# HPC and Data: When Two Becomes One

Christophe Calvin[(✉)] and France Boillod-Cerneux

CEA Saclay, DRF, Gif Sur Yvette, France
{christophe.calvin,france.boillod-cerneux}@cea.fr

**Abstract.** As claimed for many years, High Performance Computing (HPC) and high performance numerical simulation are necessary tools for fundamental science and engineering. Big data and artificial intelligence are some newcomers in the landscape, but not that new, especially in science. Finally, open data and open science are becoming now mandatory for trustable and reproducible science.

This paper presents the recent evolution of HPC with the spectacular arising of AI. HPC and AI share at least one common point: Data. Many HPC communities are struggling with data, whether they are coming from simulation and wait to be analyzed, or coming from large instruments (experiments, observatories) and wait to be treated.

Data was not a major focus in the last decades for HPC community but it reshapes HPC paradigms by introducing data as a "scientific pillar".

We will first present the current HPC context and how AI changed the current HPC landscape. We will then focus about data use in HPC and how AI can improve HPC simulations. We will also present the concept of FAIR data and why this concern shall be treated soon and embraced by HPC and AI community. We will finally conclude on the data issue and present our point of view regarding the future evolution of HPC market.

**Keywords:** HPC · AI · Open Science · Data

## 1 Introduction

With exascale challenge, Top500 [1] has really mutated, either regarding the supercomputers listed and regarding its weight to estimate a supercomputer's compute power.

The first mutation comes from the exascale challenge itself. Power wall is hard to cross and so far, (safe) decision has been made to promote low consumption processors or high multicore processors while trying to keep a reasonable power consumption of the overall machine.

Top500 is heckled by new rankings, such as Green500 [3] that proposes to classify supercomputers by Flops per Watts rather than just Flops/s capacities. This rank strategy deeply reshapes the Top500 ranking, as very few of the top 10 coming from Top500 can reach Green500 first ranks. Green500 also emphasizes

Japan's effort to build low energy consumption machines, as 6 over 10 of the best Green500 machines come from Japan.

HPCG500 (High Performance Conjugate Gradients) [4] is either getting more and more attention from HPC community: Indeed, Linpack benchmark for Top500 provides an efficiency from 50% to 80% from peak compute capacity, while real users applications have about 10% efficiency from peak performance. Few of real applications can address the compute bound profile of High Performance Linpack (HPL): In this context, HPCG is raised as a "more representative" benchmark, closer to "real HPC applications" profile. It is not surprising to see that HPCG500 and Top500 present very different ranking, but the difference is less striking than Green500. However, if we compare the HPL peak performance and HPCG peak performance, then the reality of supercomputer compute capacity becomes closer to reality.

Even though the change is unavoidable, it is either not quite smooth. Last petaflopic homogeneous supercomputers propose mostly massive parallel processors, and users must take care of their application parallelism (especially multi-level parallelism) if they really want to see a performance gain on the last petaflopic computers.

On the other side, Top500 is fastly tainted by hybrid architecture (classic processors with accelerators such as Graphic Process Units), which is a direct consequence of exascale power wall. NVIDIA® GPUs strike hard the 2018 Top500: Half of top 10 most powerful supercomputers (November 2018 Top500 list [2]) contains nodes with NVIDIA® GPUs. However, taking full advantage of hybrid nodes efficiency requires large effort for HPC community. Choosing GPUs acceleration is mostly a tradeoff portability versus efficiency (in terms of compute power and energy consumption). This trend is justified by the excellent ratio of Flopss per Watts of NVIDIA® GPUs but also by the ongoing evolution (that is actually not far from a revolution) of HPC landscape.

"Hybrid" word must be used with caution: indeed, we traditionally referred to hybrid architecture as compute nodes with accelerators (such as GPUs), but hybrid may now, depending on the context, refer to the community addressed by the supercomputer. Looking closely at the top 10 of the Top500 November 2018 list, we observe that 2 of them, Summit (USA) and ABCI (Japan) clearly expose their membership to AI research: That is to say, Summit and ABCI are not HPC dedicated but at the best, AI dedicated and at the worst, HPC and AI dedicated. This recent trend is likely to be spread to the other Top500 computers and might contaminate very quickly the full Top500.

The fusion of HPC and AI is clearly illustrated by many HPC centers that are enlarging their scope of action: As an example, Argonne Leadership Computing Facility (ALCF) from Argonne National Laboratory (ANL) has officially announced the integration of Machine and Deep Learning as one of their scientific pillars. Most of the historic and HPC actors made similar declarations.

