




An AI Factory Digital Twin Deployed Within a High Performance Edge Architecture

1st Sean Ahearne 
Dell Research
Dell Technologies
Cork, Ireland
sean.ahearne@dell.com

2nd Ahmed Khalid 
Dell Research
Dell Technologies
Cork, Ireland
ahmed.khalid@dell.com

3rd Martin Ron 
Machine Intelligence Dept.
Factorio Solutions
Prague, Czechia
martin.ron@factorio.cz

4th Pavel Burget 
Testbed for Industry 4.0
CIIRC, Czech Technical University
Prague, Czechia
pavel.burget@cvut.cz

Abstract—The exponential proliferation of big data and computation-intensive tasks, such as Artificial Intelligence (AI) applications in factories, poses a significant challenge for the current datacenter-focused technological architecture. The "Big data pRocessing and Artificial Intelligence at the Network Edge" (BRAINE) project addresses this problem by introducing an innovative system architecture designed explicitly for compute-intensive edge deployments. BRAINE focuses on decentralizing the computation tasks, enabling a significant reduction in latency, and optimizing the placement of applications within a cloud-edge continuum to ensure optimal operational efficiency. This paper presents the design, implementation, and testing of our novel system architecture in the context of an AI digital twin for factory robotics. Our empirical results indicate substantial improvements in performance metrics such as processing speed and latency compared to traditional architectures and approaches.

Index Terms—Edge, AI/ML, Latency, Motif Discovery, multi-agent systems, digital twin

I. INTRODUCTION

The rapid digitization of various sectors has caused an exponential growth in data generation, coupled with an escalating demand for computation-intensive tasks, especially in the field of Artificial Intelligence (AI) [3]. This data explosion is particularly apparent within industrial sectors such as manufacturing, where AI-driven applications and technologies are increasingly being integrated into everyday operations. For instance, the concept of an AI digital twin, a real-time digital counterpart of physical factory settings, is becoming pivotal for Industry 4.0 [17]. However, the current centralized, datacenter-focused technological architecture struggles to support the requirements posed by these demanding applications.

Edge computing is emerging as a promising solution to the limitations of traditional datacenter-focused architectures [15]. By moving computation tasks closer to the source of data generation, edge computing provides an efficient way to reduce latency, alleviate network congestion, and improve data privacy [14]. However, to fully leverage the benefits of edge computing for computation-intensive AI applications, it is

This work was funded through Electronics Components and Systems for European Leadership (ECSEL) Joint Undertaking (JU), under Grant Agreement 876967.

crucial to develop an innovative system architecture designed explicitly for this purpose.

In response to this challenge, the "Big data pRocessing and Artificial Intelligence at the Network Edge" (BRAINE) project introduces a new system architecture to enhance the capability of edge computing systems for accommodating the rigorous demands of advanced AI applications, deployed and evaluated in this case in an AI digital twin factory setting. This architecture aims to decentralize computation tasks, reduce latency significantly, and optimize application placement within a cloud-edge continuum for optimal operational efficiency.

This paper presents the design, implementation, and testing of this new system architecture under the framework of the BRAINE project. This architecture significantly enhances the capability of edge computing systems to accommodate the rigorous demands of advanced AI applications, such as pattern detection and motif discovery, within an AI digital twin factory setting.

Our architecture focuses on decentralizing the computation tasks, enabling a significant reduction in latency, and optimizing the placement of applications within a cloud-edge continuum to ensure optimal operational efficiency. It employs a variety of novel techniques, including specialized algorithms for workload scheduling, that facilitate real-time data processing and decision-making for future high-power and highly distributed edge use cases.

To demonstrate the efficacy of the proposed architecture, we examine a use case involving its deployment in the context of an AI digital twin for factory robotics. Our empirical results indicate substantial improvements in performance metrics such as processing speed and latency compared to traditional architectures and approaches.

This paper aims to provide a stepping stone towards an architecture which realizes the full potential of cloud to edge computing not just in Industry 4.0, but for many use cases with vastly different requirements. The findings and the insights generated through the BRAINE project could catalyze further innovations, shaping the future of high-powered edge deployments and digital twin technologies.

II. RELATED WORK

The necessity of dealing with burgeoning data and computation-intensive tasks in AI has triggered significant research in alternative computational paradigms. Prominently, edge computing has emerged as a potent solution, emphasizing decentralization to address latency issues and to optimize data processing [19]. This shift from the conventional cloud-centric model fosters a new approach where data processing is carried out closer to the source, enhancing performance and security [8].

