low cost.

### 2.1 Data Acquisition

The contactless system for HR monitoring consists of several parts: the first one is a data acquisition block which contains the mechanical part (sensor holder and suspension) and the direct data acquisition part. Each of these blocks is described below.

As noted above, the mechanical part should consist of two units—a mounting to the bed frame and a spring steel plate (hanger) on which the accelerometer sensor is fitted. In [13], the procedure regarding the material choice and its parameters to obtain acceptable results for respiratory measurements is presented. This block has been successfully used also for unobtrusive sleep apnea detection [14].

In this study, based on the research on the detection of vital signs using this sensor [15], we selected the ADXL355z (Analog Devices) 3-axis accelerometer sensor for data acquisition. Furthermore, this sensor model and the previous model of this product line were successfully used for cardiorespiratory measurements [16]. The sampling rate is 62 Hz, which is sufficient for the study's aims. We also used the SOMNO HD eco PSG electrocardiography system as a ground truth system to obtain reference data with a sampling frequency of 256 Hz for the cardiac signal.

Due to its compact dimensions and compliance with the low-cost requirement, an ESP32 module was used as the computing unit. This module has the option of storing data on a micro SD module, as well as the possibility of transferring the obtained data using Wi-Fi.

Another important aspect is the definition of system position or sensor arrangement. There are several positions for our system on the bed frame related to the particular subject area—from head to legs [17]. Based on the research results and the studies on accelerometer application for cardiorespiratory measurements [17], we have selected the sensor position at the chest level as the preferable position for our measurements.

## 2.2 Signal Processing

The signal processing algorithm is the second essential part of the contactless HR monitoring system described below. There are several steps in the signal processing algorithm for HR estimation. In the beginning, the first and the last 10 s of the data in each position have been removed to avoid signal epochs where processing is complicated by motion artefacts caused by sleep position changes. The mean signal amplitude was subtracted from the raw data to reduce the offset. Calculating the signal amplitude was the next step helping to consider all measurement axes. The BCG signal was then derived using the Butterworth bandpass filter with a 6th order in the range [0.7; 3.25] Hz to limit HR detection between 42 bpm and 195 bpm. In addition, we applied a Savitzky-Golay filter to enhance signal peaks, which has a 3rd order and a window size of 11 [18]. HR estimation was performed by peak detection function.

## 2.3 Experiment Design

We used a regular single bed and a standard mattress for our study. All materials were wooden and readily available for common use. Before the experiment, the subjects were instructed to lie down on the bed in four typical sleeping positions: on their stomach, right side, back, and left side. We referred to them as P1 to P4 as a consequence. Furthermore, the experiment started in position P1 and ended in position P4. In addition, all subjects had a relaxation period of at least three minutes before data collection. The data measurement lasted 140 s in each position, and the subjects were instructed to behave normally with minimal movement. They were informed that the experiment would be stopped if they became uncomfortable.

## 3 Results

The data were collected from a total of 10 subjects. Specifically, we recruited 5 males and 5 females with average age, height, and weight of $31.8 \pm 8.6$ years, $175.0 \pm 6.0$ cm, and $79.1 \pm 12.9$ kg, respectively, to participate in the experiment. All subjects read and signed the informed consent form before the experiment and were informed of the specifics of the research. The subjects didn't have or report any known disorders and diseases and weren't receiving any treatment or medication.

The contactless HR monitoring system was evaluated by estimating the obtained HR values in comparison with reference data. This research shows a mean absolute error (MAE) as the evaluation metric. It is important to mention that each signal epoch was divided into 20 s chunks for further analysis of the HR and the calculation of the MAE. Approximately $240 \times 20$ s signal windows with a total length of 4800 s were considered

for further analysis based on the experimental design and study details presented in previous sections and Subsections (60 for each sleep position). Table 1 shows the MAE values, which can help us to evaluate the performance of the HR monitoring system.

**Table 1.** The heart rate estimation of the contactless system by MAE (bpm).

| | Sleeping posture | | | |
|---|---|---|---|---|
| | Prone (P1) | Right lateral (P2) | Supine (P3) | Left lateral (P4) |
| Male | 4.81 | 5.18 | 5.00 | 4.84 |
| Female | 4.37 | 6.12 | 4.82 | 5.61 |
| Average (all) | 4.59 | 5.65 | 4.91 | 5.23 |

To evaluate contactless system performance and compare it with ground truth, the Bland-Altman analysis was conducted for each posture separately presented in Fig. 1.



**Fig. 1.** Results of Bland-Altman analysis for HR values obtained by ground truth and contactless system.

Based on the Bland-Altman analysis results, we notice the quite accurate measurement of HR in all considered sleeping postures. However, this good precision is reached within a specific HR range—[45; 80] bpm. In cases where the HR is greater than 90 bpm, there starts to be a significant difference in performance between the reference and contactless systems. At the same time, we didn't collect enough data in this range of HR to state such a trend unequivocally.

# 4   Conclusion and Outlook

Based on the findings, this contactless system based on the accelerometer sensor has good potential for HR monitoring. The level of accuracy is similar to some of the existing techniques and methods that have been considered [19]. Simultaneously, the presented system is related to novel modern approaches with significant advantages over other technologies. The system under consideration comprises hardware and software blocks that could be installed under the bed in a user-friendly way. The system works in autonomous mode, unobtrusively performing the measurements. As a result, we can state that this system can be relevant in hospitals, rehabilitation centres, and even the home. In conjunction with the processing algorithm, the proposed system has several advantages, such as a simple design combined with effective and low-cost hardware and the possibility of non-contact measurement.

For example, one of the options to improve the hardware part is to upgrade the mechanical holder due to the drift shift of the hanger. This shift could be reduced by cutting the hanger part, which induces the additional vibrations during the patient's movements. So, this upgrade could increase the quality of the obtained signal and the signal-to-noise ratio. In addition, we should expand the group of subjects with different ages, heights, etc. Work on this extension is currently underway, as is improving the hardware. Furthermore, the proposed system and approach could be part of a comprehensive sleep monitoring system in combination with the algorithms presented in [20].

# References

1. Conti M, Orcioni S, Martínez Madrid N, Gaiduk M, Seepold R (2018) A review of health monitoring systems using sensors on bed or cushion. In: Rojas I, Ortuño F (eds) IWBBIO 2018. Lecture Notes in Computer Science (2018), vol 10814. Springer, Cham. https://doi.org/10.1007/978-3-319-78759-6_32

2. Rodríguez de Trujillo E, Seepold R, Gaiduk M, Martínez Madrid N, Orcioni S, Conti M (2019) Embedded system to recognize movement and breathing in assisted living environments. In: Saponara S, De Gloria A (eds) Applications in electronics pervading industry, environment and society (2019). ApplePies 2018. Lecture Notes in Electrical Engineering, vol 573. Springer, Cham. https://doi.org/10.1007/978-3-030-11973-7_46

3. Stamatakis E, de Rezende LFM, Rey-López JP (2018) Sedentary behaviour and cardiovascular disease. In: Leitzmann M, Jochem C, Schmid D (eds) Sedentary behaviour epidemiology. Springer Series on Epidemiology and Public Health. Springer, Cham. https://doi.org/10.1007/978-3-319-61552-3_9

4. Gaiduk M, Serrano Alarcón Á, Seepold R, Martínez Madrid N (2023) Current status and prospects of automatic sleep stages scoring: review. BMEL. https://doi.org/10.1007/s13534-023-00299-3

5. Caples SM et al (2021) Use of polysomnography and home sleep apnea tests for the longitudinal management of obstructive sleep apnea in adults. AASM clinical guidance statement. J Clin Sleep Med 17(6):1287–1293. https://doi.org/10.5664/jcsm.9240

6. Asadov A, Seepold R, Martínez Madrid N, Ortega JA (2022) Non-invasive cardiorespiration monitoring using force resistive sensor. In: Hardware and software supporting physiological measurement (HSPM-2022); Hochschule Reutlingen, Reutlingen, Germany, pp 20–22. https://doi.org/10.34645/opus-4001

7. Fan D, Ren A, Zhao N, Haider D, Yang X, Tian J (2019) Small-scale perception in medical body area networks. IEEE J Transl Eng Health Med 7:1–11, Art no. 2700211. https://doi.org/10.1109/JTEHM.2019.2951670

8. Boiko A, Martínez Madrid N, Seepold R (2023) Contactless technologies, sensors, and systems for cardiac and respiratory measurement during sleep: a systematic review. Sensors 23(11):5038. https://doi.org/10.3390/s23115038

9. Malik AR, Boger J (2021) Zero-effort ambient heart rate monitoring using ballistocardiography detected through a seat cushion: prototype development and preliminary study. JMIR Rehabil Assist Technol 8:e25996. https://doi.org/10.2196/25996

10. Xu W, Yu C, Dong B, Wang Y, Zhao W (2022) Thin piezoelectric sheet assisted PGC demodulation of fiber-optic integrated MZI and its application in under mattress vital signs monitoring. IEEE Sens J 22:2151–2159. https://doi.org/10.1109/JSEN.2021.3128601

11. Gomez-Clapers J, Serra-Rocamora A, Casanella R, Pallas-Areny R (2014) Towards the standardization of ballistocardiography systems for J-peak timing measurement. Measurement 58:310–316. https://doi.org/10.1016/j.measurement.2014.09.003

12. D'Mello Y et al (2019) Real-time cardiac beat detection and heart rate monitoring from combined seismocardiography and gyrocardiography. Sensors 19:3472. https://doi.org/10.3390/s19163472

13. Boiko A, Scherz WD, Gaiduk M, Gentili A, Conti M, Orcioni S, Seepold R, Martínez Madrid N (2022) Sleep respiration rate detection using an accelerometer sensor with special holder setup. In: Proceedings of the 2022 E-Health and bioengineering conference (EHB), pp 1–4. https://doi.org/10.1109/EHB55594.2022.9991578

14. Boiko A, Gaiduk M, Martínez Madrid N, Seepold R (2023) Evaluation of a prototype for early active patient mobilization. Procedia Computer Science (Accepted)

15. Conti M, Aironi C, Orcioni S, Seepold R, Gaiduk M, Martínez Madrid N (2020) Heart rate detection with accelerometric sensors under the mattress. In: 42nd annual international conference of the IEEE engineering in medicine & biology society (EMBC). pp 4063–4066. https://doi.org/10.1109/EMBC44109.2020.9175735

16. Hersek S, Semiz B, Shandhi MMH, Orlandic L, Inan OT (2020) A globalized model for mapping wearable seismocardiogram signals to whole-body ballistocardiogram signals based on deep learning. IEEE J Biomed Health Inf 24(5):1296–1309. https://doi.org/10.1109/JBHI.2019.2931872

17. Skoric J et al (2022) Respiratory modulation of sternal motion in the context of seismocardiography. IEEE Sens J 22(13):13055–13066. https://doi.org/10.1109/JSEN.2022.3173205

18. Sadek I, Seet E, Biswas J, Abdulrazak B, Mokhtari M (2018) Nonintrusive vital signs monitoring for sleep Apnea patients: a preliminary study. IEEE Access 6:2506–2514. https://doi.org/10.1109/ACCESS.2017.2783939

19. Ullal et al (2021) Non-invasive monitoring of vital signs for older adults using recliner chairs. Health Technol 11:169–184. https://doi.org/10.1007/s12553-020-00503-9

20. Gaiduk M et al (2022) Estimation of sleep stages analyzing respiratory and movement signals. IEEE J Biomed Health Inf 26(2):505–514. https://doi.org/10.1109/JBHI.2021.3099295

# Definition of Emotional States Interval for Application of Artificial Intelligence and Stress Estimation

Wilhelm Daniel Scherz[1,2(✉)] , Juan J. Perea[2] , Ralf Seepold[1] ,
and Juan Antonio Ortega[3]

[1] HTWG Konstanz, Alfred-Wachtel-Str. 8, 78462 Konstanz, Germany
`wscherz@htwg-konstanz.de`
[2] Universidad de Loyola Andalucía, Av. de Las Universidades, 2, 41704 Dos Hermanas, Sevilla, Spain
[3] University of Seville, Avda. Reina Mercedes S/N, Seville, Spain

**Abstract.** The perception of the amount of stress is subjective to every person, and the perception of it changes depending on many factors. One of the factors that has an impact on perceived stress is the emotional state. In this work, we compare the emotional state of 40 German driving students and present different partitions that can be advantageous for using artificial intelligence and classification. Like this, we evaluate the data quality and prepare for the specific use. The Stress Perceived Questionnaire (PSQ20) was employed to assess the level of stress experienced by individuals while participating in a driving simulation for 5 and 25 min. As a result of our analysis, we present a categorisation of various emotional states into intervals, comparing different classifications and facilitating a more straightforward implementation of artificial intelligence for classification purposes.

**Keywords:** Stress · Stress perceived questionnaire · Driving

## 1 Introduction

Stress has emerged as a topic of immense significance and interest in contemporary society [1], and it has been recognised by organisations like the WHO even as a disease [1]. Stress significantly influences driving behaviour, with a weighty influence on security and safety during driving and how people drive [2]. Drivers may feel stressed when they feel like they are not in control of their vehicle or the driving situation. Stress can be caused by various factors, such as traffic, bad weather, aggressive drivers or emotional state. Additional factors that also have an essential influence on rod safety are stress, fatigue and sleepiness [3, 4].

The German Federal Statistical Office (Destatis) reported 2,569 fatalities due to traffic accidents in Germany in 2021. This number decreased to 2,719 in 2020 but increased to 3,046 in 2019, and according to the World Health Organization, approximately 1.3 million people die each year due to traffic accidents worldwide [5]. The role of stress in

risk evaluation is significant, making it an essential factor to understand. Studies show that stress can impair cognitive functions [6], it can lead to risky driving behaviour [7] and make it challenging to judge risks [8, 9].

Stress can be defined as anything that disturbs or affects the internal equilibrium of the body, and a stressor is an external stimuli, perceived threads or demands that can create stress. This balance is called homeostasis. In addition to external stressors, there are also internal stressors. Internal stressors such as worries can also trigger the stress response of the body and disturb the internal equilibrium [10]. Among the widely known adverse effects of stress on the body are: degeneration of the cardiovascular system [11, 12], the inappropriate response of the immune system [13, 14] and reduced capacity to cope with physical and mental demands [15].

Generally, stress is often categorised into the following characteristics: Insensitivity, duration, type and perception.

This research aims to introduce a method for dividing perceived stress obtained from the PSQ20 into discrete intervals, thereby optimising the classification process across diverse sections enabling the use of advanced machine learning and artificial intelligence algorithms in the future for the classification of stress depending on the emotional state. The research should help to support a more straightforward implementation of artificial intelligence.

## 2  State of the Art

In todays society, health awareness is increasingly important, and the collection of vital biological data is becoming more ordinary. This awareness is also true in the context of stress management, as people are seeking new ways to track and manage their stress levels [16]. There are mainly different methods for measuring stress. Those methods can be resumed in these categories.

Questioners and self-report. This method asks the participants to rate their stress levels. Usually, it is done using questionnaires or a journal. Some examples of this include the Perceived Stress Questionnaire (PSQ20) [17] and the Perceived Stress Scale (PSS) [18].

Measurement by biochemical hormones: This method uses the advantage that stress releases hormones, i.e. cortisol and adrenalin. These hormones can be measured by taking blood, urine or saliva [19].

Sensorics measurements. Like the last method, stress also influences the behaviour of different physiological systems. Examples include heart rate variability, skin conductivity [20] and more physiological responses.

Behavioural measurement. This method uses the subtle differences in behaviour between when a person is stressed or not is stressed. A good example of behavioural measurement is stress estimation by analysing the typing style [21].

Nonetheless, using a restricted set of parameters can pose challenges in distinguishing between healthy physical activity and adverse factors like stress. The PSQ20 questionnaire offers a distinct advantage in its simplicity when used for stress estimation and can be added to every data collection.

## 3    Methodology

To obtain the needed data to categorise various emotional states into intervals, a group of German students with driving experience was selected (n = 40). To collect biovital data, a polar H10 chest strap was given to each participant. The focus was on the Heart Rate and the R-R Intervals (interbeat intervals). Before starting the experiment, the procedure was explained in detail to the participants and the voluntary agreement was collected. At the beginning of the experiment, the PSQ20 questionnaire was given to the participants, and they had to answer it. There was no time limit to answer the questionnaire. After the questionnaire was filled, the participant could sit on the simulator and listen to 5 min of relaxing music. After this, the participant had to start driving for 25 min. During this phase, different traffic challenges were used to act as stressors.

A total of 35 male and 5 female students in Germany participated in this study. The mean age of the participants was 25.84 ± 4.83. All participants had experience driving and a valid driving license. The participants were recruited voluntarily and had no financial compensation or benefit. The participants had to fulfil the following criteria: Subjects had to be healthy, have no history of heart problems, have good vision or use glasses to correct the vision. Additionally, the participants must be sober and did not consume drugs or medicaments before the experiment.

### 3.1    Driving Simulator

The City Car Driving simulator was used for this experiment. The main advantages of this simulator are the advanced and realistic car physics that provides a reasonable degree of realism. Additionally, the traffic conditions can be configured.

For this experiment, the traffic was set to normal, but events like car driving the opposite line, accidents, and randomly crossing pedestrians were activated. Additionally, when the driver is not driving according to the driving rules, a signal is released informing him about his mistake.

### 3.2    Perceived Stress Questionnaire

The Perceived Stress Questionnaire (PSQ20) [17] is a frequently used questionnaire for estimating perceived stress. There are two forms of the questionnaire: the regular one with about 30 questions and the shorter one with 20 Questions. The shortened version of the questionnaire evaluates the following four emotional states: Worries, Tension, Joy and Demand. Higher scores in the PSQ20 mean that this emotional state is stronger. Dividing this in the quantiles helps change the score from a numerical value to intervals that are more suited for machine learning (ML) approaches.

## 4   Results

The results intend to present a division of the perceived stress into different intervals to enable a more efficient classification into different regions using machine learning / artificial intelligence algorithms. Only complete questionnaires, participants with signed agreements and recorded biovital data were analysed. Because of these requirements, three datasets were discarded, reducing the sample to 40 participants.

Figure 1 divides the participants into intervals of 5 quantiles. Figure 2 divides the participants into four quantiles, and Fig. 3 into three quantiles. As we can observe, the differences between the worries quantiles are more significant than those in Figs. 2 and 3.



| | 0-20% | 20.01-40.00% | 40.01-60.00% | 60.01-80.00% | 80.01-100.00% |
|---|---|---|---|---|---|
| Worries | 8 | 19 | 10 | 2 | 1 |
| overall score | 2 | 21 | 10 | 7 | 0 |

**Fig. 1.** Intervals of five quantiles

A similar pattern is evident for the overall score in Figs. 1, 2 and 3. This suggests that the discernible differences between the various intervals may suffice for training artificial intelligence to perform classification tasks.

| | 0-25% | 25.01-50.00% | 50.01-75.00% | 75.01-100.00% |
|---|---|---|---|---|
| Worries | 14 | 16 | 8 | 2 |
| overall score | 7 | 24 | 8 | 1 |

**Fig. 2.** Intervals of four quantiles



| | 0-33.33 | 33.34-66.67 | 66.68-100 |
|---|---|---|---|
| Worries | 22 | 15 | 3 |
| overall score | 13 | 24 | 3 |

**Fig. 3.** Intervals of tree quantiles

## 5  Conclusion

We have conducted a thorough study to quantify the stress levels of each patient in different domains, such as worries, demands, tension, joy, and overall state. Furthermore, we have developed a classification of the values in these domains for each patient and compared different divisions of these values: three, four, and five divisions. Our analysis

has shown that the most relevant domains are worries and overall state. Based on the hypothesis that more significant variability provides more information, we have shown that the optimal classification is achieved with five divisions, as it maximises the information due to its higher variability. This behaviour and this division could be used for the age range of $25.84 \pm 4.83$ of healthy persons. A validation with other ranges of ages could prove this separation as applicable.

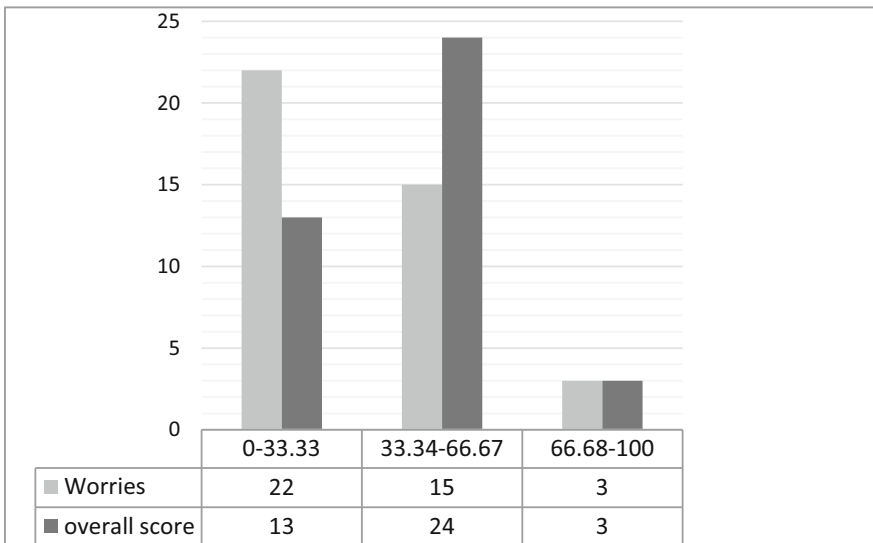Looking at Fig. 1 we can observe that more participants are located in the second quantile, telling us there is a bigger group of drivers (21 drivers) with a medium lower amount of stress. If we look at Fig. 3 we can observe that most divers (24 drives) fit in the medium quantile, telling us that most drivers are in the middle. Comparing Figs. 1 and 3 we can conclude that the smaller division in more quantiles can be advantageous for a more price description of the drivers.

## 6 Future Works

Once we have established and justified this arrangement for classifying the stress levels of patients, a future line of research could involve using this classification to train artificial intelligence systems. The process would involve training the system with the proposed classification and, based on its variability, analysing data to classify patients. The final goal would be to develop a system capable of automatically classifying patients into one of the five categories, thereby determining their stress levels without requiring direct human intervention. This would restructure the analysis of patients stress levels and could be complemented with other analytical techniques, such as the Poincaré Plot.

In addition to the collected biovital parameters and the information obtained from the questionnaires, additional information like demographic factors could deepen and help to understand deeper the nuances and characteristics of stress and enable a better classification.

## References

1. Andrade L, Caraveo-Anduaga JJ, Berglund P, Bijl R, Kessler RC, Demler O, Walters E, Kylyc C, Offord D, Ustun TB, Wittchen HU (2000) Cross-national comparisons of the prevalences and correlates of mental disorders. WHO Int Consort Psychiatr Epidemiol. Bull World Health Organiz 78:413–426
2. Yu Z, Yu T, Ge Y, Qu W (2023) The effect of perceived global stress and altruism on prosocial driving behavior, yielding behavior, and yielding attitude. Traffic Inj Prev 24:402–408. https://doi.org/10.1080/15389588.2023.2191765
3. Radun I et al (2022) Sleepy drivers on a slippery road: a pilot study using a driving simulator. J Sleep Res 31:e13488. https://doi.org/10.1111/jsr.13488
4. Taylor AH, Dorn L (2006) Stress, fatigue, health, and risk of road traffic accidents among professional drivers: the contribution of physical inactivity. Annu Rev Public Health 27:371–391. https://doi.org/10.1146/annurev.publhealth.27.021405.102117
5. Useche SA, Ortiz VG, Cendales BE (2017) Stress-related psychosocial factors at work, fatigue, and risky driving behavior in bus rapid transport (BRT) drivers. Accident Anal Prevent 104:106–114. https://doi.org/10.1016/j.aap.2017.04.023

6. Kulshreshtha A, Alonso A, McClure LA, Hajjar I, Manly JJ, Judd S (2023) Association of stress with cognitive function among older black and white US adults. JAMA Netw Open 6:e231860. https://doi.org/10.1001/jamanetworkopen.2023.1860

7. Delhomme P, Gheorghiu A (2021) Perceived stress, mental health, organisational factors, and self-reported risky driving behaviors among truck drivers circulating in France. J Safety Res 79:341–351. https://doi.org/10.1016/j.jsr.2021.10.001

8. Pooladvand S, Hasanzadeh S (2023) Impacts of stress on workers' risk-taking behaviors: cognitive tunneling and impaired selective attention. J Constr Eng Manag 149. https://doi.org/10.1061/JCEMD4.COENG-13339

9. Reynolds EK, Schreiber WM, Geisel K, MacPherson L, Ernst M, Lejuez CW (2013) Influence of social stress on risk-taking behavior in adolescents. J Anxiety Disord 27:272–277. https://doi.org/10.1016/j.janxdis.2013.02.010

10. Morrissey C et al (2023) Fear of cancer recurrence inventory scores and their correlation with quality of life and stress levels in breast cancer survivors. JCO 41:e24073–e24073. https://doi.org/10.1200/JCO.2023.41.16_suppl.e24073

11. Zaffalon Júnior JR, Viana AO, Melo GEL, de Angelis K (2018) The impact of sedentarism on heart rate variability (HRV) at rest and in response to mental stress in young women. Physiol Rep 6:e13873. https://doi.org/10.14814/phy2.13873

12. Mei Y, Thompson MD, Cohen RA, Tong X (2015) Autophagy and oxidative stress in cardiovascular diseases. Biochem Biophys Acta 1852:243–251. https://doi.org/10.1016/j.bbadis.2014.05.005

13. Martínez Fernández J, Augusto JC, Trombino G, Seepold R, Martinez Madrid N (2013) Self-aware trader: a new approach to safer trading. J Univ Comput Sci

14. Barrows IR, Ramezani A, Raj DS (2019) Inflammation, immunity, and oxidative stress in hypertension-partners in crime? Adv Chronic Kidney Dis 26:122–130. https://doi.org/10.1053/j.ackd.2019.03.001

15. Cody K, Scott JM, Simmer-Beck M (2022) Examining the mental health of university students: a quantitative and qualitative approach to identifying prevalence, associations, stressors, and interventions. J Am College Health: J ACH 1–11. https://doi.org/10.1080/07448481.2022.2057192

16. RaviPrakash H, Anwar SM (2023) Do it yourself. In: Byrne MF, Parsa N, Greenhill AT, Chahal D, Ahmad O, Bagci U (eds) AI in clinical medicine. Wiley, pp 94–103. https://doi.org/10.1002/9781119790686.ch10

17. Fliege H et al (2005) The perceived stress questionnaire (PSQ) reconsidered: validation and reference values from different clinical and healthy adult samples. Psychosom Med 67:78–88. https://doi.org/10.1097/01.psy.0000151491.80178.78

18. Lee E-H (2012) Review of the psychometric evidence of the perceived stress scale. Asian Nurs Res 6:121–127. https://doi.org/10.1016/j.anr.2012.08.004

19. El-Farhan N, Rees DA, Evans C (2017) Measuring cortisol in serum, urine and saliva—Are our assays good enough? Ann Clin Biochem 54:308–322. https://doi.org/10.1177/0004563216687335

20. Najström M, Jansson B (2007) Skin conductance responses as predictor of emotional responses to stressful life events. Behav Res Ther 45:2456–2463. https://doi.org/10.1016/j.brat.2007.03.001

21. Gunawardhane SDW, Silva PM, de Kulathunga DSB, Arunatileka SMKD (2013) Non invasive human stress detection using key stroke dynamics and pattern variations. In: 2013 international conference on advances in ICT for emerging regions (ICTer). IEEE, pp 240–247. https://doi.org/10.1109/ICTer.2013.6761185

# Novel Battery Parallelization Approach Using DC/DC Partial Power Converter in Micro-Grids

Gianluca Simonte, Roberto Di Rienzo$^{(\boxtimes)}$, Niccolò Nicodemo, Alessandro Verani, Federico Baronti, Roberto Roncella, and Roberto Saletti

Dipartimento Ingegneria Informazione, University of Pisa, Pisa, Italy
`roberto.dirienzo@unipi.it`

**Abstract.** The electrical energy generation is one of the main pollution causes as it is still largely based on oil and coal generators. Renewable energy sources are the most mature alternative but their production is intermittent and then require energy storage systems able to store a high quantity of energy to increase their usability. The second-life batteries of electric vehicles and low-cost battery chemistries are the best candidates to compose low-cost energy storage systems for these applications. These batteries must be parallel-connected to obtain the required capacity. A novel parallelization approach based on Input Parallel/Output Serial DC/DC Partial Power converter series-connected to each battery is presented in this paper. This approach is applied to different case studies showing that controlling only the 15% of the battery power allows the DC/DC converter to control the sharing of the grid current among the batteries. The reduction of the controlled power leads to a reduction of the cost of the energy storage system.

**Keywords:** Battery parallelization system · Micro grid · Partial Power DC converter · IPOS DC/DC · Second life battery

## 1 Introduction

Solar and wind energies are the best candidates to reduce the environmental impact of electrical energy generation. However, they are intermittent and not controllable by their nature. Smart grids combined with Energy Storage Systems (ESSs) aim to solve these limitations and to make more attractive renewable energy sources. The presence of an energy storage system, however, increases the grid cost and poses major ethical issues due to controversial supplying manufacturing processes [1].

Smart-grid costs can be reduced by using second-life batteries from electric mobility applications instead of new ones. In fact, the remaining capacity of second-life batteries can be exploited by stationary applications that are less demanding in terms of energy and power density than mobility ones. Usually, the capacity of one vehicle battery is much lower than the smart-grid ESS required

one. Therefore, some second-life vehicle batteries are parallel connected to compose the smart-grid ESS [2]. Unfortunately, second-life batteries or their elementary modules cannot be parallel connected using static hardware connection or simple switch-based circuits [3]. In fact, the high variability of second-life battery parameters, the non-uniform aging of the elements to be parallel-connected, and the possibility to compose the ESS with different battery types and manufacturers, impose the use of a smart parallelization system [4,5]. Moreover, an intelligent parallelization system allows us to compose the ESS using different battery chemistries, such as lead-acid and sodium-metal halide batteries. These kinds of batteries are generally less performing than the lithium-ion ones but are also cheaper.

The parallelization is historically obtained using the AC bus of the electric power distribution network and one AC/DC converter for each parallel-connected battery. This approach minimizes the management complexity of the smart-grid but the obtained energy efficiency is low. The AC/DC converters can be replaced with high-efficiency DC/DC converters if the battery parallelization is carried out with a DC bus [6]. This approach increases the system energy efficiency but requires high power and expensive DC/DC converters.

This paper explores the possibility to parallel connect the batteries using DC/DC Partial Power converters. This typology of converters is series-connected with the battery and is generally used to control the charging current profile [7]. Our idea, instead, aims to use one Input Parallel/Output Serial (IPOS) DC/DC Partial Power converter for each parallel connected battery to control the load/source current sharing among them. This idea generalizes the work of Iyer et al. [8], where a single battery is taken into account. A preliminary smart-grid simulation platform was developed in Matlab Simulink environment to demonstrate the capability of this innovative battery parallelization approach and quantify its advantages with respect to a classic DC/DC converter approach.

## 2   Novel Parallelization Approach

The block diagram of the proposed parallelization circuit is reported in Fig. 1. It uses an IPOS DC/DC Partial Power converter for each battery to control its current. The IPOS converter is a two-ports system. One port is parallel-connected with the battery while the other one is series-connected to it. The ports are modeled as a current generator $I_{bx}$ and a voltage one $V_{cx}$, respectively. Instead, the batteries are modeled as a voltage generator, which models the relation between the battery Open Circuit Voltage (OCV) with its State of Charge (SoC), series-connected with an impedance $Z$. The system current equations are reported in the system (1).

$$\begin{cases} I_{ESS} = \sum_{x=1}^{n} I_{cx} \\ I_{cx} = I_{batx} - I_{bx} \\ I_{batx} = \frac{V_{bus} - OCV - V_{cx}}{Z} \end{cases} \quad (1)$$
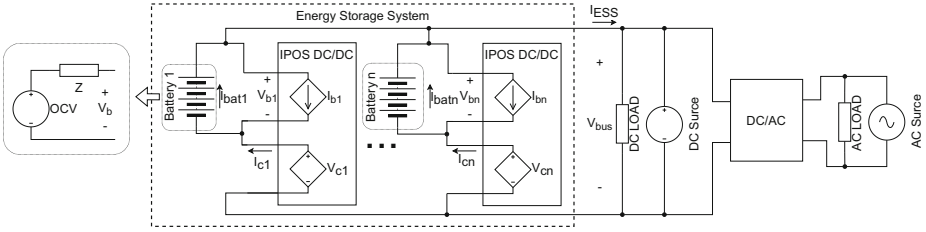
**Fig. 1.** Block diagram of a smart grid with the proposed parallelization system for the ESS.

As we can note, the current of each battery is controlled by the relative $V_c$. The battery current is equal to zero if $V_c + OCV$ is equal to $V_{bus}$. Instead, the battery is charged or discharged if $V_c + OCV$ is lower or higher than $V_{bus}$, respectively. In other words, the IPOS DC/DC can control the power exchange between the network and the battery with a very low power if the system is sized with a bus voltage $V_{bus}$ barely higher than the maximum OCV of the batteries. This is the main advantage of the proposed approach compared with the classic parallelization one. In fact, a DC/DC converter with two ports, one parallel connected to the battery and the other parallel connected to the grid, is used to control the battery power in classic parallelization approaches. In this case, the DC/DC converter needs to be able to control all the battery power and then it is more expensive than IPOS one. Higher is the DC/DC converter power, higher is its cost. Moreover, IPOS topologies are simpler than a bidirectional high-efficiency DC/DC converter [7,9].

Finally, the DC/DC efficiency can be taken into account using the equation system (2) where $V_{bx}I_{bx}$ and $V_{cx}I_{cx}$ are the input/output power of the DC/DC ports and $\eta_{bc}$ and $\eta_{cb}$ are the DC/DC efficiency for energy exchange from the battery port to the control port and vice versa, respectively.

$$\begin{cases} \eta_{bc} V_{bx} I_{bx} = V_{cx} I_{cx} & for \quad I_{cx} \geq 0 \\ V_{bx} I_{bx} = \eta_{cb} V_{cx} I_{cx} & for \quad I_{cx} < 0 \end{cases} \tag{2}$$

## 3   Simulation Platform

A simulation platform of the model described in Fig. 1 is developed in Matlab Simulink environment to verify the feasibility of the presented parallelization approach with realistic case studies. The platform is kept as simple as possible in this phase to reduce the complexity of the analysis. For example, the battery impedances are considered as purely resistive, and the grid sources and loads are emulated with a constant current profile during both charge and discharge phases. Three batteries and IPOS converters are modeled. The converters are controlled to keep the SoC of the batteries as equal as possible. Two different case studies are considered. The former is focused on second-life batteries while the latter is focused on the coexistence of different battery chemistries.
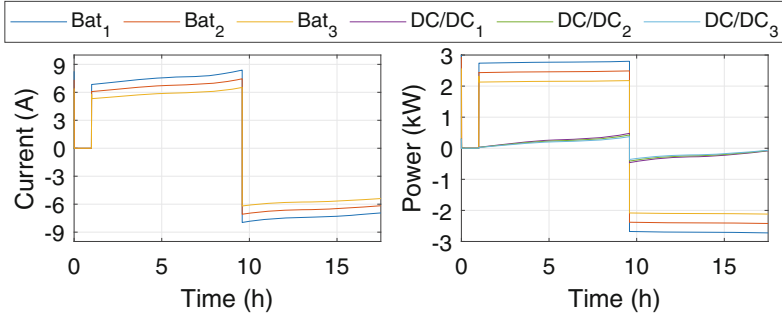
**Fig. 2.** Simulation results of the first case study in which the ESS is composed of three second-life NMC batteries with 90 ($Bat_1$), 80 ($Bat_2$), and 70%($Bat_3$) of the nominal capacity.

In the first case study, three second-life batteries composed of 96 Nickel Manganese Cobalt (NMC) series-connected cells with a nominal capacity ($Q_N$) of 60 Ah are used. The actual capacities of the three batteries are set to 70, 80, and 90% of the nominal one to emulate different battery aging levels. The resistance values are instead increased keeping the product of battery resistance and capacity constant. The OCV and resistance curves as a function of SoC for the NMC cell are obtained from [10].

The three different battery chemistries, NMC, lithium iron phosphate (LFP), and Sodium-Metal Halide (ZEBRA), are instead considered in the second case study. For the sake of simplicity, the three batteries are sized to have comparable capacity and nominal voltage of 50 Ah and 380 V, respectively. This case study highlights the ability of the parallelization approach to manage batteries with large differences in the OCV curves.

## 4   Simulation Results

The developed simulation platform was used to simulate the two presented case studies. The batteries start from a fully charged condition and are discharged with constant current profiles until they reach 10% of SoC. Then, the batteries are fully charged with a constant current profile. The grid power values are 7.3 and 5.3 kW for the first and second case studies, respectively. The DC/DC converter $\eta_{cb}$ and $\eta_{cb}$ were set to 0.88, which is a reasonable value for the IPOS converters as reported in [11].

Figure 2 shows the current and the power profiles of the three batteries and the power profiles of the three DC/DC converter control ports in the first case study. As we can note, the current values in the charge and discharge phases are directly related to the battery capacities. In fact, their values are controlled by the DC/DC converters to keep the battery SoC as balanced as possible.

Moreover, the figure highlights that the mean DC/DC output power is only the 8.4% of the battery output power. As mentioned before, this is the
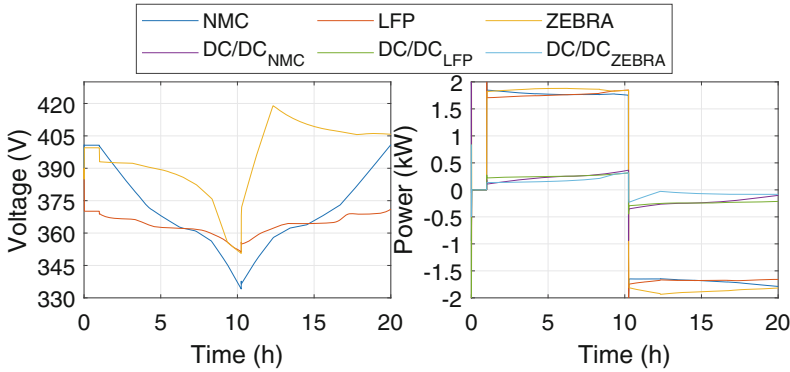
**Fig. 3.** Simulation results of the second case study in which the ESS is composed of three batteries based on NMC, LFP, and ZEBRA chemistry.

main advantage of the proposed parallelization approach which uses low-power DC/DC converter to control a much higher battery power.

Higher battery voltage differences appear in the second case study during both the charge and discharge phase as reported in Fig. 3. The voltage differences strongly change with the battery SoC and are due to the large difference in OCV and parameter curves of the three different battery chemistries.

In this case study, the output power of the DC/DC converters is about 12.8, 14.3, and 7.1% of the battery ones with a maximum control voltage of 21.4, 17.6, and 17.4% of the bus voltage (425 V) for the NMC, LFP, and ZEBRA batteries, respectively.

## 5    Conclusion

A novel battery parallelization approach based on Input Parallel/Output Serial DC/DC Partial Power converter series-connected to the battery has been presented in this paper. This technique allows us to control the current that flows through the batteries with DC/DC output powers that are lower than 15% of the battery one as demonstrated with a simulation platform in two different case studies. The very low DC/DC power requirement is a significant advantage of the proposed solution compared to the classic parallelization one. In fact, the classic solution uses parallel connected DC/DC converters that must be able to control all the battery power and are therefore highly expensive. On the other hand, the proposed parallelization approach allows us to control very high power with low-power DC/DC, resulting in a considerable reduction of the system costs.

Future works will aim to increase the simulation platform accuracy by including additional battery behaviors such as temperature and aging. Moreover, more complex and exhaustive case-studies should be considered. Furthermore, we will carry out a comparison between the Partial Power DC/DC converter topologies in order to find the most promising one for implementing an experimental setup.

# References

1. Nature: Lithium-ion batteries need to be greener and more ethical. Nature 595(7865):7 (2021)
2. Gohla-Neudecker B, Bowler M, Mohr S (2015) Battery 2nd life: leveraging the sustainability potential of EVs and renewable energy grid integration. In: 2015 international conference on clean electrical power (ICCEP), pp 311–318
3. Di Rienzo R, Baronti F, Roncella R, Morello R, Saletti R (2017) Simulation platform for analyzing battery parallelization. In: SMACD 2017–14th international conference on synthesis, modeling, analysis and simulation methods and applications to circuit design. Institute of Electrical and Electronics Engineers Inc
4. Goetz SM, Li Z, Liang X, Zhang C, Lukic SM, Peterchev AV (2017) Control of modular multilevel converter with parallel connectivity-application to battery systems. IEEE Trans Power Electron 32(11):8381–8392
5. Bruen T, Marco J, Gama M (2015) Current variation in parallelized energy storage systems. In: 2014 IEEE vehicle power and propulsion conference, VPPC 2014. Institute of Electrical and Electronics Engineers Inc
6. Boi M, Mastromauro RA, Floris A, Damiano A (2023) Integration of sodium metal halide energy storage systems in telecommunication microgrids: performance analysis of DC-DC converter topologies. Energies 16(5)
7. Lopusina I, Grbovic P (2021) Comparative analysis of input-series-output-series partial power rated DC to DC converters. In: Proceedings of 2021 21st international symposium on power electronics, Ee 2021. Institute of Electrical and Electronics Engineers Inc
8. Iyer VM, Gulur S, Bhattacharya S, Ramabhadran R (2019) A partial power converter interface for battery energy storage integration with a DC microgrid. In: 2019 IEEE energy conversion congress and exposition, ECCE 2019, pp 5783–5790. Institute of Electrical and Electronics Engineers Inc. (2019)
9. Zanatta N, Caldognetto T, Biadene D, Spiazzi G, Mattavelli P (2023) A two-stage DC-DC isolated converter for battery-charging applications. IEEE Open J Power Electron 4:343–356. https://doi.org/10.1109/OJPEL.2023.3271227
10. Morello R, Di Rienzo R, Roncella R, Saletti R, Baronti F (2018) Hardware-in-the-loop platform for assessing battery state estimators in electric vehicles. IEEE Access 6:68210–68220
11. Cao Y, Ngo M, Yan N, Dong D, Burgos R, Ismail A (2022) Design and implementation of an 18-kW 500-kHz 98.8% efficiency high-density battery charger with partial power processing. IEEE J Emerg Sel Top Power Electron 10(6), 7963–7975

# Automated Parking in CARLA: A Deep Reinforcement Learning-Based Approach

Luca Lazzaroni[✉] , Alessandro Pighetti , Francesco Bellotti ,
Alessio Capello , Marianna Cossu , and Riccardo Berta

Department of Electrical, Electronic and Telecommunication Engineering (DITEN), University
of Genoa, Via Opera Pia 11a, 16145 Genoa, Italy
{luca.lazzaroni,alessandro.pighetti,alessio.capello,
marianna.cossu}@edu.unige.it, {francesco.bellotti,
riccardo.berta}@unige.it

**Abstract.** This paper focuses on developing a Deep Reinforcement Learning (DRL)—based agent for real-time trajectory planning and tracking in a simulated parking environment, specifically low-speed maneuvers in a parking area with comb-shaped spaces and a random distribution of non-player vehicles. We rely on CARLA as a virtual driving simulator due to its realistic graphics and physics simulation capabilities, and on the Gymnasium and Stable-Baselines3 toolkits for training the agent. We show that the agent is able to achieve a success rate of 97% in reaching the target parking lot without collisions. However, integrating CARLA with DRL frameworks poses challenges, such as determining suitable environment and neural network update frequencies. Despite these issues, the results demonstrate the potential of DRL agents in developing automated driving functions.

**Keywords:** Automated driving functions · CARLA driving simulator · Deep reinforcement learning · Gymnasium · Automated parking · Proximal-policy optimization

## 1 Introduction

A growing trend in the development of automated driving functions (ADFs) is the use of DRL agents that learn specific tasks by interacting with the target environment following a trial-and-error mechanism and supported by a deep neural network brain [1]. Due to the inherent nature of this approach, it is useful to train them within virtual driving simulators that are as true to reality as possible. This allows for safe operating environments where the agent can learn the desired behavior, shaped through the awarding of rewards and penalties based on the status of the agent within the environment. Indeed, the objective of a DRL agent is to maximize the sum of expected rewards (i.e., the cumulative reward) over its lifetime. By exploiting the information learned about the expected utility (i.e., the sum of expected future rewards discounted by a factor), the agent can increase its long-term reward. This implies that the contributions that make up the reward function (RF) have to be specifically designed to enable the agent to achieve the target policy.

We are interested in developing a DRL agent able to perform real-time trajectory planning and tracking in a high-fidelity simulated parking environment, taking as input sensors data and consequently acting on the steering wheel and the throttle and brake pedals, emulating the actions of a human driver. The focus of our research is low-speed maneuvers in a parking area with comb parking spaces and random distribution of the number of parked non-player vehicles (NPVs).

The remainder of the paper is structured as follows. Section 2 presents the environment, then Sect. 3 shows the performed experiment, while Sect. 4 the obtained results. Finally, conclusions are drawn in Sect. 5.

## 2   Environment

To train our model, we took a step toward more realistic simulations than previously done in [1, 2] for parking (using Unity) and highway driving (using highway-env [3]), respectively. To increase the level of realism in terms of graphics and physics simulation, we opted for CARLA as a driving simulator [4], an open-source software created to support ADFs development. It includes several town maps and different vehicle models to which sensors such as radar, lidar, camera, etc. can be attached. In addition, CARLA exposes an API that allows users to control all aspects related to the simulation, including traffic, pedestrian behaviors, weather, sensors, etc. The official DRL-related section lacks recent updates but some recent works have recently been proposed [5, 6], showing promising results in trajectory tracking applications.

Starting from the built-in Town 5, where a parking area is already available, we removed all the unnecessary surrounding entities from the scene, so to obtain a lightweight scenario to reduce the computational cost of rendering and thus training time. The parking area, visible in Fig. 1, measures $50 \times 50$ m and is composed of 60 comb parking spaces with a 90° layout, $5 \times 2.5$ m each. Brick walls cordon off the area, preventing the agent from falling during training. The weather includes 15 different conditions of sun, rain, fog, and clouds and 41 vehicle models are available, from motorbikes to trucks, with customizable colors.



**Fig. 1.** Bird view of the parking area. The red square marks the target parking lot.

# 3  Experiment

To implement our DRL agent, we opted for Gymnasium (formerly known as OpenAI Gym [7]), an open-source Python toolkit that exposes an API to facilitate the design and implementation of DRL environments which has become a de facto standard for the communication between agents and simulation environments. In addition, Gymnasium offers the ability to define custom environments, so, we wrapped the CARLA parking setting into a Gymnasium environment to get a DRL-compatible interface. Furthermore, Gymnasium is compatible with Stable-Baselines3 (SB3) [8], an open-source Python library that offers a variety of DRL algorithms to set up the training algorithm and the neural network architecture and configuration quickly and intuitively. The schematic of the tools used and their interfacing are shown in Fig. 2.
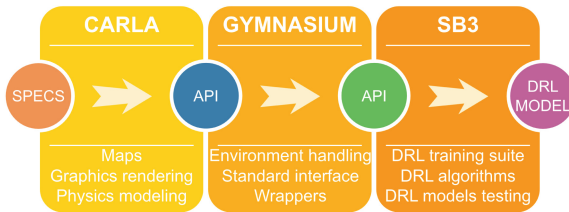


**Fig. 2.** Outline of the tools used for the implementation of the DRL agent.

At each episode, the vehicle starting point and orientation are chosen randomly within the drivable area. The target to be reached is one of the 60 lots, randomly chosen at each episode. To train our agent, we used the Audi e-tron vehicle model, available in CARLA. The car measures $4.90 \times 1.93 \times 1.62$ m, and has a total mass of 2490 kg. The vehicle dynamics are modeled by CARLA with NVIDIA PhysX [9], one of the most popular physics simulation engines. The ego vehicle (EV) is equipped with a lidar sensor placed in the center of the vehicle with a 360° horizontal field of view.