In what follows we will focus on how AI reshapes the HPC community and can help to consider data as a key point in their simulation process. We will then

focus on the data issue, and address the FAIR data concept, which may lead in the future both HPC and AI scientific research.

## 2  HPC (R)evolution

### 2.1  When AI Shakes HPC

Since the end of 2018, many HPC media relay the announcements of supercomputers dedicated to AI as well as strong positions to promote AI science and development. The recent appearance and success of Deep500 [5], the ranking of most powerful computers dedicated to AI (or HPC and AI) is a very good illustration of Machine and Deep Learning success in industry and scientific fields.

**AI Adoption.** The Deep Learning (DL) success can be attributed to enterprises and industries who have massively used it to improve their process, or wisely address their (future) customers. Deep learning is currently widely adopted by public at large scale. This might be seen as a very "late recognition" as AI discipline was born many decades ago, around 1950. Success of AI and more specifically Deep Learning is largely due to the GAFAM (Google, Apple, Amazon, Facebook and Microsoft) but not only. At a more global scale, in the last decades, every industry, every organization and everyone has gathered massive amount of data. The convergence of a large amount of data and the maturity of available computing hardware drove AI from "theoretical" discipline to "executable" discipline, and currently, "unavoidable" discipline.

**AI and HPC Fusion.** Deep Learning is driven by data and HPC is generating huge amount of data. Therefore, it is not very surprising to observe the collision of these two disciplines. On another level, HPC community is quite mature now and organized. It has access to a large panel of HPC platforms, addressing many different needs in terms of hardware and compute capacity needs. The progression of AI inside the HPC community is coming from four needs:

– AI needs large platforms to run specific experimentations, and HPC can offer such platforms,
– HPC is struggling with data and their analyses,
– Many HPC applications cannot cross petaflopic or exaflopic scale due to their algorithms and/or the physics model and numeric used,
– Governments have identified AI as a strategic point, at least as important as HPC.

**AI Supercomputers Might Cross Exascale First.** Machine Learning and Deep Learning completely reshape exascale challenge thanks to their ability to reduce the compute precision. Though many HPC simulations require a fine precision, Machine Learning and Deep Learning applications can be satisfied with very low precision.

For many years, reduced and/or mixed precision compute kernels have been considered as a strategic key to overcome the exascale challenge. Even if this strategy has been adopted for some applications, we must admit its success remains very limited and did not drastically change the performance as expected: Whether reduced/mixed precision kernels are too small to increase scalability significantly, whether their numerical impacts imply to keep them very limited. For now, reduced or mixed precision in HPC simulations is far from sufficient to cross a new computing scale.

New cards are given with the rise of Machine Learning and Deep Learning and their execution on hybrid processors that support as well classic HPC and Deep Learning simulations. One may see that "AI" dedicated or "mixed" processors and GPU addressing the market are able to deal with low precision compute. This is one of the key points when these processors are presented, showing very attractive compute capacity especially with low precision. It would not be surprising if first exascale application were a Deep Learning simulation running in low precision. In that case, Exascale should not be late for 2020 as it is expected since a few years, and might even be "greener" as expected if very low precision is used.

**Economic Market Drives Technologies.** Deep learning market is more sensitive to the economy than HPC. HPC remains a strategic point for many industries, but HPC use is restricted most of the time to very large enterprises who can afford the infrastructure costs and human resources with a specific knowledge to address large simulations (computer science, algorithms and IT). In contrast, Deep Learning is widely accessible thanks to:

– The amount of data each industry has accumulated over time,
– The "black-box" frameworks, requiring few IT knowledge to build an AI application,
– AI platform proposed by GAFAM, with reasonable costs.

A recent report coming from Market Research Future estimates that deep learning economic market is about to reach $18 billion in 2023. This is no surprise that governments provide massive efforts to support AI research and community, and fund the acquisition of proprietary infrastructure to host AI research community applications.

On the other side, vendors have greatly addressed this market by proposing either dedicated AI solution or HPC & AI compliant solutions.

Consequently, many HPC vendors either get into the AI market with "pure" AI processors or "mixed" processors, which can address the HPC or Machine Learning and Deep Learning issues.

– Intel$^{®}$ has recently proposed its AI dedicated processor, the Intel$^{®}$ Nervana processor. Intel$^{®}$ either integrates, in its "HPC" processors, new instructions set, dedicated to Deep Learning, making these processors a really good target for "converged" architecture.

