



squad.ai: A Multi-agent System Built on LLMs, Incorporating Specialized Embeddings and Sociocultural Diversity

André Neves¹  , Silvio Meira^{2,3} , Filipe Calegário¹ , Rui Belfort² ,
and Marcello Bressan³ 

¹ UFPE - Universidade Federal de Pernambuco, Recife, PE 50670-901, Brazil
Andre.neves@ufpe.br

² TDS Company, Recife, PE 50030210, Brazil

³ CESAR - Centro de Estudos e Sistemas Avançados do Recife, Recife, PE 50030-390, Brazil

Abstract. The squad.ai system emerges as an innovative proposition in the landscape of multi-agent systems, building upon the robustness of Large Language Models (LLMs). Recognizing the potentialities and limitations of LLMs, the system integrates specialized embeddings, allowing for a deepening and specialization of agent knowledge in specific domains. A distinctive feature of squad.ai is the incorporation of rich identity and educational attributes, reflecting sociocultural diversity. This diversity, coupled with behavioral archetypes, aims to facilitate more contextualized and humanized interactions, both among agents and between agents and humans. Squad.ai, therefore, represents a stride forward in the pursuit of multi-agent systems that combine vast knowledge, deep specialization, and meaningful sociocultural representation.

Keywords: Multi-Agent System · Large Language Models (LLMs) · Specialized Embeddings · Sociocultural Diversity · Collaborative AI · Real-time Adaptability · Personalized Learning · Decision-making Support · Contextualized Collaboration · squad.ai · Agent-Human Interaction · Knowledge Enhancement · Dynamic Generation · AI in Healthcare · Legal Applications of AI

1 Introduction

The rapid pace of technological evolution over the past few decades has brought to the forefront an array of computational solutions designed to mimic, enhance, and even surpass human capabilities in various domains. Among these, multi-agent systems have garnered significant attention due to their inherent ability to model and manage complex tasks by breaking them down into smaller, more manageable sub-tasks handled by individual agents. These agents, capable of autonomous operation and interaction, jointly work towards achieving overarching system goals, often emulating the collective intelligence found in natural systems.

Historically, multi-agent systems have been rooted in decentralized problem-solving. From coordinating tasks in industrial settings to optimizing traffic flows in smart cities, these systems have proven to be versatile tools. Their strength lies in the premise that many hands (or in this case, agents) make light work. By distributing tasks among multiple agents, the system can achieve parallelism, fault tolerance, and often enhanced efficiency, especially in scenarios where centralized systems are inadequate or impractical (Wooldridge, 2009).

However, as multi-agent systems found applications in increasingly diverse domains, a pertinent question arose: How does one ensure that these agents possess the depth of knowledge and expertise required to handle specialized tasks? Enter Large Language Models (LLMs). With their ability to understand, generate, and interact using human-like language, LLMs opened avenues for creating agents with vast general knowledge. Platforms such as GPT-3 and its successors demonstrated a prowess in linguistic tasks that was, until recently, considered the exclusive domain of human cognition (Brown et al., 2020). Integrating LLMs into multi-agent systems seemed like a natural progression, offering the promise of agents capable of sophisticated linguistic interactions and knowledge processing.

Yet, while LLMs brought breadth of knowledge, the challenge of depth in specific domains persisted. LLMs, despite their impressive capabilities, are not inherently specialized for every niche domain. For instance, an LLM might know about quantum mechanics, but might it possess the deep expertise of a quantum physicist? This limitation underscores the need for specialization mechanisms that allow agents to delve deep into specific knowledge areas, ensuring they are not just jack-of-all-trades, but also masters of some.

This is where the squad.ai innovation comes into play. Recognizing the potential of LLMs and understanding their limitations, squad.ai introduces a framework where agents, based on LLMs, are enhanced with specialized embeddings. These embeddings, tailored for specific domains, allow agents to possess deep expertise, extending the capabilities of the foundational LLM. In essence, while the LLM provides a wide canvas of general knowledge, the embeddings paint detailed strokes on this canvas, bringing forth intricate patterns of specialized knowledge.