The key part in the development of a DRL agent is RF. We split it into several contributions, some dense (given at each simulation step) and some sparse (awarded only when some events occur). Table 1 offers an overview the components of the RF. Distance and heading are the two dense contributions to induce EV to approach the target. Collision is sparse, needed to teach the agent to avoid other NPVs and the walls delimiting the area. Two contributions are related to reaching the target parking lot, goal and alignment. The goal reward is awarded when EV reaches the lot, while alignment rewards the agent if it parks aligned well.

For the training of the model, we used the PPO algorithm [10]. PPO aims to optimize policies for sequential decision-making tasks. It balances exploration and exploitation by iteratively updating policy parameters based on the clipped surrogate objective, ensuring stability and reliable performance. As the backbone neural network, we opted for a multi-layer perceptron configuration, composed of 2 hidden layers of 512 neurons, whose input and output values are summed up in Table 2. To give the agent the ability to better perceive its motion, at each network update we stacked the last four input values on top

**Table 1.** RF components.

| Name | Description | Range | Type |
|------|-------------|-------|------|
| Distance | Euclidean distance from EV to target | $[-0.1, 0]$ | Dense |
| Heading | Angle between EV forward vector and EV-target vector | $[-0.1, 0]$ | Dense |
| Collision | Collision with walls or NPVs event | $[-1, 0]$ | Sparse |
| Goal | Target achieved event | $[0, 15]$ | Sparse |
| Alignment | Angle between EV forward vector and parking lot orientation | $[0, 10]$ | Sparse |

of the current one. This allows the agent to capture the dynamics and changes in the environment over time.

**Table 2.** PPO network input and output.

| Type | Quantity | Dimension | Resulting size |
|------|----------|-----------|----------------|
| Input | Lidar data | 61 | 340 (5 stacked inputs) |
| | Distance from target (x, y) | 2 | |
| | Speed (x, y) | 2 | |
| | Acceleration (x, y) | 2 | |
| | Angular velocity (yaw) | 1 | |
| Output | Throttle | 1 | 4 |
| | Steering | 1 | |
| | Brake | 1 | |
| | Reverse | 1 | |

Given the results previously obtained in an environment similar to ours, although significantly less complex, such as [1], we adopted a curriculum learning (CL) strategy for training our agent [11]. In this case, we split the training into three distinct phases. The first phase involves only EV, no NPV is present and an episode terminates if the target is reached or if it goes over 240 steps. Collision is not penalized to incentivize the agent to explore the environment. In the second phase, the number of NPVs is chosen randomly between 0 and 20 at each episode and, again, only the goal reaching or the timeout can end the episode. Collisions are now slightly penalized with a weight equal to $-0.1$. Finally, in the third phase, conditions are similar to phase #2 but now a collision (penalty weight $-1$) terminates the episode.

## 4  Results

We trained our PPO agent for about 60M steps (Fig. 3). The first CL phase took 25M steps, after which the agent was able to reach the goal 100% of the time (Fig. 3a). Once this is achieved with no NPVs present in the scene, the second phase started, lasting about 15M steps in which the agent learned to park with NPVs parked near the target parking lot. The last phase, in which collisions represented a terminal state, lasted 20M steps, achieving a final success rate of 97%.
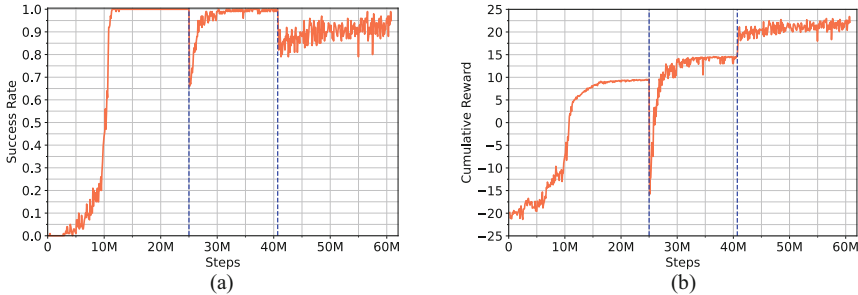


**Fig. 3.** The PPO agent training over ~60M steps. **a** is the success rate and **b** the cumulative reward. Dashed blue lines distinguish the three CL phases.

Although the success rate is satisfactory, some problems related to the use of CARLA arose during code development. In particular, the lack of direct support for DRL made the implementation of the environment and synchronization with Gymnasium and SB3 particularly complex. Too high values of the CARLA environment update frequency (greater than 20 Hz) greatly lengthened the training time due to the higher computational load but, conversely, frequencies that were too low (lower than 10 Hz) did not guarantee proper observation of the environment by the agent (e.g., it happened that the agent stepped on the goal without being detected). In addition, it was necessary to adopt the decision period (i.e., the frequency at which the agent takes an action) since too high network update frequencies (greater than 10 Hz) did not allow the agent to learn to move, since too short a press on the accelerator pedal is not sufficient to allow movement. A good trade-off we found is to use a CARLA environment update rate of 20 Hz and a network update rate of 5 Hz (equivalent to a decision period of 4) but this required a lot of trial and error.

## 5  Conclusions

This paper presented a DRL agent for real-time trajectory planning and tracking in a CARLA-simulated environment where low-speed maneuvers are performed in a parking area with comb-shaped spaces. The agent is able to achieve a 97% success rate in reaching the target parking lot without collisions. The training process involved a CL strategy with three distinct phases that gradually increased the complexity of the environment by introducing NPVs and penalizing collisions. However, the use of CARLA presented

challenges in integrating DRL into the environment. Issues related to the environment update frequency and decision period had to be carefully handled to ensure proper learning and observation by the agent but the overall success rate and performance of the trained agent confirm the capability of DRL approaches for the development of ADFs. Future work will focus on different sensor configurations and extend the range of action to more complex parking scenarios or other ADFs.

# References

1. Lazzaroni L, Bellotti F, Capello A, Cossu M, De Gloria A, Berta R (2023) Deep reinforcement learning for automated car parking. In: Berta R, De Gloria A (eds) Applications in electronics pervading industry, environment and society. Springer Nature Switzerland, Cham, pp 125–130. https://doi.org/10.1007/978-3-031-30333-3_16
2. Capello A, Forneris L, Pighetti A, Bellotti F, Lazzaroni L, Cossu M, De Gloria A, Berta R (2023) Investigating high-level decision making for automated driving. In: Berta R, De Gloria A (eds) Applications in electronics pervading industry, environment and society. Springer Nature Switzerland, Cham, pp 307–311. https://doi.org/10.1007/978-3-031-30333-3_41
3. Leurent E (2018) An environment for autonomous driving decision-making. https://github.com/eleurent/highway-env
4. Dosovitskiy A, Ros G, Codevilla F, Lopez A, Koltun V (2017) CARLA: an open urban driving simulator. http://arxiv.org/abs/1711.03938
5. Pérez-Gil Ó et al (2022) Deep reinforcement learning based control for autonomous vehicles in CARLA. Multimed Tools Appl. 81:3553–3576. https://doi.org/10.1007/s11042-021-11437-3
6. Gutiérrez-Moreno R, Barea R, López-Guillén E, Araluce J, Bergasa LM (2022) Reinforcement learning-based autonomous driving at intersections in CARLA simulator. Sensors 22:8373. https://doi.org/10.3390/s22218373
7. Brockman G, Cheung V, Pettersson L, Schneider J, Schulman J, Tang J, Zaremba W (2016) OpenAI Gym. http://arxiv.org/abs/1606.01540
8. Raffin A, Hill A, Gleave A, Kanervisto A, Ernestus M, Dormann N (2021) Stable-Baselines3: reliable reinforcement learning implementations. J Mach Learn Res 22:1–8
9. NVIDIA PhysX (2023). https://github.com/NVIDIA-Omniverse/PhysX
10. Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O (2017) Proximal policy optimization algorithms, http://arxiv.org/abs/1707.06347
11. Bengio Y, Louradour J, Collobert R, Weston J (2009) Curriculum learning. In: proceedings of the 26th Annual International Conference on Machine Learning pp 41–48. Association for Computing Machinery, New York, USA. https://doi.org/10.1145/1553374.1553380

# Design of a Portable Water Pollutants Detector Exploiting ML Techniques Suitable for IoT Devices Integration

Antonio Fotia[1]([✉]) [ID], Antonella Macheda[2] [ID], Mohamed Riad Sebti[2] [ID], Chiara Nunnari[1] [ID], and Massimo Merenda[2] [ID]

[1] DICEAM Department, University of Reggio Calabria, 89124 Reggio Calabria, RC, Italy
antonio.fotia@unirc.it

[2] DIIES Department, University of Reggio Calabria, 89124 Reggio Calabria, RC, Italy
riad.sebti@unirc.it

**Abstract.** Analyses of the presence of heavy metals are particularly important in assessing water quality. The monitoring of environmental parameters leads to the collection and analysis of a large set of data that may be used to reduce polluting actions, suggesting that Internet of Things (IoT) technology may provide an efficient contribution to the mitigation of environmental issues. The objective of this paper is to choose the best Machine Learning (ML) model, which can be loaded into a low-power microcontroller, using the sensed data obtained through the Electrochemical Impedance Spectroscopy (EIS) technique as a dataset. The real and imaginary parts of the impedance for each frequency value occurring within the range of 0.18 and 1 Hz are used as features. Five different lead concentrations are used as outputs. The features and the output constitute the dataset. Sixteen distinct scenarios were examined to train the different models in order to investigate ways to reduce the number of features. With the aim of creating a low-energy device that can conduct measurements and predict outcomes locally, the emphasis is on training a variety of models, including qualitative (classification) and quantitative (regression) approaches.

**Keywords:** Edge ML · Impedance sensor · Environmental monitoring

## 1 Introduction

The monitoring of environmental parameters leads to the accumulation and analysis of a large amount of data that can be used to mitigate/reduce polluting actions. To that extent, IoT technology is a promising candidate to contribute effectively to the mitigation of environmental issues, enabling the acquisition of high-resolution spatiotemporal data. Moreover, the recent results of embedding Artificial Intelligence (AI) onto resource-constrained devices [1] such as ultra-low-power, low-cost IoT platform, enable much more detailed and powerful data analysis closer to where sensor measurements are gathered. In assessing water quality, analyses of the presence of heavy metals are of particular

importance. Electrochemical Impedance Spectroscopy (EIS) is a frequently used technique for evaluating the presence of metal ions in various environmental and biological specimens [2]. Numerous electrospun materials for the detection of Heavy Metal Ions (HMI) have been created and evaluated using a vast array of techniques [3]. Electrospun nanofibers possess numerous advantageous properties, such as a high surface-to-volume ratio, porosity, tunability of nanofiber properties and the ability to be readily functionalized to improve material performance. The functional sites that can be conferred on fibers through various functionalization or doping techniques facilitate the activation of HMI anchor points. Moreover, the high surface area to volume ratio of fibers imparts a wider dispersion of active sites, which, in conjunction with the porous structure of fibrous membranes, determines the electrospun's ability to remove heavy metals [4]. Consequently, it has been demonstrated [5] that the incorporation of nanomaterials as electrodes can enhance the sensitivity and reproducibility of electrochemical methodologies for the detection of HMI. When considering polymeric nanofibers, polyaniline conductive polymer was discovered to be one of the most prevalent sources [6]. Typically, sensor calibration is required periodically and performed in a controlled laboratory environment, which does not always accurately reflect the conditions sensors confront when used in the field. Implementation of on-board Machine Learning (ML) techniques can enhance the quality of data and guarantee the precision of the gathered data. The purpose of this study is to develop an ML model inferencing from sensing data presence and concentration of water pollution contaminants, that can be executed on a low-power microcontroller integrated into a portable device. The dataset utilized in this study was acquired through the implementation of the EIS technique, which involved the use of a sensor featuring an active layer composed of microfibers generated via the electrospinning method. The focus is on the training of diverse models, including qualitative (classification) and quantitative (regression) approaches, with the goal of developing a low-energy apparatus that can execute measurements and forecast responses on-site.

## 2    Material and Methods

### 2.1    Sensor Preparation and Data Acquisition

Impedance measurements were carried out using AMEL 7050 galvanostat/potentiostat coupled with AMEL 7200 Frequency Response Analyzer (FRA) with a three-electrode configuration. The tests utilized Platinum Interdigitated Electrodes (PtIDEs) that were covered with micro-fibrous sensing layers (Fig. 1), serving as both the working and counter electrodes. The reference electrode utilized in the experiment was a platinum wire measuring 100 mm in length and 0.5 mm in diameter.

The active layer of the sensor was produced using the electrospinning method and the polymeric solution was prepared using Poly-Styrene (PS) as the polymer and EDTA disodium salt ($Na_2EDTA$) as the chelating agent. The electrolyte solutions utilized in the sensing experiments were prepared by different lead concentrations, specifically 10, 100 ug/l, 1, and 100 mg/l, with the buffer solution as a base. The EIS measurements were performed with a half-wave amplitude of 10 mVpK and a frequency range of 0.1–100.000 Hz. No potential bias was applied. The collected data were interpreted using
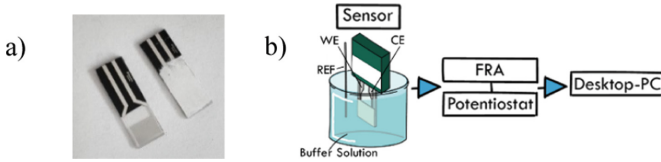
**Fig. 1. a** Bare PtIDE and PtIDE-PS-Na$_2$EDTA sensor, **b** block diagram of measurement setup.

the Nyquist plot, which plots impedance as a complex number with the imaginary part ($Z_{Im}$) along the y-axis and the real part ($Z_{Re}$) along the x-axis.

## 2.2 Dataset

**The dataset was assembled with the purpose of training various regression and classification models**. The frequency and the related real and imaginary parts of the impedance have been used as features. The frequency spectrum under analysis covers 0.18–50,000 Hz, encompassing a total of 110 frequency points. Five different lead concentrations, in terms of logarithm, were used as output (Fig. 2). Given that the difference in percentage between the results of the measurements is below 2%, the original dataset was synthetically augmented in magnitude (the number of observations) using the following equation: $Vm = random\ [-1,\ 1] * 0.02 * Vo + Vo$, where Vo represents the experimental data and $V_m$ the synthetic data. To preserve the original dataset's proportions while increasing the number of observations, the lead concentration measurements were artificially augmented by accordingly amounts. The dataset comprises 330 features and 930 observations, with 186 observations for each concentration that was analyzed.



**Fig. 2.** Structure of the dataset.

By employing two distinct forms of supervised learning, classification and regression, our objective was to train a ML model suitable for the integration on a resource-constrained microcontroller memory. For the first one, to obtain a series of classes related to different concentrations of pollutants, and for the second one, to predict a specific concentration value. Sixteen distinct scenarios were investigated to analyze the impact of different features and achieve the goal of reducing the number of employed features while maintaining model accuracy and decreasing training time. Each scenario will be described in the Table 1, which is labeled with the letters *a* to *r*. The Classification Learner App and Regression Learner App, which are included in the Statistics and Machine Learning Toolbox for MATLAB R2021 [7], were utilized to select the best models. These applications can indeed be utilized to train semi-supervised and supervised learning algorithms for binary and multi-class problems. The classification models

include Decision Trees, Discriminant Analysis, Naive Bayes, Support Vector Machine, Nearest Neighbors, and Ensemble. Meanwhile, the regression models comprise Linear Regression, Tree, Support Vector Machine, Ensemble, and Gaussian Process Regression. In the R2021a version, each group includes one or more models, for a total of 24 for classification and 17 for regression. The K-Fold Cross Validation (K-FCV) technique was employed, with a fold count of 5. This approach results in enhanced predictive accuracy by avoiding the problem of overfitting.

**Table 1.** Summary of scenarios used.

| Scenario | Number of features | Description |
|---|---|---|
| a | 330 | All features considered |
| b | 220 | Only Re and Im parts of impedance considered |
| c | 2 | PCA of scenario: 95% of the variance |
| d | 5 | PCA of scenario b: 97% of the variance |
| e | 10 | PCA of scenario b: 99% of the variance |
| f | 32 | Re and Im parts of Impedance in the range 0.18–1 Hz |
| g | 80 | Re and Im parts of Impedance in the range 1.2–100 Hz |
| h | 108 | Re and Im parts of Impedance in the range 0.12–50 kHz |
| i | 110 | Only Re part of Impedance considered |
| l | 1 | PCA of scenario i: 95% of the variance |
| m | 2 | PCA of scenario i: 97% of the variance |
| n | 5 | PCA of scenario i: 99% of the variance |
| o | 110 | Only Im part of impedance considered |
| p | 2 | PCA of scenario o: 95% of the variance |
| q | 4 | PCA of scenario o: 97% of the variance |
| r | 9 | PCA of scenario o: 99% of the variance |

## 3 Results

### 3.1 Classification Results

Twelve models have been identified and shown in Table 2 for classification purposes, based on a training time of less than 0.55 s and an accuracy of over 97%. The top-performing models, with respect to precision, duration of training, and quantity of features, are exclusively associated with scenario f. This scenario encompasses impedance values within the frequency range of 0.18–1 Hz. The top-performing models, as determined by the analysis, were Weighted KNN, Cosine KNN, and Quadratic Discriminant. Therefore, the above-mentioned models related to the f scenario turned out to be the best option, as compared to the others, it has a smaller number of features while maintaining excellent accuracy, training time, and prediction speed.

**Table 2.**  Selection of models trained for classification approach.

| Scenario | Models | Accuracy (%) | Prediction speed (Obj/s) | Training time (s) |
|---|---|---|---|---|
| f | Medium tree | 97.1 | 38000 | 0.459 |
| f | Quadratic discriminant | 100.0 | 30000 | 0.518 |
| f | Medium KNN | 100.0 | 13000 | 0.503 |
| f | Cosine KNN | 99.9 | 19000 | 0.330 |
| f | Weighted KNN | 100.0 | 19000 | 0.318 |
| g | Medium KNN | 100.0 | 11000 | 0.544 |
| g | Cosine KNN | 98.7 | 12000 | 0.495 |
| g | Weighted KNN | 100.0 | 12000 | 0.457 |
| h | Weighted KNN | 99.5 | 10000 | 0.548 |
| o | Linear discriminant | 100.0 | 21000 | 0.392 |
| o | Quadratic discriminant | 100.0 | 21000 | 0.437 |
| o | Weighted KNN | 100.0 | 10000 | 0.514 |

## 3.2  Regression Results

For the regression approach, seven models remain, considering a Root Mean Square Error (RMSE) less than 0.35 and an r-square ($R^2$) more than 0.99 (Table 3). The three models that show the best qualities, in terms of RMSE, $R^2$ and number of features, are all associated with scenario f.

**Table 3.**  Selection of models trained for regression approach.

| Scenario | Models | RMSE | $R^2$ | Prediction speed (Obj/s) | Training time (s) |
|---|---|---|---|---|---|
| a | Matern 5/2 GPR | 0.3373 | 0.99 | 8100 | 17.205 |
| a | Rational quadratic GPR | 0.3377 | 0.99 | 7700 | 35.228 |
| b | Matern 5/2 GPR | 0.3384 | 0.99 | 10000 | 14.224 |
| b | Rational quadratic GPR | 0.3390 | 0.99 | 9700 | 32.844 |
| f | Matern 5/2 GPR | 0.2796 | 1.00 | 17000 | 11.368 |
| f | Exponential GPR | 0.3237 | 0.99 | 16000 | 19.791 |
| f | Rational quadratic GPR | 0.2785 | 1.00 | 18000 | 21.462 |

The study revealed that the Rational Quadratic GPR, Exponential GPR, and Matern 5/2 GPR models exhibited the best levels of performance.

The response plot shows the number of observations in the abscissa and the responses in the ordinate, in terms of the logarithm of the various lead concentrations analyzed. The vertical red lines between the prediction and the true value indicate the prediction error committed by the model. The residuals plot of the regression models shows the difference between predicted and true responses. The analysis shows an error higher for the response at 2.30 than at 6.90, i.e., for concentrations of 10 µg/L and 1 mg/L of lead. Among the three models, Matern 5/2 GPR is the best-performing one (Fig. 3).



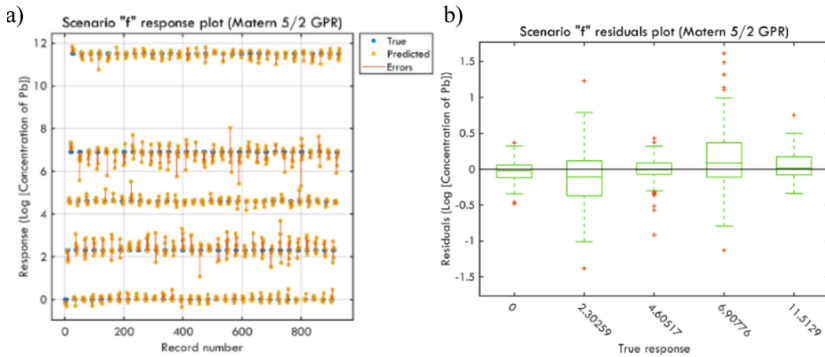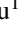**Fig. 3.** Response (**a**) and residual (**b**) plots of the best model trained with the scenario f.

An adaptation and selection of the trained models have been tested using the online tool STM32Cube. AI Developer Cloud to verify the suitability of the integration on an ARM Cortex M33 STM32U585AI microcontroller. Respectively, inference times for classification and regression are 1 and 0.99 ms. The memory occupation is less than 100 Kb for both the models.

# References

1. Merenda M, Porcaro C, Iero D (2020) Edge machine learning for ai-enabled IoT devices: a review. Sensors 20:2533
2. Wang S-L, Hsieh C-Y, Sukesan R, Chen J-C, Wang Y-L (2020) Highly sensitive lead ion detection in one drop of human whole blood using impedance-modulated field-effect transistors and a portable measurement device. ECS J Solid State Sci Technol 9:055020
3. Malara A, Fotia A, Paone E, Serrano G (2021) Electrospun nanofibers and electrochemical techniques for the detection of heavy metal Ions. Materials 14. https://doi.org/10.3390/ma14113000
4. Huang Y, Miao YE, Liu T (2014) Electrospun fibrous membranes for efficient heavy metal removal. J Appl Polym Sci 131:1–12. https://doi.org/10.1002/app.40864
5. Liu Y, Deng Y, Dong H, Liu K, He N (2017) Progress on sensors based on nanomaterials for rapid detection of heavy metal ions. Sci China Chem 60:329–337. https://doi.org/10.1007/s11426-016-0253-2
6. Fotia A et al (2021) Self standing mats of blended polyaniline produced by electrospinning. Nanomaterials 11:1–15. https://doi.org/10.3390/nano11051269
7. Lazzaro A, D'Addona DM, Merenda M (2022) Comparison of machine learning models for predictive maintenance applications. In: Advances in system-integrated intelligence. SYSINT 2022. Lecture notes in networks and systems, pp 657–666

# A Synthetic Dataset Generator for Automotive Overtaking Maneuver Detection

Luca Forneris[1(✉)] , Riccardo Berta[1] , Alessio Capello[1] , Marianna Cossu[1] , Matteo Fresta[1] , Fabio Tango[2] , and Francesco Bellotti[1]

[1] Department of Electrical, Electronic and Telecommunication Engineering (DITEN), University of Genoa, Via Opera Pia 11a, 16145 Genova, Italy
{luca.forneris,alessio.capello,marianna.cossu,
matteo.fresta}@edu.unige.it, {riccardo.berta,
francesco.bellotti}@unige.it

[2] Centro Ricerche Fiat (CRF), Orbassano, Italy
fabio.tango@crf.it

**Abstract.** This paper presents a novel tool for generating driving scenario datasets, that are a key asset to advance research and development in automated driving and driver assistance systems. The tool relies on the MATLAB. Automated Driving Toolbox and focuses on the overtaking maneuver. It uses simulated vehicular data, without relying on camera-equipped real-world vehicles, thus providing a low-cost solution, while allowing to abstract the main action features, that are very important for the pre-training of machine learning models. The tool has been designed to target customization (in terms, e.g., of road curvature radii), in order to allow meeting specific requirements, while its interoperability (e.g., multiple-format export) supports integration with other development environments. A preliminary analysis of the first scenarios generated with the tool confirms the validity of the system under development.

**Keywords:** Automated driving · Driving scenario · Synthetic datasets · Vehicular data · Driving scenario classification · Driving scenario detection

## 1 Introduction

Automated driving systems require reliable and high-quality scenario datasets for testing and validation. Scenario-based approaches are increasingly important for verifying automated driving systems compared to traditional distance-based methods. Overtaking (OV) and Lane Change (LC) maneuvers are particularly crucial due to their complex interactions between the Ego Vehicle (EV) and the Leading Vehicle (LV). To address dataset challenges, this paper proposes a tool for generating synthetic OV maneuver scenarios, highlighting its versatility, variability, adaptability, and interoperability. By offering a comprehensive dataset, this tool advances research and development in automated driving and driver assistance systems, ultimately promoting safer and more efficient autonomous vehicles.

Existing literature lacks an analysis of OV maneuver classification that doesn't rely on camera-equipped vehicles but instead utilizes only vehicular data. Our dataset creation serves two main purposes: (i) to develop an analysis tool for identifying and classifying the addressed scenario in an acquisition file, and (ii) to form the training set for a scenario classifier (e.g., left/right LC and possibly OV maneuvers). This paper focuses on dataset design and implementation, providing significant opportunities for research and development in automated driving and driver assistance systems. By focusing on vehicular data like relative position, distance, angle, and velocity between vehicles, researchers can safely develop innovative and practical solutions for automated vehicles to perform LC and OV maneuver. This approach may reduce complexity and costs compared to camera-based systems that rely on sophisticated computer vision algorithms and hardware.

Section 2 reviews existing literature related to our study. In Sect. 3, we outline our systematic approach, data collection methods, experimental setup, and tools used, emphasizing our work's key strengths. Section 4 demonstrates the usefulness of our experimental results, while Sect. 5 summarizes key findings and contributions. We highlight the importance of our methodology and dataset generation tool, along with discussing future research directions.

## 2 Related Works

Wachenfeld et al. [1] propose a scenario-based approach for validating automated driving functions, as relying solely on distance-based validation would be impractical due to the extensive kilometers required. This scenario-based method is becoming increasingly important for verifying and validating automated driving systems [2], necessitating high-quality and reliable scenario datasets.

Geyer et al. [3] compare scenarios to dynamic storylines, incorporating anticipated driver actions without rigidly specifying every detail. Scenario implementation considers factors like driver free will, allowing for flexible and context-dependent realization. Accurate scenario detection and classification involve identifying and interpreting various road situations and events, such as Car Following (CF), Lane Changes (LC) and Overtaking (OV) maneuvers. This detection enables automated driving systems to make contextually appropriate decisions.

Scenario detection is crucial in real-time for determining vehicle control operational modes and in offline analysis for assessing system performance [4]. High-quality datasets from real-world sensors are essential for training and testing autonomous driving systems, but creating them can be costly and time-consuming. Virtual simulations, like the one proposed by Cossu et al. [5], offer a complementary solution to address dataset challenges. They allow the generation of diverse scenarios through CARLA simulator [6] based on user-specified parameters, providing customized and tailored content for research purposes.

# 3  Methodology

We propose a synthetic OV maneuver scenario pipeline implemented using MATLAB [7] Automated Driving Toolbox [8]. The tool generates synthetic data extracted from simulations, focusing solely on vehicular data signals without using real sensors or cameras. Data is collected from vehicle dynamics, emulating sensors providing speed, yaw-rate, steering angle, etc. The current version allows the EV and LV to drive on a one-way road with two 3.5 m wide lanes, simulating a common highway scenario. Future improvements may include different road types like two-way roads or intersections. The main features of our tool-chain are:

- Reproducibility: scenarios are completely deterministic, guaranteeing a complete control of the experiments.
- Variability: the case histories considered in the study were generated using real statistical distributions, ensuring a diverse and variable dataset. Moreover, the range of customizable parameters enables the creation of highly distinct scenarios.
- Versatility: the tool offers complete customization according to the user's specific requirements. Users have the possibility to parameterize multiple factors, enabling the creation of varied and/or highly specific scenarios. Moreover, the tool facilitates the realization of diverse situations, including intersections, traffic roundabouts, and more.
- Interoperability: the work environment facilitates the export of scenarios generated in multiple extensions, enabling their import into other simulation and development environments. This capability empowers users to make additional modifications to the scenarios as needed.

## 3.1  Specifications

As anticipated before, one of the strengths of this tool is the possibility to customize simulated scenarios seamlessly, as users can define all the meaningful values by editing a simple text file. In particular, the user can define EV and LV speeds, the duration of the LC maneuver, the curve direction (left or right) and radius. Those parameters are required to define the generate a LC. In accordance with the user's requirements, the specifications of the scenarios generator can be documented within an external file (e.g., in a JSON format). The versatility of this tool allows for the inclusion of additional parameters to accommodate various other requirements. The chosen parameters for varying the scenarios include:

- EV speed [km/h]
- LV speed [km/h]
- Lane change maneuver duration [s]
- Road curvature radii [m], 0 for straight roads
- Curve direction, left or right.

The software utilizes the set parameters to interpret and define the scenarios.

## 3.2   Creation of Road Network and Vehicles

The road is generated, including straight and curved segments, with realistic infrastructure to represent a lifelike road layout. After road network generation, vehicles are placed according to predefined criteria and distribution to represent various traffic scenarios. To mimic real driver behavior, noise is added to the trajectories of EV and LV, simulating skids. An example of the road structure is depicted in Fig. 1.
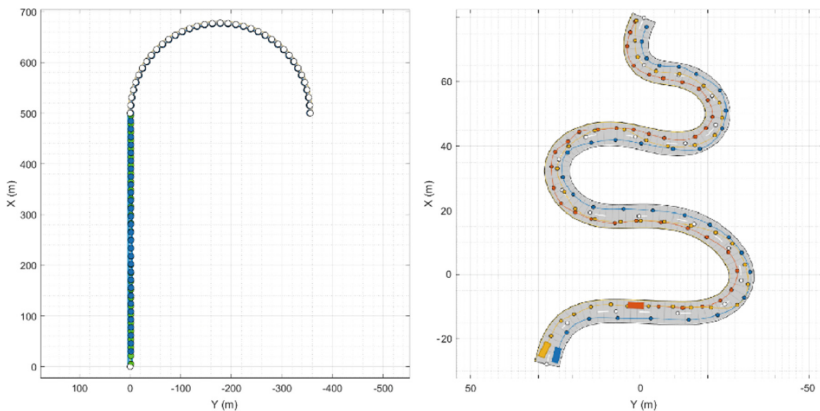


**Fig. 1.** Example of two road networks from the interface of Driving Scenario Designer application. On the left, the former portrays a road composed of a 500 m straight section followed by a curve with a radius of 178 m. Meanwhile, the latter illustrates a racetrack featuring three vehicles, each with a designated trajectory to follow.

## 3.3   Overtaking Scenario Instances

Two versions were established for each overtaking scenario instance. The first version comprises solely the EV driving along the road. The second version involves the EV following the LV without overtaking it. In the latter case, both vehicles maintain the same speed, and to ensure safety, a constant time headway of 1.3 s is maintained between them. In case of overtake, the difference between the speed of EV and LV is added to the list of parameters, so that the generated maneuvers can be as more realistic as possible.

## 3.4   Scenario Simulation and Export

The scenario is simulated and data is stored at each timesteps, which can be defined a priori according to the user needs. Users should determine an appropriate timestep value striking a balance between simulation accuracy and computational burden. Once all pertinent vehicular data have been acquired, the data relevant to the newly simulated scenario is now ready to be stored and used. Scenarios may be saved in different formats, such as ASAM OPENDRIVE® or OPENSCENARIO®.

## 4  Experimental Results

As a preliminary functional test, we defined a combination of 1000 parameter values. The corresponding baseline instances were generated and simulated, with a sample creation average time of 0.99 s. This process has been performed on a PC using 64-bit Windows 11 as operative system and equipped with an Intel core i7-12700H 2.30 GHz CPU, 16 GB RAM, and NVIDIA GeForce RTX 3060 GPU with 6 GB of memory. The generated instances have been verified by experts from the Hi-Drive project, especially to ensure physical consistency between values (e.g., speeds and radii of curvature). This demonstrates that by exploiting the proposed tool, users can effortlessly generate a dataset of OV maneuver scenarios by simply setting the desired values of some relevant parameters, with no need to write any script/code nor master a simulation environments.

## 5  Conclusions and Future Works

As efficient creation of driving scenario data is key to advance development of advanced driving assistance systems (ADAS) and automated driving functions (ADF), we have designed a tool, based on MATLAB Automated Driving Toolbox for implementing various instances of overtaking scenarios, in a variety of conditions. The relevant vehicular signals are recorded, so that can be later employed as part of machine learning datasets. This approach offers significant benefits, as it enables computationally efficient simulations and obviates the necessity of employing complex computer vision techniques to process camera-recorded data. The adaptability of the tool allows for extensive customization to suit specific user/company requirements. Authenticity of the scenario settings benefits from the incorporation in the tool of real maps from OpenStreetMap®. From serving as a diagnostic method for existing scenario classifiers with high accuracy to acting as a benchmark for comparing various classification methods. Furthermore, the flexibility of the tool goes beyond mere OV or LC scenarios. The combination of these attributes streamlines development, reducing computational complexity compared to numerous state-of-the-art tools, thereby unlocking a variety of potential applications.

A first already planned for future work will concern the implementation of other traffic vehicles beside EV and LV. On the other hand, we are interested also in involving of real users in data collection through the use of a steering wheel device connected to the tool. This approach aims to enhance the realism and reduce determinism in the data collection process, as a real driver exhibits non-deterministic behavior. Additionally, users can execute different driving styles (e.g., comfort, normal, aggressive), thus integrating the dataset with variety and human aspects (e.g., expertise, knowledge, emotions) into the dataset.

# References

1. Wachenfeld W, Winner H (2016) The Release of autonomous vehicles. In: Maurer M, Gerdes JC, Lenz B, Winner H (eds) Autonomous driving: technical, legal and social aspects. Springer, Berlin, Heidelberg, pp 425–449. https://doi.org/10.1007/978-3-662-48847-8_21
2. Nalic D, Mihalj T, Bäumler M, Lehmann M, Eichberger A (2021) Scenario based testing of automated driving systems: a literature survey. In: FISITA world congress 2021—Technical Programme. FISITA. https://doi.org/10.46720/f2020-acm-096
3. Geyer S et al (2014) Concept and development of a unified ontology for generating test and use-case catalogues for assisted and automated vehicle guidance. IET Intell Transp Syst 8:183–189. https://doi.org/10.1049/iet-its.2012.0188
4. Bellotti F et al (2020) Managing big data for addressing research questions in a collaborative project on automated driving impact assessment. Sensors 20:6773. https://doi.org/10.3390/s20236773
5. Cossu M, Berta R, Capello A, De Gloria A, Lazzaroni L, Bellotti F (2023) Developing a toolchain for synthetic driving scenario datasets. In: Berta R, De Gloria A (eds) Applications in electronics pervading industry, environment and society. Springer Nature Switzerland, Cham, pp 222–228. https://doi.org/10.1007/978-3-031-30333-3_29
6. Dosovitskiy A, Ros G, Codevilla F, Lopez A, Koltun V (2017) CARLA: an open urban driving simulator. In: Proceedings of the 1st annual conference on robot learning. PMLR, pp 1–16
7. Inc, T.M.: MATLAB version: 9.14.0 (R2023a) (2023). https://www.mathworks.com
8. Inc, T.M.: Automated Driving Toolbox version: 3.7 (R2023a) (2023). https://www.mathworks.com

# TEMET: <u>T</u>runcated <u>RE</u>configurable <u>M</u>ultiplier with <u>E</u>rror <u>T</u>uning

Flavia Guella[✉], Emanuele Valpreda, Michele Caon, Guido Masera,
and Maurizio Martina[✉]

Dipartimento di Elettronica e Telecomunicazioni, Politecnico di Torino, Turin, Italy
{flavia.guella,maurizio.martina}@polito.it

**Abstract.** Approximate computing is a well-established technique to mitigate power consumption in error-tolerant domains such as image processing and machine learning. When paired with reconfigurable hardware, it enables dynamic adaptability to each specific task with improved power-accuracy trade-offs. In this work, we present a design methodology to enhance the energy and error metrics of a signed multiplier. This novel approach reduces the approximation error by leveraging a statistic-based truncation strategy. Our multiplier features 256 dynamically configurable approximation levels and run-time selection of the result precision. Our technique improves the mean-relative error by up to 34% compared to the zero truncation mechanism. Compared with an exact design, we achieve a maximum of 60.1% power saving for a PSNR of 10.3dB on a $5 \times 5$ Sobel filter. Moreover, we reduce the computation energy of LeNet by 31.5%, retaining 89.4% of the original accuracy on FashionMNIST.

**Keywords:** Multiplier · Approximate computing · Reconfigurable computing · Image processing · VLSI

## 1 Introduction

Multiplication is the fundamental building block in most computer vision, digital signal processing, and artificial intelligence applications. Such tasks could involve millions or billions of multiplications, which significantly contribute to the overall complexity and energy demand. These high power requirements could become unbearable for mobile and edge devices, raising the urge to find viable alternatives. Approximate computing exploits the intrinsic robustness of computer vision algorithms to numeric errors to lower power consumption while preserving acceptable results. Several inexact multipliers designs have been described in the literature, leveraging different techniques, with the main aim of enhancing power performance with negligible loss in accuracy [6–9]. However, few works deal with the ever-increasing need for an architecture capable of dynamically adjusting based on specific tasks and requirements [12]. Reconfigurability allows hardware to adapt run-time to the workload requisites, potentially increasing energy saving at the cost of additional area. We propose a 9-bit input, signed

multiplier with dynamic truncation for the run-time selection of the precision (i.e., bit-width) and approximation level of the final product. This is necessary to support mixed-precision and layer-wise quantization, as it is fundamental to reduce the computational cost, the memory footprint, and the overall data movement during the execution [2]. The contributions of this paper are summarized as follows: (i) enabling simultaneous run-time selection of precision and approximation of the operands and result, (ii) fine-grain regulation of the error, with the possibility of choosing among 256 available levels, including one generating the correct result, providing several energy-quality trade-offs, (iii) introduction of a simple error compensation technique to lessen the effect of bit truncation.

## 2  Related Works

Beyond dynamic voltage-scaling [1], which generally produces an error that is complex to control and requires additional hardware to tweak the voltage, and Cartesian Genetic Programming [3], many works explore parallel multipliers. Their architecture comprises three parts on which approximation can be performed: (i) partial product generation, (ii) reduction tree, (iii) two-operands adder. The first technique simplifies the logic required to generate partial products [4]. Approximating the tree structure means substituting exact compressors with inexact ones [5] or truncating the least significant part of the matrix. Approach (iii) is quite common since there are many available approximate adders designs ready to deploy. As approximate computing gained renewed appeal for several IoT applications, including Deep Neural Network (DNN) for monitoring people health, precision agriculture, and food quality, the need for reconfigurable and inexact hardware emerged [15]. Most studies on approximate multipliers do not include dynamic reconfigurability, as it generally comes at the expense of increased power and area. Works such as [6–8] exploit previously described techniques to implement inexact unsigned multipliers with some error trimming. None of the cited works supports run-time selection of the input precision, and the configuration capability is limited; these designs sacrifice the possibility of having the correct result, thus flexibility, in exchange for higher energy saving. Dynamic truncation is introduced in [9], implemented with a signal that can set to zero selected columns of the partial product matrix. This multiplier is designed for signed operands and can execute exact multiplication besides inexact ones. A similar underlying principle is used in our architecture. As a final remark, while most previous works proposed to implement either $8 \times 8$ or $16 \times 16$ unsigned multipliers, our work suggests a $9 \times 9$ signed multiplier to support all combinations of 8-bit dot products: signed-signed, unsigned-signed, signed-unsigned, and unsigned-unsigned.

## 3    Methodology

This section covers the design phase of the signed $9 \times 9$ multiplier. From [10], it is known that the Dadda tree has a lower delay and complexity than the Wallace one. Therefore, the former structure is selected to allocate the compressors. As our multiplier should support one level producing the correct result, a final exact two-operands adder is instantiated and the tree-reduction uses only exact 3 to 2 or 2 to 1 compressors. The multiplication algorithm is implemented with the Modified Baugh-Wooley (MBW) due to its lower power consumption and area compared to Modified Booth Encoding [11]. Once the architectural structure of the exact multiplier has been fixed, the possibility of run-time configuration is achieved by using two masks, setting the precision and the approximation levels fed to the multiplier. Since reducing the dynamic power is the main objective of the optimization procedure, the choice is to perform data gating. The part of the partial product matrix (PPM) that is not involved in the operation is kept constant to reduce the switching activity and, thus, dynamic power consumption. A signal called `res_mask` is in charge of masking to zero the portion of the PPM not required for the computation, with a bitwise *and*. The `res_mask` signal covers the most significant fourteen bits of the result, as it is assumed that the minimum precision of the inputs is two bits. As the result has to be signed, additional logic is required to guarantee that its left part is correctly sign extended.



**Fig. 1.** Result sign extension



**Fig. 2.** Approximation mask

As depicted in Fig. 1, the sign bit, obtained as the *xor* between the most significant bit of the multiplicand and that of the multiplier, is used as the selection signal of a multiplexer producing the leftmost part of the output. The precision mask is used to correctly sign-extend the product. Two zero comparators are implemented to cover the corner case where one of the two inputs is zero, and the other is negative. Although the two 9-bit zero comparators are an additional cost in terms of area and power, they enable the multiplier to produce the correct result whenever at least one operand is zero. Accurate zero multiplication is essential in DNN inference or image processing, so it is an acceptable trade-off. Moreover, the sensitivity to numeric errors of different DNNs layers has high

variability [2]. Consequently, designing a multiplier with different approximation levels is crucial to ensure flexibility and adaptability. Run-time reconfigurability requires an area and power overhead; thus, the most convenient trade-off has to be evaluated. Following a holistic approach, a variation of truncation is identified as the target method for approximation, as it easily enables dynamic configuration without requiring the design of programmable approximate compressors and their insertion in specific points of the PPM. The selection of the level of approximation is achieved using a second mask signal, called `appr_mask`. While for the precision configuration the cut is on the most significant bits (MSBs) of the matrix, here the truncation is on the right part. Once an approximation method has been established, the issue of the number of reconfigurability levels has to be investigated. A first consideration, also argued by other works such as [7,8], is that there is no advantage in pushing approximation to the most significant half of the result, as the error becomes unacceptable. The design choice is thus to limit the truncation to the least significant 8 bits of the multiplier. Given a configuration, for each bit at zero in the `appr_mask`, all the corresponding columns of bits of the PPM are set to a fixed value. Configurations with more zeros in the mask are more effective than others for data gating. Although some levels are Pareto-dominated in terms of power saving and error metrics, they might be the optimal choice for a specific layer in a DNN, as this strictly depends on its inputs and weights distribution. Consequently, the choice is to keep all 256 levels, even if some are less power efficient. Figure 2 shows the implemented mechanism to decrease the generated error. The first row has the bits from 2 to 7 fixed at one, while all others are data-gated to zero. An in-depth analysis is carried out to find optimal configurations of the data-gated matrix minimizing the error.

**Table 1.** Sum and carry probabilities

| $Bit_i$ | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|
| $P_{sum}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |
| $P_{C_o}$ | $\frac{119}{256}$ | $\frac{1}{2}$ | $\frac{9}{16}$ | $\frac{5}{8}$ | $\frac{5}{8}$ | $\frac{1}{2}$ | $\frac{1}{4}$ | 0 |
| $P_{C_{o+1}}$ | $\frac{41}{64}$ | $\frac{1}{2}$ | $\frac{11}{32}$ | $\frac{3}{16}$ | $\frac{1}{16}$ | 0 | 0 | 0 |
| $P_{O_i \geq 1}$ | $\frac{107}{128}$ | $\frac{421}{512}$ | $\frac{211}{256}$ | $\frac{13}{16}$ | $\frac{3}{4}$ | $\frac{5}{8}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |

**Table 2.** Metrics of $9 \times 9$ signed exact and reconfigurable multipliers

| Design | Area [$\mu m^2$] | Arrival time [ns] | Approx level | Power [$\mu W$] |
|---|---|---|---|---|
| Exact | 607.3 | 1.7 | – | 414.0 |
| Proposed | 822.6 | 1.8 | 0 | 240.6 |
| | | | 255 | 164.8 |
| Dyn Trunc | 816.5 | 1.8 | 0 | 242.4 |
| | | | 255 | 165.0 |

Under the assumption of uniformly distributed inputs and of 9-bit full-precision inputs some considerations can be made. The purpose of the Dadda tree is to compress each of the columns of the PPM up to a single bit. Given a column, it is supposed that its bits can take all possible values with the same probability. Consequently, for each bit of the result we can evaluate its likelihood of being one or greater (i.e. generating a one in the position at its left). The probability of the sum or the carry produced by a column resulting in a one is thus evaluated.

As the sum bit is an odd function, its probability of being one is always 0.5. All the probabilities for the carries are evaluated and reported in Table 1. The probability of the bit $i$th of the result, from position 0 to 7, being one or greater, $P_{O_i \geq 1}$ in Table 1, is so estimated:

$$P_{(O_i \geq 1)} = P_{sum_{(i)}} \cup P_{C_{o_{(i-1)}}} \cup P_{C_{(o+1)_{(i-2)}}} \tag{1}$$

Results in Table 1 demonstrate that from position 2, as expected, the probability of the output bit being greater than 0 is higher than 50%. This high likelihood justifies the logic shown in Fig. 2.

## 4   Experimental Results

The synthesis is performed with the UMC 65nm technology library and a timing constraint of $2ns$, using Synopsys Design Compiler (DC). The golden model for the exact multiplier is a behavioral $9 \times 9$ bit signed multiplier whose architecture selection and optimization are left to Synopsys DC, with the library Synopsys DesignWare. Furthermore, an architecture similar to ours but with the traditional dynamic truncation to zero for the approximation logic, based on [9], is considered as a comparison. Power estimation is performed through back-annotation of the post-synthesis netlist using a testbench producing 1000 random stimuli. The procedure is repeated for every configuration of the precision and approximation masks of our multiplier. The main metrics are compared in Table 2.



**Fig. 3.** MRED comparison between Dyn Trunc and Proposed, 9-bit
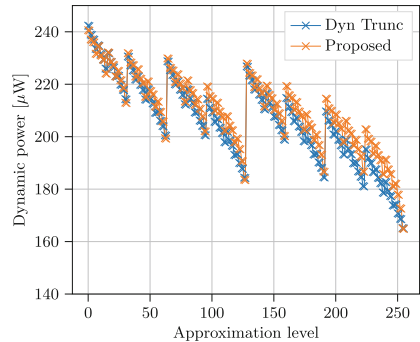
**Fig. 4.** Power comparison between Dyn Trunc and Proposed, 9-bit

The complexity of the exact architecture is lower as reconfigurability has some overhead on the occupied area and delay. Our multiplier, in the full-precision, exact configuration, saves 41.4% of power compared to the one by Synopsys, while the highest approximate level can save up to 60.1%. These results

are remarkable and demonstrate the effectiveness of approximation in reducing power. Furthermore, they prove that the reconfigurability overhead is tolerable and does not significantly impact the critical part of the design. As a final remark, modifying the approximation logic does not impact the timing of the architecture and power increase is always below 5%. Figure 3 shows the mean-relative error distance (MRED) variation with the approximation level, comparing our multiplier with traditional dynamic truncation for the full-precision case. As can be observed, the proposed strategy allows an improvement of the MRED for every approximation level, with a maximum gain of 34%. Experiments are carried out on a typical image processing task for edge detection using a $5 \times 5$ Sobel filter, as in [5]. Figure 5 reports the PSNR values, evaluated on the gray-scale peppers benchmark image[1] for changing approximation levels for both traditional dynamic truncation and ours. The proposed truncation methodology improves PSNR by up to 81.6%. Considering an acceptable threshold of 25dB for the PSNR, 64 configurations of our multiplier, against 25 for traditional truncation, are usable. Moreover, we tested the 256 approximation levels with LeNet [13] on the FashionMNIST dataset.[2]



**Fig. 5.** PSNR comparison between Proposed and Dyn Trunc multiplier for $5 \times 5$ Sobel filter

**Fig. 6.** Top 1% test accuracy variation with LeNet trained on FashionMNIST

We used PyTorch[3] and AdaPT [14] to implement the neural network, training the model for 32 epochs using stochastic gradient descent, a learning rate of $10^{-3}$, and a decay of $5 \cdot 10^{-4}$, with weights and activations quantized to 8-bit. Then we changed the approximation level of our multiplier and retrained the bias of each convolutional and dense layer for one epoch, without updating the other parameters. Figure 6 reports the results before and after the re-training. The accuracy

---

[1] https://sipi.usc.edu/database/database.php?volume=misc&image=11#top.