– NVIDIA® proposed AI supercomputers, with a "ready to go" box named
  DGX®. On the other side, NVIDIA® was able to penetrate the market thanks
  to its GPU Volta® V100 and P100, which propose very interesting features
  for AI but still remains extremely good candidates for HPC applications.

We did not mention above the Google® TPU (Tensor Process Unit) as they
are not for sale, and we do believe that Google and Facebook are not likely to
commercialize their solutions.

The large economic growth of Machine Learning and Deep Learning market
are either due to the vendors themselves: Classic vendors such as Intel® and
NVIDIA® largely and actively contribute to AI framework development and
optimization. AI success is partly due to an active and efficient contribution of
vendors to the most used Machine Learning and Deep Learning frameworks,
allowing AI researchers and enterprises to focus only on the algorithm side, and
not on the development and performance issues.

This is also one major advantage of Machine Learning and Deep Learning
communities compared to HPC, as they directly made a massive use of open
source frameworks, and started their work on, mostly, private cloud platforms.
GAFAM very quickly offered to the AI community very advanced and smooth
platforms, with free access first to attract people, and low prices when addiction
or need is established.

HPC and AI communities have very different "habits" when it comes to
production and science: HPC community executes its simulation on proprietary
or governmental platforms, develops its own frameworks and is reluctant to share
its data, as most of them are sensitive data: Therefore, platforms of execution
run in a closed and secured environment that is hardly compatible with "open-
science". The path to open-science for HPC community is harder than it is for
AI community. AI community mostly started to run on cloud platforms, non-
proprietary, with simulations based on open source frameworks, using data that
are not sensitive or whose security is not a major concern. AI community does
not need to migrate towards open-science: it directly started with open-science
mind.

## 3    Rethink HPC: Think Data

Regarding the recent market evolution and government decisions, HPC will have
to deal with AI community. This must not appear as a constraint, as many HPC
simulations can take advantage of Deep Learning or Machine Learning algo-
rithms and improve either their performance or data analysis process. Mostly,
these two scientific worlds collision is an opportunity to improve their simulation
and scientific research.

### 3.1    Data Struggling

Many scientific domains, such as astrophysics, materials, fusion... that are mas-
sively using HPC to execute their simulations are currently struggling with data.

It can be data coming from real experimentations or HPC simulations and both can be either output or input data.

Analyzing these data is vital for some scientific domains, as they might improve drastically research and therefore theory. However, analysis part is often not up to the HPC simulation level: Whether the analysis part is done manually and is excessively long and fastidious to do, including lack of resources, whether the platforms to do post treatment are inexistent or not adapted.

We recall that for very large set of data, we can hardly move Terabytes or Petabytes of data from one HPC center where data were generated to another where data could be analyzed. HPC centers hosting data consumer applications will need to find solutions to the following issues:

– For long term and interactive storage: the storage capacities and functions (which was, so far, limited to archive treatment) that we encouraged so far cannot be considered anymore. The major risk is a loss of generated data and inhibit research discoveries.
– The use of Machine Learning and/or Deep Learning solutions for analysis (pre or post-treatment) of HPC data. Data analysis can be done in situ, depending on the application workflow and needs.

It is often observed that data management is treated with a lower priority compared to the compute capacity issue. This bad habit is now changing, because of the fast and drastic rise of Deep Learning and Machine Learning in the HPC landscape.

Big-data-driven models featuring machine learning and deep learning can incredibly improve scientific research, as presented in [10]. In the Fusion energy domain, such techniques enable key discoveries and/or provide considerable time saving to produce results and simulations. The project EuroFusion Joint European Torus (JET) is faced to very large scale perturbations in modelling tokamak systems.

The JET team realized great improvements regarding the predictive capability (which currently is 80%) for disruptions that happen before to deterioration incident: such success largely overcome the "classic" HPC simulations. This work is still ongoing as upcoming International Thermonuclear Experimental Reactor (ITER) aims to provide a predictive capability of $\tilde{9}5\%$, which requires a large refinement of physical models, and this challenge can be crossed with machine learning methods.

## 3.2   HPC Community Must Think Data

Unfortunately, the data issue has not been sufficiently addressed both in terms of infrastructures and in terms of development frameworks. Many applications generate massive amounts of data, but the post treatment and storage management is not sharp enough to deal with the amount of data. Many codes do not embed the post treatment procedures and data compression is either not enough, not adapted or inexistent.