But squad.ai's innovation does not stop at knowledge enhancement. Realizing that effective multi-agent interaction is not just about knowledge but also about relatability and context, the system incorporates rich attributes of identity and education in its agents. These attributes, spanning aspects such as age, gender, nationality, and educational background, introduce a sociocultural layer to the agents. In doing so, they ensure that interactions, whether agent-agent or agent-human, are more contextual, relatable, and humanized. This addition, rooted in behavioral science and sociology, transforms the agents from mere knowledge processors to entities with simulated backgrounds, enhancing the richness of interactions within the multi-agent system.

In conclusion, as we stand on the cusp of a new era in artificial intelligence, the need for systems that are both broad in knowledge and deep in expertise becomes ever more apparent. Systems that not only compute but also relate; systems that do not just know but understand. The squad.ai, with its innovative approach, seeks to bridge these requirements, proposing a multi-agent framework that promises vast knowledge, deep

specialization, and a touch of humanity. In the ensuing sections, we will delve deeper into the theoretical foundations of this system, exploring the technologies and concepts that make squad.ai a beacon in the realm of advanced multi-agent systems.

2 Theoretical Foundation

2.1 Large Language Models (LLMs): Their Role in AI and Application in squad.ai

The emergence of Large Language Models (LLMs) like GPT-3 and its successors brought a notable transformation to the domain of Artificial Intelligence (AI). These models are trained on vast datasets, enabling them to process, understand, and generate language in a remarkably human-like manner. The capability of these LLMs to respond to a wide variety of prompts with relevant and often creative information suggests a comprehensive depth and breadth of knowledge (Brown et al., 2020).

Within the context of squad.ai, LLMs serve as a backbone, granting agents a vast base of knowledge. However, while LLMs excel in covering a broad spectrum of topics, their efficacy may be limited in specific niches or questions that require deep specialization. Thus, the need to supplement and enhance these models becomes evident. **Sample Heading (Third Level)**. Only two levels of headings should be numbered. Lower level headings remain unnumbered; they are formatted as run-in headings.

2.2 Specialized Embeddings: Introduction, Benefits, and Their Integration with LLMs

Embeddings, in the realm of machine learning, refer to vector representations of data. These representations capture semantics and relationships between data, making them powerful instruments for knowledge representation. Specialized embeddings, on the other hand, are fine-tuned or designed specifically for certain domains or topics (Mikolov et al., 2013).

By integrating specialized embeddings with LLMs, squad.ai seeks to overcome the inherent limitations of LLMs in terms of specialization. These embeddings, when integrated, serve to enhance, and specialize the knowledge of agents in specific areas, offering a depth the foundational LLM might not achieve on its own.

Sample Heading (Third Level). Only two levels of headings should be numbered. Lower level headings remain unnumbered; they are formatted as run-in headings.

2.3 Jungian Archetypes: Foundation and Their Relevance in Modeling Agent Behaviors

Jungian archetypes are concepts derived from Carl Gustav Jung's analytical psychology and refer to innate images or universal patterns that reside in the collective unconscious. These archetypes, such as the Hero, the Wise Old Man, or the Mother, play a pivotal role in shaping human behavior and perceptions (Jung, 1964).

In modelling the behavior of agents in multi-agent systems like squad.ai, Jungian archetypes can provide a foundation for developing agents that exhibit behavior patterns more aligned with human perceptions and expectations. By incorporating these archetypes, agents not only process information but do so in a manner that is perceived as more natural or relatable to human users.

2.4 Identity and Educational Attributes: Amplifying the Sociocultural Diversity of Agents

For agents in multi-agent systems to be genuinely effective and relatable, it is crucial that they mirror the rich tapestry of identities and backgrounds found in human society. Identity, encompassing characteristics like age, gender, nationality, and social class, bestows upon each agent a unique nuance that can be used to enhance relatability and efficacy of interactions. Education, on the other hand, refers to the agent's domain expertise and level of proficiency, which might include academic qualifications, technical training, or any other specialized learning.

By incorporating these attributes in squad.ai, agents can be modelled to mirror a sociocultural diversity. This diversity, when adeptly applied, has the potential to significantly enhance the collaboration and efficacy of a multi-agent network.

2.5 Prompt Engineering: Its Significance in Agent Specialization and Interactivity

Prompt Engineering refers to the art and science of formulating prompts or instructions in such a way that they guide a language model's output in a specific direction. Given the malleable nature of LLMs, careful formulation of prompts can steer, specialize, or even refine the responses provided by the model (Gao et al., 2020).