AI applications within industrial settings have further amplified the need for efficient computational models. The concept of AI digital twins, for instance, has transformed the industrial landscape, offering unparalleled real-time insights and predictive capabilities [10]. Nonetheless, handling these computation-intensive applications necessitates innovative system architectures and platforms to effectively harness the potential of edge computing.

Multi-agent systems are a well-known concept, whose revival in industry has started in recent years as shown in [7]. It relates on having independent agents representing individual hardware components or services, which act more or less autonomously based on mutual agreement within the system of interconnected agents. This concept is tightly related to the description of their capabilities as outlined in [6]. However, the ontology-based approach often encounters problems with computational demands, which are needed to reason about ontology-based relations. [5] shows how to solve the computational issue on an example from production domain. [11] deals with centralized planning approach for flexible production, which uses the Digital Twin concept to allow for changes in the production plan starting from the current state of the production system.

Motif discovery plays a pivotal role in our use case, where we leverage the innovative BRAINE platform to demonstrate its advantages. To gain a comprehensive understanding of motif discovery, extensive research has been conducted, as highlighted in [18]. Our focus on motif discovery stems from its crucial role in automating dataset preparation, a fundamental task in machine-learning applications.

In particular, preparing datasets from continuous time series data poses significant challenges as it involves manual-intensive efforts without a straightforward rule for automated segmentation and annotation. The versatility of motif discovery becomes evident as it finds applications in various domains, such as monitoring the health of machines, detecting anomalies in financial market behavior, and more.

A common limitation in existing motif discovery research lies in the assumption that the pattern's length, or even the pattern itself, is known beforehand. However, real-world scenarios often demand the ability to identify *any* repeating pattern present in the time series generated by the analyzed process. Addressing this issue, recent efforts, like those proposed in [9], have been made to handle variable-length motifs. In our work, we developed a different method of variable-

length motif discovery, but its details go beyond the scope of this paper.

Moreover, we acknowledge that applications like digital twin development and online detection require significant computational resources when higher precision is required, as discussed in [12]. The challenge arises also when the training data are scarce, which is the case for the beginning of data collection and early results of motif discovery algorithms. Authors in [13] propose a solution for learning models of time series from small datasets, which can be utilized for online detection. Problematics of digital twin are studied for several years now, and this field still lacks consensus about unified definition of digital twin, but the work [1] provide a comprehensive analysis of current state of the art. As both digital twin technology and motif discovery are computationally demanding tasks, our emphasis on an efficient platform like BRAINE becomes even more critical in facilitating their successful execution.

The BRAINE project is one such endeavor that aims to address these challenges, by designing a system architecture dedicated to edge computing use cases. This is aligned with many similar efforts in the field with works such as [16] and [4] presenting edge-based architectures and frameworks, highlighting their applicability for real-time AI tasks.

Yet, these approaches might not fully realize the rigorous demands of AI applications in industrial settings, with many works being particularly being focused on low-power edge computing. The proposed BRAINE architecture aims to bridge this gap, being designed for edge devices and locations that still contain a high level of compute power. This paper builds upon these previous efforts in edge computing, but specifically tailors its approach for high resource requiring use cases, such as AI digital twin factory settings, showcasing the substantial potential for high-powered edge use cases. Ultimately this system and architecture is intended to intertwine with traditional datacenter and cloud-based systems, forming the possible future where system architectures encompass the entirety of both the edge and cloud, forming a continuum [2].

III. REFERENCE ARCHITECTURE

The architecture proposed in the BRAINE project is a comprehensive, modular solution designed to handle the extensive demands of AI applications in edge deployments. The architecture is fundamentally based on Kubernetes (K8s), a leading open-source platform for automating deployment, scaling, and management of containerized applications. The use of K8s as a base allows the architecture to be highly flexible, with the ability to add or remove modules as per the requirements of a specific use case.