HPC community really needs to rethink their applications, especially the community with data issues. For this community, the power compute capacity is currently their focus, while data should be the first focus as it will be -soon enough- the bottleneck for their applications scalability.

Some HPC communities have understood that Machine Learning and/or Deep Learning can be a very good alternative to classic HPC compute paradigms.

Current HPC applications are overwhelmed with:

– Data captured by observations/production instruments
– Data generated by simulations that must be treated and/or analyzed
– Data generated by sensors networks.

Some HPC applications are therefore in needs of Deep Learning methods to automatize and accelerate the data treatment. This treatment can be used either as pure analysis or as a quality comparison between data coming from the simulation and data coming from real experimentations.

This can help to value the efficiency and accuracy of simulations, eventually its limits and potential. The paper [11] gives a good overview of the data issue in the context of turbulence modeling. Decades ago, the compute power was the driving metric while now the data-driven model is favored. The amount of experimental data are now used to estimate a model relevance or to improve the model itself. A data-driven model is based on making a program that is able to give results, based on large data set. Such approach is adequate for any scientific model struggling with data.

Many numerical simulations in which turbulent flows appear use Large-Eddy Simulation (LES). In [14], the turbulent flow in LES are decomposed into Grid Scale (GS) flow field and filtered with SubGrid-Scale (SGS): The scientists have used Artificial Neural Network (ANN) to find a new subgrid model of the SubGrid-Scale (SGS) applied to LES. The team trained an ANN with a backward propagation using direct numerical simulation data coming from turbulent channel flow. Such configuration improved the correlation between the GS flow field and the SGS stress tensor in LES.

As an illustration, the climate community is directly concerned by data challenge. An international consortium was built to design a climate model and make predictions. Climate community is consuming and generating massive amount of data. The recent use of Deep Learning in climate models can allow comparing simulation results with observational ones, and therefore evaluating the accuracy of simulations. Machine Learning can also help to deal with forecasting problems: Data accumulated about the climate can be effectively used to build predictors for inferring dependencies between past and short-term future values of observed values [7].

Many scientific applications are inundated with data, and for some framework, using Hadoop is a nice option as Hadoop is one of the most used framework to manage big data. Its success comes from his resilience, ease of use and very good portability. In Computational Fluid Domain (CFD), Hadoop is used to manage data analyses as CFD applications can generate large files to stock information on the systems and associated time-steps. In [12], the authors used

a one-dimensional MagnetoHydroDynamics simulation to show that Hadoop is a good approach for CFD applications struggling with large data volume. HaDoop File System (HDFS) success is also described in [13], especially because HDFS proposes any plugins that apply to each problem depending on the issue (large files, number of files, MapReduce issues...).

HPC community needs to rethink its old paradigms, where data was not an issue. Many simulations have evolved, without considering the data as a key point in both algorithms and code developments. For a long time it was assumed that HPC was based on three pillars, namely, theory, experimentation and simulation. This paradigm is not true anymore for any HPC simulation struggling with data. Data is another pillar of science and should get at least as much attention as the three others [8].

It will be tough to rethink the current HPC paradigms, especially introducing a new "key parameter" in the center of the simulations. Many HPC applications will need a complete strategy regarding data management and get educated to Data Science. There is an urge to establish, for every data consumer/producer code, data strategy: establishing a fast and efficient data treatment, that can be executed with machine learning or deep learning. Though Data science is somehow not new, and a large community has addressed this problem, the HPC community is not yet fully aware of their needs of Data Science.

### 3.3   FAIR Data

Most of the current HPC platform treats the storage as an "archive" process, while this cannot be true anymore with the current amount of data and the converged platform "AI-HPC". In the context of converged architectures, data are the key of the simulation rather than a simple input and/or output. Such application profile requires to get architecture able to feed large volume of data to the CPU with a high bandwidth and very low latency. This large data transfer is neither a pre/post process nor an input/output step, it is the "main compute process". This raises some questions as HPC centers can currently provide large compute capabilities but are not ready to sustain dynamic and open software stacks as well as dynamic and large storage capacities. One must always remember that (large amount and quality) data is the center of all AI applications. Recently, a consortium seriously treats the data-tsunami issue. This consortium proposed the concept of FAIR data [9] and [16], meaning that they must be:

– Findable: One must be able to find a data easily and quickly. This implies a standardization and strict processes for data storage and more precisely archive whatever the infrastructure is. Metadata associated to archived data must be rich enough to precisely describe the stored data themselves and respect a unique identifier.
– Accessible: Data must be accessible by scientists responsible of the data or any scientists wishing to use these data for scientific purpose. Data storage infrastructures must propose a standard protocol for data access: security and clear conditions of access (license access, use agreement...) at least.