[2] https://www.kaggle.com/datasets/zalando-research/fashionmnist.

[3] https://pytorch.org/.

loss can be recovered by updating the bias, accounting for the computation error introduced by the approximation. We evaluate the total computation energy as the product of a single multiplication and the number of operations required to execute the inference. With our approximate multiplier, it is possible to reduce by 31.5% the computation energy of LeNet, with an absolute top 1 accuracy degradation of 9.67%.

## 5    Conclusion and Future Works

This work demonstrates the effectiveness of approximate computing in reducing the power consumption in error-resilient computationally complex tasks. Furthermore, it points out the importance of reconfigurability for providing flexibility. In the future, we will investigate the effect of inexact computation on different applications, such as DNNs with layerwise approximation, and we will explore in depth the trade-offs that the concurrent adoption of approximation and reduced precision can yield. We will also consider approximate and resilient computing in different applications such as agritech [15].[4]

## References

1. Moons B, Verhelst M (2015) DVAS: Dynamic voltage accuracy scaling for increased energy-efficiency in approximate computing. In: 2015 IEEE/ACM international symposium on low power electronics and design (ISLPED), pp 237–242
2. Fasfous N, et al (2021) Hw-flowq: a multi-abstraction level hw-cnn co-design quantization methodology. ACM Trans Embed Comput Syst (TECS) 20:5s, Article 66
3. Mrazek V, Sekanina L, Vasicek Z (2020) Libraries of approximate circuits: automated design and application in CNN accelerators. IEEE J Emerg Sel Top Circuits Syst 10(4):406–418
4. Yin P et al (2018) Designs of approximate floating-point multipliers with variable accuracy for error-tolerant applications. J Signal Process Syst 90(4):641–654
5. Strollo AGM et al (2020) Comparison and extension of approximate 4–2 compressors for low-power approximate multipliers. IEEE Trans Circuits Syst I: Regul Pap 67(9):3021–3024
6. Jiang H et al (2019) Low-power approximate unsigned multipliers with configurable error recovery. IEEE Trans Circuits Syst I: Regul Pap 66(1):189–202
7. Yang T, Ukezono T, Sato T (2018) A low-power high-speed accuracy-controllable approximate multiplier design. In: 2018 23rd Asia and South Pacific design automation conference (ASP-DAC), pp 605–610
8. Gu F-Y, Lin I-C, Lin J-W (2022) A low-power and high-accuracy approximate multiplier with reconfigurable truncation. IEEE Access 10:60447–60458
9. de la Guia Solaz M, Han W, Conway R (2012) A flexible low power DSP with a programmable truncated multiplier. IEEE Trans Circuits Syst I: Regul Pap 59(11):2555–2568

10. Parhami B (2010) Computer arithmetic–algorithms and hardware designs, 2nd edn. Oxford University Press, New York
11. Sjalander M, Larsson-Edefors P (2008) High-speed and low-power multipliers using the Baugh-Wooley algorithm and HPM reduction tree. In: 2008 15th IEEE international conference on electronics, circuits and systems, pp 33–36
12. Fasfous N et al (2022) AnaCoNGA: analytical HW-CNN Co-design using nested genetic algorithms. In: 2022 design, automation & test in Europe conference & exhibition (DATE), pp 238–243
13. Lecun Y et al (1998) Gradient-based learning applied to document recognition. Proc IEEE 86(11):2278–2324
14. Danopoulos D et al (2022) Adapt: fast emulation of approximate dnn accelerators in pytorch. IEEE Trans Comput-Aided Des Integr Circuits Syst 42(6):2074–2078
15. Barezzi M et al (2022) On the impact of the stem electrical impedance in neural network algorithms for plant monitoring applications. In: 2022 IEEE workshop on metrology for agriculture and forestry (MetroAgriFor), pp 131–135

# Cycle-Accurate Verification of the Cryptographic Co-Processor for the European Processor Initiative

Pietro Nannipieri[(✉)] , Stefano Di Matteo , Luca Crocetti ,
Luca Zulberti , Luca Fanucci , and Sergio Saponara

Deparment of Information Engineering, University of Pisa, Via G. Caruso 16, Pisa,
Italy
**pietro.nannipieri@unipi.it**

**Abstract.** This paper presents a cycle-accurate verification environment for the Crypto-Tile, a cryptographic accelerator integrated into the EPI General Purpose Processor. The focus of this work is to provide a robust methodology for validating the functionality and performance of the Crypto-Tile. The verification environment includes an in-depth examination of the internal architecture and operational aspects of the Crypto-Tile, allowing for accurate modelling of hardware components and emulation of Direct Memory Access (DMA) operations. Developers can leverage this environment to simulate and verify their C-Code implementations, utilizing the functions available in the Crypto-Tile library or creating custom libraries. The verification process involves using the 32-bit AXI4 interface for communication between the processor and the Crypto-Tile while emulating DMA operations to ensure accurate testing.

**Keywords:** AES · ECC · RNG · SHA · RISC-V · EPI ·
Cryptoprocessor · Hardware · Verification · Cycle-accurate

## 1 Introduction

In today's interconnected world, where digital information flows seamlessly across networks, the importance of cryptography cannot be overstated. From securing online transactions and safeguarding personal information to protecting national security interests, cryptography plays a pivotal role in upholding the trust and privacy of individuals, organizations, and governments. As the reliance on computing systems continues to grow exponentially, so does the need for robust and efficient cryptographic processors [1,2,5,7,12]. In this context, the European Processor Initiative (EPI) [6] represents a groundbreaking collaborative effort among European Union member states, research institutions and industry partners to develop a cutting-edge high-performance computing ecosystem. The EPI aims to develop a new generation of energy-efficient and high-performance processors. With digital transformation permeating every aspect

of society, achieving self-reliance in processor technology has become imperative for Europe's strategic autonomy. By fostering homegrown expertise and innovation in processor design, the EPI seeks to reduce dependence on non-European technology providers, enhance Europe's technological competitiveness, and strengthen its position in the global market. Developing a secure and efficient cryptographic processor is an integral part of the EPI's broader vision, as it addresses the growing need for robust encryption capabilities to safeguard sensitive data and critical information infrastructures against ever-evolving cyber threats. When dealing with cryptographic algorithms, the computational complexity poses challenges for traditional software execution on Central Processing Units (CPUs). Cryptographic operations involve data manipulations that demand substantial processing power. The execution of these algorithms purely through software implementations can result in significant performance bottlenecks and increased execution times. To overcome these limitations, hardware acceleration emerges as a compelling solution. This approach improves the overall performance of cryptographic operations and ensures the secure processing and protection of sensitive data, making it a preferred choice in scenarios where efficient and high-performance cryptography is paramount. As a solution, the EPI system integrates the Crypto-Tile IP core, which is a dedicated hardware component for cryptography. It plays a crucial role in providing hardware acceleration for cryptographic algorithms, enabling efficient and secure encryption and decryption operations, and robust protection mechanisms for security-critical assets, such as keys. Cycle-accurate verification plays a pivotal role in the development and validation of complex hardware designs, such as the Crypto-Tile. Simulating the hardware at the cycle level enables a more granular and comprehensive analysis of the design's functionality, performance, and compliance with desired specifications. By precisely modelling the hardware behaviour, timing constraints, and interactions with software components, it becomes possible to identify and rectify potential design flaws or bugs. Additionally, cycle-accurate verification enables detailed measurement and analysis of power consumption and latency in software-accelerated hardware primitives. This work aims to provide an overview of the cycle-accurate verification environment developed to validate the Crypto-Tile co-processor. By presenting the verification approach followed using a cycle-accurate verification environment, we strive to contribute to the broader understanding of the rigorous verification processes undertaken for critical hardware components within the EPI framework.

## 2   The Crypto-Tile within the European Processor Initiative

The Crypto-Tile [11] is a cryptographic accelerator designed to provide a comprehensive and versatile range of cybersecurity services with advanced security features. It incorporates various engines to support different cryptographic functions. The Advanced Encryption Standard (AES) engine [10] supports both AES-128 and AES-256 ciphers, offering a minimum security strength of 128 bits

in terms of both classical and post-quantum security. The Elliptic-Curve Cryptography (ECC) engine [4] enables operations on elliptic curves with widths of 256 and 521 bits, providing symmetric-key equivalent security strength for classical security. The Secure Hash Algorithm (SHA) engine [8] generates digests of at least 256 bits and 384 bits through SHA2 and SHA-3 functions, meeting the minimum strength requirement of 128 bits. Additionally, the Crypto-Tile features a True Random Number Generator (TRNG) engine [9] that meets the security requirements for cryptographic applications. The internal architecture of the Crypto-Tile (Fig. 1) consists of several key components. These include a 32-bit AXI4 interface, which serves as a Slave Memory-Mapped interface for accessing the registers of the Crypto-Tile through the Configuration bus. The Global Management Unit handles the global configuration, control, and status of the Crypto-Tile. Each of the four independent crypto-processors is dedicated to a specific class of cryptographic algorithms: AES [10], ECC [4], SHA [8], and Random Number Generator (RNG) [9]. These cryptoprocessors function as coprocessing units for the main processor they are connected to, such as the Secure Micro-Controller Unit (MCU) or the Secure Element (SE, a compact micro-controller that handles the first stage of the secure boot routine). Each cryptoprocessor includes local registers for configuration, control, and status and a Finite State Machine (FSM) for managing cryptographic operations. They also feature engines for hardware acceleration of cryptographic algorithms and functions. The AES and ECC cryptoprocessors incorporate local resources for key storage and management. To facilitate high-bandwidth transfers, four indepen-
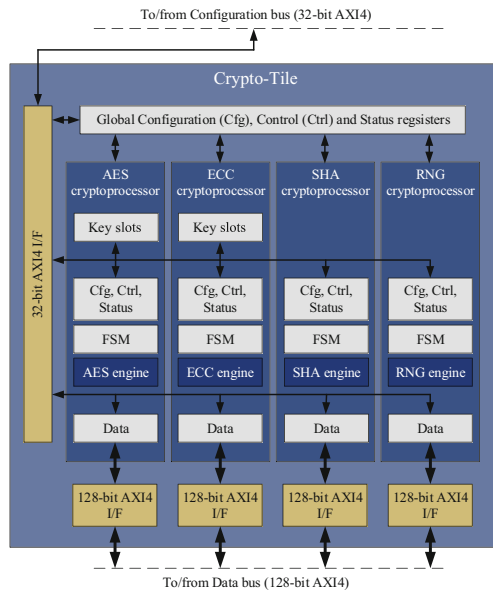


**Fig. 1.** Crypto-tile blocks scheme

dent 128-bit AXI4 interfaces, one for each cryptoprocessor, provide access to the data registers. Performances and complexity of the Cryptotile are reported in Table 1.

**Table 1.** Implementation results for the CryptoTile

| Tech | Entity | Max Freq | Complexity |
|------|--------|----------|------------|
| | Crypto_Tile | 150 MHz | 27507 CLB; 144892 CLB LUTs; 93503 CLB Reg; 64 DSPs |
| *Xilinx* | RNG engine | 260 MHz | 2294 CLB; 10154 CLB LUTs; 7122 CLB Reg; 0 DSPs |
| *VU37P* | SHA engine | 190 MHz | 3433 CLB; 10290 CLB LUTs; 10787 CLB Reg; 0 DSPs |
| *FPGA* | ECC engine | 95 MHz | 15151 CLB; 79219 CLB LUTs; 37626 CLB Reg; 64 DSPs |
| | AES Engine | 170 MHz | 1253 CLB; 6036 CLB LUTs; 2460 CLB Reg; 0 DSPs |
| | Crypto_Tile | 3.7 GHz | 1325.16 kGE |
| *7nm* | RNG engine | 4.325 GHz | 127.16 kGE |
| *Std-Cell* | SHA engine | 3.725 GHz | 128.32 kGE |
| *Technology* | ECC engine | 1.525 GHz | 658.90 kGE |
| | AES engine | 2.425 GHz | 56.01 kGE |

# 3   The Crypto-Tile RISC-V Environment



**Fig. 2.** Crypto-tile cycle-accurate verification environment

The verification environment set-up is based on the one presented in [13], which exploits an automated framework for accelerating the design space exploration of hardware/software co-designs in heterogeneous digital systems. It simplifies the customization of design tools, improves designer productivity, and allows the evaluation of different hardware/software choices. The framework integrates the RISC-V toolchain and focuses on post-synthesis analyses for accurate power

consumption evaluations. It helps design robust systems against Side-Channel Attacks [3], ultra-low-power and heterogeneous architectures, and space-grade Systems on Chip (SoCs) with varying technology outcomes. The environment architecture (Fig. 2) relies on a makefile approach that requires three different inputs: (1) a C project (i.e. the software intended to be executed by the RISC-V processor); (2) a SystemVerilog (SV) testbench (to provide the necessary inputs to the hardware and simulate its interface with the rest of the world); (3) a Recipe, to configure the work environment. These inputs are fed directly to the framework which processes them by compiling the C code and initializing a simulated RAM containing the compiled executable, then the SV testbench is executed to initialise the cycle-accurate verification. The testbench needs to provide service signals (e.g. clock and reset), and it may send data to the Crypto-Tile via the dedicated 128-bit DMA AXI4 emulator, achieving a data rate much higher than the one provided by the 32-bit AXI4 Internal interconnect. The test environment is then respfor collecting all the generated information, writing a log file of the simulation, and printing the UART output of the RISC-V processor into a dedicated file. To use the verification environment user shall perform the following steps:

(1) Write the C code to be executed, by exploiting the provided templates and referring to the Crypto-Tile documentation for details on drivers, register addresses, operation modes et al. to use the functions available in the Crypto-Tile library or by creating an own library. It is noted that the C code provides access only to the 32-bit AXI4 MCU interface on the processor. The DMAs are emulated in System Verilog and cannot be accessed via the C code. However, this limitation will be addressed with the Field Programmable Gate Array (FPGA) system. All the C source files must be added as new targets in the makefile, indicating the main C file ($<c\_test>.c$) as the main target ($<c\_test>$).

(2) Write the SV testbench, which must include at least the clock ($clk$) and reset ($rst$) signals, Also DMA operations can be included by exploiting the System Verilog Secure DMA (SDMA) functions for the AXI4 Master emulation. In this no, no default synchronization mechanisms between the RISC-V processor and the SDMA interfaces are provided, hence they can be used in the Crypto-Tile's internal signals, such as Interrupt Request (IRQ) signals. The name of the SV testbench file (e.g., $<tc\_crypto\_tile\_ID>.sv$) constitutes the top-level and must be specified in the recipe (next step).

(3) Run the simulation using the corresponding command. The logs of the compilation process and the simulation are made available in a dedicated folder ($build$), and they can be used to check errors and information.

## 4   Simulations

Thanks to the simulation environment provided, we were able to establish a comprehensive test plan that stimulates the entire system by focusing on each developed accelerator and service system. Below, we provide a concise overview of the tests performed in each test category:

- **AXI Interfaces**: This test category consists of 20 diverse tests that aim to thoroughly test the AXI communication infrastructure. The focus is mainly on the MCU AXI interconnect and each of the AXI DMA (one for each Cryptoprocessor). All supported AXI operations undergo comprehensive testing.
- **Global Management Unit**: This test category comprises 58 distinct tests that aim to execute write and read attempts on the register file. The Crypto-Tile security specification governs the register file. These tests highlight possible error situations and evaluate the security of the system against unauthorized attempts to access protected data (both read and write).
- **AES Cryptoprocessor**: This category of tests focuses on the AES cryptoprocessor and includes 138 different tests that need to be executed. These tests cover all the different AES operative modes such as ECB, OFB, CFB, CTR, XTS, CMAC, GCM, and CCM. Additionally, they involve the writing and reading of the cryptographically secured configuration and status register.
- **ECC Cryptoprocessor**: This test category focuses on the ECC cryptographic processor with 62 different tests to be executed. The writing and reading of secured configuration and status registers are tested intensively, together with all the supported ECC engine operational modes.
- **SHA Cryptoprocessor**: This category of tests is centred around the SHA cryptographic processor and includes 63 different tests that need to be performed. There is an extensive focus on testing the secure configuration and status registers' reading and writing, along with both SHA3 and SHA-2 operational modes, each with all the supported key sizes.
- **RNG Cryptoprocessor**: The final category of tests is centred around the RNG cryptoprocessor and includes 45 different tests that need to be performed. These tests heavily evaluate the writing and reading of secured configuration and status registers, as well as the various operational modes, such as changing the entropy source and number generation modes.

## 5   Conclusions

This work presented a comprehensive cycle-accurate verification environment for the Crypto-Tile, a state-of-the-art cryptographic accelerator integrated into the EPI system. The verification environment provides a systematic approach for validating functionality and performance, by accurately modelling the behaviour of the hardware components and emulating the DMA operations in System Verilog. Developers can leverage this verification environment to simulate and verify their C-Code implementations in conjunction with the Crypto-Tile. The AXI4

MCU interface is the communication channel between the processor and the Crypto-Tile, while DMA operations are emulated to ensure accurate testing. Simulations can be performed by executing the designated command within the verification environment. The generated simulation logs, accessible in the specified *build* folder, facilitate result analysis and troubleshooting in case of compilation or simulation errors.

# References

1. Intel Software Guard Extensions (Intel SGX)—Key Management on the 3rd Generation Intel Xeon Scalable Processor. Technical report, Intel (2019)
2. Coppolino L, D'Antonio S, Mazzeo G, Romano L (2019) A comprehensive survey of hardware-assisted security: from the edge to the cloud. Internet Things 6:100055
3. Crocetti L, Baldanzi L, Bertolucci M, Sarti L, Carnevale B, Fanucci L (2019) A simulated approach to evaluate side-channel attack countermeasures for the advanced encryption standard. Integration 68:80–86 September
4. Di Matteo S, Baldanzi L, Crocetti L, Nannipieri P, Fanucci L, Saponara S (2021) Secure elliptic curve crypto-processor for real-time iot applications. Energies 14(15)
5. Gupta S (2023) An edge-computing based Industrial Gateway for Industry 4.0 using ARM TrustZone technology. J Ind Inf Integr 33:100441
6. Kovač M et al (2022) European processor initiative: Europe's approach to exascale computing
7. McKeen F, Alexandrovich I, Berenzon A, Rozas CV, Shafi H, Shanbhogue V, Savagaonkar UR (2013) Innovative instructions and software model for isolated execution, vol 10
8. Nannipieri P, Bertolucci M, Baldanzi L, Crocetti L, Di Matteo S, Falaschi F, Fanucci L, Saponara S (2021) SHA2 and SHA-3 accelerator design in a 7 nm technology within the European processor initiative. Microprocess Microsyst 87
9. Nannipieri P, Di Matteo S, Baldanzi L, Crocetti L, Belli J, Fanucci L, Saponara S (2021) True random number generator based on fibonacci-galois ring oscillators for FPGA. Appl Sci (Switzerland) 11(8)
10. Nannipieri P, Matteo S, Baldanz L, Crocetti L, Zulberti L, Saponara S, Fanucci L (2022) VLSI design of advanced-features AES crypto processor in the framework of the european processor initiative. IEEE Trans Very Large Scale Integr (VLSI) Syst 30(2):177–186
11. Nannipieri P, Crocetti L, Matteo SD, Fanucci L, Saponara S (2023) Hardware design of an advanced-feature cryptographic tile within the european processor initiative. IEEE Trans Comput 1–14
12. Pinto S, Santos N (2019) Demystifying arm trust zone: a comprehensive survey. ACM Comput Surv (CSUR) 51(6):1–36

13. Zulberti L, Di Matteo S, Nannipieri P, Saponara S, Fanucci L (2022) A script-based cycle-true verification framework to speed-up hardware and software co-design: performance evaluation on ECC accelerator use-case. Electronics (Switzerland) 11(22)

# Machine Learning for SOC Estimation in Li-Ion Batteries

Di Dio Riccardo[2](✉), Aurilio Gianluca[2], Di Rienzo Roberto[1], and Saletti Roberto[1]

[1] Dipartimento Ingegneria Dell'informazione, University of Pisa, Pisa, Italy
[2] Marelli Europe Spa, Bologna, Italy
`r.didio1@studenti.unipi.it`

**Abstract.** State of Charge estimation is very important to deliver essential information about battery charge and aging level of Li-ion batteries in Electric Vehicles. This paper applies the Deep Leaning and Machine Learning approaches comparing decision tree and Long Short-Term Memory for estimating the State of Charge. The datasets for the training and the evaluation have been generated with a Digital Twin model applying driving cycles at different ambient temperature. The proposed Digital Twin model includes non-linear phenomena.

**Keywords:** State of charge (SOC) · Deep learning · Digital twin

## 1 Introduction

Accurate estimation of the battery State of Charge (SOC) is very important to deliver information about charge and aging level [1, 2].

This paper is organized in two sections: a Digital Twin of the battery system is described in the first section, whereas the machine learning approach for SOC estimation is described in the second one.

## 2 Digital Twin of the Battery System

"A Digital Twin is an up-to-date representation, a model, of an actual physical asset in operation. It reflects the current asset condition and includes relevant historical data about the asset. Digital Twins can be used to evaluate the current condition of the asset, and more importantly, predict future behavior, refine the control, or optimize operation" [3].

The advantage of using a Digital Twin is replicate everything in the physical world in the digital space and provide engineers with feedback from the virtual world. This technology allows companies to quickly analyze and solve physical problems, design and build better products, and realize value and benefits faster than previously possible.

The proposed Digital Twin is shown in Fig. 1 and includes the following models:

- battery cell model;
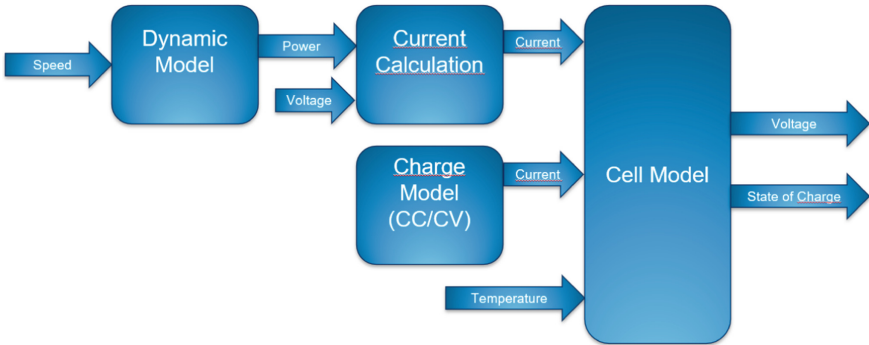- vehicle dynamic model;
- current and voltage charge strategy.



**Fig. 1.** Proposed digital twin to develop battery management system algorithms

## 2.1  Battery Cell Model

The proposed battery cell model is the second order circuit shown in Fig. 2 and includes hysteresis phenomenon.

$$V_{cell} = V_{oc} - \left[ R_{batt} * i_{cell} + \frac{1}{C_D} \int \left( i_{cell} - \frac{v_{DL}}{R_D} \right) dt + \frac{1}{C_{LF}} \int \left( i_{cell} - \frac{v_{LF}}{R_{LF}} \right) dt \right] \tag{1}$$

where:

- $V_{oc}$ is the electromotive force voltage based on OCV-SOC characteristic;
- $v_{cell}$ is the cell voltage;
- $i_{cell}$ is the cell current;
- $T_c$ is the cell temperature;
- $R_D$ and $C_D$ are the resistance and the capacity of the first RC parallel circuit;
- $R_{LF}$ and $C_{LF}$ are the resistance and the capacity of the second RC parallel circuit;
- $R_{batt}$ is the battery resistance.

All parameters are function of temperature $T_c$, SOC, and sign of $i_{cell}$. The cell current is considered positive during the discharge mode and negative during the charge mode.

The hysteresis phenomenon is modeled by adding a third RC group with resistance $R_h$ and capacitance $C_h$, where:

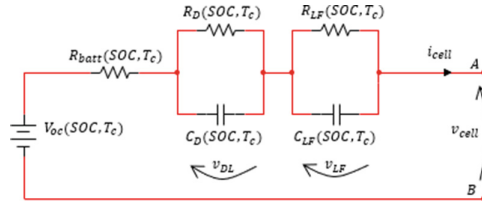$$R_h = \frac{m}{\lceil i \rceil} \tag{2}$$

**Fig. 2.** Equivalent circuit of the battery cell

$$C_h = \frac{1}{\gamma * m} \tag{3}$$

where m is the width of hysteresis and $\gamma$ is the transition speed of hysteresis.

The SOC is calculated with the Coulomb counting equation:

$$SOC = SOC_o + \frac{1}{SOH * C_N} \int (i_{cell})dt \tag{4}$$

where:

- $SOC_o$ is the initial State of Charge;
- $C_N$ is the nominal capacity of the cell;
- SOH is the State of Health.

### 2.2 Dynamic Model

The input of the dynamic model is the speed profile of a driving cycle. The dynamic model calculates the mechanical and electric power.

### 2.3 Charge Mode: Constant Current (CC) and Constant Voltage (CV)

The goal of the recharge system is to perform a charge when the SOC level is below 25%.

The implemented charge model is shown in Fig. 3 and can be described with the following steps:

- during the driving cycle the SOC level is monitored;
- when the cell SOC decreases below the threshold of 25%, the CC charge phase starts with a current value of 1 C-rate, i.e. a current value that fully charge a fully discharged cell in 1 h. This phase is considered completed as soon as the cell voltage reaches the value of the $V_{oc}$ evaluated at 80% of the SOC. Then, the CV charge phase could start.
- CV charge phase is performed, so the cell voltage is kept constant, until the cell SOC reaches the value of 80%.
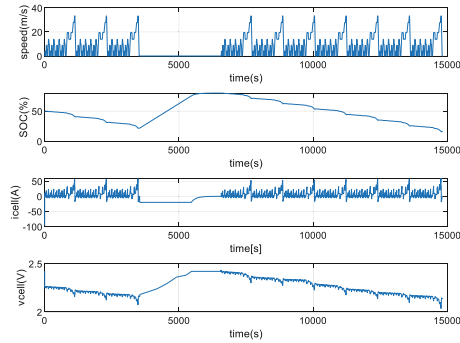
**Fig. 3.** Voltage, current and SOC calculated by the digital twin model with the NEDC speed drive cycle profile at 25 °C

## 3  State of Charge Estimation

This section presents the estimation of battery SOC using machine learning methods.

The features of the machine learning algorithms are the temperature, the current and the voltage calculated by the Digital Twin during the drive cycle operation and during the charge mode; the target of the machine learning algorithm is the cell SOC of the Digital Twin computed through the Coulomb counting equation.

The drive cycles used for the training, validation, and test datasets are reported in the following Table 1.

**Table 1.**  Drive cycle used for training, validation, and test dataset

| Drive cycle | |
|---|---|
| LA92short | Train |
| J1015 | Train |
| HWFET | Train |
| ArtMw150 | Train |
| ArtUrban | Train |
| FTP | Validation |
| ArtRoad | Validation |
| NEDC | Test |

Two different machine learning algorithms are developed and compared in this work. They are the Long Short-Term Memory (LSTM) and the Decision Tree.

The performance of the machine learning algorithms has been evaluated with the followings metrics: Mean Absolute Percentage Error (MAPE), Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE).

The LSTM is a deep learning network with memory capable of learning order dependence in sequence prediction problems.

The best implemented LSTM is composed by 5 layers with the following hyper parameters:

1. Sequence input with 3 dimensions. Each sequence is splitted into segments of 100 samples in order to improve the numerical behavior of the training.
2. LSTM with 32 hidden units.
3. 50% Dropout.
4. 1 fully connected layer.
5. Regression output.

The LSTM has been trained with a learning rate of 0.01 and 32 mini batchsize. The optmizer chosen is Adam.

The implemented Decision Tree is a regression one decision tree with a minimum leaf size of 4.

## 4   Results and Conclusions

The estimated State of Charge by the LSTM network and the decision tree with the NEDC drive cycle has been compared with the SOC calculated by the Coulomb counting equation and shown in Figs. 4 and 5. The performace of the two SOC estimation algorithms are reported in Table 2.

The results show that both methods provide very good estimates of the SOC. In particular, the LSTM shows a slightly better performance in all the metrics adopted.

The results obtained with LSTM and Decision Tree are comparable with the actual state of the art [2].
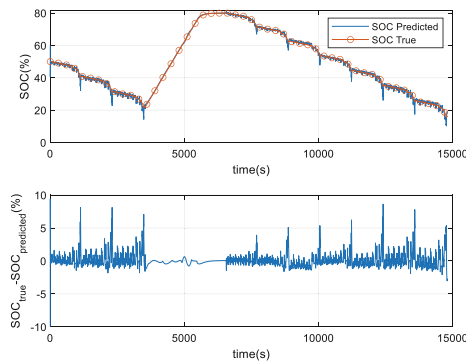


**Fig. 4.** The first plot shows the comparison between the SOC true and the SOC predicted by the LSTM; the second plot shows the difference between the SOC true and the SOC predicted by the LSTM
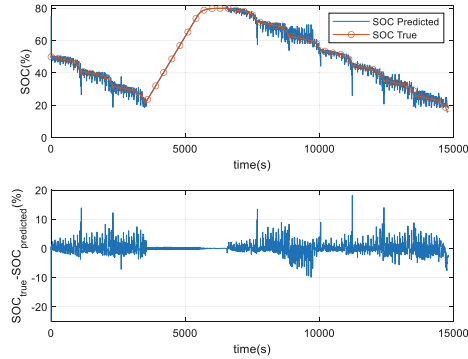
**Fig. 5.** The first plot shows the comparison between the SOC true and the SOC predicted by the decision tree; the second plot shows the difference between the SOC true and the SOC predicted by the decision tree

The RMSE, MAPE and MAE are reported in the following table:

**Table 2.** RMSE, MAPE and MAE for the LSTM deep learning network and the decision tree

|               | RMSE (%) | MAPE (–) | MAE (%) |
| ------------- | -------- | -------- | ------- |
| LSTM          | 0.85503  | 1.4828   | 0.51116 |
| Decision tree | 1.061057 | 1.582    | 0.61057 |

# References

1. Alamin KSS, Chen Y, Macii E, Poncino M, VincoS (2022) A machine learning-based digital twin for electric vehicle battery modeling. In: 2022 IEEE international conference on omni-layer intelligent systems (COINS), Barcelona, Spain, pp 1–6. https://www.researchgate.net/publication/362801333_A_Machine_Learning-based_Digital_Twin_for_Electric_Vehicle_Battery_Modeling
2. Li H, Kaleem MB, Chiu IJ, Gao D, Peng J (2021) A digital twin model for the battery management systems of electric vehicles. In: 2021 IEEE 23rd international conference on high performance computing & communications; international conference on data science & systems; 19th international conference on smart city; 7th international conference on dependability in sensor, cloud & big data systems & application (HPCC/DSS/SmartCity/DependSys), Haikou, Hainan, China, pp 1100–1107. https://doi.org/10.1109/HPCC-DSS-SmartCity-DependSys53884.2021.00171
3. https://it.mathworks.com/discovery/digital-twin.html

# HUSTLE: A Hardware Unit
# for Self-test-Libraries Efficient Execution

Nicola Ferrante[1,2]([✉]) [ID], Francesco Terrosi[3] [ID], Luca Maruccio[1], Francesco Rossi[1],
Luca Fanucci[2] [ID], and Andrea Bondavalli[3] [ID]

[1] Resiltech, 56025 Pontedera, PI, Italy
nicola.ferrante@resiltech.com
[2] University of Pisa, 56126 Pisa, PI, Italy
[3] University of Florence, 50121 Firenze, FI, Italy

**Abstract.** Online testing of computer systems is crucial in contexts such as the safety-critical domain, where the software is usually made of functional code, which is the code implementing the application-specific functionalities, and non-functional code, which implements auxiliary functionalities, e.g., test routines. By periodically running a test routine it is possible to satisfy the high dependability requirements mandated by regulators, and defined in safety standards such as ISO26262, IEC61508, and CENELEC EN 5012X. Self-Test Libraries (STLs) are a form of software-based self-test, widely used in safety-related applications. The main drawback of this safety mechanism is the overhead imposed on the execution of the functional code, and reducing this overhead is a well-known challenge in research. We propose here HUSTLE, a Hardware Unit for STL Efficient execution, which can be integrated into the chip design with no modification to the CPU's internal logic. We also propose a scheduling mechanism that allows HUSTLE to efficiently execute self-tests, by exploiting the CPU's idle time. This is achieved by storing test code in a separate memory and sending instructions to the CPU, bypassing the Instruction Cache, thus allowing to reduce the overall execution time and the cache interference of STL, while CPU utilization increases.

**Keywords:** Software-based self-test · Safety-critical systems · Embedded systems testing

## 1 Introduction

In safety-critical systems, protecting the CPU from hardware faults in the field is a fundamental requirement [1, 2]. Many protection techniques were proposed, based on both hardware (HW) and software (SW) mechanisms, each of them providing different levels of protection. HW-based techniques are faster but require modifications to the original design, whereas SW-based techniques have little to no impact on the device area, but significantly reduce performance [3, 4]. To reduce the impact of safety mechanisms on performance, we propose an approach that exploits the idle time of the CPU to execute test routines, i.e., Self-Test Libraries (STL). Well-established Functional Safety standards, such as ISO 26262, require high coverage of random hardware faults, for which

STLs have been identified as an effective safety mechanism, providing high coverage without any impact on the device's design [5]. However, to achieve the target protection level STLs need to stimulate as much as possible the internal logic of the component, imposing a significant overhead on the nominal execution [6]. We present here HUSTLE, a hardware mechanism to improve the execution efficiency of STLs. HUSTLE provides STL instructions to the core without accessing the Instruction Cache (IC), reducing the overall response time and the IC pollution (cache misses for STL instructions are eliminated), since the cache will keep in memory only instructions relative to the functional workload. To demonstrate the possibilities offered by HUSTLE we also implemented an efficient event-triggered scheduling mechanism by exploiting architectural signals to detect CPU's idle time, and use this time to execute STL instructions. HUSTLE allows to reduce the overhead on the execution time and the interference with the cache, increasing core utilization.

## 2 Related Works

The use of instruction-based self-test, i.e., STL, is a widely used technique [7], and improving their efficiency is an active research field. Most designers of such test libraries try to leverage target architecture's resources to achieve a higher protection degree while meeting constraints on the imposed overhead [8, 9], while others try to develop finer algorithms that are both general and efficient. E.g., in [10] the authors merge different SW techniques combining random program generation and signature-based self-checking; in [11] the authors focus on improving the automatic generation of test programs, while [12] proposes a technique that relies solely on available processor resources. The scheduling of such tests is crucial to their efficiency and effectiveness [9, 13]. We found no study about scheduling mechanisms that exploit idle times in the CPU by executing fragments of STL code. Using dedicated HW support to speed up software operations is a common approach [14], however, we found only one work that adopts this approach towards self-testing, by storing STL instructions in a dedicated memory [15]. In this work, the authors focus on implementing hardware support that can store test code and data. The STL is implemented as an Interrupt Service Routine, but no specific scheduling strategy is proposed.

## 3 HUSTLE Overview

### 3.1 Mechanism Description

HUSTLE is placed between the CPU and the IC, to provide STL instructions to the Core as fast as possible. This is achieved as follows: during the execution of functional code, HUSTLE, is in IDLE state, and it just forwards requests received from the CPU to the IC, and the subsequent responses from the IC to the CPU. When the CPU starts executing the STL it requests addresses in the address range of HUSTLE's ROM, HUSTLE transitions into the ACTIVE state. In this state, it blocks requests from the CPU to the IC and responds to the CPU with the requested instructions in place of the IC. Internally, HUSTLE has a Read Only Memory (ROM) containing the STL payload, i.e., the set of

STL instructions that must be executed. A high-level schema of HUSTLE's interface is depicted in Fig. 1a, while the transition diagram is shown in Fig. 1b. This addition to the design allows to: (i) reduce STL instruction latency, as STL instructions are not fetched from the IC or lower levels of cache, (ii) reduce IC pollution, since a request to an address relative to an STL instruction will be handled by a separate memory, and (iii) increase overall core's resource usage, as result of (i) and (ii).
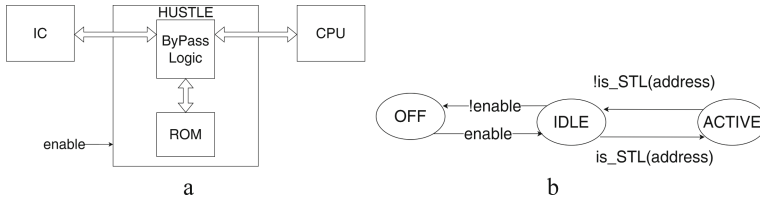


**Fig. 1.** **a** HUSTLE block diagram, **b** HUSTLE transition diagram.

### 3.2 HUSTLE Scheduling Strategy

HUSTLE's scheduling strategy is based on monitoring architectural signals to detect events that may cause instruction starvation and leave the core idle. One of the most frequent causes of instruction starvation are IC misses [16]. When a miss in the IC happens, the core must wait for its resolution to execute the next instruction. HUSTLE can exploit this time to execute STL instructions stored in its ROM.

The proposed scheduling strategy triggers a control-flow redirection sending a hardware interrupt as soon as a cache miss is detected. In our prototype, the IC miss signal is routed to HUSTLE and connected to the interrupt controller, using the STL as Interrupt Service Routine (ISR). Note that, without ad-hoc HW support it wouldn't be possible to exploit idle times caused by cache misses, because the control flow redirection to handle the interrupt could cause additional cache misses itself, thus invalidating all the benefits. In Fig. 2 we illustrate a simple scheduling example showing how idle times could be leveraged by HUSTLE. It is important to note that the STL code shall take into account proper mechanisms to avoid interference with other running applications. For example, context-switching techniques, granularity of sequences of consecutive STL tests to execute, and their priority. In general, all the aspects that can affect the impact of this scheduling strategy on the execution time of other applications should be analyzed, based on the requirements of the target system. The analysis of these aspects is not reported in this work for brevity.

## 4 Experimental Activity

To evaluate HUSTLE, a prototype has been developed using the Chipyard [17] framework. We chose as target CPU architecture the RISC-V [18] Berkley Out-of-Order Machine (BOOM) Core [19]. To evaluate the impact of the developed solution on different hardware configurations, we selected two configurations of the BOOM core:
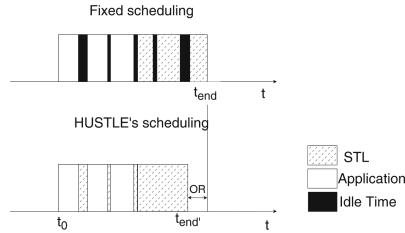
**Fig. 2.** Visual representation of HUSTLE's scheduling strategy. OR stands for overhead reduction.

SmallBoom, a single pipeline core, and the MediumBoom, a two-wide pipeline core. The Verilog RTL description of the design is generated using Chipyard's toolchain, which was then compiled and simulated using Synopsys VCS.

HUSTLE is placed between BOOM's Frontend (Containing the IC), and the Core. A Control Status Register (CSR) was created to allow enabling HUSTLE via SW. HUSTLE's ROM is initialized with STL instructions, which are mapped at a specific address that is configured at runtime as the ISR address for HUSTLE's interrupt.

The SW used for evaluation is made of three components: the workload, the STL, and a scheduler function. The workload is made of integer and floating-point operations e.g., addition, multiplication, and bit shifts, and it is large enough to fill the IC, to guarantee IC misses to happen during its whole execution. The scheduler function consists of a for-loop with a variable number of iterations. Inside the for-loop, the workload is scheduled once for each iteration. The STL is composed of 10 signature-based arithmetic and logic tests developed in assembly, and it includes a context switch routine to preserve application(s) context. We developed three variations of the code which we call: test_1, test_2, and test_3. In test_1 STL instructions are scheduled after each workload execution. Here we intend to model the worst-case for STL execution, where the whole IC is filled with functional-code instructions. In test_2, the STL is scheduled periodically, with higher priority than the workload. This test is used to evaluate the effect of HUSTLE in a common-case scenario, in which some STL instructions may be already in the IC when the STL starts its execution. Finally, in test_3 we evaluate HUSTLE's scheduling strategy, against a fixed periodic scheduling strategy. Each software is tested on the BOOM Core as is, and the BOOM Core with HUSTLE. Hereafter we define a set of metrics to evaluate HUSTLE. Values measured on HW configuration with HUSTLE are reported with the $h$ subscript.

**Execution overhead of STL (OR)**: To evaluate the reduction of the overhead caused by the STL, we define $C$ as the number of additional clock cycles required to execute the STL w.r.t. the baseline. Overhead reduction is computed as $OR = 1 - (C_h/C)$.

**CPU usage during STL ($\Delta$IPC)**: To evaluate HUSTLE's effect on the CPU resource usage we measure the Instructions Per Cycle (IPC) on the overhead imposed by the STL. We define $N$ as the number of executed STL instructions, thus $IPC = N/C$. Then we measure the difference on HW configurations with and without HUSTLE as $\Delta IPC = IPC_h - IPC$.

**Cache Interference of STL (IR)**: To measure the cache interference of the STL we define $M$ as the number of additional cache misses caused by STL instructions w.r.t the baseline. We compute the interference reduction as $IR = 1 - (M_h/M)$.

## 5   Results

In this section, we discuss the results obtained from the experimental campaign. Data are reported in Table 1.

**Table 1.**   Metrics for *test_1* and *test_2*.

| SW | HW | OR | $\Delta$IPC | IR |
|----|-----|------|------|------|
| *test_1* | Small | 0.355 | 0.278 | 1.026 |
|  | Medium | 0.496 | 0.477 | 0.907 |
| *test_2* | Small | 0.032 | 0.027 | 0.953 |
|  | Medium | 0.058 | 0.069 | 1.035 |

We see that HUSTLE provides an OR up to 40% on the worst-case scenario. Whereas in the common case, i.e., *test_2*, the benefits are reduced. This may be because STL instructions are still present in the cache, hence the use of HUSTLE has not a considerable impact as in *test_1*. We note also that the improvement provided by HUSTLE is higher with the MediumBoom configuration than with the SmallBoom in both tests. The increase in resource usage of the CPU (column $\Delta$IPC) is much higher in *test_1*; this is expected since the idle time due to cache misses is reduced and instructions are served directly by HUSTLE. The MediumBoom approximately doubles HUSTLE's benefits w.r.t. SmallBoom configuration. The results on cache interference (column IR) present a different trend than the others. As can be noted, the SmallBoom has a higher benefit in this case: in *test_1* the number of cache misses is less than the baseline (IR > 1), whereas in the MediumBoom is always lower than 0.91. We argue that this may be due to branch predictor and speculative execution.

Having confirmed the benefits of HUSTLE, we now compare two different scheduling strategies: *periodic* at a fixed time interval and *event-driven* in correspondence of cache misses. Values are reported in Table 2. Looking at the rows "HUSTLE" and "Periodic", results show that the proposed scheduling strategy outperforms the fixed scheduling strategy in both Small and Medium configurations. We can see that HUSTLE's scheduling strategy almost doubles $\Delta$IPC and OR for the Small configuration, while more than doubling these values for the Medium configuration. The Medium configuration shows the larger benefits, approximately doubling the increase in IPC w.r.t. a fixed scheduling strategy, and the same applies to the OR.

**Table 2.** *test_3* results, the baseline execution considered is *test_2*.

| HW | Scheduling | $\Delta$IPC | OR |
|---|---|---|---|
| Small | HUSTLE | 0.049 | 0.057 |
| | Periodic | 0.027 | 0.032 |
| Medium | HUSTLE | 0.133 | 0.107 |
| | Periodic | 0.069 | 0.058 |

## 6   Conclusions

In this work, we presented HUSTLE, a HW module that allows efficient execution of STL code. We showed how with this module it is possible to reduce the overhead and cache interference of STL execution while increasing the CPU utilization. We also provide an implementation of an efficient event-triggered STL scheduling mechanism to show the benefits provided by HUSTLE. In our experiments, we demonstrated that this mechanism increases CPU utilization (IPC) while reducing the overhead on execution time. Moreover, we found that by using HUSTLE it's possible to reduce the interference between functional and non-functional code, in particular the allocation of cache lines for non-functional code, and in some cases reducing also the cache misses in functional code.

## References

1. Peleska J, Siegel M (1996) Test automation of safety-critical reactive systems
2. Parnas DL et al (1990) Evaluation of safety-critical software. ACM
3. Osinski L, Langer T, Mottok J (2017) A survey of fault tolerance approaches on different architecture levels. ARCS
4. Abid A, Khan MT, Iqbal J (2021) A review on fault detection and diagnosis techniques: basics and beyond. Artif Intell Rev 54
5. Pratas F et al (2018) Measuring the effectiveness of ISO26262 compliant self-test library. In: 2018 ISQED. IEEE
6. Malaiya YK et al. The relationship between test coverage and reliability. In: Proceedings of 1994 IEEE international symposium on software reliability engineering
7. Psarakis M et al (2010) Microprocessor software-based self-testing. IEEE Des Test Comput 27(3):4–19
8. Bernardi P et al (2015) Development flow for online core self-test of automotive microcontrollers. IEEE Trans Comput 65(3):744–754
9. Paschalis A, Gizopoulos D (2004) Effective software-based self-test strategies for on-line periodic testing of embedded processors. IEEE Trans Comput Aided Des Integr Circuits Syst 24(1):88–99
10. Kranitis N et al (2008) Hybrid-SBST methodology for efficient testing of processor cores. IEEE Des Test Comput 25(1):64–75
11. Riefert A et al (2016) A flexible framework for the automatic generation of SBST programs. IEEE Trans on Very Large Scale Integr (VLSI) Syst 24(10)
12. Hatzimihail M et al (2007) A methodology for detecting performance faults in microprocessors via performance monitoring hardware. In: 2007 IEEE international test conference

13. Li Y et al (2009) Operating system scheduling for efficient online self-test in robust systems. In: International conference on computer-aided design
14. Dally WJ, Turakhia Y, Han S (2020) Domain-specific hardware accelerators. Commun ACM 63(7):48–57
15. Bernardi P et al (2013) MIHST: a hardware technique for embedded microprocessor functional on-line self-test. IEEE Trans Comput 63(11):2760–2771
16. Kanev S et al (2015) Profiling a warehouse-scale computer. In: Proceedings of the 42nd annual international symposium on computer architecture
17. Chipyard. https://chipyard.readthedocs.io/en/stable/. Accessed 28 Mar 2023
18. RISC-V ISA. https://riscv.org/. Accessed 28 Mar 2023
19. BOOM Core. https://boom-core.org/. Accessed 28 Mar 2023

# Towards Context-Aware Classrooms: Lessons Learnt from the ACTUA Project

Edgar Batista, Antoni Martínez-Ballesté, Joan Rosell-Llompart, and Agusti Solanas<sup>(✉)</sup>

Universitat Rovira i Virgili, Av. Països Catalans 26, 43007 Tarragona, Spain
{edgar.batista,antoni.martinez,joan.rosell,agusti.solanas}@urv.cat

**Abstract.** Context-aware classrooms augment traditional classrooms with unprecedented sensing and communication capabilities to monitor and analyse large amounts of contextual data. Enriching classrooms with these data opens up new opportunities to enhance the learning experience in education centres. However, implementing and deploying context-aware classrooms is not straightforward. This article describes the challenges addressed by the ACTUA project, which aims to advance the development of context-aware classrooms in primary schools.

**Keywords:** Context-Aware Classrooms · Ubiquitous Computing · Internet of Things · Education · Schools

## 1 Introduction

Technological advances are adopted to transform our cities [1], governments [2], and healthcare infrastructures [3] (to name a few) into their smart and cognitive counterparts. Schools, and classrooms in particular, are essential in today's society. These environments must be adequately set up to promote the teaching and learning processes for successful knowledge acquisition. Environmental comfort is critical to increasing productivity in workplaces and educational centres [4]. To this end, ubiquitous technologies, IoT devices and high-speed data networks enable real-time monitoring of environmental parameters at low cost. Unfortunately, this kind of monitoring is uncommon in school classrooms, despite their adverse conditions [5]. Equipping classrooms with sensing, computing and communication capabilities opens the door to context-aware classrooms, a next-generation model capable of processing vast amounts of data to adjust their behaviour (*e.g.,* devices, systems) accordingly. Building context-aware environments is not trivial. While many context-aware approximations are theoretical, deploying them in real-life scenarios is complex.