Within squad.ai, prompt engineering plays a critical role in ensuring that agents, grounded in LLMs, respond, and interact in ways that align with the system's and users' expectations. In addition to guiding the model's output, prompts can also be used to access and mobilize the specialized embeddings, ensuring agents deliver specialized responses when needed.

3 Methodology

3.1 squad.ai's LLM Base: Discussion on the Underlying Language Model

At the core of any language-based system lies its proficiency to process, understand, and generate language. For squad.ai, this linguistic prowess is provided by a Large Language Model (LLM). Trained extensively on expansive datasets, these models form the foundational generalized knowledge and linguistic capability.

What sets squad.ai apart is its ability to shape this underlying model adaptively and dynamically. Rather than relying on a pre-set agent, the system generates agents on-demand, dynamically adjusting to the user's needs and specifications.

3.2 Integration of Specialized Embeddings: Real-Time Process and Technique

Within the squad.ai realm, knowledge specialization is taken to a new echelon, courtesy of the dynamic incorporation of specialized embeddings. Unlike traditional systems that rely on pre-trained embeddings, squad.ai constructs embeddings in real-time based on user-provided data.

Through advanced machine learning techniques and semantic analysis, the system processes, scrutinizes, and translates user data into vector embeddings. These embeddings are then integrated into the LLM, granting the generated agent specialized and contextualized knowledge. This on-the-fly approach allows users to craft agents perfectly tailored to specific situations, ensuring relevance and accuracy in interactions.

3.3 Definition of Identity and Training: Real-Time Construction with Diversity and Depth

One of the pillars of squad.ai is its ability to humanize agents, making them relatable and contextualized. However, rather than relying on static or predefined identities, the system adopts a groundbreaking approach, allowing dynamic definition of identity and training.

Users can specify various dimensions of identity such as age, gender, nationality, and social class. Alongside this, users have the freedom to define the agent's training, choosing from various knowledge domains and levels of expertise. This user-supplied information is then processed in real-time, spawning an agent that mirrors the desired identity and training perfectly.

3.4 Behavior Modelling: Dynamic Generation Using Jungian Archetypes and Prompt Engineering

An agent's behavior is not just a function of its knowledge but also its personality and interaction style. In squad.ai, behavioral modeling is dynamically tailored based on user inputs and contextual necessities.

Using Jungian archetypes as a foundation, the system generates behavior patterns that align with user expectations and requirements. Whether it is the contemplative behavior of the "Wise Old Man" or the determination of the "Hero", squad.ai can model a spectrum of archetype-based behaviors.

Moreover, real-time prompt engineering plays a pivotal role. Through the meticulous formulation of prompts, the underlying LLM is steered, ensuring responses that align with the desired behavior and the situation at hand.

3.5 On the Responsibility Ascription and Ethical Concerns of Artificial Agents

Knowing that the training of LLMs and application of AI is not impervious to bias, the social and human intelligence contribution becomes paramount.

Regarding the LLMs as Artificial Moral Agents (Allen et al. 2000) allows us to ponder upon the demarcation of responsibility and involvement of such agents in co-creative processes and decision-making.

The debate on whether AI and LLMs should be regarded as AMAs is still far from maturity, however the responsibility gap should not be overlooked, and a normative ethical stance must be adopted on this matter (Behdadi and Munthe 2020).

This means that, from a practical standpoint, squad.ai's human agents systematically scrutinize all content generated by the Artificial Agents, may it be by rigorous curatorship (human intelligence) or through meticulous debate with stakeholders of the project (social intelligence).

Methodology Conclusion. Squad.ai redefines the traditional approach to multi-agent systems, emphasizing flexibility, customization, and adaptability. Through its novel methodology, the system not only spawns agents and specialized knowledge in real-time but also ensures these agents are perfectly aligned with user needs and contexts. This user-centric and dynamic approach establishes squad.ai as a disruptive force in the realm of artificial intelligence and multi-agent systems.

4 Operation of squad.ai

4.1 Formation and Interactivity of Squads: Mechanism of Grouping and Collaboration

The innovative essence of squad.ai is embedded in its dynamic and adaptive creation of 'squads' - clusters of artificial agents poised to work in tandem. When a user interacts with squad.ai, they do not merely engage with a singular, static agent but can invoke a multitude of agents, each tailored to specific tasks and domains.