– Interoperability: Data must be usable on any infrastructure and compatible with any platforms and other data. This will be a direct consequence of the "Findable" aspect in the FAIR data conception. The interoperability aspect can be achieved if and only if standardization around metadata and archive policies are defined.

– Re-usable: The re-usable access is also nested with the "Findable" and "Accessible" characteristics presented above. A perishable data should not be integrated in the FAIR process, in this part, only re-usable data matters, data that can enrich several research activities.
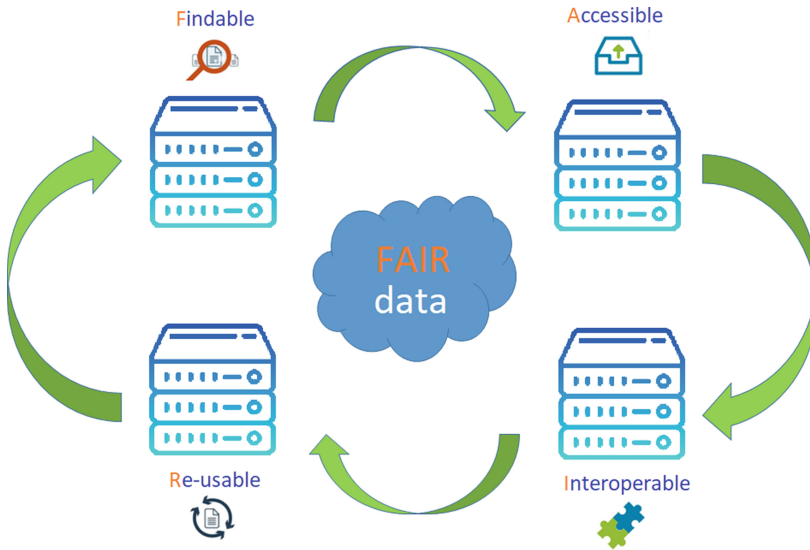


**Fig. 1.** FAIR data

Data shall not be considered as passive, stored and pending for a treatment, but active, whether we talk about "rough" or "treated" data. This implies a deep adaptation of HPC community to satisfy these requirements, but the effort coming from research community cannot happen before the establishment of FAIR standard. A direct consequence is to protect and ensure the usability/viability of scientific data.

The "DATA FAIR" consortium aims to provide guides and specifications to describe the minimal required protocols, formats and interface that will help to satisfy the Findable, Accessible, Interoperable and Reusable characteristics. The consortium gathers scientists that are involved in data issue to develop into a generic data principle, with different implementations in communities and organizations. The aims of this existing consortium is to extend and share the concept of FAIR data in order to ensure data management issue in each application design.

One danger is that data is a major economic market that jeopardizes the scientific community: Indeed, technologies for data management will very likely not be driven by science. Data issued from scientific applications (whether we talk about HPC or AI applications) are highly valuable in terms of scientific knowledge but for some security and ethical issues, those data can't be used for economic purpose, or under very strict and identified conditions.

**Towards Open Science.** FAIR data presented above is not exactly the same as open data, where the accessibility is mostly the only concern for open-data. However, FAIR data implicitly leads HPC community towards the open science concept: indeed, FAIR data is one of the corner-stone of open science.

AI scientific community has a major difference with HPC: AI success is (in part) due to the massive open-source culture. Indeed, the massive use of Deep Learning applications is a consequence of release as open source from major frameworks coming from industries.

Globally, AI community started directly with the "open source" culture while HPC went to open source more by necessity than by design.

Getting benefit from AI also implies that HPC must reconsider its "partition" design and get toward open science: namely using and contributing to open source codes, with an open access, and most of all, allowing FAIR data. Of course these paradigms do not address the defense and army sector, but many HPC scientific domains do not require secret defense or closed development and could largely get benefit from open science.