In this article, we present the main experiences, challenges, and barriers that we have faced while implementing and deploying a first approximation of context-aware classrooms within the scope of the ACTUA project. Our solution runs in a hundred classrooms in Catalonia (Spain). The rest of the article is organised as follows: Sect. 2 contextualises the goals of the ACTUA project, Sect. 3 describes

the main challenges that we have experienced during the implementation and deployment of these context-aware classrooms, Sect. 4 elaborates on the opportunities highlighted by the teaching staff, and Sect. 5 closes the article with some concluding remarks and future research lines.

## 2 The ACTUA Project

The public health emergency triggered by the COVID-19 pandemic has highlighted the need for environmental sensing tools to fight airborne transmission of respiratory diseases [6]. Classrooms have many risk factors that can increase the likelihood of transmission: enclosed spaces that are typically poorly ventilated and where high levels of aerosol formation occur over extended periods of time. To make classrooms safer and more robust environments, the ACTUA project [7] aims to understand the relationship between health effects on the school population and classroom environment variables and characteristics [8]. Primary schools are the main target of this project.
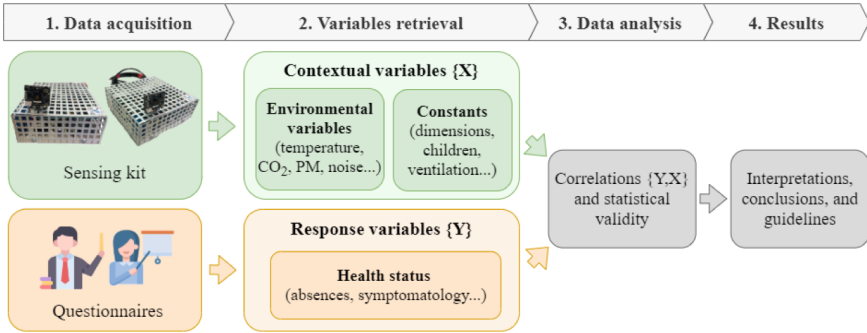


**Fig. 1.** Overview of the ACTUA project from a data perspective.

In keeping with the smart healthcare paradigm [9], the context of the classrooms is monitored automatically using low-cost sensing kits [10]. In particular, several classroom context variables are periodically collected, namely air temperature, relative humidity, barometric pressure, $CO_2$ concentration, particulate matter, light, UV radiation, ventilation, noise, and people's movement. Moreover, since classrooms can vary greatly (especially between schools), other constants must also be taken into account for a proper analysis, such as the size, floor and orientation of each classroom, the number and age of children in that classroom, and the number and dimensions of doors and windows per classroom, among others (see Fig. 1). Teachers fill out weekly questionnaires to measure the health status of the classroom's population. There, information about how many students were absent (disaggregated by reason—illness, non-medical, unknown—and temporality) or how many students had symptoms of respiratory disease is asked. During the 2022/23 academic course, our solution has been steadily installed in a hundred classrooms throughout Catalonia.

# 3   Challenges

Developing context-aware classrooms raises challenges, not only technical but societal and regulatory. Next, we highlight the main challenges addressed.

## 3.1   Technical

Sensing capabilities are the cornerstones of a contextual classroom. In this sense, nanotechnology, pervasive electronics, ubiquitous computing, and IoT devices are vital enablers. Sensing devices with self-awareness, self-reconfiguration, and self-healing are desirable to achieve absolute resiliency. This would enable devices to prevent, detect and respond autonomously to any event or failure that could affect their operation. However, these features are constrained by sensing devices' limited computational capabilities, memory and power limitations.

Our sensing kit comprises a Raspberry Pi, a single-board computer acting as a controller node, connected to several sensors that periodically collect contextual variables. To keep costs down, low-cost sensors are used, such as the SCD30 sensor to measure air temperature, relative humidity and $CO_2$ concentration, and the PMS5003 sensor to measure particulate matter [10]. These sensors are suitable for integration into context-aware environments because of their simplicity, efficiency, cost, and ease of programming [11]. However, their accuracy is inferior to that of professional metres, which are usually very expensive, require calibration and allow offline measurements [8].

Beyond sensing capabilities, communication and information sharing are additional required capacities. In particular, these are fundamental to properly coordinating and managing the whole infrastructure. In our solution, sensing kits have different communication channels. For instance, the sensing kits use the *aliveness channel* to inform external entities that they are up and running. This channel detects malfunctioning devices when they do not send signals periodically. Besides, there is a *data channel* wherein devices can transmit the collected data to external parties if needed. Also, sensing kits can react to messages received from external entities through the *command-and-control (C&C) channel*. The C&C functionality permits controlling devices remotely with predefined commands, such as modifying configuration parameters or executing specific operations. Hence, this channel allows devices to coordinate among themselves upon failures or external events. All these communication channels run over the HTTPS protocol to keep communications secure.

## 3.2   Societal

Augmenting classrooms with technology must go hand in hand with the social acceptance of the school population. The functionalities of context-aware classrooms must be transparent to students and teachers and must not interfere with teaching. For this reason, devices should be small, discreet, minimally intrusive, plug-and-play and well-integrated into classrooms. In addition, any interfaces

that require user interaction should be intuitive and as simple as possible to enhance the user experience.

However, these technological advances might cause scepticism among teachers because they have a negative perception of surveillance, resembling "Big Brother" scenarios. Concerns could magnify when augmenting classrooms with microphones or cameras for measuring ambient noise, voice recognition, locating people, or tracking people's movements. Informing and educating teachers, students and families about what context-aware classrooms can and cannot do is paramount. Societal readiness will be critical to the success of this paradigm.

### 3.3   Regulatory and Ethical

The ubiquity nature of context-aware classrooms makes user awareness complex. Current data protection regulations, such as the EU General Data Protection Regulation, make consent and awareness mandatory before data processing. As data can be collected opportunistically in context-aware classrooms, balancing meeting legal obligations and developing innovative, ubiquitous services is challenging. In this line, but primarily because most of the classroom population is underage, all systems must be built upon privacy-by-design principles, encompassing an ethical dimension in technological design. For this reason, transiting from regular classrooms to context-sensitive classrooms will not happen overnight.

Ethical committees analysed the ACTUA solution before its deployment in real school classrooms. In particular, we highlight the three most sensitive aspects that we had to address to be legally and ethically compliant: (i) measuring the classroom's noise using a microphone in the sensing kit, (ii) measuring people's movements using a camera in the sensing kit, and (iii) the collection of health data. Regarding noise measurements, the sensing kit computes the average sound level of the classroom (in decibels) without using voice analysis or speech recognition. The main drawback is that this value is susceptible to external noises (*e.g.,* vehicles, pedestrians). Regarding people's movement, we designed a numerical algorithm that compares people's places frame-by-frame using the YOLOv3's object recognition algorithm [12]. Avoiding taking images and recording videos relaxes privacy threats.

Regarding health data, the collection of medical records from students and teachers is bureaucratically complex. As an alternative approach, we have designed questionnaires that quantify the general health of the classroom population rather than asking for individual health data. These questionnaires are answered every week by the teachers in the classroom.

### 3.4   Deployment

Classrooms vary widely in terms of dimensions, equipment and furniture distribution, electrical installation and Internet availability, among other things. These differences are accentuated among schools, so there is no standard way to deploy our solution. For instance, our sensing kits are located in different places:

some are on top of tall cabinets, on shelves at medium height, or the teacher's table, and either closer to the windows, the door or at the back of the classroom. This decision depends on the availability of wall jacks, the Internet connectivity through RJ45 sockets (if WiFi is not possible), and the teaching staff's approval. Hence, installations had to adapt to each classroom's configuration.

From a technical perspective, network configurations vary among schools. We observed that schools with specific IT departments generally place more emphasis on network security (*e.g.,* strict-policy firewalls, MAC address filtering, ports disabled such as FTP...). Besides, when connecting devices to the school's WiFi, it is worth noting that any change in its configuration (*e.g.,* changes in the SSID or password) will require re-configuring the devices. These aspects must be considered to enhance the system's maintainability.

## 4   The Perspective of Teachers

Our solution has been running in 101 classrooms distributed across 38 primary schools during the 2022/23 school year. A questionnaire was distributed to the teachers and principals at the end of the course to assess our solution's acceptance and evaluate the possibilities for contextualised classrooms.

From the teachers' perspective, the technological components installed in their classrooms are not a problem. 84% report that the sensing kit has not caused any distraction to students, and 86% agree that it has not interfered with their teaching tasks. Regarding health questionnaires, around 85% of the teachers were committed and highlighted that filling them out was not an issue. Reporting this information once a week seems acceptable since 77% of the teachers could not report it more assiduously, *e.g.,* every day. Also, most teachers (57%) think that the family members of the students should report this health information. Overall, the teachers' opinion on our solution is very satisfactory (59%) and satisfactory (34%).

According to principals, 76% think that learning processes could be improved by exploiting and monitoring contextual data, and 66% will report these data to improve workplace conditions to the extent possible. Indeed, almost all principals (91%) agree on integrating monitoring systems, such as the ACTUA solution, in future schools. Regarding virus transmission, 94% of the principals believe classrooms play a crucial role in spreading disease. On this basis, 85% of them would find useful our solution to monitor and predict infections in the case of a new pandemic. All things considered, the overall perception of our solution among principals is very satisfactory (61%) and satisfactory (33%).

## 5   Conclusions

The addition of technology will enable context-aware classrooms, a groundbreaking paradigm with many new educational opportunities. However, this transformation will be gradual due to the many challenges that must be addressed.

Thanks to the experience gained during the ACTUA project, we have experienced first-hand the technical, societal, regulatory, ethical and deployment barriers to making context-aware classrooms a reality. Although our solution is only a first approximation, the opportunities to improve quality of life and teaching processes for more successful learning are enormous. With the consolidation of context-aware classrooms, networks of classrooms and schools could emerge to communicate with each other and share their experiences and knowledge. This would lead to a new paradigm of distributed, contextualised learning institutions.

# References

1. Machin J et al (2021) Privacy and security in cognitive cities: a systematic review. Appl Sci 11(10):4471
2. Jimenez CE et al (2016) Smart government: opportunities and challenges in smart cities development. In: I. Management Association (ed), Civil and environmental engineering: concepts, methodologies, tools, and applications, pp 1454–1472
3. Batista E et al (2023) Smart health in the 6G era: bringing security to future smart health services. IEEE Commun Mag 1–7. 10.1109/MCOM.019.2300122
4. Brink HW et al (2021) Classrooms' indoor environmental conditions affecting the academic achievement of students and teachers in higher education: a systematic literature review. Indoor Air 31(2):405–425
5. Barrett P et al (2019) The impact of school infrastructure on learning: a synthesis of the evidence. Technical report The World Bank Group, Washington DC, USA
6. Saini J et al (2020) A comprehensive review on indoor air quality monitoring systems for enhanced public health. Sustain Environ Res 30(1):1–12
7. ACTUA: anàlisi Contextual dels factors de mitigació de la Transmissió de la COVID19 a l'aUlA. https://actua-urv.cat/, (Catalan). Accessed 25 June 2023
8. Villanova O et al (2022) Sensing kit for the study of respiratory disease transmission in school classrooms. In: 11th international aerosol conference (IAC), p 825
9. Solanas A et al (2014) Smart health: a context-aware health paradigm within smart cities. IEEE Commun Mag 52(8):74–81
10. Batista E et al (2022) On the deployment of low-cost sensors to enable context-aware smart classrooms. In: Berta R, De Gloria A (eds) Applications in electronics pervading industry, environment and society (ApplePies), vol 1036. Lecture notes in electrical engineering. Springer, Cham, pp 333–338
11. Batista E et al (2021) Sensors for context-aware smart healthcare: a security perspective. Sensors 21(20):6886
12. Redmon J et al (2018) YOLOv3: an incremental improvement. arXiv:1804.02767

# Comparison of Lithium-Ion Battery SoC Estimation Accuracy of LSTM Neural Network Trained with Experimental and Synthetic Datasets

Luca Amyn Hattouti, Roberto Di Rienzo[✉], Niccolò Nicodemo, Alessandro Verani, Federico Baronti, Roberto Roncella, and Roberto Saletti

Dipartimento di Ingegneria dell'Informazione, Università di Pisa, Pisa, Italy
`roberto.dirienzo@unipi.it`

**Abstract.** Data-driven algorithms, such as the neural network ones, seem very appealing and accurate solutions to estimate the lithium-ion battery's State of Charge. Their accuracy is strongly related to the amount of data used in their training phase. Therefore, huge experimental campaigns are needed to effectively train the neural network used to State of Charge estimation. The main idea behind this paper is to mitigate this drawback by training the algorithm with synthetic datasets generated from simulations of a model of the battery, instead of experimentally collected data. Two instances of the same Long-Short-Term-Memory neural network architecture designed for battery State of Charge estimation are trained, one with an experimental dataset, and the other with a synthetic one. The two neural network instances are then evaluated with the same test dataset derived from experimental data and their estimation accuracies are compared. Results show that the performances of the two networks are comparable. The experimental trained neural network scored a RMSE of only 0.3 % lower than the RMSE of the synthetic trained one. These results suggest the possibility of fruitfully using a synthetic training dataset to speed up and reduce the complexity and cost of the training phase of neural network algorithm for battery state of charge estimation.

**Keywords:** State of Charge estimation · Neural network · Model-based training dataset

## 1 Introduction

Recent years have seen a steady growth of the electric vehicles market thanks to the Lithium-ion Batteries (LIBs) characteristics [1]. However, LIB technology still presents some open points concerning safety issues and battery state estimation. One of the most important challenges in LIB management is the accurate estimation of the State of Charge (SoC), which is defined as the charge currently stored in the cell normalized by its actual capacity.

Several types of SoC estimation methods are presented in the literature. They are usually divided in: Direct, Model-Based, and Data-Driven methods. Direct methods use the direct relationship between SoC and measurable battery quantities, such as the cell terminal voltage, and the exchanged current [2]. These approaches are typically more straightforward than the other methods but they require specific test procedures which are not suitable for real-time applications. Instead, the Model-Based approaches describe the behavior of the battery cells with a model that is used to estimate the cell SoC [3]. In particular, estimation algorithms based on Kalman filtering usually offer a very good compromise between complexity and accuracy [4,5]. Finally, Data-Driven methods use machine learning techniques to estimate the cell SoC starting from the measurable quantities and without using a specific cell model and any knowledge of the specific electrochemical properties of the cells [6]. Many machine learning architectures have been explored in the last years to implement SoC estimation algorithms, such as Feedforward Neural Network [7], Recurrent Neural Network [8], and Long-Short-Term-Memory (LSTM) [9], reaching very good accuracy with low computational complexity. However, these methods require a huge quantity of experimental data to execute an efficient training phase [10] and are therefore heavily hampered by the lack of publicly available data and by the high cost of LIB test campaigns.

This paper investigates the effects of using synthetic data rather than experimental ones during the training phase of Data-Driven SoC estimation algorithms. The synthetic dataset is generated starting from a battery model. The time needed to cycle a modeled cell is much less than the one needed to cycle an actual cell, even if a complex model is used. Therefore, the proposed approach allows a consistent reduction of the amount of experimental data, and thus a quicker and less expensive training process. In this work, two instances of a LSTM neural network are used to estimate the cell SoC using different training datasets. The first is trained using experimental data obtained from real cell measurements while the other is trained on a synthetic dataset. Finally, the two LSTM neural networks are tested using the same experimental dataset and compared in terms of accuracy to verify if the use of simulated data could reduce the accuracy of the SoC estimation algorithms.

## 2   Design of the Neural Network and Training Datasets

The neural network used in this work is composed of two hidden layers: an LSTM layer with 95 computational nodes, and a Fully connected layer (FCL) with 60 nodes. The input of the LSTM layer is the sequence of the battery voltage and current samples, while the output of the neural network is the estimated SoC value. The cell voltage and current input sequences are provided to the network a single sample at a time. The network elaborates the couple of samples and provides the estimation cell SOC value of that sample time. Then the process is repeated for all the samples of the input sequences.

The neural network design, training, and test are carried out using the Matlab Deep learning tool. The ADAM optimization method is used in the training with

an initial learning rate of 0.01, and a learning rate drop factor of 0.1 every 800 epochs[11]. The training procedure concludes when 2000 epochs are reached.

As mentioned before, the aim of this work is to compare the SoC estimation accuracy of a LSTM neural network trained with an experimental dataset and a synthetic one. A TSWB-LYP60AHA lithium iron phosphate (LFP) cell manufactured by ThunderSky-Winston with a nominal capacity of 60 Ah is used as case-study. The cell is cycled with specific test profiles to create the experimental dataset. Seven different standard driving cycles are chosen to emulate the use of the LFP cell in an electric vehicle: HWFET, LA92, ArtUrban, ArtMw150, J1015, UDDS, and NEDC. In particular, the LFP cell is first fully charged, then it is discharged repeating one driving cycle until its voltage reaches the cut-off value of 2.8 V. The reference cell SoC for each test is obtained using the Coulomb counting technique with the cell capacity identified in the cell characterization phase [12]. All the experimental tests were performed using a Chroma 17020 battery tester.

The synthetic dataset is instead obtained using the 2-RC equivalent Electrical Circuit Model (ECM) reported in Fig. 1 [13].
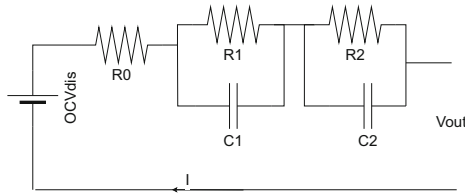


**Fig. 1.** Equivalent electrical circuit model

The model is composed of a resistor $R0$ that is the sum of all the resistive contributions of the cell, two RC groups, $R1$–$C1$ and $R2$–$C2$ to model cell relaxation phenomena, and a voltage generator $OCV_{dis}$ that models the open circuit voltage in the discharge phase. It must be noted that open circuit voltages of LFP cells show a strong hysteresis between the charge and discharge phases. As the experimental dataset considered does not contain charge phases, only the discharge curve of the open circuit voltage is considered in the model. The model parameters are estimated starting from the Pulsed Current Tests (PCTs) performed on the considered cell, according to the procedure reported in [14]. The characterization phase and the execution of the driving cycles are performed only at room temperature. Therefore, the temperature dependency of the model parameters is neglected.

The ECM model is implemented in Matlab Simulink environment and is used to generate the synthetic dataset. This dataset is composed of the model response to the same driving cycle current profiles carried out on the actual cell. It is important to assess whether the model can achieve voltage profiles similar to those of the actual cell. The RMSE obtained by comparing the real

and simulated cell voltages is $16.36\,\text{mV}$, $14.58\,\text{mV}$, $15.6\,\text{mV}$, $13.52\,\text{mV}$, $24.1\,\text{mV}$, $15.7\,\text{mV}$, and $22.78\,\text{mV}$ for ArtUrban, J1015, LA92, NEDC, HWFET, UDDS, and ArtMw150, respectively. The results show a very good accordance between the real and simulated cell voltages, because the relative RMS error goes from 0.4 to 0.7 % of the nominal LFP cell voltage of $3.4\,\text{V}$.

## 3    Comparison Between Experimental and Model-Based Training Datasets

Both the experimental and synthetic datasets are divided in two parts: training and test ones. The training part is composed of the data obtained with HWFET, LA92, ArtUrban, ArtMw150, and J1015 driving cycles. The data of UDDS, and NEDC cycles, instead, compose the test dataset.

Two instances of the same SoC estimation LSTM neural network are separately trained using the experimental dataset and the synthetic one, respectively. For the sake of simplicity, we refer to the network trained with the experimental dataset as $LSTM_{exp}$ while $LSTM_{synth}$ is the network trained with simulated data. The test part of the experimental dataset is used to test the SoC estimation of the two network instances whose results are shown in Fig. 2. The test part of the synthetic dataset is discarded.
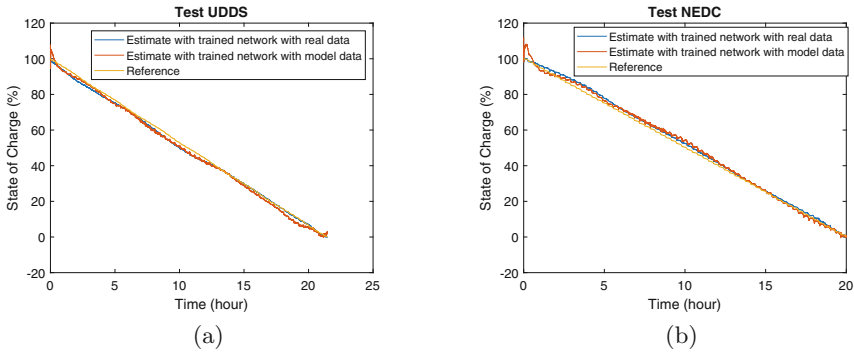


**Fig. 2.** Comparison among reference cell SoC and estimated SoC values by $LSTM_{exp}$ and $LSTM_{synth}$ for UDDS (**a**) and NEDC (**b**) driving cycles

A SoC RMS error of 1.65 and 1.92 % is obtained for the UDDS driving cycle with the $LSTM_{exp}$ and $LSTM_{synth}$, respectively. Similar results are also obtained for the NEDC test. In this case, a SoC RMS error of 1.81 and 2.1 % is obtained with the $LSTM_{exp}$ and $LSTM_{synth}$, respectively. Finally, the obtained results highlight that the $LSTM_{synth}$ network performs worse than $LSTM_{exp}$ especially for SoC values higher than 98 % and lower than 10 %. These errors can be attributed to the reduced accuracy of the ECM in replicating the cell behavior in these SoC ranges, as usually found in the literature [15]. On the other hand, EV batteries rarely operate in these SoC ranges.

# 4    Conclusion

The effects of training a SoC estimation Data-driven algorithm using simulated data instead of experimental data are investigated in this paper. Two datasets, one composed of real experimental data and the other of synthetic data, were used to train and evaluate a SoC estimation Long-Short-Term-Memory neural network. In particular, the experimental dataset was obtained with laboratory tests carried out on a 60 Ah lithium iron phosphate cell. Instead, a 2-RC equivalent circuit model is used to model the cell and generate the synthetic dataset. The same current profiles derived from standard driving cycles were used to discharge the cell and its model in both datasets. Two instances of the same LSTMs were trained, one with the experimental dataset and the other with the synthetic dataset. Both LSTM instances were tested using data obtained from experimental measurement. Results show that the SoC estimation accuracies of the two LSTMs are comparable. Indeed, the LSTM error trained with the synthetic dataset reported an RMS error of only 0.3 % higher than the one of the LSTM trained with the experimental dataset. These findings suggest that simulated data could fruitfully be used as training data in SoC estimation Data-Driven algorithms, resulting in a significant speedup and cost reduction in the algorithm development and validation. Future works will extend the comparison by improving the cell model used, by taking into account additional phenomena, such as hysteresis, temperature, aging, and measurement noise.

# References

1. Bibra EM, Connelly E, Dhir S, Drtil M, Henriot P, Hwang I, Le Marois JB, McBain S, Paoli L, Teter J (2022) Global ev outlook 2022: securing supplies for an electric future
2. Chang WY (2013) The state of charge estimating methods for battery: a review. ISRN Appl Math. https://doi.org/10.1155/2013/953792
3. Morello R, Di Rienzo R, Roncella R, Saletti R, Baronti F Tuning of moving window least squares-based algorithm for online battery parameter estimation. In: 2017 14th international conference on synthesis, modeling, analysis and simulation methods and applications to circuit design (SMACD), pp 1–4. doi: 10.1109/SMACD.2017.7981558

4. Hossain M, Haque M, Arif M (2022) Kalman filtering techniques for the online model parameters and state of charge estimation of the li-ion batteries: a comparative analysis. J Energy Storage 51(104):174. https://doi.org/10.1016/j.est.2022.104174

5. Zhang S, Zhang C, Jiang S, Zhang X (2022) A comparative study of different adaptive extended/unscented kalman filters for lithium-ion battery state-of-charge estimation. Energy 246(123):423. https://doi.org/10.1016/j.energy.2022.123423

6. Chen L, Song Y, Lopes AM, Bao X, Zhang Z, Lin Y (2023) Joint estimation of state of charge and state of energy of lithium-ion batteries based on optimized bidirectional gated recurrent neural network. IEEE Trans Transp Electrification 1. https://doi.org/10.1109/TTE.2023.3291501

7. Ali O, Ishak MK, Memon F, Asaari MSM (2022) Estimation of battery state-of-charge using feedforward neural networks. In: 2022 19th international conference on electrical engineering/electronics, computer, telecommunications and information technology (ECTI-CON), pp 1–4. https://doi.org/10.1109/ECTI-CON54298.2022.9795401

8. Sherstinsky A (2020) Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network. Phys D: Nonlinear Phenom 404(132):306. https://doi.org/10.1016/j.physd.2019.132306

9. Xiao B, Liu Y, Xiao B (2019) Accurate state-of-charge estimation approach for lithium-ion batteries by gated recurrent unit with ensemble optimizer. IEEE Access 7:54192–54202. https://doi.org/10.1109/ACCESS.2019.2913078

10. Vidal C, Malysz P, Kollmeyer P, Emadi A (2020) Machine learning applied to electrified vehicle battery state of charge and state of health estimation: state-of-the-art. IEEE Access 8:52796–52814. https://doi.org/10.1109/ACCESS.2020.2980961

11. Kingma DP, Ba J (2017) Adam: a method for stochastic optimization

12. Kanchan D, Nihal, Fernandes AP (2022) Estimation of SoC for real time EV drive cycle using Kalman filter and coulomb counting. In: 2022 2nd international conference on intelligent technologies (CONIT), pp 1–6 (2022). https://doi.org/10.1109/CONIT55038.2022.9848066

13. Morello R, Di Rienzo R, Roncella R, Saletti R, Baronti F (2018) Hardware-in-the-loop platform for assessing battery state estimators in electric vehicles. IEEE Access 6:68210–68220. https://doi.org/10.1109/access.2018.2879785

14. Barzacchi L, Lagnoni M, Rienzo RD, Bertei A, Baronti F (2022) Enabling early detection of lithium-ion battery degradation by linking electrochemical properties to equivalent circuit model parameters. J Energy Storage 50(104):213. https://doi.org/10.1016/j.est.2022.104213

15. Zhang L, Peng H, Ning Z, Mu Z, Sun C (2017) Comparative research on RC equivalent circuit models for lithium-ion batteries of electric vehicles. Appl Sci 7(10). https://doi.org/10.3390/app7101002

# Smart Kinetic Floor System for Energy Harvesting and Data Acquisition in High Foot-Traffic Areas

Gabriele Ciarpi[(⊠)], Ettore Noccetti, Luca Ceragioli, Marco Mestice, Daniele Rossi, and Sergio Saponara

Department of Information Engineering, University of Pisa, 56122 Pisa, Italy
`gabriele.ciarpi@unipi.it`

**Abstract.** This work explores energy harvesting through kinetic energy capture from human steps. The proposed smart floor system, consisting of multiple smart tiles, offers a promising solution for energy generation and data acquisition in high foot-traffic areas, such as shopping centers. The smart tile incorporates an energy generation and storage system, along with a data acquisition and transmission system. The use of only an accelerometer for both step and energy data acquisition minimizes power impact. An edge computing approach processes acceleration data directly on the tile, transmitting essential information, such as event steps and generated energy, to the cloud via the tile WiFi connection. This information can be used for floor optimization and commercial uses based on customer tracking. Sustainability analysis indicates that for real-time monitoring the current smart tile system requires around 15540 steps per tile for 10 h for sustainability, but this can be reduced to 362 steps by implementing power-saving techniques if a real-time feature is not required. Further research can lead to practical and commercially viable applications, contributing to a greener future.

**Keywords:** Kinetic energy harvesting · Smart tile · People tracking

## 1 Introduction

Energy harvesting is a promising field for collecting energy from various environmental sources, like light, heat, wind, and kinetic energy. Kinetic energy from human and vehicle motions is an untapped resource that could be harnessed efficiently. Researchers and companies have introduced energy harvesting floors to capture this energy. They are based on piezoelectric materials, electromagnetic motors, triboelectric effects, and hybrid combinations [1–3].

Foot traffic is crucial for the effective use of this technology. Since kinetic floors have a higher cost than traditional pavement, placing the kinetic floor in a more crowded path leads to the generation of more energy. Therefore, proper placement of the floor allows for a faster compensation of the floor cost. Kinetic tiles could also have commercial applications in retail spaces. Energy generated by shoppers while doing their groceries or checking out the latest trends is valuable data [4] for shops and shopping centers to

know customers' preferences. Kinetic floors can extract and analyze this information while generating energy.

Currently, video surveillance in shopping centers is a common practice to monitor people and their movements. However, video surveillance is permitted by rules and regulations only for enhancing security, ensuring the safety of shoppers and staff, and preventing theft and criminal activities. Indeed, since cameras can have very high resolutions, they can recognize individuals, leading to privacy concerns [5]. People tracking through kinetic floors does not pose the same privacy issues because it cannot precisely identify individuals, thus serving as an alternative to camera recordings for commercial purposes. In this work, a smart tile for smart kinetic floors is developed and optimized for energy harvesting and data information acquisition. Section 2 presents the architecture of the smart kinetic floor composed of several tiles. The proposed smart tile system is discussed in Sect. 3, while Sect. 4 focuses on the self-sustainability of the proposed system. Conclusions are drawn in Sect. 5.

## 2 Floor System Architecture

The kinetic floor architecture is based on a combination of several smart tiles, which are mechanically and electrically connected to each other. The electrical connection consists of two main parts. The first part is related to energy accumulation and its transport to external users, such as lighting systems, information displays, and small hydroponic farms. The second part includes cabled bus systems for data communications and wireless connections for cloud data transmissions, creating an Internet of Things (IoT) device. Figure 1 depicts the two electrical systems for energy and data management.
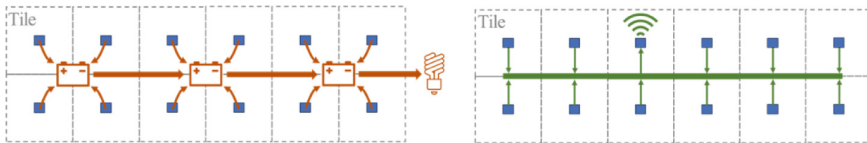


**Fig. 1.** Floor system architecture. The energy system and data management systems are depicted on the left and right sides, respectively.

The energy generated by each tile is channelled toward a battery for storage. Multiple tiles can be connected to the same battery depending on the battery size, energy generated per step, estimated daily step count, and user's power consumption. For instance, if each step generates approximately 0.5 J, with a step rate of 15 steps/min for 10 h, and assuming zero load consumption, a 12 V 0.5 Ah battery can store the energy generated by four tiles in a day. Foot traffic data is crucial for optimizing the floor design, as step rates may vary based on the installation site.

For data management, the acquired data from the tiles need to be sent to the cloud for storage and analysis. A wireless data communication is preferable to avoid the challenges of accessing physical data cables, ensuring a plug-and-play floor. One or more wireless router tiles can be used for wireless data transmission. The data network comprises a

wired bus network using the RS485 protocol, connecting each tile to the nearest master tile able to provide wireless communication to the cloud.

## 3 Smart Tile

The smart tile system proposed in this work consists of three main elements: the mechanical structure, the energy harvesting, and the data extraction and transfer.

The mechanical assembly of the tile is provided by a commercial partner and is not addressed in this work. Instead, the focus is on describing the proposed energy harvesting and data management solutions in the following subsections.

### 3.1 Energy Harvesting

The energy harvesting system is based on converting the vertical movement of the tile during a step into an angular movement that drives a dynamo. Springs are used to return the tile to its original level after the vertical force from the step is finished. The energy generated by this system is impulsive and depends on the force generated by the step. Higher force leads to faster dynamo speed and greater generated energy. Experimental measurements of the dynamo's output signal for steps and jumps performed by individuals of different weights show that the dynamo produces alternating signals with a short duration and an amplitude greater than 15 V.

While maximum power tracking systems are typically used for energy harvesting applications [6, 7], in this case, the impulsive nature of the energy and its low level would lead to complex and impractical solutions, as they would consume more power than the generated one. Therefore, a simple and low-cost solution was adopted consisting of a cascade of the dynamo, a diode bridge rectifier, and a lead-acid battery. To power the additional circuitry for data management, a DC/DC step-down converter is used to generate 3.3 V [8, 9], as shown in Fig. 2.
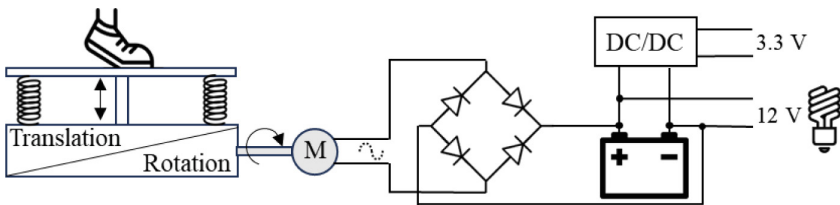


**Fig. 2.** Energy harvesting network based on a bridge rectifier and a 12 V lead-acid battery, which powers the IoT node and other users.

### 3.2 Data Information Extraction and Communications

In addition to energy harvesting, the tile aims to extract information about the energy generated per step and the step rate efficiently to minimize its impact on energy generation. Although the basic way to measure the energy generated is to measure the current

and voltage provided by the dynamo, in this work we opted to use only an accelerometer fixed to the tile. This approach reduces the tile cost and improves its energy efficiency by using just one sensor for both step characterization and energy estimation.

For IoT data measurement and processing, we utilized an Espressif ESP32 Azure IoT kit. This prototyping board was selected because it already includes an accelerometer (MPU6050) and an ESP32 module with WiFi capabilities. Additionally, it features I$^2$C I/O pins for connecting an I$^2$C/RS485 converter module.

The accelerometer data is acquired by the sensor and transmitted to the ESP32 processor via an I$^2$C bus for processing.

Figure 3 shows an example of acquired data for two steps performed by two individuals with different weights. Each step, independently from the individual, exhibits four distinct phases labeled A, B, C, and D. Phase A shows a negative acceleration peak caused by the person's weight on the tile, resulting in downward movement. In phase B, a positive peak is observed, resulting from two factors. Firstly, the internal springs slow down the tile's movement, bringing it to a stop in the low-level position. Secondly, the release of the person's weight as they continue the step accelerates the tile upward toward the top-level position. The negative peak in phase C corresponds to the tile decelerating towards its top-level position. Phase D is attributed to tile vibration and its oscillatory behavior. By analyzing the measured data and using a simple comparison with a threshold, the ESP32 core can identify each step and assign it a timestamp.
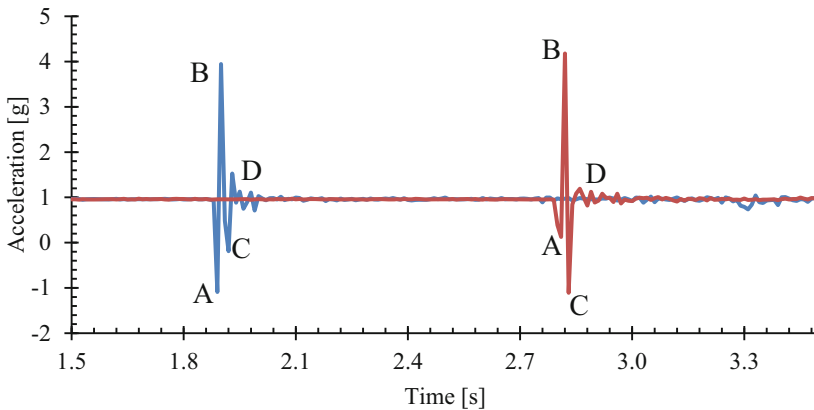


**Fig. 3.** Acceleration data measured by the accelerometer. The data are related to two steps by two individuals. The first step is in blue at 1.9 s and the second step is in red at about 2.8 s.

Regarding the estimation of the generated energy, it is proportional to the dynamo's rotation speed, which in turn is related to the tile's vertical speed. Theoretically, calculating the tile speed and the generated energy is possible by integrating the acquired acceleration data. However, this algorithm is complex for a low-power application as it requires storing all data samples, performing a moving-window integral, and dealing with saturation issues. In this work, a simpler algorithm is implemented on the ESP32 core for the generated energy estimation. It relies on measuring the acceleration peak during phase B, where the behavior is mainly influenced by the tile springs and by

the dynamo's torque. Consequently, the maximum peak in phase B remains unaffected by odd vibrations caused by the person's foot touching the tile. To verify the proposed algorithm, several measurements of the acceleration peak in phase B and the corresponding generated energy were performed. These measurements involved several steps and jumps on the tile performed by individuals with different weights. The acceleration peak is read using the MPU6050 and the ESP32 core, while the generated energy is measured by directly connecting a 5 mF capacitor to the diode bridge rectifier. The voltage across the capacitor is then read after each step, and the generated energy is calculated using the formula $E = \frac{1}{2}CV^2$.

A clear relationship between the acceleration peak of phase B and the generated energy is evident in Fig. 4. Therefore, with the use of only acceleration data, the system is capable of estimating the generated energy by the tile after a suitable prior calibration.

Concerning the communication, in this work the master tile with the wireless capability for data transfer on the cloud is developed. The communication stack for the WiFi 802.11.b/g/n protocol is developed and integrated into the ESP32 core to establish the connection with the cloud. On the cloud side, an IoT Hub in the Microsoft Azure service is configured to receive and display the data provided by the smart tile.

Following the edge computing approach, the complete acceleration data acquired through the MPU6050 are processed directly on the tile and not sent to the cloud. Indeed, to reduce the power consumption of the tile, only the time stamp related to each recorded step and the estimated energy generated by each step are transmitted to the IoT Hub. However, additional data could be transmitted if required.
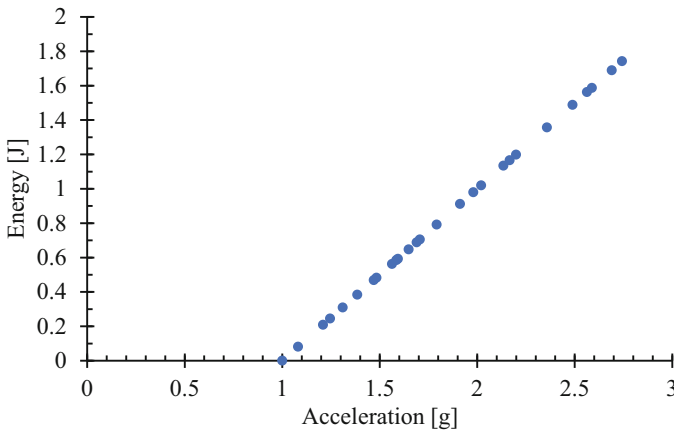


**Fig. 4.** Step generated energy as a function of the acceleration peak in phase B.

## 4   System Sustainability

The measured average energy generated by a step of a person weighing 70–80 kg is approximately 0.5 J. The measured power consumption of the basic tile, which acquires data from the accelerometer, is 0.19 W, and the master tile, with wireless capability,

consumes 0.5 W. In the worst-case scenario, with a twelve-tile floor and the basic master tiles active for ten hours for real-time application, it would require 186480 steps on the floor, averaging 15540 steps per tile, to maintain the system self-sustainable. However, these step numbers can potentially be significantly reduced by adopting duty cycle techniques in smart tile management if the real-time feature is not required. For instance, the entire basic tile system could be in sleep mode and only awakened when the accelerometer detects an acceleration peak. The ESP32 core in light sleep mode consumes 3.3 mW, and the MCU6050 accelerometer consumes 1.65 mW [10]. Although this solution is not currently implemented in this work due to the unconnected interrupt pin of the accelerometer on the selected development board, it could be considered in the future. Moreover, data transmission to the cloud could be limited to once a day, taking only about one minute, further reducing power consumption [11]. Considering all these potential conditions, the estimated number of steps required to achieve the self-sustainability of the floor is 362 steps per tile. The energy generated by additional steps can be used for the other user devices.

## 5  Conclusions

In this work, we explored energy harvesting by capturing kinetic energy from human steps. The proposed smart floor system, composed of several smart tiles, could be a valuable solution for energy generation and data acquisition in high foot-traffic areas like shopping centers.

The proposed smart tile comprises an energy generation and storage system, along with a data acquisition and transmission system. The energy generated by the tile movement is converted into electrical energy by a dynamo and stored in a battery solution.

The information related to the person's step and the related generated energy are acquired using only an accelerometer for reducing the power impact of the system.

Following an edge computing approach, acceleration data are processed directly on the tile, and only essential information is sent to the IoT Hub through the master tile WiFi connection.

The sustainability analysis indicates that for a real-time application, the current developed smart tile system requires around 15540 steps per tile for its sustainability. However, if the real-time characteristic is not required, with additional power reduction using sleep mode and duty cycling, the estimated number of steps needed for sustainability drops to 362.

Further research, such as the design of custom electronic boards, the use of more energy efficient communication protocols, and the introduction of more energy harvesting sources in the same floor, can lead to practical and commercially viable applications, contributing to a greener future.

# References

1. Visconti P et al (2022) Available technologies and commercial devices to harvest energy by human trampling in smart flooring systems: a review. Energies 15(2):432
2. He M et al (2019) Study of a piezoelectric energy harvesting floor structure with force amplification mechanism. Energies 12(18):3516
3. Kim K et al (2018) Optimized composite piezoelectric energy harvesting floor tile for smart home energy management. Energy Convers Manag 171:31–37
4. Liono J et al (2018) Predicting the city foot traffic with pedestrian sensor data. In: 14th EAI international conference on mobile and ubiquitous systems
5. Gomez IL (2022) Machine learning applications to video surveillance camera placement and medical imaging quality assessment. Illinois Institute of Technology
6. Ciarpi G et al (2023) Power optimization of systems for direct thermal to electrical energy conversion. Electronics 12(10):2163
7. Motahhir S et al (2020) The most used MPPT algorithms: review and the suitable low-cost embedded board for each algorithm. J Clean Prod 246:118983
8. Ciarpi G et al (2020) Inductorless DC/DC converter for aerospace applications with insulation features. IEEE TCAS-II 67(9):1659–1663
9. Saponara S et al (2018) Electrical, electromagnetic, and thermal measurements of 2-D and 3-D integrated DC/DC converters. IEEE TIM 67(5):1078–1090
10. Perdomo A et al (2023) ESP32 based low-power and low-cost wireless sensor network. In: The conference on Latin America control congress. Springer, Cham, pp 275–285
11. Liu D et al (2022) Concurrent low-power listening: a new design paradigm for duty-cycling communication. ACM Trans Sens Netw 19(1):1–24

# YoloP-Based Pre-processing for Driving Scenario Detection

Marianna Cossu[1]([✉]), Riccardo Berta[1], Luca Forneris[1], Matteo Fresta[1], Luca Lazzaroni[1], Jean-Louis Sauvaget[2], and Francesco Bellotti[1] [ID]

[1] Department of Electrical, Electronic and Telecommunication Engineering (DITEN), University of Genoa, Via Opera Pia 11a, 16145 Genova, Italy
{marianna.cossu,luca.forneris,matteo.fresta,
luca.lazzaroni}@edu.unige.it, {riccardo.berta,
francesco.bellotti}@unige.it
[2] STELLANTIS N.V. in Taurusavenue 1, 2132 LS Hoofddorp, The Netherlands
jeanlouis.sauvaget@stellantis.com

**Abstract.** Recognition of driving scenarios is getting ever more relevant in research, especially for assessing performance of advanced driving assistance systems (ADAS) and automated driving functions. However, the complexity of traffic situations makes this task challenging. In order to improve the detection rate achieved through state-of-the-art deep learning models, we have investigated the use of the YoloP fully convolutional neural network architecture as a pre-processing step to extract high-level features for a residual 3D convolutional neural network We observed thar this approach reduces computational complexity, resulting in optimized model performance, also in terms of generalization from training on a synthetic dataset to testing in a real-world one.

**Keywords:** Driving scenarios · Synthetic datasets · Automated driving · YoloP · Deep learning · Pre-processing · Video classification · Time-series · Convolutional neural network · Three-dimensional convolution

## 1 Introduction

Development of advanced driving assistance systems (AFAS) and automated driving functions (ADFs) needs a precise analysis of their behaviour in different operational design domains (ODD). To this end, original equipment manufacturers (OEMs) and suppliers are developing rule-based and machine learning-based systems to detect different types of driving/traffic scenarios (e.g., [1, 2]), in which the various systems and functions should be tested. Moreover, detecting driving scenarios stands as a fundamental objective in expanding the ODD of ADFs. This capability empowers the systems to anticipate and proactively adapt to dynamic external conditions. However, detection of scenarios through deep learning models is very challenging especially when, for cost reasons, the input is provided by a single camera, and research in the field has just started. In order to improve the performance of such a detector, we have focused on high-level image pre-processing.

In general, pre-processing plays a key role in training machine learning (ML) models. In the domain of computer vision, raw data pre-processing can encompass several steps, such as input scaling, dimensionality reduction [3], normalization, data cleaning [4], feature extraction and selection. In this work we investigate the use of the state-of-the-art YoloP fully convolutional neural network architecture as a high-level feature extractor to enhance performance of a state-of-the-art residual 3D convolutional neural network that we have trained to classify 1 s driving scenario single camera recordings. In this context, this article focuses on three main research questions (RQs): (1) What is the effect of the proposed high-level pre-processing on the detector's training time? (2) What is the effect on network performance, in terms of the typical classification metrics? (3) What is the effect on the model's generalization ability?

The remainder of the paper is organized as follows: Sect. 2 presents the dataset used during the experiments; Sect. 3 illustrates the pre-processing approaches implemented in the paper; Sect. 4 presents the experiment executed and the results obtained, while the final section provides the concluding remarks.

## 2 Datasets

We trained our model to classify 5 driving scenarios (Fig. 1), namely cut-in executed in front of the ego vehicle (1) and behind it (2); cut-out in front of the ego (3) and behind it (4); and ego lane change (5).
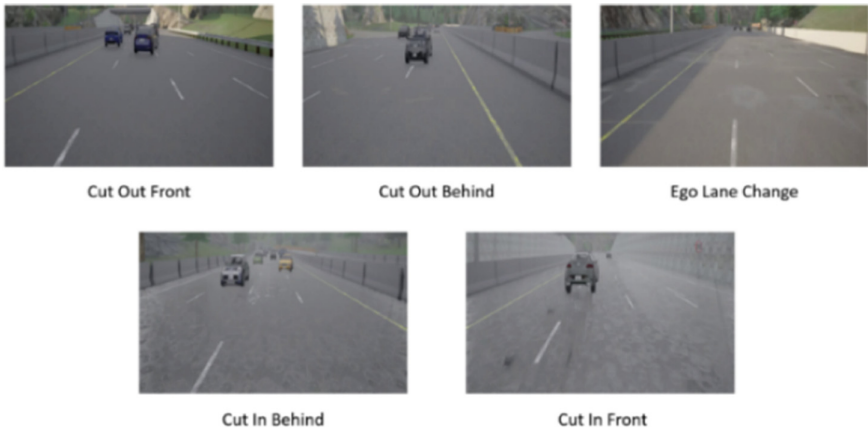


**Fig. 1.** Examples of the synthetic dataset's classes.

We used two different types of video-clip datasets: a synthetic one, exploiting a CarLA-based scenario generator [5, 6]; and a real-world one obtained after a manual labeling of the PREVENTION Dataset [7]. Manual labeling was necessary to extend Prevention to cover all our 5 types of labels. We trained and validated our system on the synthetic dataset only, while we tested it on both the datasets.

The implemented detector (described in the next section) receives as input sample a 1 s. time-window, consisting of three frames (3 frames $\times$ 3 channels $\times$ 224 px $\times$ 224

px). Each sample is extracted from a single scenario class instance (which typically last more than 1 s., e.g., 1–10 s), with a stride of 0.2 s for training and validation, 0.1 s for testing. Split of training, validation and testing set is done at level of single scenario instance. Samples for training are shuffled before their utilization.