The overall reference architecture for BRAINE can be seen in Figure 1, which shows the core software functionality of BRAINE in the centre of the diagram. These core platform functions are separated into three categorized groups known as service/resource management, data/policy management, and telemetry/monitoring. These categories simply serve to aid in

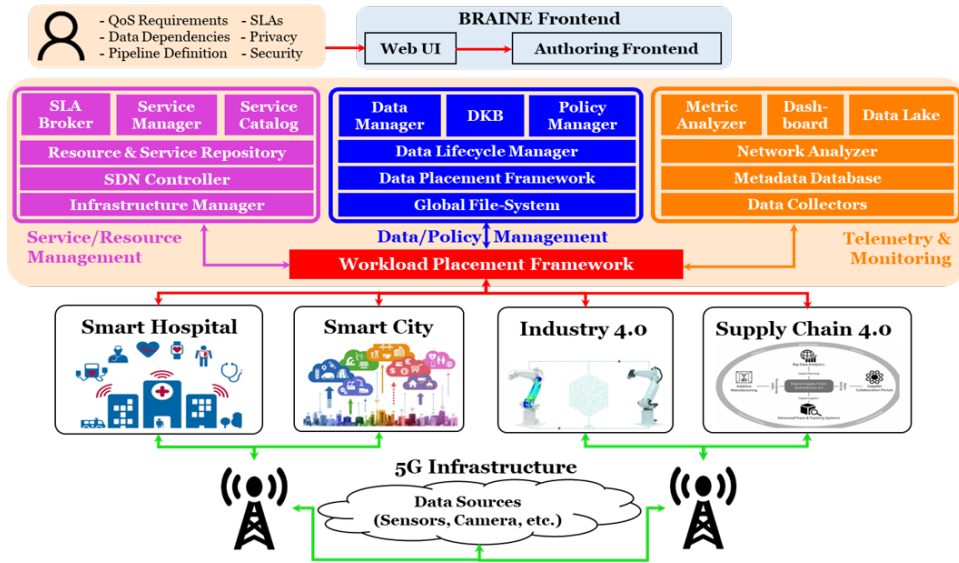


Fig. 1: BRAINE's Reference Architecture and System Components With Example Use-Cases

understanding the functionalities of the individual software modules within each group.

The workload placement framework is a core part of the overall architecture and is one of the key differentiators of the BRAINE architecture. It is an AI-driven Kubernetes scheduler based on reinforcement learning. This scheduler is designed with a specific focus on optimizing power efficiency. However, its robust design allows it to be modified and fine-tuned to prioritize other metrics such as latency, thus catering to a diverse range of edge deployment scenarios. The full details of every software module is beyond the scope of this paper, but details on modules relevant to the Industry 4.0 use case will be provided.

In the BRAINE architecture, Apache Ozone serves as the global filesystem. This robust, distributed object store is designed to handle large volumes of data and billions of objects, which are typical characteristics of big data workload. Apache Ozone's support for horizontal scalability makes it particularly useful in factory edge deployments, where the volume of data may vary dramatically based on production requirements and machinery involved. Moreover, Ozone's inherent design aims at eliminating the scaling bottleneck of a traditional Hadoop Distributed File System (HDFS), enabling more efficient data management in an edge environment.

Apache Ranger, integrated into the architecture for policy management, provides a framework for managing data security in Hadoop clusters. Industrial settings often involve dealing with sensitive data, including proprietary designs, machinery status, and production volumes. The access to such information needs to be carefully controlled to avoid security breaches and intellectual property theft. Apache Ranger's comprehensive security suite, including features like centralized policy management, dynamic row filtering, and data masking, provides granular access control, thus ensuring the security and integrity of data in the edge deployment.

The architecture's Software-Defined Networking (SDN) controller plays a critical role in managing the network's behavior dynamically. In a factory edge deployment, network conditions can change rapidly, for example, due to a sudden surge in data generation or the addition of new devices. The SDN controller can programmatically configure network behavior in real-time, ensuring optimal network performance and resource allocation. It enables dynamic adjustment of network bandwidth, priority assignment, and other parameters based on the real-time requirements of AI applications running in the factory.

The Service Level Agreement (SLA) broker within the BRAINE architecture monitors the latency of applications, triggering necessary actions if defined performance requirements are not met. The importance of meeting latency requirements cannot be overstated in industrial applications, where real-time decision-making is often critical. For example, AI applications involved in quality control or failure prediction need to provide immediate outputs to prevent production delays or equipment damage. The SLA broker ensures that these latency-sensitive applications receive the required resources and network priority to function optimally.