The open science issues concern the biology and medical scientific communities. More precisely, the genomic studies are facing to the problem of data struggling and FAIR data: They generate about 250 000 genomes per year. Nowadays genomic analysis are used for clinical R&D and under pressure regarding the confidentiality issues. The project sponsored by French government called "France genomic plan" [6] aims to combined the use of big data and HPC to implement an efficient use of genomics in healthcare pathway. Today, it could take between days to weeks to obtain a feedback from biology experts in charge of genomic analysis. Due to the huge increase of data and since we will be in healthcare pathway, the genome's analysis should be done in 2 or 3 days, and has to be reproducible. HPC and AI are two solutions to reach this objective, such as using a specific hardware to genome alignment process, which produces the differences between the sequenced genome to analyze and a references one. The use of deep learning is also a solution to genomic studies issues, and could help to analyze the genome and identify the "defaults": Annotation phase, which needs lots of cross-comparisons, correlation and data analysis with many other data (existing annotated genomes, databases, bibliography, clinical data, medical images...).

Such scientific domains can get a large benefit from open science; sharing genomic data can greatly accelerate the research. Medical topics are clearly caught between the data protection issue and the needs of open science.

## 4   Conclusion

Data is the new gold and both AI and HPC communities must take this very seriously. It is urgent to treat the data issue, and the FAIR concept and consortiums around are very promising. However, both HPC and AI community must get involved in this consortium to make sure it contributes to scientific research.

FAIR data is not easy to set up as it involves not only the scientific community but also vendors, HPC centers, legal sector... Nevertheless, this is an unavoidable step to ensure a viable and substantial improvement of scientific research (HPC and AI areas). The FAIR consortium is a good protection to protect scientific data and limit their use in an economic context. The current approach of FAIR to address the data issue enables the communities to prepare their frameworks to big data issues and exploiting/analyzing these data.

HPC community will have large benefit from FAIR data but either from AI community, as it can drastically help simulations scalability and data treatment. We believe that HPC will be strongly impacted by Deep Learning applications. It is very likely that the first exaflopic computer will be in fact a supercomputer dedicated to Deep Learning frameworks, which illustrates well the HPC strong evolution facing the AI rising.

## References

1. Top500. https://www.top500.org
2. Top500 November 2018. https://www.top500.org/lists/2018/11
3. Green500. https://www.top500.org/green500
4. HPCG500. http://www.hpcg-benchmark.org
5. Deep500. https://www.deep500.org
6. Plan France Médecine Génomique 2025. https://solidarites-sante.gouv.fr/systeme-de-sante-et-medico-social/recherche-et-innovation/france-genomique
7. Bontempi, G., Ben Taieb, S., Le Borgne, Y.A.: Machine learning strategies for time series forecasting. In: Aufaure, M.A., Zimányi, E. (eds.) Business Intelligence; eBISS 2012. Lecture Notes in Business Information Processing, vol. 138. Springer, Heidelberg (2013)
8. Tolle, K.M., Tansley, D.S.W., Hey, A.J.G.: The fourth paradigm: data-intensive scientific discovery. Proc. IEEE **99**(8) (2011). https://doi.org/10.1109/JPROC.2011.2155130
9. Wilkinson, M.D., et al.: The FAIR Guiding Principles for scientific data management and stewardship. Sci. Data https://doi.org/10.1038/sdata.2016.18
10. Big Data and Extreme-Scale Computing: Pathways to Convergence; BDEC report; Tech Report No. ICL-UT-17-08; http://www.exascale.org/bdec; January 2018. Int. J. High Performance Comput. Appl. **32**(4) (2018)
11. Duraisamy, K., Iaccarino, G., Xiao, H.: Turbulence modeling in the age of data. Annu. Rev. Fluid Mech. **51**, 357–77 (2019)
12. Kim, M., Lee, Y., Park, H.-H., Hahn, S.J., Lee, C.-G.: Computational fluid dynamics simulation based on Hadoop Ecosystem and heterogeneous computing. Comput. Fluids **115**, 1–10 (2015)

13. Lange, B., Nguyen, T.: A Hadoop distribution for engineering simulation; [Research Report] INRIA Grenoble - Rhône-Alpes (2014). ffhal-01130630; https://hal.inria.fr/hal-01130630/document
14. Gamahara, M., Hattori, Y.: Searching for turbulence models by artificial neural network. Phys. Rev. Fluids **2**, 054604 (2017)
15. Wilkinson, M., Dumontier, M., Aalbersberg, I.J., Appleton, G., Axton, M., Baak, A., Blomberg, N. and Boiten, J.W., da Silva Santos, B., Olavo, L., Bourne, P.E., Bouwman, J., Brookes, A.J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C.T., Finkers, R., Richard and Mons, B.: The FAIR Guiding Principles for scientific data management and stewardship. Sci. Data **3** (2016). https://doi.org/10.1038/sdata.2016.18
16. FAIR data consortium. http://www.datafairport.org/