## 2.1 Synthetic Dataset

The synthetic dataset contains over 800 simulated scenario instances for each class, with substantial variability. The corresponding number of samples is 73 K for the training set, 25 K for the validation set, and 6 K for the test set. The dataset includes 14 distinct weather and lighting conditions, various traffic densities (ranging from car-free to low, up to high traffic cases), and 21 car models, with different colors. A key focus of our work is on intra-scenario action variability. To this end, each scenario is parametrized (e.g., in terms of duration, speed of the ego vehicle, distance from the leading vehicle), and their distribution is sampled from data extracted from the real world ADScene dataset, which accounts for over 1M km manual and automated driving recordings [8]. As a limitation, the simulator currently includes highways with 4/3 lanes, either in a straight configuration or with road curvature degree of 90°, with two different radius values.

## 2.2 Real-World Dataset

To perform tests with a real-world dataset, we selected PREVENTION, which is one the very few publicly available recorded with cameras on highways. PREVENTION labels only front lane-change actions, thus we manually ported part of it to our target classes. During the labelling, we removed the samples that are not represented by the synthetic dataset ODDs (i.e., those with a number of lanes less than 3). However, we did not remove samples with curve radii very different from those of our synthetic dataset, which caused a reduction in performance. The real-world test set includes 10K time windows.

## 3 Pre-processing

To accomplish our objective, we preprocessed the previous presented dataset using the YOLOPv2 [9], as it reaches better performance compared to other networks sharing the same main goal. [9] allows extracting high-level information, including semantic segmentation of the drivable area and lane lines, and bounding boxes (BBs) of the vehicles.

We performed our experiments applying two different types of pre-processing, which corresponds to the utilization of three versions of input datasets (see Fig. 2):

1. Raw dataset: we directly use the frames extracted from the video-clips. The only pre-processing step involves pixel normalization.
2. Pre1 dataset: it is obtained with the first version of the high-level pre-processing method. This method extracts the most important information for each frame and stores it in the three different color channels. The red channel contains the vehicle images BBs; the green channel stores the drivable area and lane lines. The gray scale of the original image minus the drivable area is stored in the blue channel.

3. Pre2 dataset: it is obtained with the second pre-processing version. Differently from the first one, it removes all the information that is not extracted by the YoloP, leaving only (in the red channel) the grayscale inside each vehicle's BB.



**Fig. 2.** An example frame for each pre-process: **a** Raw, **b** Pre1 and **c** Pre2, respectively

## 4 Experiments

The experiments were conducted on an Ubuntu PC equipped with an NVIDIA RTX A4000 with 16 GB of VRAM. Results are reported in terms of accuracy, precision, recall, and F1 score [10]. The model used for the classification task is a Residual 3D Convolutional Neural Network (R3D), presented in [2]. The R3D was trained for 6 epochs. Figure 3 illustrates the training and validation accuracy's trend. From the plots it appears that both the pre-processing methods significantly speed up the training phase and enhance the network's performance.
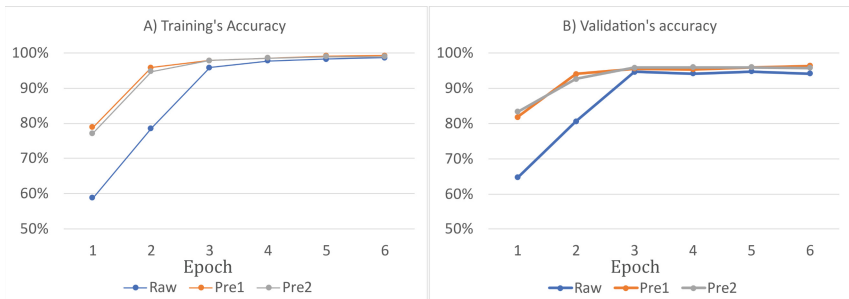


**Fig. 3.** Comparison in training and validation accuracy among the 3 investigated approaches.

The training times for the R3D model are as follows: 7.17 h for the Raw dataset, 3.81 h for Pre1, and 2.92 h for Pre2, which clearly demonstrates the impact of our pre-processing techniques. Specifically, Pre1 led to a remarkable 50% reduction in training time compared to the Raw dataset, while PreV2 achieved a 60% reduction. This hints at a significant decrease in data dimensionality and complexity achieved by the pre-processing.

The experimental results obtained on the synthetic test set show that all versions of the model achieved high performance values (Table 1). Particularly, both Pre1 and Pre2 outperform the baseline by 3% and 2%, respectively.

**Table 1.** Performance reached on the synthetic test set by the R3D trained with and without the proposed pre-processing.

| Dataset | Accuracy (%) | Precision (%) | Recall (%) | F1 Score (%) |
|---|---|---|---|---|
| Raw-synthetic | 92 | 92 | 92 | 92 |
| **Pre1-synthetic** | **95** | **95** | **95** | **95** |
| Pre2-synthetic | 94 | 95 | 94 | 94 |

In the second experiment, we utilized the real-world test set to assess the model's generalization ability. Since no real-world samples were included during the training phase and the real-world test set differed in elements such as camera resolution, position and road and lane type, compared to the synthetic dataset, the performance obtained at this stage was notably lower than on the synthetic dataset. However, this experiment highlights the robustness provided by the pre-processing techniques, that led to a substantial improvement in accuracy, up to 10% (Table 2).

**Table 2.** Performance reached on the real-world test set by the R3D model trained with and without the proposed pre-processing.

| Dataset | Accuracy (%) | Precision (%) | Recall (%) | F1 Score (%) |
|---|---|---|---|---|
| Raw-real-world | 63 | 68 | 63 | 60 |
| Pre1-real-world | 66 | 70 | 63 | 63 |
| **Pre2-real-world** | **73** | **74** | **73** | **73** |

## 5 Conclusion and Future Works

In order to improve the performance of a state-of-the-art camera-based driving scenario detector, we have explored the use of the YoloP system as a high-level feature extractor.

Concerning our first RQ, our experimental results demonstrated that the proposed pre-processing allows to cut more than 50% of the computational time needed to perform a full training. Performance measured with the typical classification metrics is increased by 3%, on the original synthetic dataset (RQ2). Finally, we addressed the RQ3 exploiting the PREVENTION dataset. Results show that the proposed pre-processing allowed a significant 10% increase in accuracy. This indicates a clear improvement in the model's ability to generalize also to real-world scenarios. Since the preliminary results of our experiments are promising, a possible next goal of this research is to expand the ODD of the synthetic training dataset as much as possible (e.g., by introducing more road environment variability such as increasing and decreasing the number of lanes, etc.) and then assess the capability to synthetically generate training datasets that allow training models able to generalize well in the real-world.

# References

1. Izquierdo R et al (2021) Vehicle lane change prediction on highways using efficient environment representation and deep learning. IEEE Access 9:119454–119465. https://doi.org/10.1109/ACCESS.2021.3106692

2. Cossu M, Villon JLQ, Bellotti F, Capello A, De Gloria A, Lazzaroni L, Berta R (2022) Classifying simulated driving scenarios from automated cars. Lect Notes Electr Eng (LNEE) 866:229–235. https://doi.org/10.1007/978-3-030-95498-7_32

3. Obaid HS, Dheyab SA, Sabry SS (2019) The impact of data pre-processing techniques and dimensionality reduction on the accuracy of machine learning. In: 2019 9th annual information technology, electromechanical engineering and microelectronics conference (IEMECON), pp 279–283. https://doi.org/10.1109/IEMECONX.2019.8877011

4. Maharana K, Mondal S, Nemade B (2022) A review: data pre-processing and data augmentation techniques. Glob Transit Proc 3:91–99. https://doi.org/10.1016/j.gltp.2022.04.020

5. Motta J, Bellotti F, Berta R, Capello A, Cossu M, De Gloria A, Lazzaroni L, Bonora S (2022) Developing a synthetic dataset for driving scenarios. Lect Notes Electr Eng (LNEE) 866:310–316. https://doi.org/10.1007/978-3-030-95498-7_43

6. Cossu M, Berta R, Capello A, De Gloria A, Lazzaroni L, Bellotti F (2023) Developing a toolchain for synthetic driving scenario datasets. Lect Notes Electr Eng (LNEE) 1036:222–228. https://doi.org/10.1007/978-3-031-30333-3_29

7. Izquierdo R, Quintanar A, Parra I, Fernández-Llorca D, Sotelo MA (2019) The prevention dataset: a novel benchmark for prediction of vehicles intentions. In: 2019 IEEE intelligent transportation systems conference (ITSC), pp 3114–3121. https://doi.org/10.1109/ITSC.2019.8917433

8. ADScene (2021) towards an industrial scenarios plateform for driving assistance systems design & validation—DSC 2021 EUROPE VR. https://dsc2021.org/adscene-towards-an-industrial-scenarios-plateform-for-driving-assistance-systems-design-validation/. Accessed 30 May 2023

9. Han C, Zhao Q, Zhang S, Chen Y, Zhang Z, Yuan J (2022) YOLOPv2: better, faster, stronger for panoptic driving perception. http://arxiv.org/abs/2208.11434

# Open Source Remote Diagnostics Platform for Custom Instrumentation in Nuclear Applications

Iván René Morales[1][✉], Maria Liz Crespo[1], and Sergio Carrato[2]

[1] Multidisciplinary Laboratory (MLab), STI Unit, The Abdus Salam International Centre for Theoretical Physics (ICTP), Strada Costiera, 11, 34151 Trieste, Italy
{imorales,mcrespo}@ictp.it

[2] Dipartimento di Ingegneria e Architettura (DIA), Universitá degli Studi di Trieste, Piazzale Europa, 1, 34127 Trieste, Italy
ivanrene.moralesargueta@phd.units.it, carrato@units.it

**Abstract.** An open source remote diagnostics platform for custom electronics in experimental acquisition setups has been developed, targeting nuclear and high-energy physics (HEP) applications. We aim at enabling remote access to instrumentation hardware prototypes located in radiation-controlled areas by using existent network infrastructures. The platform relies on two components: a graphical user interface (GUI) developed in GNURadio and a remote hardware bridge (HB). The GUI was designed using the GNURadio's optimized building blocks, ensuring a fluid visualization of real-time pulse traces and energy spectrum. No third-party libraries are used, turning our solution into an easy drop-in tool for instrumentation diagnostics in HEP and radiation-related experiments. Reliable communication between the remote instrument and the GUI is accomplished using TCP/IP, carried out by the HB in case of remote instruments without network interface. A Raspberry PI Zero-W is tested as the HB to demonstrate the few computational resources required for this end.

**Keywords:** Nuclear instrumentation · GNURadio · GUI · Remote diagnostics · Hardware prototypes

## 1 Introduction

Testing custom hardware prototypes for acquisition systems in experiments related to nuclear or particle physics instrumentation usually requires scheduling time slots for measurements during a beam run [1]. Thus, having a graphical user interface (GUI) to assess the behavior of the detectors in real time is crucial to quickly verify the system response under controlled stimuli. Since no physical access to radiation-controlled areas is allowed during the beam operation, personnel are forced to be located in a control room away from the detectors.

Commercial solutions exist for both remote oscilloscope and multichannel analyzer (MCA) applications [2,3], which are the most commonly used tools for this testing stage [4]; however, the compatibility is limited to the vendor-specific products and no APIs are available to create wrappers for custom hardware solutions. Moreover, many of them require a dedicated computer to be placed in the radiation-controlled area to act as a bridge with the instrument prototype under test to enable the connectivity via remote desktop.

Our proposal addresses both features in a single solution: a GUI for real-time diagnostics and a hardware bridge (HB) to enable remote access for instrument prototypes without network connectivity. The open source GUI was developed using the GNURadio software development kit, taking advantage of its building blocks, which are already optimized for fluid visualization and real-time multi-core processing, regardless of the operating system [5]. Hence, the intensive computation tasks are carried out in the computer where the GUI is being executed, making possible to use low-end and low-power devices as the remote HB. Moreover, no third-party libraries are required, enabling a flawless deployment to easily start testing data acquisition system (DAQ) prototypes. To ensure the reliability of the transmitted data, the communication between the DAQ and the GUI is carried out using the transport control protocol (TCP).

In terms of alternative communication protocols, some commercial solutions rely on dedicated high-level solutions built atop TCP, such as LAN eXtensions for Instrumentation (LXI) [6], which consumes significant computational resources at the DAQ side and increase the development time for diagnostics purposes. Other protocol alternatives exist to address real-time remote data transmission: Padded Jittering Operative Network (PJON) designed to exploit noisy or long transmission distances on wireless interfaces [7], or Lab Streaming Layer (LSL), designed for multiple channel brain signals [8]; however, they require third-party or custom libraries to provide GNURadio compatibility and may hinder the deployment process with missing dependency or version problems [9]. Our platform overcomes these issues by using TCP transactions managed by GNURadio built-in blocks, leveraging existing cabled or wireless network infrastructures. Moreover, the GUI is compatible with multiple operating systems and the HB can be deployed in many different models of inexpensive off-the-shelf low-power single-board computers (SBC). These features provide flexibility to use the available hardware at the laboratories without compromising the fluid oscilloscope and MCA visualization.

## 2   Methodology

We chose a custom DAQ under development as the instrument to be diagnosed. In our tests, the DAQ was connected to the proposed HB via universal serial bus (USB) emulating a universal asynchronous receiver-transmitter (UART) interface. Three different gamma detectors were tested to assess the flexibility of the platform under diverse signal types: a sodium-iodide (NaI:Tl) crystal with a silicon photomultiplier (SiPM), a NaI:Tl crystal with a photomultiplier tube (PMT), and a CLYC (Cs2LiYCl6:Ce) crystal with a SiPM.

Albeit the HB can be deployed using several Linux-compatible SBC models, we opted for a low-end and low-power Raspberry PI Zero W (version 1.1) to demonstrate the few computational requirements needed for its operation and the flexibility to power it from any available USB port or phone charger. This SBC features only 512 MB of DDR2 RAM and a single-core 1 GHz ARMv6 CPU: humble specifications compared to other models available nowadays. We measured the power consumption during one hour of heavy CPU load and continuous WiFi transmission, using a 5 W USB adapter as power supply.

The GUI was tested in a desktop computer (located in the control room) with 16 GB of DDR4 RAM and a Core i7-8700 processor running GNURadio version 3.10.5.1 on Linux Ubuntu 20.04.6. The DAQ and detectors were located about ten meters away, in the radiation-controlled area. Figure 1 shows the proposed remote diagnostics platform block diagram, including all the elements of the system. It is worth noting that the HB acts as an auxiliary element in case the DAQ under development is not provided with network connectivity; otherwise, the computer running the GUI may be directly connected to the DAQ.
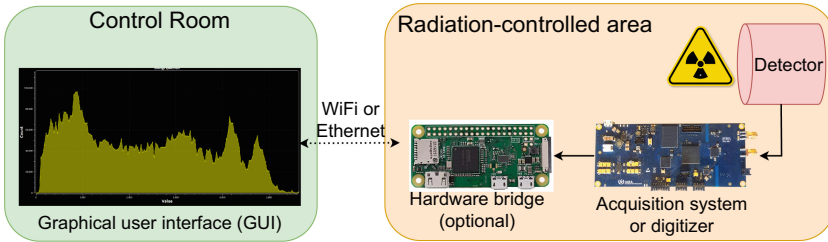


**Fig. 1.** Block diagram of the system components distribution, including the hardware bridge between the remote digitizer and the network infrastructure.

## 3   Results

The oscilloscope and MCA were tested with a remote DAQ connected to multiple detectors (one at a time) stimulated with a Cobalt-60 (Co-60) gamma source. The source was located close to the scintillators and served as a reference for pulse shape visualization and energy spectrum computation. The HB was able to transmit the data using its wireless interface between the DAQ (at the radiation-controlled area) and the computer executing the GUI (in the control room). A sustained throughput close to 2.1 Mbps was measured, while the system drew an average of 934 mW from the power supply. Moreover, rates close to 27 Mbps were reached with a cabled Ethernet connection. The open source GUI source files and the script for the HB are available at the online repository [10].

### 3.1    Graphical User Interface

Since the GUI was developed in GNURadio, it is also compatible with multiple operating systems (Linux, MacOS, Windows). Moreover, a GNURadio port for Android (currently under development) may also enable its execution on a smartphone or tablet [11]. Such flexibility provides access to the DAQ or detector under diagnostics from mostly any available device in the control room of the experiment. The GUI features some of the most common tools required for debugging and diagnostics of nuclear and particle physics detectors and digitizers, i.e. oscilloscope (real-time trace visualization) and MCA (energy spectrum) [12]. Stability of both GUI elements during long-term measurements was tested in 24-hours continuous sessions with no issues reported or detected.

The following settings are required to be defined prior to the first execution of the application: remote IP address and port number (either from the HB or the remote DAQ), sampling rate of the DAQ, and trace length (in sample count). After setting these values up the GUI can be executed directly from GNURadio, leading to a window with two tabs (oscilloscope and MCA), described next.

**Oscilloscope** Based on GNURadio's oscilloscope block, provides the user the basic tools required for the real-time diagnostics of scintillation detectors and the DAQ prototypes attached to them, such as baseline, pulse shape and pulse amplitude distribution. Figure 2 shows the oscilloscope tab featuring the pulse shape of a NaI:Tl scintillator coupled to a SiPM and digitized by the remote DAQ under test. During run-time, this tool provides the following features and settings: cross-level trigger threshold with pulse polarity selection, multiple trigger sources (signal, peak detector), baseline restorer (BLR), peak detector and pulse height analysis (PHA), and recording of individual traces in binary format. An update rate close to 15 traces per second was measured with short trace lengths (128 samples) in cabled Ethernet mode.

**Multichannel Analyzer** Based on GNURadio's histogram block, displays the estimated deposited energy spectrum in the detector. The individual energy of each event is computed as the peak amplitude (using the PHA) of the pulse minus its baseline (from the BLR value). Flexibility to the user is provided by tuning the following settings during execution: number of histogram bins, autoscale and window zooming, limited number of events or accumulative spectrum. Figure 3 shows the MCA tab with experimental data built upon a Co-60 $\gamma$ source located close to the NaI:Tl detector with SiPM, shown in ADC channels. The two characteristic photopeaks at 1.17 MeV (above channel 4000) and 1.33 MeV (below channel 5000) as well as the Compton back-scatter peak (210 keV) and the Compton edge (936 keV) evidence the expected behavior of the system.

### 3.2    Hardware Bridge with SBC

The tested HB demonstrated consistent remote data transmission using the embedded WiFi 802.11n peripheral. Configuring the SBC simply required
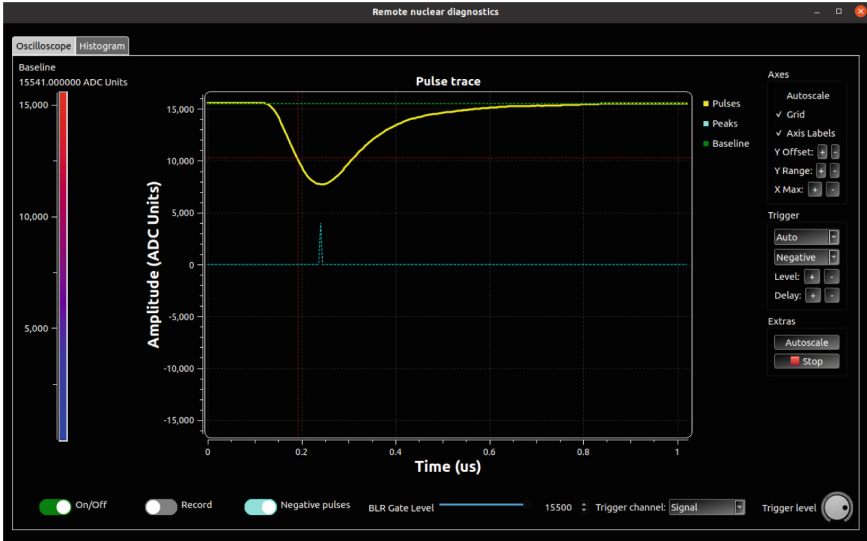
**Fig. 2.** Oscilloscope tab for real-time visualization of the triggered traces showing experimental data from a Co-60 gamma source close to a NaI:Tl SiPM detector.
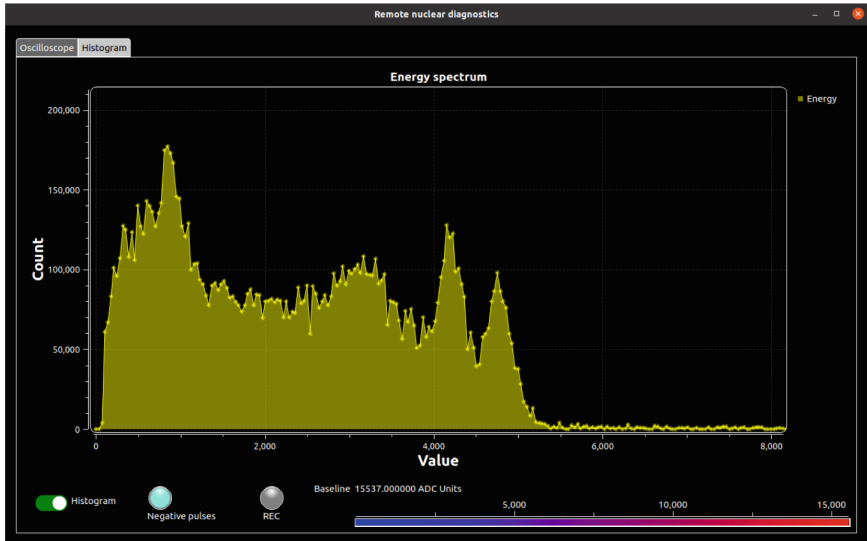


**Fig. 3.** Multichannel analyzer (MCA) tab. The histogram represents the accumulated energy spectrum of measured pulses from the NaI:Tl crystal with SiPM when placed close to the Co-60 $\gamma$ source. The spectrum matches the expected detector response.

installing the default operating system (Raspberry PI OS, release 20230503) into a 8 GB micro-SD card and setting up the WiFi access-point configuration. Taking the hardware bridge script up and running only required to download the source code from the provided public repository and execute the Python script.

## 4  Conclusions

The idea of this platform aroused from the necessity of deploying a fast diagnostics system for custom instrumentation hardware under experimental conditions. An open source platform for remote DAQ diagnostics has been consequently developed and tested with diverse detectors. The system enables an easy setup to test instrumentation prototypes with an oscilloscope and multichannel analyzer from remote locations taking advantage of existent network infrastructures. The open source GUI is optimized for multi-core processing to provide a fluid visualization of the data transmitted from the DAQ under test and no third-party libraries are required to execute it, apart from installing GNURadio. A remote hardware bridge has also been developed to enable connectivity for remote instruments that lack of network capabilities. The proposed platform provides a drop-in solution for fast debugging and diagnostics of custom electronics in nuclear and HEP experiments.

## References

1. Diener R et al (2019) The DESY II test beam facility. Nucl Instrum Methods Phys Res Sect A: Accel Spectrometers Detect Assoc Equip 922:265–286. https://doi.org/10.1016/j.nima.2018.11.133
2. CAEN: CAENScope (2023). https://www.caen.it/products/caenscope/. Accessed 15 June 2023
3. AMETEK: MAESTRO multichannel analyzer emulation software (2023). https://www.ortec-online.com/products/application-software/maestro-mca. Accessed 15 June 2023
4. Crespo ML et al (2021) Remote laboratory for e-learning of systems on chip and their applications to nuclear and scientific instrumentation. Electronics 10:2191. https://doi.org/10.3390/electronics10182191
5. Valerio D (2008) Open source software-defined radio: a survey on GNUradio and its applications
6. Mumford R (2007) The LXI standard: past, present and future. Microwave J 50
7. Mitolo GB (2022) PJON (Padded Jittering operative network). https://github.com/gioblu/PJON. Accessed 22 June 2023
8. Stenner T et al (2022). Lab streaming layer project. https://doi.org/10.5281/zenodo.6387090, github.com/sccn/liblsl/tree/v1.16.0
9. USRP (2020) LSL toolkit for GNURadio/USRP data. https://github.com/lwa-project/usrp. Accessed 22 June 2023
10. Morales I (2023) Remote diagnostics platform for custom instrumentation in nuclear applications—code repository. https://doi.org/10.5281/zenodo.8082581. Accessed 10 July 2023

11. Bloessl B, Baumgärtner L, Hollick M (2020) Hardware-accelerated real-time stream data processing on android with GNU radio. ACM 103–109. https://doi.org/10.1145/3411276.3412184

12. Yang J, Xiong M, Zhao X (2020) Improved low power multichannel pulse amplitude analyzer. OALib 07:1–7. https://doi.org/10.4236/oalib.1106629

# A Tool for Design and Simulation of Battery Operated Trains

Luca Pugi[✉], Luca di Carlo, Aljon Kociu, Lorenzo Berzi, and Massimo Delogu

Department of Industrial Engineering, University of Florence, Via di Santa Marta 3, 50139
Florence, Italy
luca.pugi@unifi.it

**Abstract.** Multi-modal battery-operated trains are gaining an increasing favor
and diffusion, as a sustainable alternative to Diesel-Electric and Diesel-Hydraulic
Powertrain that are conventionally used on not electrified or partially electrified
lines. In this work authors examine how current improvements on battery technologies
are going to increase performances autonomy and diffusion of this technology
also proposing a model and a power management architecture that can be used to
simulate and accelerate the development of this technology.

**Keywords:** Battery operated trains · Power management of multi-modal trains ·
Power loop strategy

## 1 Introduction

BEMU (Battery-Electric-Multiple Units) are currently proposed as valid alternative to
DMU (Diesel Multiple Units) to increase sustainability of railway service on not electrified
lines [1, 2]. In previous research activities authors have developed [3] a complete
model of FCHMU (Fuel Cell Hybrid Hydrogen Multiple Unit) that has been used to
simulate a possible update to Hydrogen technology of an existing HMU (Hybrid Multiple
Unit) the Hitachi BLUES platform [4]. The same builder of the BLUES platform has
planned the construction of BEMU train whose powertrain is described in Fig. 1: traction
inverters are directly coupled to overhead line through pantograph, assuring propulsion
on electrified lines. Batteries are connected to traction DC bus through a bi-directional
DC-DC converter assuring propulsion on not electrified lines. This further conversion
stage introduces further losses during battery operated service but at the same time assure
the possibility of a decoupled lower voltage levels for accumulators.

Most of existing BEMU are currently equipped with high power lithium cells, mainly
based on LTO technology (Lithium Titanate Oxide) which assure a long and reliable operational
life sacrificing specific energy performances that are limited to about 80 Wh/kg.
Thanks to the progress of railway battery technology due to the increasing technology
transfer from the automotive sector, a new generation of accumulators for railway application
is offering the possibility of increasing specific energy to about 110–120 Wh/kg,
maintaining a high-power capability (3C) which can assure the possibility of a complete
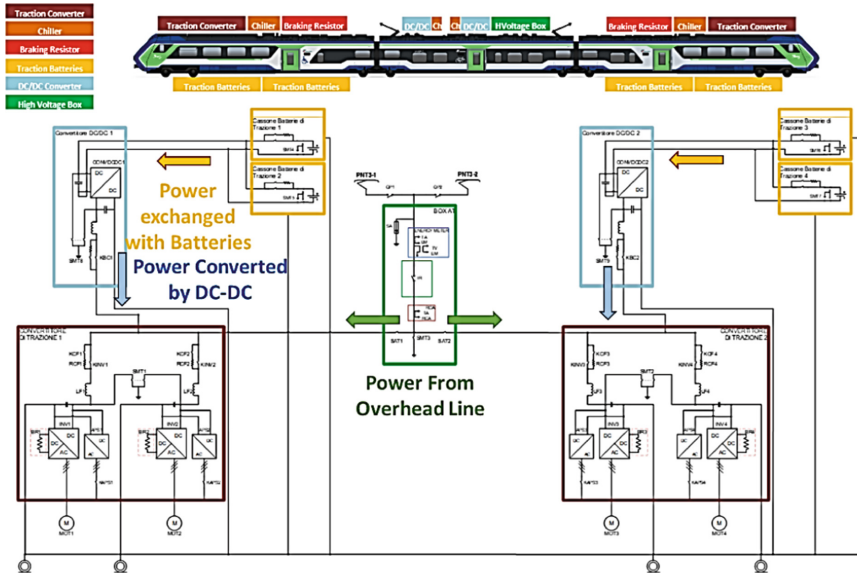recharge of the storage in about twenty minutes as shown in Fig. 2 (Table 1).

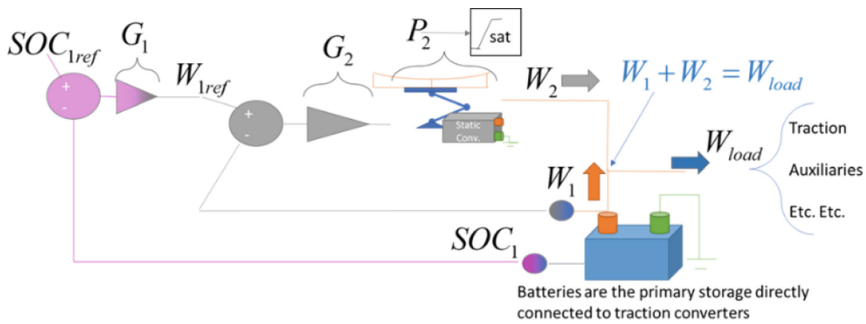**Fig. 1.** Powertrain layout proposed for BEMU platform of Hitachi as presented in 2021 [4]



**Fig. 2.** Proposed power management layout

**Table 1.** Main Feature of benchmark test case [4] and of proposed battery [5]

| Mean. Aux power | 120 [kW] | Max speed | 140 [kmh] |
|---|---|---|---|
| Traction power min max | 889–1333 [kW] | Capacity | 300 seats |
| Reg. braking | 1933 [kW] | Max long. effort | 160 kNN |
| Battery Spec. energy | 352 [W/kg] | Battery spec. power | 117 [Wh/kg] |

## 2   Proposed Solution

Respect to the original plant layout authors proposed a power management strategy that was originally developed for the power management of hybrid storage system for fast charge of buses [6, 7]. For the specific application related to BEMU authors proposed to use a high voltage battery storage directly connected to traction inverter to maximize efficiency of power flows between them. The aim of this approach is to further increase the autonomy of the system minimizing these losses. Considering installed power of the proposed BEMU and the current increase of voltage levels of storages used for automotive application, this solution can be feasible in a short-term scenario.
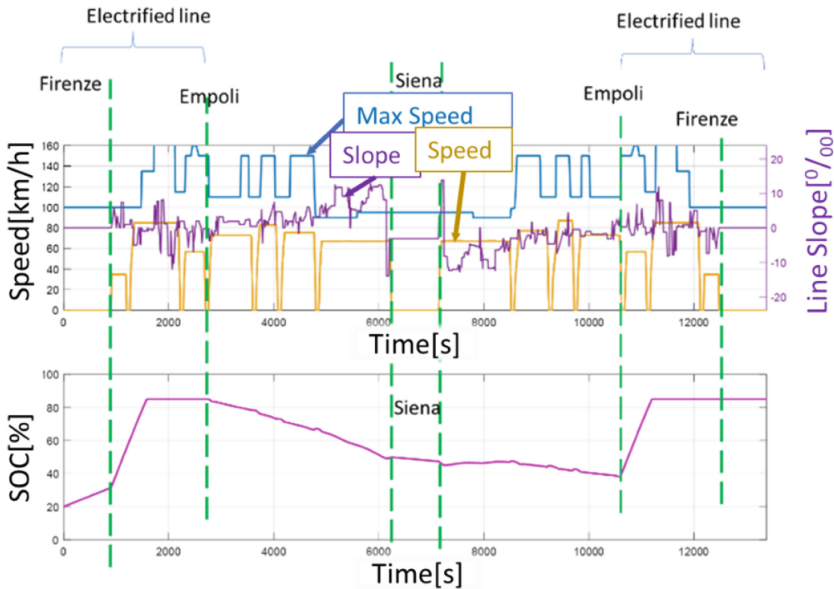


**Fig. 3.** Simulation results on Florence-Siena line (speed, slope and SOC vs. time)

It's the connection of electrified line through pantograph that is controlled by a static converter. This solution probably introduces some additional losses when the vehicle is operating in pure electric mode, but it also allows the possibility of exploiting on the same train different electrification standard, further increasing the flexibility of the system. As shown in Fig. 3 the external power connection to overhead line is treated as an additional power source that can be regulated to assure that the power exerted by the battery $W_1$ is equal to a desired power profile $W_{1ref}$ by suppling the power $W_2$.

This functionality can be easily implemented introducing a closed loop control of $W_1$, where $G_2$ and $P_2$ are respectively regulator and plant transfer functions. Saturation of system response respect to amplitude constraints is also introduced.

If the nonlinearity represented by the saturation is not excited the response of the system is described by transfer functions (1)

$$W_1 = \frac{W_{load} + W_{1ref}G_2P_2}{1 + G_2P_2}; \quad W_2 = \frac{(W_{load} - W_{1ref})G_2P_2}{1 + G_2P_2}; \tag{1}$$

Closed loop has a finite bandwidth, also the response is saturated, so there is a region of frequencies and amplitudes of traction load $W_{load}$ for which the closed loop is not effective. In these conditions energy stored in the primary storage should be stabilized within an acceptable range with an additional closed loop aiming to stabilize the primary storage by regulating the value of $W_{1ref}$. This external loop through the controller $G_1$ is able to regulate the battery $SOC_1$ to follow a desired reference value $SOC_{1ref}$. This further loop is not devoted to directly track the load but to extend the duration of the primary storage by reducing the oscillation of $SOC_1$. With an adequate loop shaping of both $G_1$ and $G_2$ is possible to obtain different behaviors of the power flow management blending a "load following" behavior and "range extender" one.

As shown in Table 2, it is considered a storage of about 1 MWh whose weight and encumbrances are allowable respect to the simulated train composition. Also, power flows are saturated to take count of limitations arising from both battery and pantograph features. Proposed solution is simulated considering the service on a partially electrified line the Firenze-Siena one which is electrified from Firenze to Empoli (32.5 km) and not electrified from Empoli to Siena (63 km). In Fig. 4 some simulation results are shown: the train starts its mission in Florence with a depleted battery (SOC = 20%) on the electrified line. The 20% of minimum SOC is considered since for a worn EOL (End Of Life) battery a capacity degradation of 20% should be expected. As the train waits in the station, SOC of the battery slowly increases since in parking mode max collected power is limited to 600 [kW] and auxiliaries are consuming about 120 [kW]. As train starts to move collected power can increase and battery is recharged fast to desired setpoint (SOC = 85%). A modest value of desired SOC is chosen mainly for two reasons:
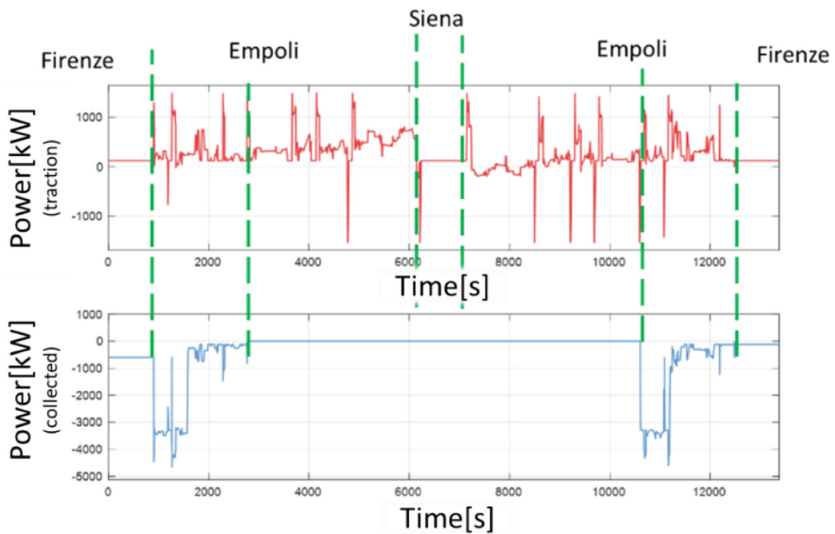
- A part of the battery capacity is kept unused to allow prolonged regenerative braking.
- By reducing the depth of discharge of the battery (chosen range is from 85 to 20%) equivalent FCE is reduced prolonging life and reliability of the storage.

After Empoli the line is not electrified so the battery-operated mission starts. Since the altitude of Siena is higher respect to Empoli, consumptions of outward and return courses are not equal. In the station of Siena, it's not supposed the presence of any recharge system, so a slow reduction of SOC is also recorded during the stop due to consumption of auxiliaries. During the battery-operated mission some small oscillations of SOC can be detected around stations where sequences of regenerative braking and accelerations produce peaks of regenerated and consumed power. As the train returns on the electrified line in Empoli the train is recharged reaching Florence with the maximum desired SOC of 85%. The total discharge depth for a complete mission is equal to about the 45% of the SOC. So, mean expected autonomy of the train in battery mode is at least 230–250 km. Expected life of the battery in terms of equivalent FCE is typically equal to 4000–5000 for slow recharges and its reduced to about 2000 FCE for fast charging profiles that are typically adopted on automotive applications.

**Table 2.** Size limitations, saturation of power flows, additional load of auxiliaries

| Parameter | Value |
|---|---|
| Size of the storage | 1060 [kWh] (9000 kg, 6685.86 dm$^3$) |
| Max power collected on a 3 kV catenary in motion | 6 [MW] (2000 [A]) |
| Max power collected on a 3 kV catenary in standstill conditions | 600 [kW] (200 [A]) |
| Max continuous charge discharge power on batteries | 3 C (about 3 [MW]) |
| Allowable voltage catenary range (For 3 kV) | 2400–3700 [V] |
| Installed traction power | 1333 [kW] |
| Max regenerative braking | 1333 [kW] |
| Power required by installed auxiliaries | 120 [kW] |

Even supposing an equivalent life of the battery limited to 2000 FCE, the expected life of the battery due to cycle aging corresponds to about 4500 complete missions (567000 km in full battery mode, 860000 km of total mileage). Since the duration of a complete mission is about 3–3.5 h, no more than three mission/day can be reasonably performed, the service life of the battery is expected to be at least about four-five years. This is only a cautious estimation of the life of the component that can be probably much longer since the mission profile is not excessively stressing for the chosen battery. Availability of experimental data from testing devices should help to further refine this preliminary rough estimation.



**Fig. 4.** Simulation results on Florence-Siena line (collected and consumed power profiles)

# 3  Conclusions and Future Developments

In this work, authors have proposed a simulation model of a BEMU to verify how technological improvements of batteries and an innovative layout of power management can improve the performances of existing battery trains. Preliminary results show that is possible to considerably improve, nearly doubling the foreseen autonomy. Authors are currently focusing their attention on a more accurate simulation (higher detail on implemented power electronics) of the system considering different operational scenarios and mission profiles. Also, authors are investigating the possible application to railway sector of innovative hybrid supercapacitors that are currently proposed for micro-mobility applications [8].

# References

1. Schenker M, Schirmer T, Dittus H (2021) Application and improvement of a direct method optimization approach for battery electric railway vehicle operation. Proc Inst Mech Eng Part F: J Rail Rapid Transit 235(7):854–865
2. Juston M, Vulturescu B, Chamaret A (2023) A statistical approach for the optimal sizing of partial electrification for battery trains. In: 2023 IEEE international conference on electrical systems for aircraft, railway, ship propulsion and road vehicles & international transportation electrification conference (ESARS-ITEC). IEEE, pp 1–6
3. Pugi L, Berzi L, Cirillo F, Vecchi A, Pagliazzi V (2023) A tool for rapid simulation and sizing of hybrid traction systems with fuel cells. Proc Inst Mech Eng Part F: J Rail Rapid Transit 237(1):104–113. https://doi.org/10.1177/09544097221092622
4. Vannucchi A, La piattaforma MASACCIO di Hitachi Rail per la decarbonizzazione dei treni regionali, LA TRANSIZIONE TECNOLOGICA DALLA TRAZIONE DIESEL AI NUOVI TRENI A BATTERIA E IDROGENO Mercoledì 29 settembre 2021 Convegno Webinar in occasione di Expo Ferroviaria
5. Technical data from hoppecke (visited last time on 14/05/2023). https://www.hoppecke.com
6. Pugi L, Alessandrini A, Barbieri R, Berzi L, Pierini M, Cignini F, Genovese A, Ortenzi F (2021) Design and testing of a supercapacitor storage system for the flash recharge of electric buses. Int J Electr Hybrid Veh 13(1):57–80. https://doi.org/10.1504/IJEHV.2021.115197
7. Cignini F, Genovese A, Ortenzi F, Alessandrini A, Berzi L, Pugi L, Barbieri R (2020) Experimental data comparison of an electric minibus equipped with different energy storage systems. Batteries 6(2):26. https://doi.org/10.3390/batteries6020026
8. Corti F, Gulino M-S, Laschi M, Lozito GM, Pugi L, Reatti A, Vangi D (2021) Time-domain circuit modelling for hybrid supercapacitors. Energies 14(20):6837. https://doi.org/10.3390/en14206837

# Lightweight Neural Networks for Affordance Segmentation: Enhancement of the Decoder Module

Simone Lugani, Edoardo Ragusa[✉], Rodolfo Zunino, and Paolo Gastaldo

Department of Electrical, Electronic, Telecommunications Engineering and Naval Architecture (DITEN), University of Genoa, Genoa, Italy
edoardo.ragusa@unige.it

**Abstract.** The deployment of deep neural networks for visual affordance segmentation on wearable robots poses may prove critical, due to some conflicting aspects of the problem. On one hand, affordance segmentation requires high-level abstraction capabilities, that typically involve large-size models. On the other hand, computing resources hosted on wearable robots prevent to run large-size models in real-time. The paper presents an analysis of the role of the segmentation head in the trade-off between generalization performance and compute cost. The obtained models outperform modern baseline solutions in well-known, real-world datasets while meeting low computing requirements.

**Keywords:** Affordance Segmentation · Embedded Computer Vision · Wearable Robots

## 1 Introduction

Wearable robots such as prostheses and exoskeletons are advanced systems with multiple degrees of freedom. The control of these devices is an open problem because the operation of all functions requires to generate generating explicit commands [1].

Semiautonomous control automatizes parts of the action by reducing the users efforts [2,3]. Toward that end, the system selects the action autonomously; this implicitly requires high-level capabilities in both semantic reasoning and fine-grained sensing. Teleceptive sensing [4], i.e. sensing without contact, enables sensing information from the environment, for example by using cameras. Visual affordance segmentation (VAS) divides an object into functional parts, thus paving the way to planning fine-grained semiautonomous actions. State-of-the-art approaches use deep neural networks (DNNs) [5–7] to tackle this problem, but the typical hardware resources embedded in wearable robots prevent real-time processing. As a result, the design of those solutions combines very complex computer vision problems with tight computing requirements. In [8] the authors proposed a simplified version of the original VAS problem supported by portable

hardware. The original solution was subsequently improved by proposing additional DNNs modules [9] and sensory information [10].

Balancing generalization performance and hardware requirements is an open problem, further complicated by many levels of abstraction. When designing DNNs, most approaches take into account the design of the backbone part of the architecture, for example by using Network Architecture Search (NAS) [11]. This choice seems reasonable when considering that the backbone plays a crucial role when compared with the segmentation head (SH), in terms of both computations and generalization performance. However, the SH plays a significant role that has been under-explored in the specific scenario considered in the research presented here.

This paper considers the impact of the Segmentation Head components on visual affordance segmentation, when using hardware-friendly DNNs. The analysis first shows that even minor adjustments in the SH architecture can affect the overall performance of the model. Secondly, a hierarchical problem formulation based on multitask learning leads to improved generalization performance. The resulting set of hardware-efficient SHs outperform existing solutions on established benchmarks in terms of accuracy, while ensuring a comparable computational cost.

## 2   Method

Small-size networks may typically exhibit reduced generalization ability as compared with large-size models [8]. As a consequence, the role of the SH becomes critical, as the backbone might fail to disentangle effectively the implicit features involved in the problem.

Segmentation Heads accomplish two basic tasks: (1) the reconstruction or rearrangement of geometrical information, based on the features extracted by the backbone levels; (2) the eventual classification of affordable parts. The literature proposes two main solutions for the reconstruction of geometrical information. Standard upsampling operations, such as bilinear upsampling, can augment image resolution, and are not involved in the solution of the learning problem. Conversely, transposed convolution layers perform a convolution operation with a stride value smaller than one, thus yielding larger output tensors. In this setup, the kernel plays a crucial role in the learning process, at the expense of a larger number of flops with respect to standard upsampling. A few convolutional layers typically support the eventual classification step. Standard convolution exhibits satisfactory expressive capabilities, whereas depth-wise separable schemes can better balance computing requirements and generalization capabilities.

Object localization is critical in VAS [12], as it is preliminary to understanding the object's affordable parts. One can formalize a hierarchical version of the learning problem, where the network includes two output heads. The first SH supports a binary Object segmentation (OS) classification task. The computational graph includes an additional SH, which tackles the actual VAS problem. This component processes an input representation that is optimized to isolate

the object, and focuses on the distinctive parts of the object, thus biasing the feature extraction process. The loss function is formalized as:

$$\mathcal{L} = \mathcal{L}_{seg} + \mathcal{L}_{aff} \tag{1}$$

where the two terms refer to OS and VAS, respectively.

## 2.1   Segmentation Head

The large networks, which yet fit the memory constraints of embedded devices, typically feature inference times of many seconds [13], which proves inadequate for AS. FLOPs have a major impact on inference time and can be computed independently from the target platform [11]. Therefore, the FLOPs count was selected as the constraint in the SH design process. When paired with the MobileNetV3 backbone for one image of size $128 \times 128$, the admissible SHs were limited to the ones with a FLOP count in the range of 700/800 M. Figure 1 shows the template architecture for a multitask scenario. In the case of a single task, one just removes the OS branch. The subset of connections with the backbone is selected based on the solution selected for the *base blocks* to control the number of FLOPs. The base blocks always contain a convolutional and an upsampling layer.

The research presented in this paper compared three instances of the base block, namely, depthwise separable convolution with nearest upsampling (U), depthwise separable convolution with transposed convolution layer (T), and standard convolution and nearest upsampling (B).
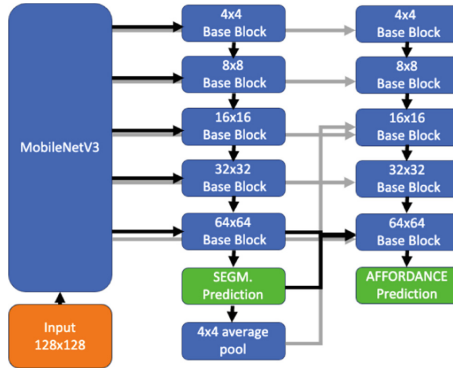


**Fig. 1.** General scheme of the proposed SHs

## 3    Experiments

The experiments were divided into three sets. The first experimental section compared the segmentation heads; the second set of tests measured the accuracy values scored for the auxiliary task. Finally, the role of the interconnection between SH and backbone was considered.

Two datasets provided the experimental benchmarks: the UMD dataset held 28.843 RGB-D images of 7 categories of objects. 5135 images formed the test set. Likewise, the IIT dataset (IIT) [5], including 8835 images at different framing and resolution settings, was used for the assessment. In compliance with the setup described in [8], all the grasping affordances were grouped into a unique class and the other affordances were grouped into a single category (i.e. "don't grasp").

All models were trained on both training sets. The IIT dataset was augmented by applying standard geometrical and color distortions. The UMD was augmented by a custom procedure, replacing the original blue background in the images with a different picture. The objects were isolated by using the segmentation mask of the labels and a color based-procedure relying on the uniform background. Finally, standard augmentation was applied.

The baseline SH relied on standard convolution and nearest upsampling, and adopted the design strategy presented in [8]. In the following, the architectures with multiple SHs will be denoted by a letter for each head. For example, decoder 'UB' will denote an architecture including Segmentation Head 'U' and Visual Affordance Segmentation 'B'. All the SHs were connected to the backbone using four connections [with reference to Fig. 1, the $64 \times 64$ connection was not implemented].

Class-wise weighted pixel-wise accuracy was computed over the test set. Three test sets were considered: the IIT test set, the original version of UMD, and the version with the replaced background UMD_B. A summary performance measure will be indicated with TOT computed as the weighted average of the three accuracies, weighting IIT 0.5 and the two versions of UMD 0.25.