IV. FACTORY USE-CASE

The adaptability and versatility required in modern production processes necessitate a novel approach that extends beyond the machines themselves and their local control systems. The scope of higher-level systems traditionally managed by IT systems, such as Manufacturing Execution Systems (MES) and Enterprise Resource Planning systems (ERP), must also be reevaluated. In this context, the focus of this use case is:

- At the MES level, aiming to propose a multi-agent platform that enables the implementation of a distributed MES. This system will be responsible for overseeing the

production process, conducting diagnostics, and providing supervision

- On a Motif Discovery tool (MOD) to detect anomalies in observable manufacturing processes, based on domain-agnostic motif/pattern discovery in multidimensional time series data.

A. Multi-agent platform for flexible modular production

The multi-agent approach enables flexibility of the planning and orchestration. The overall production plan is divided into smaller tasks which are planned separately, assuming resources are available at the time of planning. This makes the planning problem more optimized compared to the case when a full plan is calculated at the beginning of production. The multi-agent approach facilitates the management of the resources and products at runtime. New resources may be added/removed from the platform, and new products may be ordered during runtime without changing the platform's configuration or interrupting ongoing production.

The agents themselves can be divided into three fundamental groups: (1) Service agents, (2) Product agents, and (3) Hardware agents. The service agents provide the basic functionality of the platform, such as service discovery, monitoring, and communication with external product order services. Product agents represent the products to be manufactured, related to the specific product description in the ontology as described below. The product agent sequentially negotiates the execution of particular operations composing the product recipe. The core of the platform includes the service agents and the product agents. The hardware agents may be running as separate applications outside the platform core. The architecture is displayed in Figure 2.

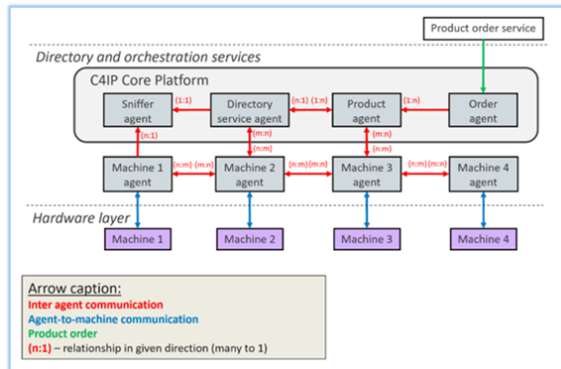


Fig. 2: Multi-Agent Platform Architecture

The physical system and the products are described semantically in the ontology together with capabilities and properties (e.g., the machines' manufacturing activities, the composition of the product, and its production recipe). The recipe concept in the ontology specifies operations that must be performed during the production and their order [5]. The overall plan splits into a sequence of jobs representing specific tasks performed by the machine. The top-level job represents a single operation of the production recipe. Within the job, the

agent submits requests to other agents to satisfy prerequisites needed for the job (e.g., the robot may need to supply material and transport the product to its station before executing its task). Receiving the request from the agent invokes a new job, which may, in turn, submit requests to other agents.

The multi-agent platform is integrated into the BRAINE platform, where its deployment takes place within Kubernetes containers. Each agent runs in a container to allow for flexible deployment at one or several EMDCs as needed. BRAINE's optimized scheduler ensures proper deployment with respect to meeting the timing requirements for the multi-agent application through the use of BRAINE's telemetry monitoring system.

B. Anomaly Detection in Manufacturing Processes using the Motif Discovery (MOD) Tool

We also introduce a novel approach to support the adaptability and efficiency of smart manufacturing by utilizing the Motif Discovery for anomaly detection and highly compressed discrete event log extraction in the context of observable manufacturing processes. The MOD tool offers a versatile solution that encompasses three core building blocks: the Motif Discovery module, Motif Learning module, and Online Detection module. The architecture overview is depicted in Fig. 3.