Upon a user's request, the system swiftly evaluates the contextual requirements. Leveraging the underlying LLM, coupled with real-time embeddings derived from user-provided data, squad.ai dynamically assembles a squad. This assemblage ensures a diverse range of skills and knowledge, closely mirroring the user's contextual needs.

The interactivity within these squads is orchestrated meticulously. Agents within a squad are interconnected, enabling seamless data exchange and collaborative decision-making. Each agent, built with specific expertise and behavior, plays its role - from analyzing complex data to drafting responses or making predictions. The collective intelligence of the squad, therefore, is greater than the sum of its parts.

4.2 Agent-Agent Collaboration: Protocols and Collaboration Mechanisms

The bedrock of squad.ai's efficacy lies in its ability to foster seamless collaboration between agents. Such Agent-Agent collaboration is crucial for the holistic and efficient operation of the system.

At the onset of a task, protocols define the hierarchy and sequence of operations among agents. While one agent might specialize in initial data processing, another could focus on in-depth analysis, and yet another might be responsible for synthesizing a comprehensible output.

The communication protocol adopted by squad.ai ensures low latency and high accuracy. When one agent completes its segment of the task, it dispatches the results to

the next agent in the sequence, often with meta-information or suggestions for further processing. This iterative and collaborative approach ensures optimal outputs.

Furthermore, agents within a squad possess the capability to seek assistance or validation from their peers. If an agent encounters an anomaly or uncertainty, it can request peer agents for insights or alternate perspectives, ensuring that the squad's final output is refined and accurate.

4.3 Agent-Human Collaboration: Interfaces, Protocols, and Feedback

Squad.ai is not just designed to foster collaboration among agents but is equally adept at facilitating Agent-Human interactions. This symbiotic relationship augments the capabilities of both parties, ensuring superior outcomes.

The user interface is intuitive, designed to cater to a wide range of users. Upon initiating a task, users can specify their requirements, which are then interpreted by the system to select or generate the most apt squad. As agents work on the task, users can monitor progress, interact with individual agents, or adjust parameters in real-time.

The protocols governing Agent-Human interactions prioritize transparency and flexibility. Users receive updates, insights, and even suggestions from agents, ensuring they remain in the loop. In scenarios where user intervention or decision-making is required, agents proactively seek input.

A unique feature of squad.ai is its feedback mechanism. Post-task, users can provide feedback on the squad's performance. This feedback is processed and integrated, allowing the system to learn and refine its operations. Over time, this continuous feedback loop ensures that squad.ai becomes increasingly adept at tailoring its squads to users' needs.

Conclusion on squad.ai's Operation. Squad.ai stands as a testament to the evolution of AI systems, emphasizing collaboration both internally among agents and externally with human users. Its dynamic, adaptive, and user-centric approach reimagines how we perceive and interact with AI. The system's ability to amalgamate diverse agents into cohesive squads, each tailored to specific tasks, coupled with its robust collaboration protocols, sets a new benchmark for AI-driven solutions. The emphasis on continuous learning, driven by user feedback, ensures that squad.ai is not just a static tool but a constantly evolving partner, poised to meet the ever-changing demands of its user base.

5 Potential Use Cases and Applications

The inherent flexibility and adaptability of squad.ai, which generates specialized agents in real-time to form collaborative squads, positions it as a priceless tool across a broad range of contexts. Here, we will delve into four pivotal areas, shedding light on how squad.ai can revolutionize decision-making, learning, health, and the legal environment, with special emphasis on the enriching interaction between agents and humans.

Decision-Making. Today's corporate environments are replete with intricate challenges where decisions often intersect various vectors of information.

Real-Time Data Analysis. In a volatile corporate setting, decisions need to be swift and informed. An executive contemplating changes in the supply chain might benefit from a squad with agents specializing in logistics, global economics, market analysis, and finance. The agents' interactivity with professionals amplifies their capability to evaluate scenarios, enabling more precise and contextualized decision-making.

Strategic Simulations. When evaluating market strategies or product launches, squads can craft simulations involving various variables, always in collaboration with strategy teams, ensuring multiple perspectives are contemplated.

Learning Environments. Learning, whether formal or informal, is ever-evolving, and personalization has become crucial.

Personalized Student Support. A student grappling with a particular concept could benefit from a multifaceted squad offering theoretical, practical, and interactive approaches. Collaboration between agents and students enhances comprehension and engagement.