Table 1 gives the weighted accuracy of the networks on the three versions of the test sets, also including the number of parameters. All the tested SH outperformed the baseline approach. The baseline implementation relied on a different subset of connections and prioritized high-level features, thus leading to an increased number of parameters. The networks with 'B' heads benefited from the use of standard convolution, leading to the reported good results. The two architectures based on the single head with depthwise convolution both obtained excellent scores; this proved that, even without multitasking, depthwise segmentation could support the segmentation process effectively.

Table 2 presents the results of multitask architecture for the object segmentation task. In this case, the result on the augmented UMD_B dataset has been added. The Table uses the same convention of Table 1. The results confirmed that all the SHs managed to discriminate the object pixels from the background pixels with an accuracy higher than 90% on average.

**Table 1.** Affordance segmentation performance

|  | TOT | IIT | UMD | PARAMS. |
|---|---|---|---|---|
| Decoder B | 91.8 | 88.6 | 94.6 | 2.2 M |
| Decoder U | 92.6 | **91.0** | 94.6 | **1.3 M** |
| Decoder T | 92.6 | 90.5 | 95.1 | 2.0 M |
| Decoder BB | 92.5 | 89.9 | **95.6** | 3.6 M |
| Decoder UU | 92.1 | 89.6 | 94.9 | 1.6 M |
| Decoder TT | 92.5 | 90.3 | 95.2 | 2.1 M |
| Decoder TB | **92.8** | 90.8 | 95.2 | 3.6 M |
| Decoder UB | 91.7 | 90.3 | 95.2 | 3.3 M |
| Baseline [8] | 91.2 | 88.6 | 94.6 | 5.9 M |

**Table 2.** Object segmentation performance

|  | TOT | IIT | UMD | UMD_B |
|---|---|---|---|---|
| Decoder BB | 92.7 | 90.2 | 95.4 | 95.0 |
| Decoder UU | 93.1 | 90.2 | 96.0 | 95.9 |
| Decoder TT | 93.7 | 91.3 | **96.2** | **96.1** |
| Decoder TB | **93.9** | **91.8** | 96.2 | 96.0 |

Table 3 illustrates the results of the experiments about the connections between the SH and the backbone, comparing a new version of the SHs that used the $64 \times 64$ layer as a connection. The $8 \times 8$ connection was removed in this setup while maintaining four levels of feature maps. Table 3 follows the conventions adopted in Table 1 to report the performance of the new versions of the SHs; the differences with respect to the previous versions are reported in brackets. The results confirmed that exploiting connection with low-level features did offer some improvements, although the previous feature set actually yielded the best overall average result. This result confirmed the importance of connections, and highlighted that they should represent a crucial factor in automated design strategies.

**Table 3.** VAS performance with low-level feature maps

|  | TOT | IIT | UMD | PARAMS |
|---|---|---|---|---|
| Decoder UT-X | **92.7** (+1.5) | 90.5 (+2.4) | **95.4** (+0.7) | **1.1M** (−0.6M) |
| Decoder TU-X | 92.6 (+0.5) | **90.7** (+1.5) | 94.9 (−0.6) | 1.4M (−0.6M) |
| Decoder TT-X | 92.6 (+0.1) | 90.5 (+0.2) | 95.1 (−0.1) | 1.5M (−0.6M) |
| Decoder TB-X | 92.6 (−0.2) | 90.5 (−0.3) | 95.3 (+0.1) | 1.8M (−1.8M) |
| Decoder BB-X | 92.5 (+0.0) | 90.3 (+0.4) | 95.1 (−0.5) | 1.8M (−1.8M) |

## 4   Conclusion

The paper presented an analysis of the role of the segmentation head for lightweight neural networks in applications of visual affordance segmentation. The paper analyzed the problems that may arise when using constrained networks, and presented some design criteria to limit the issues introduced by the reduced generalization capabilities of the models. Empirical results confirmed that a careful design of the segmentation head, with a slight reformulation of the learning problem, could improve over the baseline solutions proposed in recent works.

## References

1. Salminger S, Stino H, Pichler LH, Gstoettner C, Sturma A, Mayer JA, Szivak M, Aszmann OC (2022) Current rates of prosthetic usage in upper-limb amputees-have innovations had an impact on device acceptance? Disabil Rehabil 44(14):3708–3713
2. Tang Z, Zhang L, Chen X, Ying J, Wang X, Wang H (2022) Wearable supernumerary robotic limb system using a hybrid control approach based on motor imagery and object detection. IEEE Trans Neural Syst Rehabil Eng 30:1298–1309
3. Sun Y, Fei T, Li X, Warnecke A, Warsitz E, Pohl N (2020) Real-time radar-based gesture detection and recognition built in an edge-computing platform. IEEE Sens J 20(18):10706–10716
4. Krausz NE, Hargrove LJ (2019) A survey of teleceptive sensing for wearable assistive robotic devices. Sensors 19(23):5238
5. Nguyen A, Kanoulas D, Caldwell DG, Tsagarakis NG (2017) Object-based affordances detection with convolutional neural networks and dense conditional random fields. In: 2017 IEEE/RSJ international conference on intelligent robots and systems (IROS). IEEE, pp 5908–5915
6. Jiang Z, Zhu Y, Svetlik M, Fang K, Zhu Y (2021) Synergies between affordance and geometry: 6-dof grasp detection via implicit representations. arXiv:2104.01542
7. Khalifa ZO, Shah SAA (2022) Towards visual affordance learning: a benchmark for affordance segmentation and recognition. arXiv:2203.14092
8. Ragusa E, Gianoglio C, Dosen S, Gastaldo P (2021) Hardware-aware affordance detection for application in portable embedded systems. IEEE Access 9:123178–123193
9. Apicella T, Cavallaro A, Berta R, Gastaldo P, Bellotti F, Ragusa E (2021) An affordance detection pipeline for resource-constrained devices. In: 2021 28th IEEE international conference on electronics, circuits, and systems (ICECS), IEEE, pp 1–6
10. Ragusa E, Ghezzi MP, Zunino R, Gastaldo P (2023) Affordance segmentation using RGB-d sensors for application in portable embedded systems. In: Applications in electronics pervading industry, environment and society: APPLEPIES 2022. Springer, pp 109–116
11. Benmeziane H, Maghraoui KE, Ouarnoughi H, Niar S, Wistuba M, Wang N (2021) A comprehensive survey on hardware-aware neural architecture search. arXiv:2101.09336

12. Nguyen A, Kanoulas D, Caldwell DG, Tsagarakis NG (2016) Detecting object affordances with convolutional neural networks. In: 2016 IEEE/RSJ international conference on intelligent robots and systems (IROS). IEEE, pp 2765–2770

13. Canepa A, Ragusa E, Zunino R, Gastaldo P (2022) Detection-based video surveillance using deep neural networks on stm32 microcontroller. In: 2022 29th IEEE international conference on electronics, circuits and systems (ICECS). IEEE, pp 1–4

# Low-Cost, Edge-Cloud, End-to-End System Architecture for Human Activity Data Collection

Matteo Fresta<sup>(✉)</sup> , Ali Dabbous , Francesco Bellotti , Alessio Capello , Luca Lazzaroni , Alessandro Pighetti , and Riccardo Berta

Department of Electrical, Electronic and Telecommunication Engineering (DITEN), University of Genoa, Via Opera Pia 11a, 16145 Genova, Italy

{matteo.fresta,ali.dabbous,alessio.capello,luca.lazzaroni,
alessandro.pighetti}@edu.unige.it, {francesco.bellotti,
riccardo.berta}@unige.it

**Abstract.** Research in the Internet of Things (IoT) have paved the way to a new generation of applications and services that collect huge quantities of data from the field and do a significant part of the processing on the edge. This requires availability of efficient and effective methodologies and tools for a workflow spanning from the edge to the cloud. This paper presents a generic, complete workflow and relevant system architecture for field data collection and analysis with a focus on the human physical activities. The data source is given by a low-cost embedded system that can be placed on the user body to collect heterogeneous data on the performed movements. The system features a 9 DoF IMU sensor, to ensure a high level of configurability, connected to a custom board equipped with a rechargeable battery for wireless data collection. Data are transmitted via Bluetooth Low Energy (BLE) to a smartphone/tablet app, which manages the data transfer to Measurify, a cloud-based open-source framework designed for building measurement-oriented applications. Results from a preliminary functional experiment confirm the ability of the proposed end-to-end system architecture to efficiently implement the whole targeted edge-cloud workflow.

**Keywords:** Field data collection · Human activity recognition · Internet of things · Sensor-based classification · Wearable sensor · Edge-cloud hardware/software architectures · Embedded systems

## 1 Introduction

In the Internet of Things (IoT) scenario, collecting large amounts of data from sensors and storing them within a cloud database has become a major challenge [1]. The need for data is constantly increasing with the growing development and spread of ever new machine learning (ML) technologies that rely on supervised learning techniques, thus requiring datasets.

For this reason, obtaining large amounts of data quickly and easily, while maintaining high accuracy during acquisition, is crucial for training accurate ML models. For this

purpose, various runtime applications have been developed to handle the data generated by edge devices (e.g. [2]). The challenge is also to build a "smart" database that adapts as closely as possible to the resources being sent [3] and that allows data to be managed using a simple, lightweight and fast sending method, while still maintaining a high level of security during data exchange. One of the possible applications of this paradigm, based on an embedded system and requiring a flexible database due to the diversity of data to be collected, is the classification of specific actions during a human physical activity. Over the past decades, motion classification has been a constantly growing research area. Sensors are applied to the human body to accurately represent movements [4] and collect as much data as possible to best support the training of state-of-the-art ML algorithms in a smart controlled environment [5].

In this paper, a data collection workflow is implemented from the physical sensor to a database to manage the measurements. This embedded system consists of an Arduino Nano 33 BLE Sense mounted on a custom board with a battery installed for autonomous energy support. The IMU sensor on the device provides accelerometer, gyroscope and magnetometer data that will be collected and sent to the database. We were looking for a flexible and data-oriented framework that best fit our use-cases, so the choice fell on Measurify, formerly Atmosphere [6], as it is an open-source, cloud-based, measurement-oriented API Framework, which is connected to MongoDB [7] as database.

A Flutter [8] application, that can be installed on any tablet or mobile phone device, is used to store in memory a great number of data received from Arduino without stopping the data collecting phase, this maximizes the amount of data obtainable per minute and simplify the connection between the embedded system and Measurify.

The design choices aim to ensure that this workflow remains easily accessible by adopting an open-source approach, with all components available for download on Github [9]. The instrumentation is intentionally low-cost, using mainly Arduino as the only physical device, while all other components are free-to-use and can be hosted locally. This configuration allows for widespread adoption and easy replicability, making it feasible for a broad audience to participate in and benefit from the workflow.

## 2   Workflow

The proposed workflow can be decomposed into three main sectors as shown in Fig. 1: Edge, Fog, and Cloud. These represent the spectrum of distributed computing, from immediate data processing at the source (Edge), intermediate processing in local networks (Fog), to centralized processing and storage in remote servers (Cloud). The Edge consists of a versatile wearable embedded system designed to be attached to any part of the human body. Its primary function is to collect data while the wearer performs specific actions or activities. On the other hand, the Fog sector comprises a Flutter application and a personal device, which together facilitate the visualization of the data obtained from the wearable embedded system. Lastly, the Cloud sector consist of the Measurify Framework, which acts as a receiver for data sent from the Flutter application and saves them into a MongoDB database, enabling secure storage for future use and analysis.

The embedded system used consists of an Arduino NANO 33 BLE Sense soldered to a board as shown in Fig. 2, which allows a rechargeable battery to be installed in order to power the Arduino and data to be collected wirelessly.
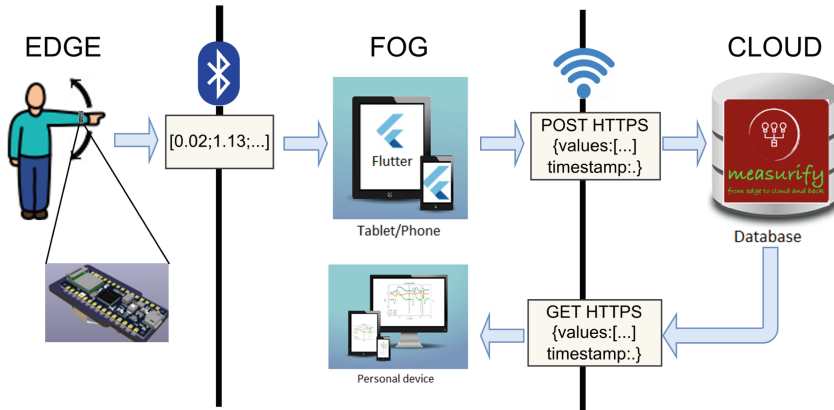
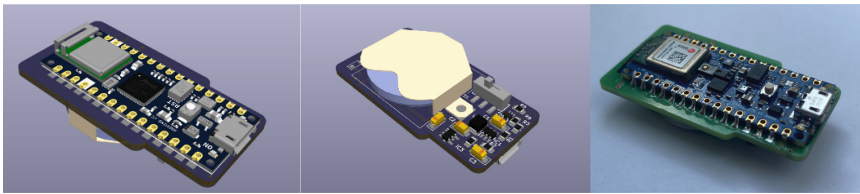**Fig. 1.** Workflow enabled by the proposed system architecture.



**Fig. 2.** Arduino board with battery: on the left, up-side prototype design; in the center, bottom-side prototype design; on the right the board used.

The board features the 9-axis IMU LSM9DS1 sensor: a 3D accelerometer with a default range $[-4, +4]$ g $-/+0.122$ mg, 3D gyroscope with a default range $[-2000, +2000]$ dps $\pm70$ mdps and 3D magnetometer with a default range $[-4, +4]$ gauss $\pm0.14$ mgauss. With an LIR 2450 battery with a capacity of 120 mAh it is estimated that data can be collected for 6 consecutive hours.

Sensor data are collected by the Arduino via a script with a minimum sampling period of 5 ms and transmitted via Bluetooth Low Energy (BLE). The script sets up the Arduino to expose Bluetooth services with features that a device can subscribe to and be notified whenever new data are available.

To receive these data, a custom application was developed for the subscription of the characteristics exposed by the Arduino and to assign a label to the measurement, and then send it via an HTTPS POST RESTful API to the Measurify framework in the format of timeseries, a type of data that contains the numerical values and the timestamp in which it was measured.

This application can be installed on a tablet or mobile phone and was developed using the open-source Flutter framework [8] and it connects to the embedded device via BLE and stores the received numerical data in memory. Once the user stops the data recording, the stored dataset is sent to the Measurify's timeseries route through a HTTPS POST. While the process is running, values are sent in blocks of 1000 samples

as they are collected. Request's body encapsulate values organized in JSON format. An example of a timeseries sample is as follows:

```
{
"timestamp": "1684833177652",
"values": [−0.50488, −0.40173, −0.75158, 32.959, −66.284, 123.474, 0, 0, 0]
}
```

This protocol minimizes the amount of space used during the calls and speeds up the data transmission. For the data visualization, it is possible to get the values via a HTTPS GET request from a personal device. The timeseries route, secured through authentication, allows user to retrieve previously inserted values in common formats: JSON, Pandas Dataframe, and CSV [10]. User can also filter measurements to obtain only samples in a specific period, or to retrieve only values exceeding a certain threshold.

## 3 Results

We performed a functional test for our system with a very simple preliminary data collection experiment. For the test, the sampling period of the device was set to 250 ms and the default sensitivity values of the IMU sensors of the Arduino Nano 33 BLE Sense were used (Table 1). The dataset collected consists of the samples taken during the action of repeatedly raising and lowering an arm progressively increasing the movement speed. To perform the test, we attached the Arduino near the wrist of the hand and started recording values. After 38 samplings, we stopped collecting data and the dataset was correctly uploaded to the database. As the speed of the movement increased, we expected a reduction in the number of timesteps required to complete the movement and also an increase of the acceleration vector. Using a Python script, we plotted the values of the accelerometer, gyroscope and magnetometer (Fig. 3).
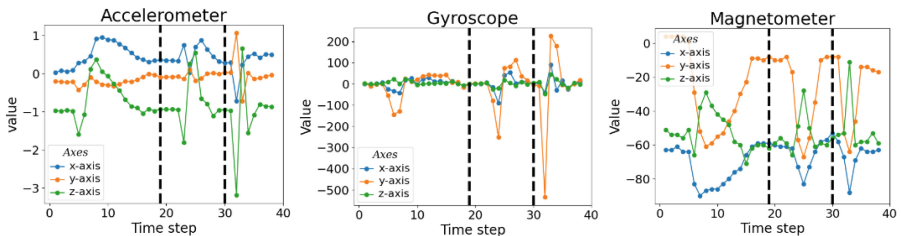


**Fig. 3.** Accelerometer, gyroscope and magnetometer plot of the movement replicated three times with incremental speed.

As expected, the plots show a movement reproduced three times at incremental speed, with a progressive decrease in the number of steps required to complete each action and an increase in the intensity of the acceleration vector.

Given the target of making the supported workflow flexible and generic for different types of tracking, we defined a set of configurable parameters (e.g., to increase the

amount of data collected per time unit or to increase the sensitivity of the sensors), as summarized in Table 1.

**Table 1.** Configurable parameters.

| Variable | Default | Range | Sensibility |
|---|---|---|---|
| Sample period | 250 ms | >5 ms | – |
| Accelerometer | ±4 g, ±0.122 mg | ±2, ±8, ±16 g | 0.061, 0.244, 0.732 mg |
| Gyroscope | ±2000 dps, ±70 mdps | ±245, ±500 dps | ±8.75, ±17.50 mdps |
| Magnetometer | ±4 gauss, 0.14 mgauss | ±8, ±12, ±16 gauss | 0.29, 0.43, 0.58 mgauss |

## 4 Conclusion and Future Works

Data collection and management have become an essential part in the development of ML models based on supervised learning. Therefore, it is crucial to obtain large amounts of data easily, quickly and accurately. One of the main challenges lies in building a workflow that starts with the edge device until it interfaces with a "smart" database that adapts to the type of data being sent, ensuring also a secure connection.

The proposed workflow offers a fully accessible, low-cost and user-friendly infrastructure for the collection and management of data from embedded systems. This approach can be applied to various applications, including classification of human physical activities, providing high quality data for training advanced machine learning models.

As future work, it is planned to integrate the study of ML techniques from the collected data, in order to generate models that can be imported into the same Arduino used previously, obtaining a complete pipeline, from edge to cloud and back. This will increase the autonomy and efficiency of the system, enabling a rapid classification of the detected actions.

## References

1. Ali I, Ahmedy I, Gani A, Munir MU, Anisi MH (2022) Data collection in studies on internet of things (IoT), wireless sensor networks (WSNs), and sensor cloud (SC): similarities and differences. IEEE Access 10:33909–33931. https://doi.org/10.1109/ACCESS.2022.3161929
2. Berta R, Mazzara A, Bellotti F, De Gloria A, Lazzaroni L (2021) Edgine, a runtime system for IoT edge applications. In: Saponara S, De Gloria A (eds) Applications in electronics pervading industry, environment and society. Springer International Publishing, Cham, pp 261–266. https://doi.org/10.1007/978-3-030-66729-0_31
3. Berta R, Bellotti F, De Gloria A, Lazzaroni L (2022) Assessing versatility of a generic end-to-end platform for IoT ecosystem applications. Sensors 22. https://doi.org/10.3390/s22030713
4. Kos A, Umek A (2019) Wearable sensor devices for prevention and rehabilitation in healthcare: swimming exercise with real-time therapist feedback. IEEE Internet Things J 6:1331–1341. https://doi.org/10.1109/JIOT.2018.2850664

5. Alemayoh TT, Lee JH, Okamoto S (2020) A new motion data structuring for human activity recognition using convolutional neural network. In: 2020 8th IEEE RAS/EMBS international conference for biomedical robotics and biomechatronics (BioRob), pp 187–192. https://doi.org/10.1109/BioRob49111.2020.9224310

6. Berta R, Kobeissi A, Bellotti F, De Gloria A (2021) Atmosphere, an open source measurement-oriented data framework for IoT. IEEE Trans Ind Inf 17:1927–1936. https://doi.org/10.1109/TII.2020.2994414

7. MongoDB: the developer data platform, https://www.mongodb.com. Accessed 17 July 2023

8. Flutter—Build apps for any screen. //flutter.dev/. Accessed 17 July 2023

9. Measurify. https://github.com/measurify. Accessed 18 July 2023

10. Fresta M, Bellotti F, Capello A, Cossu M, Lazzaroni L, De Gloria A, Berta R (2023) Efficient uploading of .Csv datasets into a non-relational database management system. In: Berta R, De Gloria A (eds) Applications in electronics pervading industry, environment and society. Springer Nature Switzerland, Cham, pp 9–15. https://doi.org/10.1007/978-3-031-30333-3_2

# Analysis of the Divider Control Policy for a Fractional Low-Power Time Synchronization Algorithm

Giuseppe Coviello[✉], Antonello Florio, Giuseppe Brunetti,
Caterina Ciminelli, and Gianfranco Avitabile

Department of Electrical and Information Engineering, Polytechnic University of
Bari, Bari BA, 70125, Italy
Giuseppe.Coviello@poliba.it, Antonello.Florio@poliba.it,
Giuseppe.Brunetti@poliba.it, Caterina.Ciminelli@poliba.it,
Gianfranco.Avitabile@poliba.it

**Abstract.** In the Internet of Medical Things (IoMT) applications, time synchronization assumes a relevant role as in every distributed system. Various protocols and algorithms have been developed for this purpose. Recently, Fractional Low-Power Synchronization Algorithm (FLSA) was proposed as a potential solution for low-power multi-board time synchronization within a Wireless Body Area Network (WBAN). FLSA is based on a fine timer resolution control thanks to a dual-modulus divider. In this paper, we aim to study the impact of the modulus control policy on time synchronization accuracy.

**Keywords:** Time synchronization · Dual-modulus prescaler · Internet of Medical Things (IoMT)

## 1 Introduction

Time synchronization constitutes an integral component of all distributed systems. In the case of the IoMT scenario, multiple boards equipped with sensors and actuators are positioned on the body to monitor health parameters or to dispense medications. Two distinct levels of synchronization are possible. Strong Time Synchronization (STS) calls for a higher degree of precision, as even minor discrepancies in time-stamping may result in significant errors during the data merging phase. An example application requiring STS is gait analysis based on surface electromyography (sEMG) and Inertial Measurement Units (IMUs). This increased precision requires an increase in the number of exchanged synchronization messages. In contrast, Weak Time Synchronization (WTS) allows for a lower degree of accuracy, thereby reducing the complexity of the synchronization process. WTS is appropriate for monitoring slowly varying phenomena such as fluctuations in body and environmental temperature.

Over the years, various protocols and algorithms have been developed to synchronize network time. Some protocols, such as the Reference Broadcast Algorithm (RBS) [1], Timing synchronization Protocol for Sensor Networks (TPSN) [2], and Flooding Time Synchronization Protocol (FTSP) [3], are widely used as references. Recently, *Coviello et al.* introduced the FLSA as a potential solution for low-power multi-board time synchronization within a WBAN configured in a star topology and employing a master/slave (M/S) architecture [4,5]. The algorithm demonstrates exceptional performance in terms of both accuracy and power consumption. The algorithm has already been employed in the low-cost synchronized platform for offline acquisition of heartbeat rate (HR), peripheral oxygen saturation (SpO2), gait analysis, posture analysis, and muscle activity related to gait analysis, as presented in the [6]. The proposed approach utilizes a two-modulus divider to enhance the granularity of timer synchronization. This study focuses on examining the influence of the divider control policy on the System Timer (ST) synchronization accuracy.

The paper is organized as follows. Section 2 reviews the fundamental operations performed by FLSA. Section 3 discusses the issues arising from the dual-modulus approach and considers potential solutions. Section 4 presents the results of a simulation campaign that compares different correction approaches and their achieved accuracy. The conclusions close this work and discuss possible future improvements.

## 2  Recalls on the Algorithm

FLSA discretizes time into synchronization slots of length $T_M$. Hence, every $T_M$ seconds, the master transmits a synchronization packet containing its ST to the slaves, which is processed through routines ensuring the highest processing priority for the synchronization packet, and computing the accumulated time shift compared to the master time so that it can be proactively updated. FLSA also deploys routines in charge of deciding whether or not to keep the radio section of the device to be on for the synchronization packet reception, thus providing a notable energy-saving feature and making it eligible for implementation in battery-powered devices aiming for low-power consumption (Fig. 1).
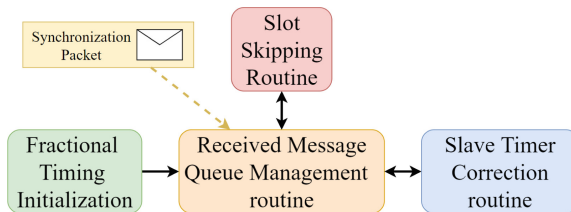


**Fig. 1.** FLSA routine interaction diagram

FLSA employs a pair of $N/N+1$ type comparators/counters to execute dual modulus division. The comparators are denoted as $\{A, B\}$. The number of ticks $N_M$ needed to form a synchronization slot is computed as:

$$N_M = \left\lfloor \frac{T_M}{t_{RTC}} \right\rceil \tag{1}$$

where $\lfloor \cdot \rceil$ is the rounding at the closest integer number operator, and $t_{RTC}$ is the granularity of the Real-Time Clock (RTC). The ST granularity $t_{ST}$ can be retrieved as:

$$\begin{cases} R_C^A = \left\lfloor \dfrac{t_{ST}}{t_{RTC}} \right\rfloor \\ R_C^B = R_C^A + 1 \end{cases} \tag{2}$$

where $R_C^i$ is the number of ticks needed to form a $t_{ST}$ for comparator $i$. The variation of the comparator values determines the synchronization slot, as shown by:

$$N_M = R_C^A \cdot I_A + R_C^B \cdot I_B \tag{3}$$

where $I_i$ represents the number of times comparator $i$ must be utilized.

## 3    The Impact of the Modulus Control Policy

The implementation of a dual-modulus divider is based on the principles of fractional-PLL theory [5]. According to the literature, this approach enables precise control over the generated output frequency by dynamically selecting one of two division values in the frequency divider. However, this advantage is accompanied by the presence of frequency spurs [7]. In the context of FLSA, a similar effect happens in the time domain. Even if the steady state ST on average is the same, by a close look at the TS values it is clear how the transient value is dependent on the modulus selection strategy, which directly impacts the final synchronization error.

Several techniques can be employed for the modulus selection. The most trivial one is referred to as "default", and is performed by using $R_C^A$ and $R_C^B$ for $I_A$ and $I_B$ times, respectively, in a cyclical manner. This method ensures consistency and predictability in TS generation and serves as a fundamental benchmark for evaluating the performance and effectiveness of alternative implementations. However, the timer curve exhibits a double slope during the synchronization period, as depicted in Fig. 2. Consequently, although perfect synchronicity is achieved at each synchronization slot, errors on the order of hundreds of microseconds occur in between synchronization points.

A first variant consists in randomly choosing the comparator value, still keeping the same $I_A$ and $I_B$ for each synchronization slot. This approach is referred to as "random". Another possible variation requires an analysis of the algorithm initialization phase. More specifically, a fixed pattern is derived from a prediction of the TSs trend obtained through theoretical computations and applied accordingly. The main advantage of this implementation relies on the fact the

values are not computed at run-time. This approach is referred to as "proportional". Finally, another technique proactively adjusts the $R_C^A$ and $R_C^B$ choices in real time. Once the comparator values have been computed in the initialization phase, those are chosen through an optimization process aiming to minimize errors arising from the RTC approximations. This approach is referred to as the "predictive" method.
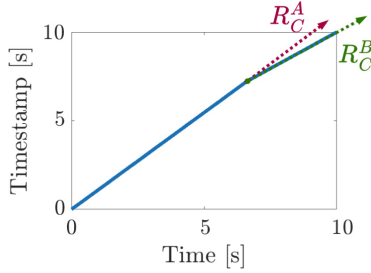


**Fig. 2.** Dual-modulus ST setting using the trivial approach. The slope is fragmented into two pieces, one due to the use of $R_C^A$, and the other related to $R_C^B$

## 4  Simulation Results

To assess the performance of each of the described modulus control policies, we performed a simulation campaign implementing FLSA on MATLAB. The slaves are driven by a clock signal $f_{clk} = 32.768$ kHz, thus leading to $t_{RTC} \approx 30.52$ μs. The ST resolution was set to $t_{ST} = 100$ μs. By the (2) the initial values for the dividers were found to be $R_C^A = 3$ and $R_C^B = 4$. The impact of the modulus control policy was studied during a time synchronization slot, whose duration was set to $T_M = 10$ s. We considered the error as the metric to evaluate the policy impact, defined as:

$$err_i[k] = \text{GT}[k] - ST_i[k] \tag{4}$$

where $\text{GT}_i[k]$ and $ST[k]$ represent the expected (Ground Truth, GT) TS and the system timer for the discrete time sample $k$ using the modulus control policy $i$.

The results of the ST mismatch are depicted in Fig. 3a for all the analyzed techniques, with a zoom on a time interval of 0.3 s in Fig. 3b.

By looking at Fig. 3c we can analyze the ST error using the different policies. With the default policy, as discussed, we obtain the highest error. The random approach introduces the first mitigation. Compared to the previous method, errors are halved, yet they remain in the range of hundreds of microseconds. With the proportional approach, errors are mitigated up to a maximum error of 200 ms. Optimal results are obtained through the predictive approach, where there is no predefined pattern, but rather the value to be utilized is chosen iteratively and
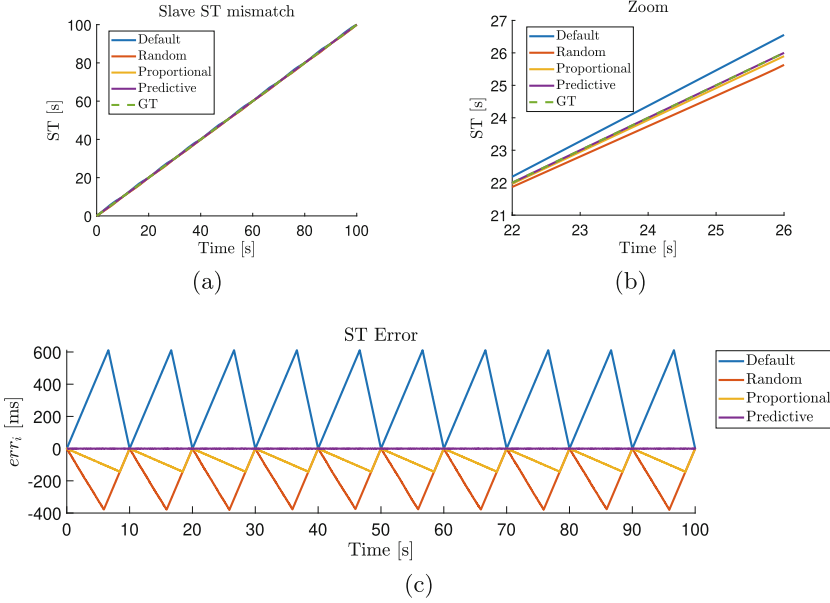
Fig. 3. Analysis of different modulus control policies: **a** the slave ST mismatch and **b** its zoom. In **c** it is shown the ST error according to the (4)

in real-time. This selection process involves identifying the value that exhibits the least deviation from the ideal value. Consequently, this technique yields a noteworthy degree of precision, surpassing the minimum threshold set by $t_{RTC}$. The accuracy is visible in Fig. 3b. In Fig. 4 we visually represented the use of each prescaler modulus value during the synchronization slot. With reference also to Fig. 3b. It can be observed that synchronization errors improve in correlation with the alternating of divisors in a uniform manner.



Fig. 4. Use of the modulus values in the four analyzed implementations: **a** Default, **b** Pure random, **c** Proportional and **d** Predictive technique

It's important to underline that each of the described methods nominally guarantees a zero error at the end of each synchronization slot, ensuring the

proper functioning of the FLSA algorithm as shown in Fig. 3c. Figure 5 shows an example of the default method applied to M/S synchronization.



**Fig. 5. a** M/S ST synchronization over a 100 s time frame and **b** zoom on a smaller interval

## 5  Conclusions

In this study, we examined the effect of the dual-modulus control policy on both in-slot and overall time synchronization accuracy. The choice of policy is driven by the goals, requirements, and characteristics of the specific application. If the primary concern is synchronization at the end of each synchronization round, the trivial technique may be suitable due to its low computational complexity. However, larger errors may occur in ST behavior between sync packets. In rapidly evolving systems with significant events requiring precise timing, in-slot accuracy is crucial. In such cases, the predictive approach is preferred as it ensures a maximum error equal to half the resolution of the RTC.

## References

1. Elson J, Girod L, Estrin D (2002) Fine-grained network time synchronization using reference broadcasts. ACM SIGOPS Oper Syst Rev 36(SI):147–163
2. Ganeriwal S, Kumar R, Srivastava MB (2003) Timing-sync protocol for sensor networks. In: Proceedings of the 1st international conference on Embedded networked sensor systems, pp 138–149
3. Maróti M, Kusy B, Simon G, Lédeczi, Á (2004) The flooding time synchronization protocol. In: Proceedings of the 2nd international conference on Embedded networked sensor systems, pp 39–49
4. Coviello G, Avitabile G, Talarico C, Roveda J, Florio A (2021) A fractional approach to time synchronization in wireless body area networks. In: 2021 IEEE international midwest symposium on circuits and systems (MWSCAS), pp 933–936

5. Coviello G, Avitabile G, Florio A, Talarico C, Wang-Roveda JM (2021) A novel low-power time synchronization algorithm based on a fractional approach for wireless body area networks. IEEE Access 9:134916–134928
6. Coviello G, Florio A, Avitabile G, Talarico C, Roveda JM (2022) Distributed full synchronized system for global health monitoring based on FLSA. IEEE Trans Biomed Circuits Syst (2022)
7. Razavi B (2012) RF microelectronics. Prentice Hall, Los Angeles, California

# Investigating and Importance of Fetal Monitoring Methods and Presenting a New Method According to Convolutional Deep Learning Based on Image Processing to Separate Fetal Heart Signal from Mother

Morteza Zilaie[1(✉)], Zohreh Mohammadkhani[1], Keyvan Azimi Asrari[2], and Sadaf Noghabi[3]

[1] Islamic Azad University Science and Research Branch, Khorasan Razavi, Mashhad, Iran
`mortezazilaie@gmail.com`
[2] The University of Pavia, Mashhad, Italy
[3] Imam Reza International University, Khorasan Razavi, Mashhad, Iran

**Abstract.** This Fetal electrocardiogram is a standard method to identify and diagnose fetal diseases. Therefore, effective techniques are needed to monitor fetal conditions during pregnancy and delivery. Meanwhile, obtaining the fetal electrocardiogram (FECG) signal, which contains the electrical activity of the fetal heart, has great importance, Of course, this signal is contaminated with many noises and disturbances and the most important of them is the mother's electrocardiogram signal. Inherently, the NI-ECG signal contains the maternal ECG signal, which has a larger amplitude than the fetal ECG signal. Therefore, it is not easy to detect the fetal QRS complex in order to control the condition of the fetus and prevent congenital defects. According to the above explanations, in this article, a deep learning approach based on a convolutional neural network is proposed to separate the electrocardiogram signals of the mother from the fetus without separating the mother's ECG signal. The proposed algorithm is able to reliably detect the fetal QRS complex. Also, in addition to not needing feature extraction steps, it has been able to show more suitable performance than the best methods proposed in previous research in terms of detection accuracy.

**Keywords:** Fetal heart signal · Deep learning · Electrocardiogram · Convolutional neural network · QRS complex

## 1 Introduction

Today, one of the prominent and challenging issues in the field of medicine is to check the condition of the heart of the fetus, so that according to the available statistics, one out of every 125 babies is born with heart failure. It may look healthy or it can be so severe that the child will have problems at birth [1]. There are also some possible complications such as Hypoxia caused by the closure of the umbilical cord around the neck of the fetus,

which causes heart failure, or congenital heart disorders that may affect the growth of the fetus after years.

Ways and methods of monitoring the fetal heart and checking its performance can prevent damage to the fetus and increase the hope of the fetus' survival [2]. The fetal heart rate (FHR)1 provides valuable information regarding the physiological conditions of the fetus, therefore, many studies have been conducted to determine the fetal heart rate and monitor and check the function of the fetal heart automatically [3]. Doppler ultrasound methods, fetal phonocardiography (PCG), fetal magnetocardiography (MCG), cardiotocography (CTG), cranial ECG (SECG), invasive electrocardiogram, non-invasive electrocardiogram, are among the suggested methods [4].

Doppler ultrasonography and phonocardiography are only able to determine the fetal heart rate (FHR) and do not provide the doctor with information about the morphology of the electrical signal of the fetal heart, so it is considered ineffective in many cases. The method of fetal magnetocardiography, which uses sensors located near the mother's abdomen to identify the magnetic field and the heart of the fetus, in addition to high costs and patient complaints about the signal collection process, is not recommended for long-term use due to the need to not move during recording. At the same time, this method shows the waveform of the fetal electrocardiogram signal with high accuracy and has a relatively good signal-to-noise ratio.

Cardiotocography (CTG) uses an ultrasound transducer and a pressure-sensitive transducer of uterine contractions to determine the fetal heart rate. Cardiotocography requires extensive training and high costs, as well as the safety and safety of the fetus in this method is doubtful. The SECG method is considered an invasive method and is very risky because it causes infection, so it is not suitable for long-term fetal monitoring. Doppler ultrasound and phonocardiography methods are not capable of showing the morphology of the electrical signal of the fetal heart.

Among all the methods, the non-invasive electrocardiogram method, in addition to providing an accurate estimate of the FHR, has complete monitoring of the morphology of the fetal heart rate waveform, etc.

According to the explanations given, the fetal electrocardiogram method (non-invasive) has priority over other methods, but there are limitations such as; Noises that have both biological and non-biological origin and make it difficult to record the signal limit of this method [5]. In order to increase the efficiency of this method, various types of research have been proposed to extract and separate the maternal electrocardiogram (MECG) from the fetal (FECG) signals more efficiently with this method, which will be briefly mentioned below and their disadvantages will be stated.

In research by Ping Gao et al., the blind source separation (BSS)2 method was used to extract fetal signals from a single-channel signal [6]. The results showed that the use of SVD alone is inefficient and the combination of ICA + SVD is more efficient in separation. The research was done by Fanelli et al., which used the combination of principal components extraction method, decision rules, and Matched Filter to extract the ORS complex of the mother signal [7]. In this method, researchers have used the Kalman method as well as pattern matching to estimate the mother signal. The research was done by Kwang Jin Lee and Boreom Lee in 2016 using the full filter of successive changes, which first reduces the noise components of the abdominal ECG to effectively remove

the maternal ECG through TSPCA and then separates the fetal ECG [8]. Xuwilson and colleagues have used principal component analysis and vertical imaging to reduce the effects of the maternal signal on the fetal signal [9].

They used the PCA clustering method to identify the fetal QRS complex. Research was done by Karimi Rahmati and colleagues in 2016 using PCA, and ICA techniques to extract FECG and detect R peak [10]. In another study, Feng tried to reduce the non-linear effects of maternal heart rate in fetal heart extraction by filtering [11]. Another research was conducted by Nannan Zhang et al. in 2017, which used ANC to remove MECG signals from recorded abdominal signals and combined SVD + SW techniques for estimation.

This method requires a reference signal that is morphologically similar to the maternal abdominal signal waveform [12]. Two other significant investigations are [5, 13]. In the work [13] on the compression and then recovery of the fetal signal with deep learning, no direct signal recovery has been done. In fact, the focus of [13] is based on the assumption that the fetal ECG signal is available, and we need a proper compression and recovery algorithm to send and receive it and recover the original signal from the signal received from the channel and remove its noise, which is the same as deep learning. In [5], which is closer to the proposed method of the article, the main features of the fetal ECG signal have been directly extracted using the convolutional neural network. The accuracy of this work in the best case has reached 77%.

The database cited in this study was Phiziont [14–16] have all used signal-based deep learning, which has the limitations of one-dimensional and noisy deep learning systems in all three. None of these three references have a separate step for noise and artifact removal. In the methods based on adaptive filters, due to the need for reference MECG signal (maternal electrocardiogram signal), the absence of a pure signal, and its contamination with noise, efficiency decreases. There are many problems in non-linear methods of signal separation due to the complexity of calculations. Linear decomposition is stationary and therefore not a good estimator for non-stationary FECG signals. The proposed method of this research uses fetal and maternal heart rate signal images and deep learning to detect the fetal QRS complex, which is expected to be more effective than conventional methods due to the innovative combination of deep learning and image approach in signal separation.

## 2 Deep Learning

Below is an implicit classification of the set of deep learning methods available [14].

An overview of a convolutional neural network architecture is shown in Fig. 1. In general, a CNN network consists of three main layers, which are: the convolution layer, pooling layer, and fully connected layer.

### 2.1 Convolution Layer

In these layers, the CNN network uses different kernels to convolution the input image and also create a map of intermediate features, and this map creates different features. Convolution operation has several benefits:

**Table 1.** Implicit classification of the set of existing deep learning methods

| | | |
|---|---|---|
| Alex Net | CNN-based Methods | Deep Learning Methods |
| Clarifai | | |
| VGG | | |
| Google Net | | |
| SPP | | |
| Deep Belief Networks | RBM-based Methods | |
| Deep Boltzmann Machines | | |
| Deep Energy Models | | |
| Sparse Autoencoder | Autoencoder-based Methods | |
| Denoising Autoencoder | | |
| Contractive Autoencoder | | |
| Sparse Coding SPM | Sparse Coding-based Methods | |
| Laplacian Sparse Coding | | |
| Local Coordinate Coding | | |
| Super-Vector Coding | | |



**Fig. 1.** A convolutional neural network [14]

The weight-sharing mechanism in each feature map reduces the number of parameters. Local connectivity learns the relationship between signal samples. It causes immutability and stability against noise and shift.

## 2.2 Pooling Layer

A pooling layer is usually placed after a convolution layer and it can be used to reduce the size of the feature map and network parameters. Like convolutional layers, pooling layers are stable to noise and limited spatial variations due to considering the relationships between signal samples in their calculations.

## 2.3  Fully Connected Layer

After the last pooling layer, there are fully connected layers that transform the feature map into a unified vector to continue the recognition process. Fully connected layers work like their counterparts in traditional artificial neural networks and include almost 90% of the parameters of a CNN network.

# 3  Suggested Method

Considering that in the past works, based on the review, all the articles have taken help from the deep neural network with the input of fetal and maternal heart rate signals, and with this approach, they have tried to determine the QRS complex of the fetus and separate it from the mother. The designed approach in this research, instead of examining the signal in the form considered in previous works, uses the image recognition approach. In the proposed approach, which is completely designed based on the behavioral pattern of the doctor and specialist, by drawing the heart signal of the fetus and the mother in terms of time, the deep neural network is trained, which separates the mother's heart rate as a long peak from the fetal heart rate as a peak. In short, label the right place of the QRS complex of the fetus.

This idea, that is, dealing with an image with a complex medical time-varying signal, can be more efficient than conventional methods, because there are some deep neural networks with very comprehensive initial training on servers such as Google or Amazon, which are used to separate specific areas of the image from other areas have been trained and this issue will greatly help to increase the accuracy of the work. After the text edit has been completed, the paper is ready for the template. Duplicate the template file by using the Save As command, and use the naming convention prescribed by your conference for the name of your paper. In this newly created file, highlight all of the contents and import your prepared text file. You are now ready to style your paper; use the scroll-down window on the left of the MS Word Formatting toolbar.

## 3.1  Database

Considering that the research is related to FECG and the database [17] is used, according to the given explanations, the database is simulated in 10 situations related to pregnancy based on real data and simulated for each pregnancy. 7 different physiological parameters have been measured. Taken tests, which are specified by Base Line definitions and numbers 0 to 5, include tissue, abdominal, baseline with noise, fetal movement with noise, and uterine contraction with noise. These parameters are measured and measured along with different levels of noise for each situation in 5 repetitions. Therefore, according to 10 pregnancy conditions, 7 measured gauges, 5 different noises, and 5 repetitions, a total of 1750 signals will be produced, which can be seen in 145.8 h of time efficiency and 1.1 million information peaks related to the fetus. Each sub-branch shows the pregnancy status in the form of a name, and the next two digits show the noise level. The noise levels are assumed to be 0, 3, 6, 9, and 12 dB respectively in the simulation, and the pregnancy status includes the numbers 0 to 20. In addition, the test repetition number

is the third number that is placed in the file. This database can be added to the current MATLAB path, and it should be noted that in the default state, the data cannot be read by Windows software and has a Linux extension. Figure 2 shows an example of database images. The horizontal axis is the sample number based on the sampling frequency of 1000 Hz and the vertical axis is the amplitude of the sampled signal in millivolts after amplification.
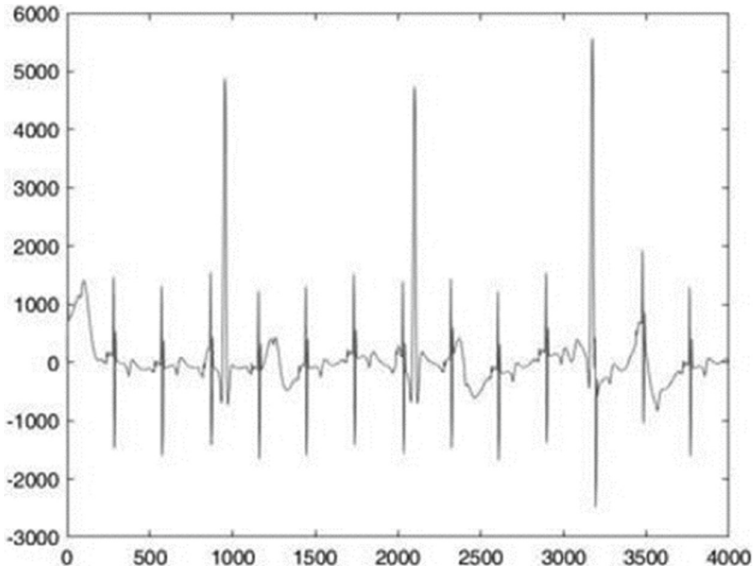


**Fig. 2.** Sample database image

## 3.2   Labeling Images According to the Input Pattern of Deep Neural Network

In the first step of the proposed method, fetal and maternal heart rate signal images should be drawn and labeled according to the input pattern of the trained deep neural network available in MATLAB software. Considering that the collected database had a different labeling than what we intended, this was done by using the information of that database and manual labeling, and finally, for each signal, a large number of separated signal images with labels on the heart rate of the fetus was determined. The input labels in the previously trained deep neural network were in the form of rectangular windows from the beginning to the end of the desired signal in the image, which were stored separately in a text file, and this file was later read by the software and converted into information related to the network as the appropriate label was added. This stage took a lot of time due to the high number of images and the need for high accuracy, but in the end, one hundred and forty-one images, each of which included seven to eight fetal heartbeats on average, were labeled. For this part of the work, a separate code was designed which, after drawing each image, allowed the user to specify the left and bottom right corners of the

rectangle with a mouse click, After this step, the labeled coordinates were automatically displayed and in the user's confirmation form was added to the database.



**Fig. 3.** Sample labeled image

In this way, about 140 images were obtained for deep neural network training. With the identification of the manual labeling of the database, the second step in the proposed method was to find a deep neural network that could properly separate and show the specified parts in the image. Image segmentation is determined and after implementation, the neural network with the best correct percentage was selected as the selected network.

### 3.3  Deep Neural Network

The pre-trained network of MATLAB software, which was used with some minor changes in the parameters in the proposed method, to isolate and recognize the selected part of the image, includes the following 15 layers.