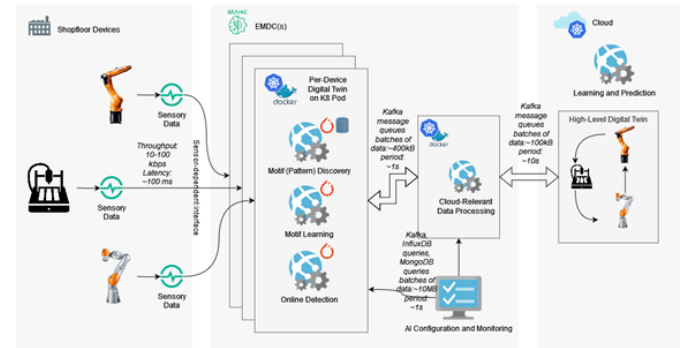


Fig. 3: Motif Discovery Architecture Overview

Motif Discovery Module: The Motif Discovery module serves as the foundation of the MOD tool's data processing pipeline. Its purpose is to autonomously discover repeating patterns in data streams without any prior information. It consumes batches of sensory data collected from industrial robots operating in the CIIRC testbed at CTU. The data comprises measurements such as speed, position, and current, recorded at a mean sampling rate of 73 milliseconds (ms) for each of the six axes of each robot.

Significant repetitive patterns emerge in the behavior of the robots, interspersed with periods of resting time. To ensure the accuracy of anomaly detection, the pre-processing step removes these resting periods, focusing the KPI evaluation solely on meaningful active-time segments.

Motif Learning Module: Building upon the data generated by the Motif Discovery module, the Motif Learning module embarks on the resource-intensive task of training a set of probabilistic detection models. These models, derived from

annotated data of time series segments, enable this module to learn and recognize patterns indicative of anomalies within the manufacturing processes. The detection models are stored both on-premises and in the cloud, facilitating efficient access and utilization. The learning process pipeline is depicted in Fig. 4.

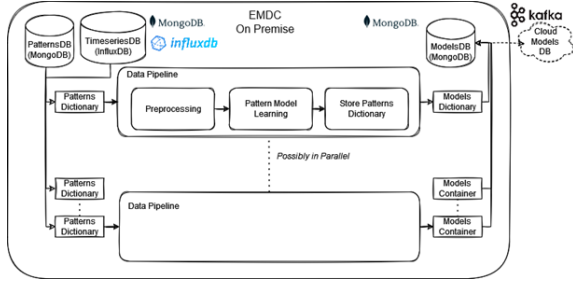


Fig. 4: Learning Pipeline on BRAINE EMDC

Online Detection Module: The Online Detection module utilizes a pre-trained detection model at the commencement of its operation. It subscribes to streams of sensory data from the modeled logical devices, which represent collections of interconnected sensory data streams, such as temperature, speed, and position data from each robot. As the Online Detection module continuously receives live data streams through Apache Kafka topics, it undergoes a series of detection steps, including pre-processing, segmentation, and pattern detection. The module then presents the resulting detections, communicating its findings to a Digital Twin (DTwin) deployed on a local BRAINE EMDC. The communication is efficiently compressed into a stream of single discrete events per detected pattern and then uploaded to a cloud-based DTwin, greatly reducing the amount of data transfer. The detection pipeline overview is depicted in Fig. 5.

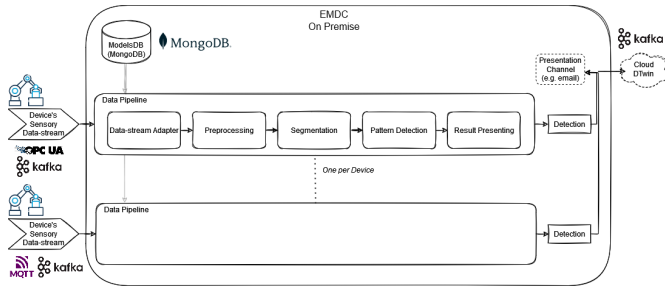


Fig. 5: Detection Pipeline on BRAINE EMDC

The deployment of these modules takes place within containers on the BRAINE platform. Each module operates within a separate container, ensuring modularity and scalability. Leveraging BRAINE's optimized scheduler and SLA Broker, these applications gain a crucial advantage in a latency guarantee, enabling online detection of operational states of machines with precision. BRAINE's telemetry system is harnessed to monitor and present data on the data processing pipelines during Motif Discovery and Model Learning, as well as to provide health monitoring for the online detection process.