Multidisciplinary Content Creation: For educators and researchers, multidisciplinary squads can assist in integrating varied perspectives, always in tandem with the educators themselves, ensuring the generated content is holistic and pertinent.

Health Environments. Healthcare is an arena where timely and accurate decisions are vital.

Assisted Diagnosis. In situations where a patient presents intricate symptoms, a squad of agents specializing across various medical specialties can assist healthcare professionals, providing a holistic view and collaborating for a more accurate diagnosis.

Treatment Suggestions. Once diagnosed, a squad can work alongside doctors to explore treatment options, considering everything from current medical literature to emerging practices, ensuring the patient receives the most apt and up-to-date treatment.

Legal Environment. The intricacies of law require a multifaceted approach.

Case Analysis. When facing a complex legal case, lawyers and jurists can benefit from the assistance of a squad specialized in different branches of law, pertinent jurisprudences, and legislations. The agents, collaborating with professionals, can facilitate a more comprehensive and in-depth analysis.

Document Preparation. In the drafting of legal documents, precision is paramount. A squad can aid in the process, ensuring all legal facets are contemplated, always in collaboration with the responsible professional.

Benefits of squad.ai Across Different Contexts. The potential of squad.ai extends beyond its technical functionality, bringing tangible benefits across all scenarios:

Adaptability. The dynamic squad generation ensures relevant and contextualized solutions.

Deep Collaboration. The interaction between agents and humans amplifies individual capabilities, fostering enriching collaboration.

Amplified Engagement. In all contexts, user engagement is enhanced by personalized and contextual interaction.

Increased Efficiency. The ability of squads to bring specialized expertise accelerates processes and heightens accuracy.

Use Cases Conclusion. Squad.ai is more than just a tool; it is a partnership. Across all scenarios, it acts as an enhancer, amplifying the individual capacities of the humans it interacts with. Through its adaptive and collaborative capabilities, it redefines paradigms across a broad spectrum of areas, promising to revolutionize how we approach and tackle challenges.

6 Experimental Application

6.1 Experiment Description

An instance of squad.ai was constructed, comprising 36 agents (Fig. 1), with the aim of assisting consultants at TDS.company, a strategic transformation consultancy firm based in Digital Port, Recife, Brazil. The agents were intentionally designed with a range of diversity covering different ages, genders, races, educational backgrounds, and behavioral profiles, based on Jungian archetypes. These agents were deployed for participation both in debates among themselves and in debates involving human beings. Discussions took place within the *Strateegia.digital* platform (Fig. 2), and their interventions were overseen by human consultants.

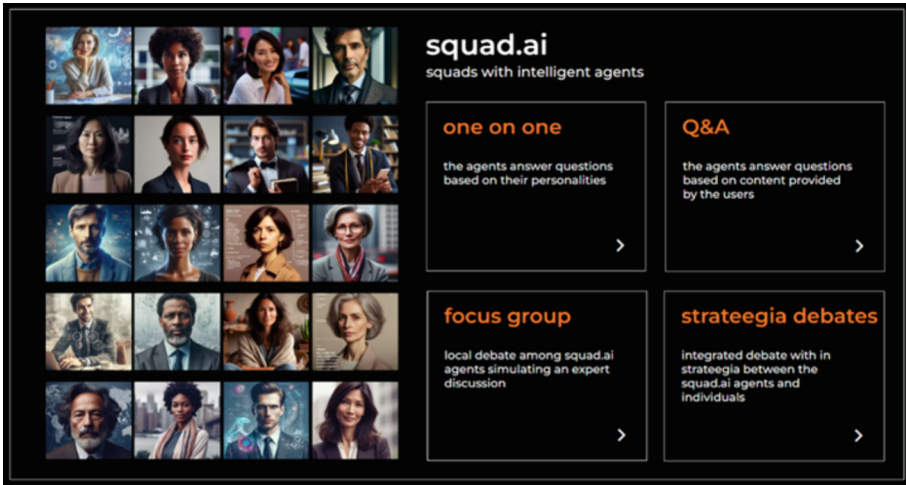


Fig. 1. Interface of the experimental squad.ai application

6.2 Discussion of Results

The findings indicate that the agents maintained performance that was similar to or even superior to human beings in strategic debates. During debates among the agents, effective self-regulation was observed, which mitigated the phenomenon of hallucinations



Fig. 2. Results of the experimental squad.ai debate with human beings in strateegia.digital

commonly seen in single-agent systems. This aspect suggests a significant improvement in the quality of the responses generated.