As can be seen from the structure of the network in Tables 1 and 2, the first and last layers are considered input and output, respectively, and the rest of the layers are a combination of five different types of layers. The input of this network is the black and white image of the signal. Layers No. 2, 5, and 8 are considered convolutional layers that perform two-dimensional input convolution. The number of rows and columns padded around the base matrix is assumed to be 2 to maintain the dimensions of the output. The next layer is the Pooling layer, in this layer the inputs are down sampled and divided into square areas based on the settings given in the software, and finally the average of each square and five images for each area is calculated. In the proposed neural network, layers 3, 7, and 10 are pooling layers, respectively, and in these layers, the amount of pad performed was equal to zero. The next layer is the ReLU layer, which includes layers 4, 6, 9, and 12. In these layers, the data is threshold and values less than the threshold level

**Table 2.** Details of deep neural network used

| Layer type | Layer name | Layer number |
|---|---|---|
| Image Input | image input | 1 |
| Convolution | Conv | 2 |
| Max Pooling | max pool | 3 |
| ReLU | Relu | 4 |
| Convolution | Conv_1 | 5 |
| RelU | Conv_2 | 6 |
| Average Pooling | Avgpool | **7** |
| Convolution | Conv_2 | **8** |
| ReLU | Relu_2 | **9** |
| Average Pooling | Avgpool_1 | **10** |
| Fully Connected | Fc | **11** |
| ReLU | Relu_3 | **12** |
| Fully Connected | fc_rcnn | **13** |
| Softmax | Softmax | **14** |
| Classification Output | Classoutput | **15** |

**Table 3.** Description of each layer

| | |
|---|---|
| $32 \times 32 \times l$ images with zero center normalization | **1** |
| $32\ 5 \times 5 \times 3$ convolutions with stride [1 1] and padding [2 2 2 2] | **2** |
| $3 \times 3$ max pooling with stride [2 2] and padding [0 0 0 0] | **3** |
| ReLU | **4** |
| $32\ 5 \times 5 \times 32$ convolutions with stride [1 1] and padding [2 2 2 2] | **5** |
| ReLU | **6** |
| $3 \times 3$ average pooling with stride [2 2] and padding [0 0 0 0] | **7** |
| $64\ 5 \times 5 \times 32$convolutions with stride [1 1] and padding [2 2 2 2] | **8** |
| ReLU | **9** |
| $3 \times 3$ average pooling with stride [2 2] and padding [0 0 0 0] | **10** |
| 64fully connected layer | **11** |
| ReLU | **12** |
| 2 fully connected layer | **13** |
| Softma $\times$ | **14** |
| Cross entropy e $\times$ with classes stopSign and Background | **15** |

**Table 4.** Test results of the trained neural network

| Acc | PPV | Se |
| --- | --- | --- |
| 0.97 | 0.93 | 0.96 |

are equal to zero. The 13th and 11th layers are all connections, the details of which were mentioned in the theory related to neural networks, and finally, the penultimate layer is the 14th SoftMax layer, which is the final division based on values for the output layer. In this way, the deep neural network used to detect the image area is a neural network with fifteen layers, in which different layers including convolutions, pooling, ReLU, Fully Connected, and Softmax are used.

This network has been the foundation of deep learning in this work. After this step, the pre-trained deep neural network is re-trained based on the existing database data, and its correct percentage is estimated in the labeled images of fetal and maternal heart signals.

## 4   Results

According to the structure of the deep neural network, the data was divided into two categories, training and testing, First, it was retrained with 70% of the data, which means a total of 100 images of the network. After this stage, the deep neural network was evaluated using the test data, which included a total of 41 fetal and maternal heartbeat signals (287 fetal heartbeats). For analysis similar to previous works, in addition to the fetal heart rate in these signals, about 821 peaks (including noise and mother's heart rate) were detected with MATLAB software, which could be wrongly detected as fetal heart rate.) and not specifying the false peaks as the location of the beat, correct negative diagnosis (TN), the following results were obtained. The result of the implementation is given in Table 3.

| Tn | Fn | Fp | Tp |
| --- | --- | --- | --- |
| 800 | 11 | 21 | 276 |

where in

$$Se = \frac{TP}{TP + FN} PPV = \frac{TP}{TP + FP}$$

Due to the lack of a similar approach to this work, the results were compared with signal-based deep neural network [14–16]. These results, which took place in the years 2018 to 2020, are very similar to the current research in terms of the database. For example, the deep neural network used in [14] including input and output has 12 layers and is trained directly. According to the results given in [14–16], the best average results obtained in these methods with the same conditions as the present research were Se = 89.06, and PPV = 92.77 in terms of percentage, which is clearly weaker than the results obtained in the proposed method.

## 5 Conclusion

In this article, a method for detecting fetal heart rate signals and separating them from maternal heart rate signals based on a deep neural network is proposed. Considering the importance of automatically separating the heart rate of the fetus from the mother in diagnosing the deficiencies in the fetus and their appropriate treatment, providing an accurate estimate of the state of the heart rate of the fetus and analyzing its morphology can be very useful.

Despite the limitations such as the mother's strong heartbeat and other noises in the environment, the proposed method based on a deep neural network was able to combine the fetal signal with the mother based on the drawn image without the need for feature extraction or pre-processing steps. Acceptability compared to the previous methods to determine the fetus's heart rate in different areas of the image.

In addition to not needing feature extraction steps, the proposed method has shown a very good performance in terms of correct percentage. Finally, the work that was done compared with relatively new articles in this field and its proper efficiency was confirmed. It is hoped that due to the increasing progress of medical engineering sciences and wireless networks, systems can be produced that can be easily carried and used by pregnant mothers, and the information of the fetus and the mother, including heart rate and all items her vitals should be taken at the right time and sent to the medical condition control centers and provided to specialist doctors, to minimize possible risks for the fetus and the mother.

## References

1. Sahaf MA (2007) Applying blind separation of correlated signal sources to extract fetal electrocardiogram based on second-order statistics. SID 1
2. Wrobel J et al (2015) Pregnancy telemonitoring with smart control of algorithms for signal analysis. Journal of Medical Imaging and Health Informatics 5(6):1302–1310
3. Boudet S et al (2018) An online viewer of FHR signal for research, E-learning, and tele-medicine. In: Bioinformatics and Biomedical Engineering: 6th International Work-Conference, IWBBIO 2018, Granada, Spain, April 25–27, 2018, Proceedings, Part II 6. 2018. Springer
4. Martinek R et al (2015) Enhanced processing and analysis of multi-channel non-invasive abdominal fetal ECG signals during labor and delivery. Electron Lett 51(22):1744–1746
5. Zhong W et al (2018) A deep learning approach for fetal QRS complex detection. Physiol Meas 39(4):045004
6. Gao, P., E.-C. Chang, and L. Wyse. Blind separation of fetal ECG from single mixture using SVD and ICA. in Fourth International Conference on Information, Communications and Signal Processing, 2003 and the Fourth Pacific Rim Conference on Multimedia. Proceedings of the 2003 Joint. 2003. IEEE

7. Fanelli, A., et al. A smart wearable prototype for fetal monitoring. in Proceedings of the 4th International Conference on Applied Human Factors and Ergonomics (AHFE), San Francisco, CA, USA. 2012

8. Lee KJ, Lee B (2016) Sequential total variation denoising for the extraction of fetal ECG from single-channel maternal abdominal ECG. Sensors 16(7):1020

9. Xu-Wilson, M., et al. Spatial filtering and adaptive rule-based fetal heart rate extraction from abdominal fetal ECG recordings. in Computing in Cardiology 2013. 2013. IEEE

10. Rahmati AK, Setarehdan S, Araabi B (2017) A PCA/ICA based fetal ECG extraction from mother abdominal recordings using a novel data-driven approach to fetal ECG quality assessment. Journal of biomedical physics & engineering 7(1):37

11. Dang S, Tian W, Qian F (2009) De-noising fractional noise in fiber optic gyroscopes based on lifting wavelet. Chin J Lasers 36(3):625–629

12. Zhang N et al (2017) A novel technique for fetal ECG extraction using single-channel abdominal recording. Sensors 17(3):457

13. Muduli, P.R., R.R. Gunukula, and A. Mukherjee. A deep learning approach to fetal-ECG signal reconstruction. in 2016 Twenty Second National Conference on Communication (NCC). 2016. IEEE

14. Lee JS et al (2018) Fetal QRS detection based on convolutional neural networks in noninvasive fetal electrocardiogram. in 2018 4th International Conference on Frontiers of Signal Processing (ICFSP). 2018. IEEE

15. Jagannath D, Dolly DRJ, Peter JD (2020) Composite Deep Belief Network approach for enhanced Antepartum fetal electrocardiogram signal. Cogn Syst Res 59:198–203

16. Jagannath D, Dolly DRJ, Peter JD (2020) Deep learning strategies for fetal electrocardiogram signal synthesis. Pattern Recogn Lett 136:286–292

17. Andreotti, F., et al. Comparing feature-based classifiers and convolutional neural networks to detect arrhythmia from short segments of ECG. in 2017 Computing in Cardiology (60). 2017. IEEE

# Real-Time Sea Monitoring Using FMCW Radar

Assil Chawraba[1,2(✉)], Ali Rizik[3], Andrea Randazzo[1], and Daniele Caviglia[1]

[1] Department of Electrical, Electronic, Telecommunication Engineering and Naval Architecture (DITEN), University of Genoa, Genoa, Italy
assilchawraba97@gmail.com

[2] Lebanese School of Science and Technology (MECRL Lab-EDST), Lebanese University Beirut, Beirut, Lebanon

[3] Department of Environment, Land and Infrastructure Engineering (DIATI), Politecnico Di Torino, Turin, Italy

**Abstract.** This paper presents a method for sea monitoring based on a low-cost 24 GHz off-the-shelf FMCW radar and an embedded Raspberry Pi PC. The technique relies upon the use of range-Doppler maps, which are obtained through a 2D FFT processing of the acquired scattered data. The resulting insights may offer new possibilities for monitoring and understanding sea waves, which have important applications The study's findings provide valuable insights into monitoring and understanding sea waves, which can have significant implications for predicting natural disasters and mitigating their impact based on predicted disaster scenarios.

**Keywords:** Sea monitoring · FMCW radar · Range-Doppler processing · 2D FFT · Doppler features · Feature extraction

## 1 Introduction

Sea wave detection and monitoring is an active area of research due to its potential for predicting and mitigating the impact of natural disasters and managing coastal infrastructure [1]. Indeed, understanding the properties of water waves is crucial in various applications such as the near shore region, ship operations, level measurement in tanks, and sea keeping and maneuvering. In [2], the authors present a method based on x band Radar for estimating different wave parameters such as wave height, period and velocity using Frequency-modulated continuous-wave (FMCW) radar. This highlights the importance of accurate and reliable methods for measuring and analyzing water wave behavior, which can help improve safety and efficiency in maritime operations.

The authors in [3] design a radar system for deploying on an unmanned surface craft. It includes two subsystems, the first is a long-range radar for detecting sea-surface objects a few Kilometers away, and the second one entails a short range FMCW radar.

The physical characteristics of sea waves, including their wavelength, height, and breaking patterns, are crucial for predicting their behavior and potential impact on coastal areas. Categorizing waves based on these characteristics can facilitate the development

of effective strategies for coastal management and disaster risk reduction. Understanding these different types of waves is important for various applications, including oceanography, marine engineering, and coastal management, as it can help improve safety and efficiency in maritime operations. By accurately identifying and classifying different types of waves using appropriate sensors, it is possible to develop more effective strategies for managing and mitigating the risks associated with waves in different marine environments [4].

To address the inefficiency of using separate systems for sea surface monitoring with microwave remote sensing techniques, the authors in [5] integrated the measurement functions of scatterometers, altimeters, and Doppler radars into a single observation system. This integration enhances the efficiency of system operation and enables cross-analysis of the observation data. The result is a more streamlined and effective approach to sea surface monitoring using microwave remote sensing techniques.

FMCW radar is a promising technology for sea wave monitoring, which may provide a range of information on sea waves. Indeed, FMCW radar has the added advantage of being able to detect both the water surface level and the surface velocity in a flood warning system [6].

In particular, the developed system relies upon the use of a low-cost FMCW radar board and a Raspberry Pi PC, in which information for near-shore sea monitoring are obtained through a range-Doppler map processing technique.

## 2 FMCW Radar Processing

Figure 1 illustrates the adopted data acquisition system, which utilizes the "Distance2Go" radar board by Infineon Technologies [7] for data collection. The system is housed in a box, which also contains a Raspberry PI embedded PC, a camera and a power supply (Power bank). The Distance2Go board is locally connected to the Raspberry PI mini-PC through a USB connection, and the Raspberry PI mini-PC is in wireless communication with the control PC that collect and process the radar raw data. A power bank or a power source with a USB output can be used to power both the Raspberry Pi and the Distance2Go board. The camera is connected to the raspberry PI, allowing to also collect video images of the monitored scenario.

FMCW radar is a type of radar that uses a continuous wave frequency modulation technique to measure the distance and velocity of objects. The range and velocity of the target can be obtained by analyzing the frequency and phase shifts of the transmitted and received signals.

The delay between transmitted and received signals is determined by:

$$t_d = \frac{2R}{c} \tag{1}$$

where $R$ is the target range and $c$ are the speed of light in a vacuum.

In an FMCW radar transceiver, the beat frequency ($f_{beat}$) is the frequency of the beat signal that is produced by mixing the transmitted and received signals.

**Fig. 1.** The monitoring system's hardware prototype.

The beat frequency produced by an FMCW radar can be used to calculate the range of a target object by:

$$R = \frac{cT_{chirp}}{2B_{sweep}}(f_{beat}) \tag{2}$$

$B_{sweep}$ is the sweep bandwidth, and $T_{chirp}$where is the chirp time.

The range $R$ is proportional to the beat signal frequency.

The maximum velocity is limited by ensuring that the calculated phase is clear because velocity is calculated by:

$$v = \frac{\lambda}{4T_{chirp}} \tag{3}$$

where $v$ is the target velocity and $\lambda$ is the wavelength.

Radar range resolution ($\Delta R$) refers to the ability of the radar system to distinguish or resolve closely targets in the range can be calculated by:

$$\Delta R = \frac{c}{2B_{sweep}} \tag{4}$$

The resolution of velocity can be defined as the smallest difference in velocity that a radar can detect and is an important parameter that affects the radar's ability to accurately identify and track targets.

$$\Delta v = \frac{\lambda}{2(T_f)} \tag{5}$$

where $Tf$ is time period.

The range-Doppler (R-D) map generated by the radar signal processing system provides valuable information on the behavior of the sea, including the presence of waves and velocity. The R-D map is obtained using a double FFT algorithm that samples the intermediate frequency (IF). In FMCW radar systems, target information is obtained

by performing a range-FFT along each chirp to generate range bins representing the distance to the target, and a second FFT along the chirps for a single range bin to extract Doppler information that reflects the target's velocity. These two FFT processes, known as range-FFT and Doppler-FFT, are combined to generate a 2D FFT result, which produces a Range-Doppler map.

The radar experimental simulation parameters are shown in (Table1).

**Table 1.** Experimental simulation parameters.

| Parameters | Value |
| --- | --- |
| Bandwidth | 200 MHz |
| Center Frequency | 24 GHz |
| Chirp Time | 1.5 µs |
| Maximum Range | 25 m |
| Range Resolution | m |

## 3   Experimental Results and Discussion

The experimental setup involved placing the radar device approximately 1.5 m away, with the device oriented directly towards the sea as shown in Fig. 2. The data were collected wirelessly using a MATLAB code and a Raspberry PI 3B + embedded PC. The Raspberry PI collected the raw data from the radar board via USB and transmitted them to MATLAB through a Wi-Fi network. The data were collected specifically for the purpose of detecting sea waves.

This paper examines the performance of an FMCW radar for sea monitoring using the MATLAB tool.

We set up the measurement environment using a radar front-end module and a real-time data acquisition module in order to test the effectiveness of the proposed sea wave schema based on Doppler spectrum.

Figure 3 shows the Doppler spectrum and Range-Doppler map results from an experiment on sea waves.

The Range-Doppler (R-D) map in Fig. 3 indicates a water flow, as evidenced by the high intensities observed. The camera image confirms the water flow.

Furthermore, the R-D map of the sea exhibits distinctive reflection patterns that indicate the presence of waves. These patterns are characterized by multiple reflection centers that rapidly move with an approximate velocity of –2.5 km/h. The reflection centers exhibit variations within a range of 5 m, indicating the dynamic nature of the waves. This information provides valuable insights into the sea's wave characteristics, including their speed and spatial distribution.

The red line shown in Fig. 3 represents threshold for range detection in the FMCW processing based on the 2D FFT analysis.

The negative sign in the velocity measurement also provides information about the flow direction towards the radar.

Further analysis of the Doppler spectrum for sea waves reveals numerous variations, resulting in multiple peaks at different Doppler frequencies. This information provides evidence that the variations in the Doppler spectrum and R-D map are effective tools for analyzing the sea waves.



**Fig. 2.** Experiment setup for Sea Waves detection.



**Fig. 3.** Camera Image, R-D map and Doppler spectrum at 5 m range for sea wave.

The presence of waves, the flow direction, and the Doppler spectrum are just a few of the features that can be analyzed to obtain valuable information about the sea.

## 4   Conclusion

This paper has presented the effectiveness of a low-cost 24 GHz FMCW radar system operating in multi-chirp sequence mode for sea waves and detecting both range and velocity. The results of the Doppler analysis and the Range-Doppler map demonstrate

that sea waves produce distinctive Doppler features. The use of FMCW radar to study sea waves and flow has great potential for advancing our understanding of these environments and improving our ability to manage and protect against hazards such as floods and debris flows. These insights have opened possibilities for the development of innovative monitoring systems that can enhance safety and efficiency.

## References

1. Hwang J-H, Kim D-J (2020) Sea Surface Imaging with a Shortened Delayed-Dechirp Process of Airborne FMCW SAR for Ocean Monitoring on Emergency. Sensors 20(24):7310
2. Eric A, Ericson, David R. Lyzenga and David T. Walker (1999) Radar backscatter from stationary breaking waves. Res Oceans 104.C12, pp. 29679–29695
3. Cui J et al (2015) Wave height measurement using a short-range FMCW radar for unmanned surface craft. Oceans 2015-MTS/IEEE Washington. IEEE
4. Linghu L, Wu J, Huang B, Wu Z, Shi M (2018) GPU-accelerated massively parallel computation of electromagnetic scattering of a timeevolving oceanic surface model I: Time-evolving oceanic surface generation. IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens 11, 8, 2752–2762
5. Hwang J-H, Kim D-J, Kang K-M (2022) Multifunctional scatterometer system for measuring physical oceanographic parameters using range-Doppler FMCW radar. Sensors 22(8):2890
6. Scherhaufl M, Hesch C, Sevar JM (2019) Fluid surface velocity estimation using a 77 GHz radar module. In: Proc. IEEE Topical Conf. Wireless Sensors Sensor Netw, pp 1–4
7. Infineon Technologies (2020) Distance2Go—XENSIV 24 GHz RadarDemo Kit With BGT24MTR11 and XMC4200 32-Bit ARM CortexTMM4MCU for Ranging, Movement and Presence Detection. https://www.infineon.com/cms/en/product/evaluationboards/demodistance2go/

# Multi-agent Systems for Pervasive Electronics: A Case Study in School Classrooms

Juvenal Machin, Edgar Batista, and Agusti Solanas(✉)

Universitat Rovira i Virgili, Av. Països Catalans 26, 43007 Tarragona, Spain
{juvenal.machin,edgar.batista,agusti.solanas}@urv.cat

**Abstract.** Complex electronic systems comprise heterogeneous components operating concurrently to achieve complex goals. Electronic systems (like everything else) can fail. Hence, having resiliency in mind, components must adapt to changing circumstances to meet their goals. In this context, multi-agent systems (MAS) offer a flexible approach to managing complex problems by leveraging agents' collective behaviour and sensory abilities. This article presents a multi-layered event-driven hierarchical MAS model, which enhances the resilience of pervasive electronic systems by adapting to events within the application domain. We report a case study with classrooms augmented for contextual sensing.

**Keywords:** Multi-agent systems · Hierarchy · Pervasive electronics

## 1 Introduction

Electronic systems might encounter many unexpected circumstances in dynamic and ever-changing environments that negatively impact their functioning. Some of these adverse events might be internal (e.g., device failures and dead batteries), and others might be external, i.e., originated in the environment where these systems are placed (e.g., stolen components, power outages and fires). Thus, electronic systems can hardly be programmed statically but should adapt to changing situations autonomously. To this end, systems should have computational and learning capabilities to understand their nearby environment, identify potential threats, and respond accordingly.

Multi-agent systems (MAS), an active research area in distributed artificial intelligence, are key enablers to meet these requirements. They comprise multiple autonomous computing entities, called agents, able to communicate and cooperate in a shared environment while performing complex tasks [1]. Through sensors, agents can acquire information from the environment, make local information-based decisions, and take action on the environment accordingly. While each agent decides the best action to solve a specific task, MAS leverage social abilities, such as coordination, negotiation and cooperation, to deal with complex problems that cannot be solved using the capabilities of a single agent. Since

MAS can flexibly and efficiently perform complex problem-solving, they are suitable for uncertain environments under constant change. Despite their benefits, it is not yet clear how to organise agents to maximise their adaptability and reliability capabilities, and this design decision remains problem specific.

In his article, we discuss the opportunities of adopting a hierarchical organisation for MAS-based systems in the context of pervasive electronics. More specifically, after describing the hierarchical organisation, we exemplify its application with a use-case focused on smart classrooms. The remainder of the article is organised as follows: Sect. 2 provides a quick overview of MAS organisations. Next, Sect. 3 outlines our hierarchical approach, and Sect. 4 presents a case study in school smart classrooms. The article finishes with some conclusions in Sect. 5.

## 2    Quick Background on MAS Organisations

Multi-Agent Systems (MAS) is a field within the broader discipline of Distributed Artificial Intelligence. They have been defined as systems comprising multiple autonomous computing elements—called agents—that can interact and cooperate with each other [1]. MAS use these cooperation (i.e., social) abilities and their interactions with the environment to solve problems collaboratively (e.g., by dividing them and addressing smaller tasks) in a flexible way. Thanks to their knowledge of the environment, agents in a MAS can choose and perform actions, making them suitable for dynamic and complex environments. MAS have been used in various areas, including computer science, logistics, transportation, industry, and control engineering. However, most of the work on applying MAS is still domain-specific, and there is a lack of flexibility and resilience. Moreover, there is no consensus on a single taxonomy for MAS, and several approaches, following diverse criteria can be found in the literature. Regarding how MAS organise their agents, the following approaches are mostly considered:

– Flat organisation: all agents are considered equal and collaborate and communicate at the same level. The flat organisation significantly increases communication and local processing overhead.
– Hierarchical organisation: agents follow tree-like hierarchical relations by which higher-level agents can control and assign tasks to lower-level agents.
– Holonic organisation: fractal-like organisation in which agreed commitments in nested groups bound agents together. Any group (holon) can be logically seen as an agent comprising sub-agents arranged in a self-similar structure.
– Coalitions: agents with similar goals form temporary groups that are dismantled when they are no longer needed. Cooperative and self-interested agents may form coalitions, which usually have a flat structure.
– Teams: Agents in teams pursue a commonly agreed objective. That is, agents collaborate to reach a team goal, which can differ from their own goals.
– Markets: resemble real-world market economies. Based on self-interested agents that compete to receive items (such as services, shared resources, or tasks), trying to maximise their individual utilities.

– Congregations: agents with complementary abilities form congregations, which are long-lived flat organisations, to attain requirements that they cannot achieve individually. Agents in a congregation do not have any concrete predefined goal, and they group together to achieve a stable set of capabilities.
– Matrices: are multi-hierarchy organisations in which at least one manager agent manages each agent. That is, managed agents may follow multiple chains of command and rules.
– Societies: are open systems formed by heterogeneous agents having any goal or capability. Interactions inside a society are flexible, and agents can form other types of structures or be unrelated while belonging to the society.
– Federations: grouped agents transfer some degree of autonomy to a single intermediary or broker agent representing the group. Group members interact only with their broker, which can communicate with other brokers.

It is still unclear which organisation strategy is best, and the election is made ad-hoc and problem specific.

## 3    Multi-layered Event-Driven Hierarchical MAS

Inspired by the Complex Adaptive Systems (CAS) properties [2], we propose a multi-layered event-driven hierarchical MAS model to enhance the resiliency of pervasive electronics systems applications. As introduced in the previous section, in a hierarchical organisation, agents follow tree-like hierarchical relations by which higher-level agents can control and assign tasks to lower-level agents. The tree-like structure in which agents are organised allows partitioning problems into smaller, more manageable portions that can be solved collaboratively. This facilitates the coordination of multiple behaviours and enables the resolution of conflicts that may arise among agents [3].

In our multi-layered event-driven hierarchical MAS model, agents are organised hierarchically in stacked layers according to a set of semantic categories representing different levels of abstraction. Each layer represents an increased level of complexity. Our approach is characterised by (i) **events**, i.e., relevant pieces of contextual information about the system's state sent from one agent to another, (ii) **roles**, i.e., behaviour, capabilities and constraints that agents are compelled to fulfil, acting as a representation of the observable behaviour exhibited by agents. According to a role-transition graph, agents change roles in response to the events received and the world perceived. And (iii) the very **agents**, i.e., heterogeneous software or hardware entities with diverse capabilities assigned in real-time based on the system's needs.

Regarding communications, a double flow of data and control between adjacent levels guide the MAS operation. In particular, agents operate in a continuous Observe-Orient-Decide-Act (OODA) loop [4] through which they: (i) perceive the world, (ii) enrich perception with contextual data, (iii) decide which role is enacted, and (iv) enable the proper roles and follow the associated behaviours.

The hierarchical organisation of the agents is exploited in two ways: first, lower-level data flows towards higher levels, where it gets aggregated to achieve

contextual awareness. Second, control messages are pushed towards lower-level agents to keep the MAS aligned with the needs of higher abstraction levels.

Agents activate roles according to their capabilities, the occurrence and priority of events, and the needs of the collective. Collective needs, evaluated following distributed multivariate criteria, comprise: (i) minimising the duration of adverse events, (ii) maximising role redundancy across any level, and (iii) maximising role diversity across any level. Using a repository of roles (available to all agents for installation) promotes code modularity and re-usability. In addition, enforcing role redundancy and diversity enables flexible responses to various adverse events, which enhances the preparedness and resiliency of the system.

Several properties of our model can be easily identified:

– Robustness: The distributed nature of this organisation reduces performance bottlenecks and service outages, improving robustness.
– Heterogeneity: The ability to allocate resources across a network of agents allows the use of resource-constrained devices (e.g., IoT devices) in lower-level agents. By providing a hierarchical relation among agents, lower-level agents can sit in low-capacity devices and send data to a parent agent hosted by a powerful device, where complex processing is made.
– Awareness: Agents can act as malfunction detectors for their sibling agents (e.g., by comparing the readings of several sensor agents or polling them with a heartbeat event), and take action, such as role switching, to support the continuity of operations.
– Compartmental-friendly: Our organisation allows partitioning the scenario into several levels of complexity, which aids problem conceptualisation, access control, and agent behaviour design.
– Distribution of Labour: Low-level agents can respond to environmental changes by sensing their surroundings and taking action. Higher-level agents concerned with abstract contexts can run more complex roles that imply deliberation and learning to make inferences and gather knowledge.

Since this model is very theoretical, in the next section, we put it into practice and show a real application to clarify its functioning.

## 4   Case Study: School Classrooms Monitoring

To validate the suitability of our model, we have implemented a monitoring system for contextual sensing [5]. This system, installed in a hundred classrooms from forty schools across Catalonia, captures and analyses several contextual variables, such as air temperature, relative humidity, $CO_2$ concentration, particulate matter, room ventilation, UV light, noise level, and people's movement, among others, using low-cost sensors [6,7]. Within the classrooms' context, these variables gain relevance due to their relationship with academic performance and the spread of respiratory viruses [8].

According to our model, we distinguish four agent roles (see Fig. 1). At the lowest level, the *sensor agent* is responsible for interacting with a specific sensor. This agent uses the sensor's library to collect the corresponding contextual

**Fig. 1.** Our MAS model adapted to contextual monitoring in schools

variable. Note that if a sensor is broken, its agent will malfunction, but the rest of the sensor agents will not be affected. Also, this agent can raise alarms if a given variable exceeds a threshold (e.g., air temperature higher than $28\,°C$). At the second level, there is the *controller agent*, which orchestrates the logical functioning of several sensor agents. This agent is responsible for (i) verifying that all sensor agents are up and running and (ii) forwarding the data received from the sensor agents to the higher level. At the third level, there is the *classroom agent*, which is responsible for controlling all the sensor agents within a classroom. With the advent of smart classrooms [9], several technological systems will interact within the classroom: contextual sensing, multimedia systems, teaching facilitators, and smart equipment. Classroom agents will therefore be essential to coordinate all the functionalities, operations, and data flows within each classroom. Last but not least, at the fourth level, there is the *school agent*, able to manage and coordinate all the information gathered from its classroom agents. This agent provides a graphic data visor to display and report all the school's related information to the school's principal (see Fig. 2).



(a) Data visor

(b) Catalan primary schools using our model

**Fig. 2.** Contextual monitoring in primary schools' classrooms in Catalonia (Spain)

## 5    Conclusions

Environments are becoming increasingly unpredictable, and electronic and computing systems must be ready to adapt to changing circumstances. Regarding resiliency, systems must be able to adapt to guarantee proper functioning. This property is essential in sensitive environments like healthcare. Thus, in this article, we have presented a multi-layered event-driven hierarchical MAS model, and we have exemplified its functioning in a real-world application. With our solution, agents are organised hierarchically with multiple layers of abstraction and can dynamically adapt their roles/behaviours depending on events.

This is a first step towards a more comprehensive analysis. We recognise that with this paper, we have only started to scratch the surface of a very complex and promising research line. Further work could focus on the following lines:

– Performance: There is a need for a thorough analysis of the performance of our hierarchical organisation. Protocols for the exchange of messages and events should be carefully studied.
– Roles marketplace: Creating a "marketplace" for agents to access a variety of roles is paramount for the real scalability and reusability of our solution.

## References

1. Wooldridge, M.: An Introduction to MultiAgent Systems. Wiley Publishing, 2nd edn. (2009)
2. Carmichael, T., Hadžikadić, M.: The Fundamentals of Complex Adaptive Systems, pp. 1–16. Springer International Publishing, Cham (2019)
3. Gao X, Liu R, Kaushik A (2021) Hierarchical Multi-Agent Optimization for Resource Allocation in Cloud Computing. IEEE Transactions on Parallel and Distributed Systems 32(3):692–707
4. Boyd J (1995) OODA loop. Tech. rep, Center for Defense Information
5. Solanas A, Patsakis C, Conti M, Vlachos IS, Ramos V, Falcone F et al (2014) Smart health: A context-aware health paradigm within smart cities. IEEE Communications Magazine 52(8):74–81
6. Batista E (2022) at al.: On the Deployment of Low-Cost Sensors to Enable Context-Aware Smart Classrooms. Applications in Electronics Pervading Industry, Environment and Society (ApplePies), vol 1036. LNEE. Springer, Cham, pp 333–338
7. Villanova, O. et al.: Sensing kit for the study of respiratory disease transmission in school classrooms. In: 11th International Aerosol Conference (IAC). p. 825 (2022)
8. Bluyssen PM (2017) Health, comfort and performance of children in classrooms - New directions for research. Indoor and Built Environment 26(8):1040–1050
9. Zhang, Y., Li, X., Zhu, L., Dong, X., Hao, Q.: What Is a Smart Classroom? a Literature Review. Shaping Future Schools with Digital Technology pp. 25–40 (2019)

# Where Are My Cryptos?

Miquel Calonge[1], Edgar Batista[1], Julio Henández-Castro[2],
and Agusti Solanas[1(✉)]

[1] Universitat Rovira i Virgili, Tarragona, Spain
miquel.calonge@urv.cat, edgar.batista@urv.cat, agusti.solanas@urv.cat
[2] University of Kent, School of Computing, Canterbury, UK
j.c.hernandez-castro@kent.ac.uk

**Abstract.** The financial sector has suffered a groundbreaking transformation with the advent of cryptocurrencies, shifting from centralised to decentralised schemes. Hardware wallets play an essential role in storing cryptocurrencies securely. However, these electronic devices generally have limited resources that open the door to attacks. In this article, we describe three attacks against them. Several wallets with funds or recent transactions have been discovered with these attacks.

**Keywords:** Hardware wallets · Cryptocurrencies · Entropy · Derivation paths · Seed phrases · Security

## 1 Introduction

In recent years, the financial landscape has undergone a revolutionary transformation with the emergence of cryptocurrencies [1]. One of the most significant changes this digital revolution brings is the newfound responsibility placed on users to safeguard their assets. Eliminating intermediaries, such as banks and other financial institutions, implies users have absolute control over their digital wealth. This decentralisation eradicates the single point of failure that could expose funds to potential hacks. Blockchain technology is widely used in this new digital financial ecosystem [2], with over 100 million users worldwide and a total market capitalisation of cryptocurrencies exceeding the billion dollars.

Thanks to cryptocurrencies, users hold a private cryptographic key to their financial sovereignty. To help users securely manage their private keys, electronic devices have emerged to enhance the safety of funds transactions: hardware wallets [3]. However, these purpose-built devices are often resource-constrained, with limited computational capabilities, memory capacity, and reduced entropy sources. Despite these limitations, they remain a popular choice for cryptocurrency holders.

This article aims to demonstrate that the resource-constrained nature of hardware wallets may lead to potential security attacks. In particular, three attacks against hardware wallets targeting their random wallet generation system and their limited memory are shown. Preliminary results have demonstrated the

feasibility of these attacks. The rest of the article is organised as follows: Sect. 2 presents some background notions, Sect. 3 describes the proposed security attacks, Sect. 4 provides some experimental results, and Sect. 5 closes the article.

## 2 Background

This section elaborates on the main concepts addressed in this article: hardware wallets, seed phrases and derivation paths. Figure 1 illustrates the relationship amongst these concepts.



**Fig. 1.** From bits to cryptographic keys in hardware wallets

### 2.1 Hardware Wallets

Hardware wallets have emerged as a reliable solution for safeguarding cryptocurrencies. These physical devices serve as secure offline storage, providing an impregnable fortress for digital assets against potential online threats. Unlike software wallets that run on computers or smartphones, hardware wallets are purpose-built devices designed solely for the secure management of cryptocurrencies. They operate in an isolated environment, commonly known as a "cold storage" setup, which ensures that the private keys are never exposed to the Internet. Hardware wallets offer two primary functionalities: (i) creating new wallets and (ii) restoring existing wallets. On the one hand, the hardware wallet needs to generate a new seed phrase with some randomness source for creating new wallets. On the other hand, hardware wallets allow users to restore a previously used wallet by introducing a seed phrase and, hence, regain access to the funds and transaction history. Some popular hardware wallet brands include Ledger, Trezor, and KeepKey [4,5].

Introducing randomness into a hardware wallet is essential to provide users with a truly random and secure cryptocurrency wallet. However, this implementation can pose several risks. For instance, hardware limitations can create challenges in generating sufficient levels of entropy required for cryptographic operations. When the randomness generation process is compromised or biased, it can lead to weak cryptographic keys or predictable seed phrases. This, in turn, exposes users' funds because adversaries might exploit patterns in the keys or phrases to gain unauthorised access.

## 2.2    Seed Phrases

Seed phrases were introduced to store the private keys in a secure and user-friendly fashion. A seed phrase (a mnemonic or recovery phrase) is a sequence of human-readable words generated deterministically from the private key. The relationship between the seed phrase and the private key is pivotal. When a wallet is created or restored, the seed phrase generates the private key and all the corresponding public addresses. A seed phrase may consist of 12, 15, 18, 21 or 24 words representing the underlying private key of 128, 160, 192, 224 and 256 bits, respectively.

Most hardware wallets work with seed phrases following the BIP39 (Bitcoin Improvement Proposal 39) standard, a widely adopted standard for generating seed phrases and their corresponding private keys. BIP39 aims to enhance the security and usability of mnemonic phrases by providing guidelines for generating seed phrases and converting them into deterministic seeds. Also, BIP39 defines a wordlist of 2048 unique words from which the words in the seed phrase are chosen.

## 2.3    Derivation Paths

In cryptocurrencies, derivation paths are hierarchical structures that generate and organise cryptographic keys to simplify key management and enhance compatibility. BIP32, BIP44, BIP49, BIP84, and BIP141 are some of the most common standards to generate derivation paths.

BIP32 introduces a hierarchical deterministic key generation method, allowing users to generate a master key from a random seed value. From this master key, a virtually infinite number of child keys can be derived deterministically, making it easier to manage multiple cryptocurrency addresses securely. Each derivation in the path is determined by a specific index (the values between the slashes), and the resulting keys can be used for receiving or sending funds. However, BIP32 does not define a specific structure for organising these keys, leading to the development of BIP44. BIP44 defines a multi-level derivation path that includes specific account indexes, making it easier for wallets to support multiple cryptocurrencies without mixing their keys. For example, a BIP44's derivation path for Bitcoin (BTC) might look like m/44'/0'/0', where the first index (44') denotes the BIP44 standard, the second index (0') represents the BTC account,

and the third index (0') points to a specific address. From here, the last index might change (e.g., 1', 2', and so on) to obtain other cryptographic keys.

## 3    Attacks over Hardware Wallets

In this section, we describe three attacks against hardware wallets, namely (i) attacks based on low entropy, (ii) attacks based on uncommon derivation paths, and (iii) attacks based on public domain seeds. As a result of all these attacks, a list of public keys is obtained, which is then compared to online databases of public keys (e.g., Loyce.club) that have received transactions at some time. When matches are found, the wallet's private key is known (see Fig. 2).



**Fig. 2.** Methodology of the security attacks conducted

### 3.1    Attack #1: Low Entropy Attack

Low entropy refers to a state of limited randomness or unpredictability in data. In the context of cryptography, low entropy indicates that the information is relatively easy to guess or predict, thus being more vulnerable to attacks. This lack of randomness reduces the strength of seed phrases. Malfunctioning hardware wallets could potentially result in generating seed phrases with low entropy.

There are several mechanisms to generate strings with low entropy. One of these mechanisms consists in generating binary strings containing significantly more zeros than ones, or vice versa, e.g., "10001001000010000000010000001000" for a 32-bit representation (for the sake of simplicity). Errors during randomness exploitation in hardware wallets (e.g., repetitive character sequences or pattern-based sequences) lead to biased combinations of zeros and ones. A low entropy source provides a low entropy string generator. From these binary strings, the corresponding seed phrases are computed using BIP39. Then, public keys are obtained after deriving the seed phrases following BIP32 and BIP44, as hardware wallets do.

**Algorithm 1** Generation of low entropy strings

**Require:**
$N$ = number of bits (typically, $N = \{128, 160, 192, 224, 256\}$).
$p$ = maximum number of ones ($p < N$).
**Ensure:**
wordList = list of all combinations of low entropy strings.
1: **function** STRINGGENERATOR($N, p$)
2:     wordList ← empty list
3:     **for** countOnes **from 1 to** $N$ **do**
4:         **for** i **from 0 to** $p$ **do**
5:             binary ← 0;
6:             **for** j **from 0 to** countOnes - 1 **do**
7:                 position ← (i + j) % N
8:                 binary ← binary OR (1 << position)
9:                 wordList.add(binary)
10:             **end for**
11:         **end for**
12:     **end for**
13:     **return** wordList
14: **end function**

### 3.2    Attack #2: Uncommon Derivation Path Attack

Derivation paths are commonly generated using the BIP standards described in Sect. 2.3. However, derivation paths could be generated without following any standard, i.e., customising the indexes of the derivation path.

This attack builds upon the strings with the lowest entropy obtained from the previous attack. The lower the strings' entropy, the more possibilities that might correspond to someone's public addresses. So, after generating the corresponding seed phrases from these strings, it is worth applying different kinds of derivations (i.e., considering uncommon derivation paths too).

In particular, two approaches have been followed to generate uncommon derivation paths. The first approach creates custom derivation paths by building all the combinations using the most frequent index values. As observed, these values are 0, 0', 1, 1', e.g., m/0/1/0, m/0'/1/0' and m/0'/1/0/0. This way, a new set of public keys is obtained. The second approach analyses how derivation paths are implemented in popular programming libraries, such as Python's *bip32utils*. With this approach, their funds could be compromised if users have started from strings with very low entropy for their key generation when using these libraries.

### 3.3    Attack #3: Public Domain Seed Phrases Attack

This attack leverages the functionality of hardware wallets to reset an existing wallet by inputting a previously generated seed phrase. Also, this functionality can be used by users to create a new wallet by deciding their seed phrase. Unfortunately, hardware wallets do not have any mechanism to alert users that their seed phrases are weak, e.g., publicly available on the Internet. A strategy to discover weak seed phrases involves seeking the least common words from the BIP39's wordlist on the Internet, e.g., using tools like Pastebin or PrivateBin. If seed phrases are found, public keys could be derived from them.

(a) Wallet with funds in the past        (b) Wallet with funds

**Fig. 3.** Wallets discovered

## 4 Initial Preliminary Results

Some preliminary tests have been conducted to prove the attacks' feasibility. Access has been gained to several wallets with transactions involving different cryptocurrencies (see Fig. 3). Also, although numerous empty wallets have been found, they had funds for a long time. If they receive funds from now on, they could be compromised. The most significant findings are listed as follows:

- By using Attack #1, a Ripple (XRP) wallet with 490.64932 XRP ($337 worth at this time) was discovered. This wallet is active with several weekly transactions.
- By using Attack #2, by exploiting the *bip32utils* library, some wallets with no current funds were found. However, one had 0.003 BTC ($87 worth) for three months.
- The highest success rate was achieved applying Attack #3. A Tron (TRX) wallet with 197 USDT ($197 worth) was discovered. This wallet is active, with many transactions every week.

## 5 Conclusions

Hardware wallets are essential to safeguard cryptocurrencies securely and reliably. However, the resource-constrained nature of these devices opens the door to numerous vulnerabilities. In this article, we have described three attacks against these devices: attacks against their low entropy, attacks against uncommon derivation paths, and attacks against public domain seed phrases. Some wallets with recent activity, and even with funds, have been discovered. Further work will concentrate on analysing the electronic components of hardware wallets during the generation of the private key (e.g., side-channel attacks) and evaluating the functioning of these components under stressful conditions, such as extreme temperatures.

# References

1. Fang F et al (2022) Cryptocurrency trading: a comprehensive survey. Financ Innov 8(1):1–59
2. Casino F et al (2019) A systematic literature review of blockchain-based applications: current status, classification and open issues. Telemat Inform 36:55–81
3. Arapinis M et al (2019) A formal treatment of hardware wallets. In: 23rd international conferences financial cryptography and data security, pp 426–445. Springer
4. Dabrowski A et al (2021) Better keep cash in your boots-hardware wallets are the new single point of failure. In: Proceedings of ACM CCS workshop on decentralized finance and security, pp 1–8
5. Almutairi E et al (2019) Usability and security analysis of the keepkey wallet. In: IEEE international conferences blockchain and cryptocurrency. IEEE, pp 149–153

# Real-Time Implementation of Tiny Machine Learning Models for Hand Motion Classification

Razan Khalife[1], Rawan Mrad[1], Ali Dabbous[2], and Ali Ibrahim[1]([✉])

[1] Lebanese International University, Beirut, Lebanon
ali.ibrahim@liu.edu.lb
[2] University of Genoa, Genoa, Italy

**Abstract.** This paper investigates the design, implementation, and assessment of different embedded neural network models for hand motion classification. Using an Inertial Measurement Unit (IMU) sensor, a dataset of five distinct hand motion classes has been collected from the Arduino nano BLE 33 targeting hand motion analysis. Two different machine learning models namely Convolutional Neural Network (CNN) and Multilayer Perceptron (MLP) are implemented and compared in terms of classification accuracy and hardware complexity. Experimental results show identical classification accuracy of 100% for the two models. After deployment, results show that the MLP network exhibited the lowest processing time and RAM usage, taking 4 ms and 4.0 K, respectively. In contrast, the CNN model needs 313ms and utilizes 15.9 K of RAM. Moreover, the MLP model also had the lowest flash usage with 25.3 K, while the CNN model used 67.7 K.

**Keywords:** TinyML · Embedded Neural Networks · Hand motion recognition · Convolutional Neural Network · Multilayer Perceptron

## 1 Introduction

The human hand is one of the most significant organs in the body. As the hand is essential to artistic expression and syntactical communication, the loss of hand is a destructive experience demanding significant psychological support and physical recovery. To facilitate the gesture recognition and accurately perceive the user motions, an Inertial Measurement Unit (IMU) sensor may be used to determine the motion, orientation, position and acceleration of a hand. Recently, Machine learning ML methods ensured their powerful approach to classify hand motions and gestures, enabling the interpretation of human movements. This technology has significant applications in various fields, including human-computer interaction, virtual reality, robotics, and healthcare. In order to classify five different classes of hand motions namely extension, flexion, pronation, supination and ulnar deviation. Two different ML models, MLP and CNN were implemented. The proposed networks can learn an appropriate connectivity patterns for classifying and distinguishing five different hand motion classes.

## 2   State of the Art

In recent years, different works have addressed hand motion classification using inertial measurement unit sensors. IMU sensors are composed of accelerometers, gyroscopes, and magnetometers which measure linear acceleration, angular velocity, and magnetic field strength respectively. These sensors capture information about hand movement, position, and torque [1]. In [2], a soft tactile sensing system is proposed to detect the location of the contact points on the surface of the sensor. The system is composed of silicon rubber with an IMU embedded on a deformable silicon-based surface. The reported results showed a mean accuracy superior to 80%, and peak accuracy of 97.97% (for the single point contact locations) and 97.54% (for the linear regions of contact points). An improved tactile probe that employs an accelerometer as a transducer in [3] was introduced to gather information related to the mechanical properties of object surfaces. Two methods were used: (1) supervised learning, where surface identification performance of the tactile was quantified and classification was done using Support Vector Machines (SVM); and (2) unsupervised learning classification where the Dirichlet Process Mixture Model was used to estimate the number of clusters, based on a prior distribution. A surface recognition rate of 96.7% was achieved with 1 s of data. In [4], a new sensor glove allows separate measurements of proximal and distal finger joint motions as well as position/orientation detection with an IMU. The whole system is controlled and monitored by an Arduino micro board and the software is written on Arduino IDE software.

Authors in [5] utilized Convolutional Neural Networks (CNN) to process the signal estimating ground slope and tilting of the foot. Experimental results provided that the sensor obtained complete information solely from contact deformation of foot skin of legged robots for locomotion tasks. In [6], authors described the design and implementation of an end-to-end Tactile Cyber Physical System. For the human operator end, they designed a tactile glove having IMU sensors for motion capture and vibrotactile actuators for tactile feedback. For tracking hand motion, authors used five MEMS based IMUs that have 9-axis degrees of freedom. The measurements were done by keeping the glove stationary and 3000 samples were collected.

## 3   Methodology and Experimental Setup

### 3.1   System Description

The block diagram of the proposed system is given by Fig. 1. The system is composed an IMU sensor (LSM9DS1) that determines the orientation, position, and acceleration of the hand. The sensor is integrated in the Arduino Nano BLE sense forming the electronic interface. Data from IMU sensor is collected synchronized with respect to the three sensors i.e., accelerometer, gyroscope, and magnetometer through Ardo spreadsheet library used with the Arduino IDE tool. Then, the dataset is organized to be fed to the neural network models for training and validation. The deployment process take place when the trained model provides an acceptable validation accuracy. The deployment process involves converting the trained model to TensorFlow lite before it could be integrated with the inference code that performs the on-board real-time classification of the hand motions.

**Fig. 1.** System block diagram.

## 3.2 Sensors and Dataset

The collected dataset consists of five hand motion namely: extension, flexion, pronation, supination, and ulnar deviation. The time span to complete each motion was set to be 2 s, and an overall dataset is of 250 samples. The Arduino Nano board was placed in the palm of the hand in different positions (orientations) taken randomly, and secured by a rubber band. The motions of the hand were carried out accordingly, as illustrated in Fig. 2. It is important to note that, during training and inference, if the user hold the board in the left or right hand and perform a specific motion, the model will predict it accurately. However, if the board is with the wrong orientation while performing the motion, the predictions might not be completely accurate.



**Fig. 2.** Five different hand gestures.

## 3.3 Tested Models

**Multilayer Perceptron**. To Classify Five Different Hand Motion Classes, an MLP Neural Network Was Constructed Using a Flatten Layer that Converts the Input Data (40,9) into a 1D Array, Two Dense Layers with 10 Units and ReLU Activation Function. A Final Output Layer with Five Units (One Per Each Activity) and SoftMax Activation Function is Used, as Illustrated in Fig. 3. This Model is Trained on the Training Set for 20 Epochs with a Batch Size of 5.