V. EXPERIMENTAL EVALUATION

Multi-agent Platform: Multi-agent planning embraces uncertainty as an inherent part of the system. Changing the topology of the production line when it is composed of agents is more than 10 times faster compared to a traditional centralized approach. This was validated with the CIIRC testbed at CTU. It relates not only to the topology itself but also to the capability to adjust the agents' plans after a change in production is introduced. Another scenario relates to adding a new instance of a machine (e.g. a robot), which is realized by spawning a new container with an agent representing the machine. The system reconfiguration is done automatically after the agent is registered in the agent repository. This represents a dramatic change in the time it takes compared to a traditional manufacturing approach, with the key improvements being:

- Reduction of robot topology reconstruction from 10 days to 5 hours
- Scaling up agents in an existing topology now takes a few seconds, where it previously required several hours,

In addition, robotic image recognition for picking manufacturing parts was also evaluated while running as processes on the BRAINE platform, achieving the following key results:

- Robotic image recognition achieved picking approach accuracy of above 99 %, exceeding the KPI target of 90 %,
- Robot calibration accuracy of 0.4 cm on the XY axis and 1.9 % angle deviation exceeded the results obtained using the robots built-in calibration system, which achieved 1.0 cm and 2.5 % angle deviation.

Pattern Detection: Throughout the testing and evaluation, the MOD tool demonstrated its efficiency in detecting anomalies within manufacturing processes. By harnessing the power of motif discovery, the MOD tool provides a comprehensive and real-time analysis of operational states. Our evaluation and testing setup consisted of two testbeds separated geographically by 500 kilometers connected through a VPN tunnel, forming a distributed BRAINE cluster. In the CTU testbed, the robotic manipulators were measured with a low latency on-premises container, while the Discovery and DTwin were deployed to other cluster nodes in the second testbed. The distribution of testing between two locations showcased the potential for remote evaluation, emphasizing the adaptability of the tool for various hardware setups and coping with low-latency requirements.

Fig. 6 shows an example of the data collected by the MOD tool, and the automatically discovered motifs in the time series data of a robotic movement of axis 5 during an arbitrary manufacturing operation. The dark blue and the light blue line are the first and the second occurrence of the motif in the time series. They are plotted aligned to each other to emphasize their discovered similarity. These patterns together with other patterns form a training dataset that is used for training the detection models. Details on this topic is out of scope of this paper. For more detail on detection model training, see [13].

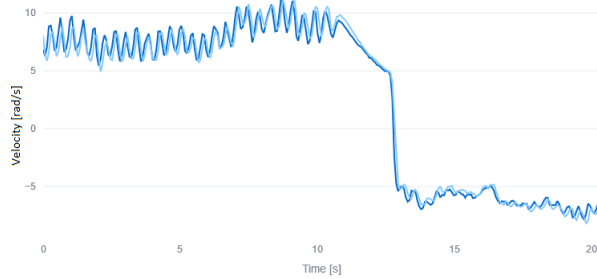


Fig. 6: Discovered motifs in robot (velocity) on axis 5

With respect to latency in detection, it depends on the particular robot movement that was learned. For one particular experiment, we created a robot operation 7.3 seconds long (sampled each 73 ms) for detection. In addition, another important focus of the MOD tools is the reduction of data-transfer. It is dependent on the dataset and defined as the ratio of amount of the raw data collected over amount of data transferred. For these metrics, we achieved the following results:

- The detection module achieved a latency of (418 ± 65) ms plus 25 % of the motif length. This result indicates that the detection models are precise enough to be able to detect an operational motif, while only seeing a short part of it.
- It also achieved a data reduction ratio of 0.17, which exceeds the targeted ratio of 0.2. This was achieved thanks to precise models, which are normally complicated to process. However, the BRAINE platform was shown to be able to support such complex models.

VI. CONCLUSION

The BRAINE project addresses the challenges posed by the rapid proliferation of big data and compute-intensive AI applications in Edge computing applications. This paper presents a comprehensive edge computing architecture developed under this project, designed specifically to meet the rigorous demands of such future Edge applications. By decentralizing computation tasks, optimizing application placement, and introducing advanced workload scheduling algorithms, this architecture enhances the capacity of edge computing systems to handle high-power, highly distributed edge use cases.

Through an empirical examination of the deployment of this architecture in an AI digital twin and Motif Discovery for factory robotics, we demonstrated substantial improvements in key performance metrics, including deployment speed, latency, robotic accuracy, and data processing. These findings underscore the potential of the proposed architecture to catalyze innovations not only within the realm of Industry 4.0 but also in various other Edge computing domains with diverse requirements.