Human consultants who interacted with the agents displayed enhanced performance compared to those who did not make use of the agents. This improvement was notable both quantitatively, in terms of the number of interventions, and qualitatively, with respect to the depth and relevance of contributions.

As consultants became familiar with the specificities of individual agents, there was a tendency to select preferred agents for subsequent debates. This choice seemed to be influenced by the agents' skills and knowledge in specific domains.

Additionally, it was found that agents with backgrounds in areas directly related to the debate topic tended to steer the discussion toward well-established pathways within that respective field. In contrast, agents with backgrounds in more distant domains provided more innovative insights that veered off traditional pathways.

6.3 Conclusion on squad.ai's Experiment

These findings underline the potential of squad.ai to enrich strategic debates by offering a diversity of perspectives and by enhancing human performance. The observed dynamics suggest that well-orchestrated collaboration between agents and humans can result in more thoroughly grounded and innovative decisions.

7 Conclusion

As the digital age unfolds, the intersection of artificial intelligence and human expertise becomes an increasingly important frontier for innovation. squad.ai stands emblematic of this confluence, illustrating the profound potential when collaborative AI is harnessed in diverse sectors, from business decision-making to the intricate corridors of healthcare and law.

The essence of *squad.ai* is not solely rooted in its advanced technical capabilities but rather in its pioneering approach to integrating AI agents dynamically, in real-time, to human needs. This adaptability ensures that solutions are not just computational but contextually rich, responsive, and deeply collaborative. It demonstrates the value of a system where AI does not replace, but rather enhances, human capabilities, serving as an augmentative force that empowers individuals across various professional spheres.

Furthermore, *squad.ai*'s emphasis on real-time, user-driven configurations, combined with the application of specialized embeddings and rich socio-cultural attributes, sets a new benchmark in personalized AI solutions. This is not just automation; it is true collaboration, where AI agents and humans engage in an iterative, co-creative process, producing outcomes that neither could achieve in isolation.

The diverse use cases explored in this work attest to the system's transformative potential. Yet, it is crucial to understand that the true impact of *squad.ai* lies ahead, as users across the globe adapt, innovate, and uncover new applications we have yet to envision.

In conclusion, *squad.ai* is more than an advanced multi-agent system; it represents a paradigm shift in how we perceive AI's role in our professional and personal lives. By fostering a harmonious environment where AI and human intelligence coalesce, *squad.ai* is not just setting a trajectory for the future of AI but is also reshaping our collaborative future.

References

- Allen, C., Varner, G., Zinser, J.: Prolegomena to any future artificial moral agent. *J. Exp. Theor. Artif. Intell.* **12**(3), 251–261 (2000)
- Behdadi, D., Munthe, C.: A normative approach to artificial moral agency. *Mind. Mach.* **30**, 195–218 (2020)
- Brown, T.B., et al.: Language models are few-shot learners. arXiv preprint [arXiv:2005.14165](https://arxiv.org/abs/2005.14165) (2020)
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. arXiv preprint [arXiv:1810.04805](https://arxiv.org/abs/1810.04805) (2018)
- Gao, L., et al.: The next generation of large-scale language models: challenges and research directions. arXiv preprint [arXiv:2012.15903](https://arxiv.org/abs/2012.15903) (2020)
- Jung, C.G., Hull, R.F.C.: *Archetypes and the Collective Unconscious*. Princeton University Press, Princeton (2014)
- Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. arXiv preprint [arXiv:1301.3781](https://arxiv.org/abs/1301.3781) (2013)
- OpenAI Blog: GPT-3: language models are few-shot learners (2020). <https://openai.com/research/language-models-are-few-shot-learners>. Accessed 15 Nov 2023
- Ribeiro, M.T., Singh, S., Guestrin, C.: Why should i trust you? Explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1135–1144 (2016)
- Russell, S.J., Norvig, P.: *Artificial Intelligence: A Modern Approach*. Pearson Education Limited, Berkeley (2010)
- Vaswani, A., et al.: Attention is all you need. In: *Advances in Neural Information Processing Systems*, pp. 5998–6008 (2017)
- Wooldridge, M.L.: *An Introduction to MultiAgent Systems*, 2nd edn. Wiley, Chichester (2009)