**Convolutional Neural Network**. The CNN Model Has Seven Layers, Two Conv1D Layers Where the Input Shape, Activation Function, Kernel Size is Specified. MaxPooling1D Layer that Applies the Max Pooling Operation Along the Time Dimension of the Input, Which is Assumed to Be a One-Dimensional Time Series Signal, a Flatten Layer and a Dense Layer with 15 Units and a ReLU Activation Function. A Final Output Layer with Five Units (One Per Each Activity) and SoftMax Activation Function is Used, as

**Fig. 3.** MLP neural network for hand motion classification.



**Fig. 4.** CNN neural network for hand motion classification.

Shown in Fig. 4. This Model is Trained on the Training Set for 20 Epochs with a Batch Size of 5.

The purpose of incorporating MLP and CNN models is because they provided the best results among all other tested models. So, this work limits the comparative analysis on these two models providing their performance while dealing with hand motions classification. To assess the performance of the proposed models, a K fold cross-validation with 10 folds was applied. The dataset is divided into 10 equal-sized folds, and shuffled before splitting, as shown in Fig. 5.

Finally, the Keras model is converted into a TensorFlow Lite (TFLite) model and then the hex data of the TensorFlow Lite model is converted into an array for C programming. The resulting C file is then utilized to deploy the TensorFlow Lite model to Arduino, enabling real-time predictions.

## 3.4 Deployment

The deployment process involves transferring the converted model to the Arduino board and integrating it with the code that runs on the board. This code will include instructions for loading the model, preprocessing input data, and making predictions using the model inference capabilities. The C model is deployed on the Arduino Nano BLE Sense After

**Fig. 5.** K-Fold cross validation technique.

performing the classification based on the weights extracted from the training, the code prints the predictions determined by the classification process. These predictions are typically displayed along with their corresponding confidence values.

## 4 Results

### 4.1 Multilayer Perceptron Model

After running the MLP model and completing 20 epochs, the model achieved a test loss of 0.0028 and a test accuracy of 100%. After deploying the model on Edge Impulse which is an end to end platform that facilitates the development and deployment of embedded ML models for edge devices, it was observed that the MLP model had a processing time of 4ms. In terms of resource usage, the MLP model exhibited the lowest RAM and flash usage, with values of 4.0K and 25.3K respectively. To perform online classification, positioning the Arduino in the palm of the hand, specific motions such as extension, flexion, etc., are performed.

### 4.2 Convolutional Neural Network

Regarding the CNN model, the evaluation demonstrated a test loss of 0.036 and a test accuracy of 100% after accomplishing 20 epochs. Using Edge Impulse platform, it was observed that the CNN model required a long processing time of 313ms. In terms of resource usage, the CNN model utilized 15.9.2K of RAM and 67.7K of flash memory. To perform online classification, the Arduino was placed in the palm of the hand and specific motions are accomplished. The CNN model accurately predicts each motion, determining the one with the highest probability.

A comparison of MLP and CNN models is presented in Table 1. The evaluation is conducted using K-fold cross validation and focuses on accuracy, processing time, RAM usage, and flash usage.

**Table 1.** Comparison of the embedded implementation results for the presented models.

| Model | Accuracy (%) | RAM usage (K) | Flash usage (K) | Processing time (ms) |
|-------|--------------|---------------|-----------------|----------------------|
| CNN   | 100          | 15.9          | 67.7            | 313                  |
| MLP   | 100          | 4.0           | 25.3            | 4                    |

## 5   Conclusions

This study investigates the development of a real-time TinyML system for hand motion recognition. The design, implementation, and evaluation of two different models for classifying hand motions were presented. The comparative analysis when performing online classification using the MLP and CNN models has shown the effectiveness of the MLP over the CNN in terms of processing time and memory usage.

## References

1. Osborn LE, Iskarous MM, Thakor NV (2020) Sensing and control for prosthetic hands in clinical and research applications. In: Wearable robotics. Academic Press
2. Iwamoto Y, Meattini R, Chiaravalli D, Palli G, Shibuya K, Melchiorri C (2020) A low cost tactile sensor for large surfaces based on deformable skin with embedded IMU. In: 2020 IEEE Conference on Industrial Cyberphysical Systems (ICPS), vol 1, pp 501–506
3. Dallaire P, Emond D, Giguere P, Chaib-Draa B (2011) Artificial tactile perception for surface identification using a triple axis accelerometer probe. In: 2011 IEEE International Symposium on Robotic and Sensors Environments (ROSE)
4. Weber P, Rueckert E, Calandra R, Peters J, Beckerle P (2016) A low-cost sensor glove with vibrotactile feedback and multiple finger joint and hand motion sensing for human-robot interaction. In: 2016 25th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN), New York, NY, USA, pp 99–104. https://doi.org/10.1109/ROMAN.2016.7745096
5. Zhang G, Du Y, Zhang Y, Wang MY (2021) A tactile sensing foot for single robot leg stabilization. In: 2021 IEEE International Conference on Robotics and Automation (ICRA)
6. Arjun N, Ashwin SM, Polachan K, Prabhakar TV, Singh C (2018) An end to end tactile cyber physical system design. In: 2018 4th International Workshop on Emerging Ideas and Trends in the Engineering of Cyber-Physical Systems (EITEC), Porto, Portugal, pp 9–16

# Detecting Patient Readiness for Colonoscopy Through Bowel Image Analysis: A Machine Learning Approach

Nour Kaouk[1(✉)], Lamis Amer[2], Tina Yaacoub[1], Youssef Bakouny[1], Chantal Hajjar[1], Flavia Khatounian[1], Joseph Amara[1,2], Rita Slim[2], Ali Mansour[3] , and Cesar Yaghi[1,2]

[1] Saint Joseph University of Beirut, Beirut, Lebanon
`nour.kaouk1@net.usj.edu.lb`
[2] Gastroenterology Service, Hotel Dieu de France, Beirut, Lebanon
[3] Lab-STICC, UMR 6285, CNRS, ENSTA Bretagne, 29806 Brest, France

**Abstract.** Bowel preparation is a crucial step in ensuring the success and accuracy of colonoscopy procedures. Adequate bowel cleansing allows for better visualization and detection of abnormalities within the colon. In this study, we present an AI tool developed to assess the quality of bowel preparation in colonoscopy procedures. The dataset used in this study consists of 350 images of toilet bowls obtained from patients at the hospital "Hôtel Dieu de France" in Beirut, Lebanon. Their images are labeled by the professionals using the Boston scores. Our methodology involves a comprehensive pre-processing phase, encompassing detection, cropping, color adjustment, and Principal Component Analysis (PCA) on the image dataset. Subsequently, we applied different machine learning (ML) models for classification, achieving a high accuracy of 92% with Gradient Boosting. This AI-based approach exhibits great potential in enhancing the efficiency and reliability of colonoscopy evaluations, ultimately leading to improved patient outcomes and early detection of gastrointestinal disorders.

**Keywords:** Colonoscopy · Bowel preparation · Artificial intelligence · Classification · Machine learning

## 1 Introduction

Colonoscopy is of utmost importance for the screening and diagnosis of colorectal cancer, which is the third-leading cause of cancer death. The available evidence suggests that inadequate bowel preparation reduces the diagnostic yield of colorectal neoplasia and increases post-colonoscopy colorectal cancer risk [1]. Suboptimal bowel preparation has been shown to prolong the overall procedure time (e.g. increased time washing, and suctioning debris, prolonged withdrawal time), decrease the cecal intubation rate, and increase the risk of missing polyps or adenomas during the colonoscopy. It furthermore entails a shorter colonoscopy follow-up interval with a higher economic burden, therefore resulting in shorter surveillance intervals and increased costs [2]. Quality of bowel preparation has been associated with adenoma detection rates (ADR) in multiple

studies [3–5]. Similarly, a prospective multicenter randomized control trial published in 2022 included 413 patients, it reveals an ADR of 45.3% on repeat colonoscopy (on a median interval with the index colonoscopy of 28 days), an advanced ADR of 10.9%, and a serrated polyp detection rate of 14.3% [6]. A meta-analysis showed that, compared with low-quality bowel preparations, ADRs were significantly greater with intermediate (odds ratio [OR], 1.4; 95% CI, 1.1−1.8) and high-quality bowel preparations (OR, 1.4; 95% CI, 1.2−1.6) [3].

Usually, physicians rely on patients' subjective observations and descriptions to be informed about bowel preparation. Information obtained include the adherence of patients to colon cleansing guidance, and the quality of the stools in terms of color, translucency, and presence of particles. This approach often falls short in practice mainly because subjective self-assessment by patients regarding the quality of bowel preparation can be inconsistent and unreliable.

To assess the quality of bowel preparation through an advanced AI CNN model, Yan-Xing Hu et al. study [7] primarily focused on the medical approach, showcasing the potential implications and benefits of their innovative model, without delving into the technical aspects of their methodology. Previous articles in this field have predominantly emphasized the development of segmentation and classification algorithms for the detection of polyps in colonoscopy [8–10]. However, they did not explicitly address the significance of bowel preparation in the context of the colonoscopy procedure, which can be considered a potential drawback in their approaches.

To address this challenge, the proposal is to employ AI tools aimed at objectively assessing bowel preparation. This technology would allow to reduce the need for repeat procedures, saving both time and the fatigue experienced by patients. Our article aims to provide a comprehensive overview that not only highlights the medical significance of bowel preparation but also delves into the intricacies of our machine-learning model.

The rest of the paper is organized as follows to provide a clear structure and comprehensive understanding of the research. Section 2 highlights the materials and methods used to conduct this study. Section 3 focuses on the data pre-processing techniques utilized to mitigate variations in the dataset. In Sect. 4, the methodology employed for classification is described in detail. Finally, the conclusion is provided in Sect. 5, discussing the implications of the research in the field of classification.

## 2 Material and Methods

### 2.1 Image Collection and Dataset

We used a comprehensive dataset consisting of 350 images of toilet bowls obtained from patients at the hospital "Hôtel Dieu de France" in Lebanon. Patients were instructed to take pictures of the toilet bowl in the second half of their colon cleansing procedure. These images were anonymized and meticulously collected and curated to encompass a wide range of bowel preparation conditions. Figure 1 illustrates some pictures of the original dataset used in this study.

**Fig. 1.** Some figures of the original dataset before and after processing.

### 2.2 Bowel Preparation

A colonoscopy preparation begins by adjusting one's diet a few days ahead of his colonoscopy. Typically, the patient eats a low-fiber diet for two or three days, followed by a clear liquid diet on the last day. The afternoon or evening before the colonoscopy, he takes a laxative formula to purge the bowels. Studies have found that inadequate bowel preparation can lead to failed detection of cancerous lesions and are associated with an increased risk of procedural adverse events [11].

### 2.3 Colon Cleansing Classification

In the article, we have employed a comprehensive approach by incorporating one distinct scoring system namely the Boston scores. The Boston Bowel Preparation Scale (BBPS) is a medical scoring system used to evaluate the quality of a patient's bowel preparation before a colonoscopy. It assesses how effectively the bowel has been cleansed, with higher scores indicating better bowel cleanliness. Each segment of the colon, the right colon, transverse colon, and left colon are assigned points from 0 to 3 with regard to the cleanliness of the colon. A score of 0 includes an unprepared colon, 1 includes those in which only a portion of the mucosa of the colon segment is visible, and 2 includes those with a minor amount of residual staining and small fragments of stool present. Lastly, 3 includes those where the entire mucosa of the colon is seen well with no residual stool. The entire colon is assigned a cumulative score [12].

The dataset contains pictures of toilet bowls, which have been linked to the patients' colonoscopy tests with preparations classified into 10 different classes based on the Boston score. For classes 0, 1, 2, and 5 we have 0 images. For classes 3, 4, 6, 7, 8, and 9 we had 1, 3, 13, 25, 74, and 235 images, respectively We split the dataset into two categories, namely 1 (good bowel preparation) for a Boston score of 9 and 0 (poor bowel preparation) corresponding to the other scores. The binarization of Boston scores helps doctors and medical professionals quickly identify the quality of preparation. Our team of doctors has validated the importance of this categorization in the medical context.

## 3 Pre-Processing and Data Preparation

The data pre-processing involved several techniques aimed at reducing the variations caused by different patient photographs taken with various camera phones, lighting conditions, and orientations.

Firstly, template matching approach to detect and crop the region of interest in the image dataset was applied to focus solely on the toilet and bowl area by computing

the maximum correlation value and its index. If the maximum correlation exceeds the empirical threshold of 0.6 the region of interest is cut off from the image.

Additionally, Histogram Stretching and Hamming's filter were employed to normalize the color profiles across the images, minimizing the impact of varying lighting conditions. The first algorithm expands or contracts the intensity values of the image to utilize the full dynamic range, enhancing the contrast and adjusting the color levels. The second method minimizes frequency leakage and distortion caused by sharp edges in the image. A filter size of 12 was chosen to be large enough to capture more details, smooth the image more, and reduce noise.

Finally, principal component analysis was used to orient the bowls in the same direction then, a rotation was applied with the calculated angle of rotation for each frame. The images were rotated by 45°. However, regardless of pictures orientation, the method consistently yields improved results.

Figure 1 also displays some images after the completion of all pre-processing steps carried out using MATLAB and Python.

## 4   Classification

Several pre-trained models, such as VGG16 and ResNet50, were tested to explore the potential use of deep learning in our specific use case. Unfortunately, a common challenge was encountered across these models: overfitting. The validation accuracy appeared to plateau, which is a known symptom of overfitting. This issue can be attributed to the relatively small sample size we have to work with, consisting of only 350 images that need to be divided into training and testing sets.

Machine learning (ML) algorithms were thus investigated. In such cases, where binary classification is appropriate, the Support Vector Machine (SVM) algorithm emerges as a good choice, as it succeeds at finding optimal decision boundaries between two classes. First, the categories are defined, and empty lists are initialized to store input and output arrays. Next, the images are resized to a $32 \times 32$ pixel size, and then the pixels are flattened into a 1D feature vector. This feature vector is used as the input to the SVM classifier. The model is trained using the training data that contains 80% of the data and its performance is evaluated on the testing data that comprises the remaining 20% of the data. The SVM classifier uses the default hyperparameter settings: The C (Regularization Parameter) is in its default value of 1. This parameter controls the trade-off between maximizing the margin and minimizing the classification error on the training data. The Kernel Type is used in its default value which is the radial basis function (RBF) kernel, which is commonly used for SVM classification. The Gamma (for RBF Kernel) parameter is 'scale', which is based on the inverse of the number of features. It controls the shape of the decision boundary.

To optimize the SVM's performance, we have considered tuning the hyperparameters by explicitly setting them and performing a hyperparameter search to find the best combination of hyperparameters for our specific dataset. To this purpose, we have used Grid Search technique in combination with cross-validation. We have obtained the best hyperparameters: {'C': 10, 'gamma': 'scale', 'kernel': 'rbf'} and achieved accuracy with best hyperparameters: 90.77%.

In our quest to find the most suitable machine learning model for our classification task, we tested four different models: Random Forest, Logistic Regression, Decision Tree, and Gradient Boosting. The Random Forest model and Logistic Regression model achieved an accuracy of approximately 87.7%, demonstrating good performance. The Decision Tree model achieved an accuracy of approximately 78%. While Decision Trees are interpretable, in this case, they had a lower accuracy compared to the other models. The Gradient Boosting model achieved the highest accuracy of approximately 92.3%.

## 5  Conclusion

The goal of this paper is to detect the colonoscopy preparation degree of the patient through image analysis, so after employing various pre-processing techniques to prepare our image dataset, we meticulously trained and tested our model and classified the images, achieving an accuracy of 92.3% with the Gradient Boosting model. This automated detection of bowel preparation quality addresses a medical need to help patients and physicians obtain the best quality during colonoscopy. Consequently, this innovative model will shorten the follow-up interval and reduce the economic burden, then, ensure a better medical service for the patients.

In the ongoing development of the bowel preparation assessment ML model, several areas for future work have been identified. First, there is the potential for the model to evolve into a clinical decision support system, aiding healthcare providers in real-time decision-making during colonoscopy procedures. Additionally, efforts may focus on enhancing patient education and engagement through user-friendly interfaces or mobile applications, which can positively impact patient compliance and preparation quality.

It is crucial to emphasize that the sample size of our experiment is limited, and the dataset contains a small number of images. This limitation arises from the fact that these images are obtained from real patients within a short time frame. As part of our ongoing efforts, we are committed to augmenting the number of images in our dataset over time.

## References

1. Hernández G, Gimeno-García AZ, Quintero E (2019) Strategies to improve inadequate bowel preparation for colonoscopy
2. Sharma P, Burke CA, Johnson DA, Cash BD (2020) The importance of colonoscopy bowel preparation for the detection of colorectal lesions and colorectal cancer prevention.
3. Clark BT, Rustagi T, Laine L (2014) What level of bowel prep quality requires early repeat colonoscopy: Systematic review and meta-analysis of the impact of preparation quality on adenoma detection rate
4. Anderson JC, Butterly L, Robinson CM (2014) Impact of fair bowel prep on adenoma and serrated polyp detection: Data from the New Hampshire Colonoscopy Registry using a standardized preparation quality rating
5. Noh Hong S, Kyung Sung I, Hwan Kim J (2012) The effect of the bowel preparation status on the risk of missing polyp and adenoma during screening colonoscopy: A tandem colonoscopic study
6. Pantaleón Sánchez M, Gimeno Garcia A.-Z, Bernad Cabredo B (2022) Prevalence of missed lesions in patients with inadequate bowel preparation through a very early repeat colonoscopy

7.  Bor Lu Y, Cun Lu S, Ning Huang Y (2022) A novel convolutional neural network model as an alternative approach to bowel preparation evaluation before colonoscopy in the COVID-19 era: A multicenter, single-blinded, randomized study

8.  Wang L, Chen L, Wang X (2022) Development of a convolutional neural network-based colonoscopy image assessment model for differentiating crohn's disease and ulcerative colitis

9.  Wen Y, Zhang L, Meng X (2022) Rethinking the transfer learning for FCN based polyp segmentation in colonoscopy

10. Wang Y, Feng Z, Song L (2021) Multiclassification of endoscopic colonoscopy images based on deep transfer learning

11. Patel N, Kashyap S, Mori A (2023) Bowel Preparation

12. Lai EJ, Calderwood AH, Doros G (2009) The Boston bowel preparation scale: a valid and reliable instrument for colonoscopy-oriented research

# Machine Learning Model for Fault Detection in Safety Critical System

Pragya Dhungana[1(✉)], Rupesh Kumar Singh[2], and Hariom Dhungana[3]

[1] Nepal Telecommunications Authority, Kathmandu, Nepal
pdhungana@nta.gov.np
[2] Università Degli Studi Di Genova, 16126 Genova, Italy
[3] Western Norway University of Applied Sciences, 5063 Bergen, Norway

**Abstract.** Common bearing failure modes (wear, contamination, corrosion, overload, misalignment, etc.) have unique characteristics, requiring diverse identification and mitigation strategies. No single definition can encompass all contributing factors. Understanding these complexities is crucial for safety critical system to implement fault detection. Machine learning offers a data-driven and intelligent approach to fault detection to improve safety, efficiency, and cost-effectiveness. In this work we propose a supervised machine learning model using naïve bayes classifier for safety critical system fault detection using time domain vibration features. Isolation forest-based anomaly detection is used for labeling faults or healthy condition. The proposed model is tested on the PRONOSTIA dataset, and the model detects fault before their failure criteria in all eleven experiments. The models hold promise for early fault detection in safety-critical systems.

**Keywords:** Machine learning · Health monitoring · Condition monitoring · Biomimetic and Bio-inspired systems

## 1 Introduction

There are basically three maintenance strategies that exist for bearing: (I) Run-to-break involves running until failure occurs, offering maximum uptime. However, it poses safety risks and may lead to costly consequences or prolonged downtime. (II) Usage-based maintenance replaces parts at set intervals, but it can be costly due to unnecessary replacements. (III) Condition-based maintenance overcomes the limitations of the other strategies by relying on reliable monitoring techniques. This strategy emerged to minimize expenses associated with overly frequent scheduled maintenance. It uses computed features from raw measurements, making it easier to track changes and identify damage types, leading to broader successful applications [1].

Vibration analysis, a widely used in condition based maintenance, has evolved from the observation that excessive vibration often signifies bearing defects. A defect in bearing is an imperfection that can be shown to cause failure by a quantitative analysis and that would not have occurred in the absence of the imperfection [2].

Machine learning has greatly improved industrial inspection quality but faces challenges when relying on single features for fault detection. Accurate predictions depend on effective bearing health indicators [3]. Typically, health indicators are constructed by combining appropriate statistical features extracted from vibration signals. However, there are two common shortcomings associated with existing bearing health indicators. First, not all statistical features carry equal weight in the construction of these indicators due to varying ranges [4]. Second, determining a failure threshold becomes challenging as health indicators differ across machines at the point of failure. Handling noisy sensor data, constrained computing power, and limited processing time can be likened to a bounded rationality problem in fault detection [5].

In the past, research focused solely on prognostics, neglecting safety, and financial considerations regarding repair/replacement time. This study aims to optimize the time frame to ensure both adequate maintenance and maximum equipment lifespan. This model helps operators prevent failures and use equipment efficiently, particularly in critical systems like aerospace and healthcare. Early fault detection before the meantime to repair/replacement is crucial for safety. When assessing machine replacement or repair, it's essential to consider factors beyond operating costs, including downtime, depreciation, opportunities, and customer goodwill [6].

The rest of the paper is structured as follows: Sect. 2 presents the methodology of fault identification, including datasets, data preprocessing steps, data labeling technique, fault detection model. Section 3 provides the result and Sect. 4 draws the concluded remarks.

## 2  Methodology

In this section, the steps of fault detection process from the vibration signal processing from the theoretical issues based on literature review are presented. Figure 1 shows the flowchart of proposed fault detection model before failure, to provide sufficient time to repair/replace bearing. The model is composed of three modules, data preprocessing, data labeling and fault detection.

This model's core is machine learning, which revolves around learning from sensor data and continuously improving over time. It heavily relies on probability and statistics to make informed decisions and predictions. It surpasses standard statistical methods in decision-making process [7]. By uncovering hidden relationships and correlations not apparent using simple rules, machine learning algorithms can identify important features or combinations of features, contributing to accurate predictions even without explicit predefined relationships.

### 2.1  Data Sets

In the evaluation of the model, we are searching for condition monitoring datasets from multiple sensors. We found PRONOSTIA, bearings' accelerated life tests provided by FEMTO-ST Institute. Two accelerometers were used to monitor vibration signals, sampled every 10 s with a frequency of 25.6 kHz. Seventeen experiments were conducted under three different operating conditions, as summarized in Table 1. Six experiments (Bearing1_1, Bearing1_2, Bearing2_1, Bearing2_2, Bearing3_1, Bearing3_2) are used

**Fig. 1.** Flowchart of the fault detection model.

for training and the remaining eleven experiments are used for testing. Detailed information is presented in [8]. Although a common practice is to split training and testing data as 70 and 30%, we maintain the configuration provided for the competition, ensuring consistency with the evaluation setup using full test set measurements. The datasets are released on an open-source basis (https://www.nasa.gov/content/prognostics-center-of-excellence-data-set-repository), with a goal to support prognostics research for condition monitoring.

**Table 1.** Operating conditions of various experiments.

|  | Operating conditions | | |
|---|---|---|---|
|  | Condition 1 | Condition 2 | Condition 3 |
| Load (Newton) | 4000 | 4200 | 5000 |
| Rotation speed (RPM) | 1800 | 1650 | 1500 |
| Training sets | Bearing1_1 Bearing1_2 | Bearing2_1 Bearing2_2 | Bearing3_1 Bearing3_2 |
| Testing sets | Bearing1_3 Bearing1_4 Bearing1_5 Bearing1_6 Bearing1_7 | Bearing2_3 Bearing2_4 Bearing2_5 Bearing2_6 Bearing2_7 | Bearing3_3 |

Failure refers to the state or condition in which a system, process, or component does not perform its intended function or meet its expected requirements. Safety critical system has strict requirement of performances; therefore, they fix their failure criteria own based on their requirements. Failure criteria for PRONOSTIA experiment was the amplitude of the vibration signal overpassed 20 g; and the datasets contain the measurement until the either of the accelerometer reaches the threshold.

## 2.2  Data Preprocessing

Vibration data contains random noise that can hide meaningful information. We use 5 point moving average filter for data smoothing to reveal underlying patterns, trends, and anomalies. The feature extraction process substantially reduces raw vibration data dimensions while retaining crucial bearing health information. Employing multiple features, as opposed to just one, is advantageous since a single feature might not encompass all degradation-related bearing insights. For each experiment, thirteen-time domain features (Maximum, Minimum, Average absolute value, Peak to peak, Variance, Standard Deviation, Root mean square, Crest factor, Clearance Factor, Impulse factor, Skewness, Kurtosis, and Shape factor) are extracted from both sensor data. Detailed definitions, physical meanings, and statistical equations are described on [9]. Feature selection is vital in data preprocessing, converting extensive computation into informed decisions. It's pivotal for efficiency. Highly correlated features can be represented by just one, while uncorrelated variables might introduce noise and bias in fault detection. The selection algorithms can be classified into four groups: similarity based, information-theoretical-based, sparse-learning-based, and statistical-based methods. In [10], the author discussed evaluation parameters of feature selection on supervised, unsupervised, and semi-supervised machine learning algorithms. In this work we use Pearsons's correlation coefficient measures to select features as shown in Eq. (1), that deals with the relationship between two features. Maximum value, Variance, Skewness, and Kurtosis features are picked as promising features for fault representation.

$$\rho_{X,Y} = \frac{COV(X,Y)}{\sigma_x \sigma_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} \tag{1}$$

## 2.3  Data Labeling

Since the datasets contain the measurement until the either of the accelerometer overpassed 20 g, there is no pre-defined label. We redefine the fault mode for safety critical system based on the anomaly. Anomalies are unique data patterns with distinct characteristics from normal instances. Detecting anomalies is vital across various domains, offering critical actionable insights. In this work, the Isolation Forest algorithm is utilized for data labeling due to its fast tree-based approach, which assigns anomaly scores based on binary search tree path lengths for each data point [11]. It can handle multiple features to enhance accuracy of data labeling and leverages contextual insights from various monitoring approaches [12]. The Isolation Forest is particularly effective in

high-dimensional problems with many irrelevant attributes and scenarios where anomalies are rare, or the training set lacks anomalies. Its computational efficiency makes it suitable for streaming data applications.

The figure in Table 2 shows the central tendency of decision score in anomaly detection, that reveals the distribution and characteristics and helps to identify potential outliers. To get the clear decision boundary between normal and anomaly we use two parameters anomaly and decision score. We label the training data based on the rule presented on the right side of table.

**Table 2.** Table captions should be placed above the tables.

|  | Criteria | Label |
|---|---|---|
| Box plot of decision score (Bearing1_1) that shows the distribution of decision score in anomaly detection. We use the lower quartile value as threshold for data labeling | Anomaly not found | Normal |
|  | Anomaly found & Anomaly score $< \; =$ Threshold | Normal |
|  | Anomaly found & Anomaly score $>$ Threshold | Anomaly |

After labeling the training data, the pairwise relationship between time domain vibration features is observed to identify patterns, that indicates the health status of bearing either normal or anomaly. Figure 2 shows the scatter plots to visualize these pairwise relationships between selected eight features. The blue points show the normal label, and the orange points show the anomaly label for training model.

## 2.4 Fault Detection by Classification

While recent classifiers like ROCKET (RandOm Convolutional KErnel Transform) offer high accuracy and rapid training [13], our focus is on interpretability rather than black box model. Thus, we opt for the Naive Bayes classifier in this study. It has shown effectiveness in various practical applications, such as text classification, medical diagnosis, predictive maintenance, and fault detection [14]. The classifier simplifies learning by assuming features are independent of a given class. While this assumption may not hold true in many cases, naive Bayes often performs well in practice, even when compared to more complex classifiers. Bayesian theory is used in Naive Bayes classifier as presented in Eq. (2). Binary classification is used for fault detection. A model can be trained from

**Fig. 2.** Pairwise relationship between eight vibration features from two sensor, the blue dots represent the normal conditions and orange dots represent the anomaly.

labeled data obtained from isolation forest-based anomaly from vibration features.

$$P(D_k|x) = \frac{P(D_k)}{P(x)}P(x|D_k) = \frac{P(D_k)}{P(x)}\prod_{i=1}^{n}P(x_i|D_k) \tag{2}$$

## 3   Results

To assess the effectiveness of the proposed method in the fault detection of rolling element bearings, the detected times are compared with those failure time mentioned by data provider. The fault identification time from the proposed method is displayed in the third column of Table 3. Utilizing data from three consecutive samples to confirm faults reduces the risk of drawing conclusions based on a single sample, ensuring that the findings are representative of the overall situation. This approach enhances generalization, mitigates bias, accommodates outliers, and ultimately bolsters the statistical power of our fault detections.

**Table 3.** Stoppage decision from immediate decision-making model.

| Experiments | Actual failure time (Sec) | Fault detection by ML based model (Sec) |
|---|---|---|
| Bearing1_3 | 23750 | 16430 |
| Bearing1_4 | 11720 | 10900 |
| Bearing1_5 | 24630 | 24510 |
| Bearing1_6 | 24480 | 16360 |
| Bearing1_7 | 22590 | 22140 |
| Bearing2_3 | 19550 | 2580 |
| Bearing2_4 | 7510 | 3440 |
| Bearing2_5 | 23110 | 4030 |
| Bearing2_6 | 7010 | 6890 |
| Bearing2_7 | 2300 | 2240 |
| Bearing3_3 | 4340 | 4250 |

In some test sets, like Bearing 2_7 and Bearing 3_3, fault detection time is close to the actual failure criteria due to short test durations (e.g., 29 min for Bearing 2_7). In the Bearing2_5 experiments, initial jerks and heightened vibrations were observed, which later stabilized into smooth operation. To eliminate false detections, faults occurring before 10 percent of the useful life are disregarded, and only those detected after this point are considered genuine. Despite small variation in different experiments, the model consistently detects faults before the actual failure criteria. Therefore, the model can be considered suitable for early fault detection in safety-critical systems.

For detail illustration of fault detection time we pick on Bearing1_3. The top two graphs in Fig. 3 show the raw vibration data from horizontal and vertical accelerometer and the last graph shows the optimal time for fault detection in safety critical system. The model gives 16410 s or 4.56 (Hours) is optimal time for maintenance. During the transition period, we identified significant variations in the variance and kurtosis values of the vertical vibration sensor by cross-checking the feature data frame.

The quick machine learning based fault detection model is simple. By leveraging this approach, technicians can make informed decisions about the fault going to occur and take appropriate actions for maintenance to improve the overall reliability and performance of the system.

## 4   Conclusion and Future Works

Failure occurs when a bearing does not perform as expected, leading to undesirable performance issues. Safety-critical systems have strict performance requirements and set their own failure criteria. For PRONOSTIA, stopping tests when vibration signal amplitude exceeds 20 g is insufficient for safety-critical systems. In this work our motivation is different than prognostic; the goal of this work is to make fault detection before the system reaches failure criteria, so that we can get the maintenance time to avoid down time. The model detects fault before the failure criteria defined by data provider in all eleven test experiments, therefore this association-based fault detection model can be valuable in making prompt and effective decisions for time-critical industrial equipment.

**Fig. 3.** Fault detection time for Bearing1_3 for safety operation is after four and half hours.

To achieve sufficient maintenance time to avoid downtime in safety-critical systems, this ML based fault detection model can be one alternative. Moreover, its performance can be enhanced by implementing proper noise reduction techniques and utilizing complex features from frequency and time-frequency domains opens new possibilities for improving the model's performance in real-world applications.

# References

1. Cawley P (2018) Structural health monitoring: Closing the gap between research and industrial deployment. Struct Health Monit 17(5):1225–1244
2. Zerbst U, Madia M, Klinger C, Bettge D, Murakami Y (2019) Defects as a root cause of fatigue failure of metallic components. I: Basic aspects. Eng Fail Anal 97:777–792
3. Guo L, Li N, Jia F, Lei Y, Lin J (2017) A recurrent neural network based health indicator for remaining useful life prediction of bearings. Neurocomputing 240:98–109
4. Sait AS, Sharaf-Eldeen YI (2011) A review of gearbox condition monitoring based on vibration analysis techniques diagnostics and prognostics. Conference Proceedings of the Society for Experimental Mechanics Series 5:307–324
5. Simon HA (1947) Administrative behavior: A study of decision making processes in business organization. Macmillan, New York
6. Yung C, Bonnett AH (2004) Repair or replace? IEEE Ind Appl Mag 10(5):48–58
7. Jayatilake SMDAC, Ganegoda GU (2021) Involvement of machine learning tools in healthcare decision making. J Healthc Eng
8. Nectoux P, Gouriveau R, Medjaher K, Ramasso E, Chebel-Morello B, Zerhouni N, Varnier C (2012) PRONOSTIA: An experimental platform for bearings accelerated degradation tests. In: IEEE International conference on prognostics and health management, PHM'12. pp 1–8

9. Wang T, Han Q, Chu F, Feng Z (2019) Vibration based condition monitoring and fault diagnosis of wind turbine planetary gearbox: A review. Mech Syst Signal Process 126:662–685
10. Cai J, Luo J, Wang S, Yang S (2018) Feature selection in machine learning: A new perspective. Neurocomputing 300:70–79
11. Liu FT, Ting KM, Zhou ZH (2008) December. Isolation forest. In: 2008 eighth ieee international conference on data mining. IEEE, pp 413–422
12. Hayes MA, Capretz MAM (2014) Contextual anomaly detection in big sensor data. In: Proc IEEE Int Congr Big Data, pp. 64–71
13. Dempster A, Petitjean F, Webb GI (2020) ROCKET: exceptionally fast and accurate time series classification using random convolutional kernels. Data Min Knowl Disc 34(5):1454–1495
14. Rish I (2001) An empirical study of the naive Bayes classifier. In: IJCAI 2001 workshop on empirical methods in artificial intelligence, vol 3, no 22. pp 41–46

# Modeling and Simulation of Optically Transparent Brain Computer Interfaces

Rayan Kheirldeen[1(✉)], Ali Shaito[2], Houssein Hajj Hassan[3], Ali Cherry[4], Ali Dabbous[5], and Mohamad Hajj-Hassan[1]

[1] Department of Biomedical Engineering, Lebanese International University, Bekaa, Lebanon
  11832116@students.liu.edu.lb, mohamad.hajjhassan@liu.edu.lb
[2] Department of Mechanical Engineering, Lebanese International University, Bekaa, Lebanon
[3] Department of Biological and Chemical Sciences, Lebanese International University, Bekaa, Lebanon
[4] Department of Biomedical Engineering, The International University of Beirut, Beirut, Lebanon
[5] Department Department of Electrical, Electronic, Telecommunication Engineering and Naval Architecture DITEN, University of Genoa, Genova, Italy

**Abstract.** The use of brain computer interfaces (BCIs) is crucial to the development of neural prosthetics and neuroscience. BCIs assist paralyzed patients by enabling them to control computers or robots using their neural activity. Optically transparent BCIs, combined with optical imaging modalities, enable simultaneous gathering of high-resolution electrophysiological signals and imaging of neural activities. In this work, we report the investigation of the mechanical behavior of an ultra-thin glass-based brain machine interface probe when subjected to various types of forces encountered during insertion and engagement into the brain tissues. Results of probe thickness optimization show that electrodes can be 25 µm thick while maintaining a factor of safety against failure > 1.

**Keywords:** Brain machine interfaces · Transparent neural microelectrodes · SolidWorks

## 1 Introduction

Loss of muscle function in a specific area of the body is known as paralysis. Any portion of the body can experience paralysis at any time during a person's lifetime [1]. Due to the damage or injury of the spinal cord, stroke, or a motor neuron disease, patients with these conditions are often referred to as having complete paralysis. Brain-machine interface (BMI) technology has been of particular interest in these cases. This technology can be used in combination with other neuromuscular stimulation techniques to help patients regain their movement [2]. Various technologies are required to create a functional brain-machine interface (BMI). One of these is the neural electrode. Neural electrodes are small structures that are inserted into the brain and act as a link between the electrically active neurons in the brain and external electronic circuits. They have two

main functions: recording the electrical impulses generated by neurons and stimulating particular regions of the brain. These electrodes are vital for developing prostheses and treatments for patients with spinal cord injuries, stroke, sensory impairments, and paralysis by allowing the stimulation of nerve tissue and the recording of neural electrical activity. So, the use of neural microelectrodes is crucial to the development of neural prosthetics and neuroscience. The electrode contains a back support area that carries recording sites to record electrical activity from the neurons, connecting traces, and bond pads for connecting the probe to the electronics [3].

In order to develop specialized treatments for neurological and psychiatric issues as well as to understand brain function, it is essential to understand how the complex neural networks of the human brain function. Recent advances in neuroimaging have aided our understanding by producing high-resolution wiring schematics of the brain. It has been demonstrated that functional optical imaging of brain tissue can reveal a wealth of information on the dynamic characteristics of many cells (more than 100) at once. The advantages of optical imaging and simultaneous electrophysiology in terms of temporal and spatial resolution could be utilized. Recent research has shown that the combination of electrophysiological and optical modalities is faced with additional difficulties and limitations when using opaque electrodes made of common metals, such as gold or platinum. Metal microelectrode arrays, which are often employed to record neural activity, are ineffective for such applications because they obstruct the field of view, create optical shadows, and are more likely to introduce light-induced artifacts into the recordings. This spatial-temporal resolution conundrum can be resolved by fully transparent microelectrodes that permit simultaneous electrophysiology and imaging from the same microcircuit. The high light transmittance rate of transparent microelectrode arrays makes them promising because more physiological signals can be observed when more light passes through the recording electrodes. Electrophysiological recording and image capture could be done simultaneously using transparent microelectrodes. Transparent electrodes, as opposed to metal-based opaque MEAs, eliminate the drawbacks of multimodal systems by allowing them to measure electrical signals from the brain in real-time without optical interference via material and structural controls [3].

## 2 Design

### 2.1 Design Model

The upcoming section describes the design of a proposed transparent neural electrode that can utilize see-through imaging which enables tissue response monitoring during optical or electrical stimulation or electrical recording. This neural electrode is made of glass quartz material and is extremely long, having a length of 1.007 cm. Its shape is tapered, which makes it easier to penetrate. The electrode's geometry is divided into three main parts: the base, the measuring region (which includes the recording sites), and the piercing region. The base measures 250 μm in width, which gradually narrows down to a width of 150 μm when the length reaches 270 μm. This design is helpful in reducing damage to the brain tissue and minimizing displacement. The measuring region is 1 cm long and has a width of 150 μm at the base, gradually decreasing to 50 μm at the other end. After the measuring region, there is the piercing region, which is 50 μm

long and is designed to be 10 μm wide at the end of the probe. After selecting the perfect dimensions and geometry, a 2D electrode model is created and then extruded into 3D where a specific thickness of 50 μm is taken. The 3D electrode design is analyzed using SolidWorks. Then, analytical analysis is used to study and understand the mechanical behavior of the probe when it is subjected to different types of forces during insertion into the brain. Different simulations will be done to predict the behavior of the probes.

Figure 1 shows the relative dimensions of the different electrode regions using SolidWorks.



**Fig. 1.** 3D Model of the electrode

## 2.2   Analytical Modeling of Neural Probes Mechanics

In an effort to anticipate the mechanical performance of the proposed probe, two approaches were utilized—a simulation approach using finite element modeling, and an analytical calculation approach.

The three main forces that affect the electrode probe while it is being handled and inserted into the brain tissue are the bending, buckling, and shear force. Bending, which occurs under an out-of-plane loading that causes parallel displacement to the tissue plane, occurs after penetration. Buckling, which occurs under axial compression and is primarily represented by the force that counteracts the normal force by the probe, occurs directly after penetration. Shear force which prevents the tip of the probe from slipping on the surface of the brain tissue. Shear forces rarely happen because when the probe is inserted, it is affixed to a motion controller that moves solely in one direction toward the brain. So, both bending and buckling forces will be focused on throughout this paper.

The maximum bending force is important since it provides the maximum out-of-plane force at which probe breakage is possible. When moving the probe's base after its tip area has been implanted into the brain tissue, this transverse bending can happen. Once the tip has experienced a bending force, the electrode deflects in a quarter-circle in this region. This causes maximum stress at the fixed bottom of the probe connected to the base holding the bonding pad and an average (normal) stress at the middle region

of the probe. These stresses are given by the following equations [4].

$$\sigma_{\max} = \frac{3Et\delta}{2L^2} \tag{1}$$

$$\sigma_{Avg} = \frac{3Et\delta}{4L^2} \tag{2}$$

The buckling force that results from axial compression is the second type of force acting on the probe. The force at which the probe tip enters the brain tissue is determined by the buckling load, which is a critical parameter. The applied force is assumed to be axial in the case of probe insertion in the brain tissue, and the electrode probes are modeled as rigidly supported on one end (base region) and free to rotate on the other end (tip region). When the axially loaded tapered probe is subjected to an external force, P, greater than the critical load, $P_{cr}$, necessary to buckle the single probe, the probe loses its initial stable shape [4]. Determining the critical load is important, where it should be compared to the maximum anticipated forces applied to the electrode during its insertion into the brain. So, to prevent the probe from breakage, the force should be within the buckling strength range.

$$P_{cr} = \left(\frac{n\pi}{L}\right)^2 EI \tag{3}$$

where n is the buckling mode, L is the length of the probe, and EI is the flexural rigidity of the probe. Moreover, in addition to the buckling failure mode, another failure mode that the electrode can be subjected to during insertion is fracture. For the analysis of brittle materials, the maximum normal stress theory is adopted. Brittle materials tend to fail suddenly by a fracture with no apparent yielding. This theory states that a brittle material will fail when the maximum tensile stress, $\sigma_1$, in the material reaches a value that is equal to the ultimate normal stress [5]. Then, the factor of safety (n) is calculated. By calculation, the factor of safety is determined to be approximately 9.9.

$$|\sigma_1| = \sigma_{ult} \tag{4}$$

$$n = \frac{\sigma_{ten}}{\sigma_1} \tag{5}$$

where $\sigma_{ten}$ is the maximum tensile strength of glass quartz material.

## 3   Results

The following four simulation types are performed: buckling analysis to determine the load factor, static analysis with buckling force, identification of the maximum bending force before breakage, and thickness optimization.

## 3.1   Buckling Analysis

After constructing the 3D model electrode, a custom material is created in SolidWorks (glass quartz) with specified properties (tensile and compression strength, young's modulus, and Poisson's ratio). After that, the base region is selected as a fixed geometry. Knowing that the maximal penetration force for a neural microelectrode is equal to 0.00242 N [4], a pressure of 4.84 MPa was applied to the tip of the electrode (pressure = force/area = 2.42 mN/500 $\mu$m$^2$). Finally, the simulation is executed using mesh analysis as shown in Fig. 2. The load factor obtained from the buckling analysis which is equal to 59 indicates that the applied load is 1/59th of the critical buckling load. This suggests that the structure is stable and has a large margin of safety against buckling failure.



**Fig. 2.**  Computing the load factor

## 3.2   Static Analysis with Buckling Force

The same procedure used for buckling analysis was applied, but the selected study was static. Upon completion of the simulation, the von Mises stress, displacement, and strain values at different regions for the electrode are obtained. The results are shown in Fig. 3a, b and c. The factor of safety is then determined, using the normal stress theory, which yields a value of 9.8 as shown in Fig. 3d, which is very close to the analytical result of 9.9. The close correlation between the simulation and analytical results further validates the accuracy and reliability of the simulation approach in evaluating the design's safety and performance. If the factor of safety is greater than one, then neural electrode design is considered safe.

## 3.3   Identification of the Maximum Bending Force Before Breakage

According to the article "Bending strength of thin fused silica membranes for optical applications" by M. Trunec, V. Kolarik, and P. Bouchal, the bending tensile strength of thin fused silica membranes with a thickness of 50 $\mu$m was measured to be approximately 370 MPa using a four-point bending test [6]. So, depending on the maximum bending tensile strength, the maximum bending force before breakage will be computed. To find this value, the optimization tool on SolidWorks was be used by varying the force until the maximum bending tensile strength is reached. After selecting the optimization study tab, the variable to be optimized is added which is the force applied to the electrode (along the z-axis). A constraint is set by selecting the maximum allowable stress to

**Fig. 3.** Static analysis results with buckling force.

be 370 MPa. As a result, a graph is generated as shown in Fig. 4, which shows the relationship between the force and stress. From the graph, the maximum force that can be reached before exceeding the maximum bending tensile strength is 0.02125 N. This finding is significant as it provides valuable information about the maximum force that can be applied to the electrode without compromising its structural integrity. By staying within this limit, the design ensures the safety and reliability of the neural electrode during its intended application.



**Fig. 4.** Computing the maximum force before breakage

The electrode bending force is much lower than its buckling force (0.00242 N) which is also less than 0.02125 N (maximum force before breakage). This indicates that the electrode is not prone to failure due to bending. However, to ensure its safety, another static study is performed, with pressure applied along the z-axis to represent the bending force. Since the exact bending force is unknown, the value selected was set to be equal to the buckling pressure as an estimate. The analysis reveals a factor of safety of 2.6,

indicating that the design has a substantial safety margin. Moreover, the actual factor of safety is expected to be even higher because in reality the bending load is lower than the buckling. This ensures the electrode's reliability and reduces the risk of failure or damage.

## 3.4  Thickness Optimization

Optimizing the thickness of the neural electrode enhances flexibility, signal quality, and spatial resolution, improving its performance and biocompatibility. This optimization enables more accurate and reliable neural recordings. To do thickness optimization for the neural electrode, bending and buckling force is considered as well as the factor of safety. An optimization tool on SolidWorks is utilized to obtain the optimal thickness for the electrode. In the initial design study focusing on thickness optimization, the variable thickness of the electrode was considered while applying only the buckling force. The constraint was set to maintain a minimum factor of safety greater than three, and the goal was to minimize the mass. The resulting graph, shown in Fig. 5, indicates that the factor of safety remains high, ensuring the safety of the electrode. As the thickness decreases, the factor of safety also decrease based on these results. While the factor of safety provides valuable insights into the electrode's response to buckling forces, it is essential to consider other factors that can influence structural behavior like bending force. By conducting additional design studies, it will be possible to determine the most suitable thickness that ensures a balance between structural integrity and functionality. These investigations will contribute to refining the electrode design and ultimately improve its performance, reliability, and safety during practical applications.



**Fig. 5.**  Thickness optimization when applying buckling force

Another design study is done when applying the bending force. As mentioned before, the force is estimated to be 0.00242 N, but in reality it is lower. The graph in Fig. 6 shows that as the thickness decreases, the factor of safety decreases. Based on the graph, it is observed that for a thickness of 45 μm and higher, the factor of safety exceeds one,

indicating that the electrode design is safe under bending forces. Therefore, optimizing the neural electrode thickness to 45 µm would be a suitable choice.



**Fig. 6.** Thickness optimization when applying bending force

## 4   Conclusion

This paper focused on the investigation of the mechanical behavior of a three-dimensional optically transparent neural microelectrode. The simulation results provided valuable insights into the behavior of the electrode under different loading conditions and showed that the probe possesses favorable mechanical characteristics, indicating its suitability for neural interface applications.

## References

1. Medically reviewed by William Morrison, M.D.— By Heaven Stubblefield, 22 march 2018. [Online]. Available https://www.healthline.com/health/paralysis
2. Lebedev M (2014) What limits the performance of current invasive brain machine interfaces? J Front Syst Neurosci
3. Cho YU, Lim SL, Hong JH, Yu KJ (2022) Transparent neural implantable devices: a comprehensive review of challenges and progress. Flexible Electronics By Numbers
4. HajjHassan M (2008) NeuroMEMS: Neural Probe Microtechnologies. Sensors (Basel)
5. Hibbeler R (2000) Mechanics of material, United States of America: Pearson Prentice Hall, 2014, 2011, 2008, 2005, 2003, 2000
6. Trunec M, Kolarik V, Bouchal P (2021) Bending strength of thin fused silica membranes for optical applications. Optical Material

# Author Index