In conclusion, by harnessing the power of the BRAINE architectures' cloud-edge continuum, advanced AI applications, and multi-agent systems, this architecture paves the way for future high-powered edge deployments and digital twin technologies. The successful results from this industry use case present a promising outlook for the application of such systems in manufacturing. The flexible and modular nature of the BRAINE architecture also allows it to be adapted to various use cases with different requirements, showcasing its potential for broader application.

REFERENCES

- [1] BOYES, H., AND WATSON, T. Digital twins: An analysis framework and open issues. *Computers in Industry* 143 (2022), 103763.
- [2] FIROUZI, F., FARAHANI, B., AND MARINŠEK, A. The convergence and interplay of edge, fog, and cloud in the ai-driven internet of things (iot). *Information Systems* 107 (2022), 101840.
- [3] GANDOMI, A. H., AND HAIDER, M. Beyond the hype: Big data concepts, methods, and analytics. *International Journal of Information Management* 35, 2 (2015), 137–144.
- [4] GÜLTEKIN, Ö., CINAR, E., ÖZKAN, K., AND YAZICI, A. Real-time fault detection and condition monitoring for industrial autonomous transfer vehicles utilizing edge artificial intelligence. *Sensors* 22, 9 (2022), 3208.
- [5] HRADECKÝ, P., JANŮ, V., BURGET, P., JOCHMAN, T., AND BECKER, T. Description and evaluation of production goals. In *2023 IFAC World Congress* (2023).
- [6] JIRKOVSKÝ, V., ŠEBEK, O., KADERA, P., BURGET, P., KNOCH, S., AND BECKER, T. Facilitation of domain-specific data models design using semantic web technologies for manufacturing. *iiWAS2019, Association for Computing Machinery*, p. 649–653.
- [7] KADERA, P., JIRKOVSKÝ, V., AND MAŘÍK, V. *Revival of MAS Technologies in Industry*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2021, pp. 131–144.
- [8] KHAN, W. Z., AHMED, E., HAKAK, S., YAQOOB, I., AND AHMED, A. Edge computing: A survey. *Future Generation Computer Systems* 97 (2019), 219–235.
- [9] LINARDI, M., ZHU, Y., PALPANAS, T., AND KEOGH, E. Matrix profile goes mad: variable-length motif and discord discovery in data series. *Data Mining and Knowledge Discovery* 34 (07 2020).
- [10] LU, Y., LIU, C., KEVIN, I., WANG, K., HUANG, H., AND XU, X. Digital twin-driven smart manufacturing: Connotation, reference model, applications and research issues. *Robotics and Computer-Integrated Manufacturing* 61 (2020), 101837.
- [11] NOVÁK, P., AND VYSKOČIL, J. Digitalized automation engineering of Industry 4.0 production systems and their tight cooperation with digital twins. *Processes* 10, 2 (2022), 404.
- [12] RON, M., AND BURGET, P. Stochastic modelling and identification of industrial robots. In *2016 IEEE International Conference on Automation Science and Engineering (CASE)* (2016), pp. 342–347.
- [13] RON, M., BURGET, P., AND HLAVÁČ, V. Parameter continuity in time-varying gauss-markov models for learning from small training data sets. *Information Sciences* 595 (2022), 197–216.
- [14] SATYANARAYANAN, M. The emergence of edge computing. *Computer* 50, 1 (2017), 30–39.
- [15] SHI, W., CAO, J., ZHANG, Q., LI, Y., AND XU, L. Edge computing: Vision and challenges. *IEEE Internet of Things Journal* 3, 5 (2016), 637–646.
- [16] SINGH, R. K., BERKVEN, R., AND WEYN, M. Agrifusion: An architecture for iot and emerging technologies based on a precision agriculture survey. *IEEE Access* 9 (2021), 136253–136283.
- [17] TAO, F., CHENG, J., QI, Q., ZHANG, M., ZHANG, H., AND SUI, F. Digital twin-driven product design, manufacturing and service with big data. *The International Journal of Advanced Manufacturing Technology* 94, 9-12 (2018), 3563–3576.
- [18] WANKHEDKAR, R., AND JAIN, S. K. A brief survey on techniques used in discovering time series motifs. *SSRN Electronic Journal* (2020).
- [19] ZHOU, B., XU, L., WANG, Y., AND TAO, F. Edge intelligence: Paving the last mile of artificial intelligence with edge computing. *Proceedings of the IEEE* 107, 8 (2019), 1738–1762.