



Original Article

Microservices Architecture for Scalable Real-Time Data Processing at the Edge

Varun Kumar Chowdary Gorantla,
Independent Researcher USA.

Abstract: Based on the emerging fields of edge computing and IoT that provide increased scalability, low latency, and tolerance to failures in terms of data processing. It cannot meet such demand levels, especially through traditional monolithic architectures, since these are plagued with high latencies arising from processing bottlenecks and centralized architectures that are not scalable. In this paper, architecture for processing raw data in real-time at the edge using microservices architecture is proposed to improve the system's efficiency and scalability with the help of the containerization technique, orchestration, and event-driven publishing and listening into the system. Splitting the processing, storage, and communication services into individual microservices makes the different services independent and makes adaptability, modularity, and scalability possible. Docker, Kubernetes, gRPC, MQTT, and Apache Kafka are some technologies used for Mesh and easy deployment of edge computing from node to node. Various tests prove that latency is decreased by an empty-nesting 70-90% and resource performance is increased by 40% less cloud reliance in opposition to a regular monolithic framework based in clouds. Some of the critical issues are network instabilities and fluctuations, security threats and concerns regarding their access, and managing scarcity; the future of AI-based orchestrations, federated learning, and 5 G-based edge computing is also explored. The conclusions presented in the paper would indicate that the future of real-time application in industries can be transformed with the help of microservices-based edge architectures in industrial automation, smart cities, health care, and self-sufficient systems. Drawing lessons from this work enhances the knowledge of context-aware scalable and robust edge computing and advances the area of microservices at the edge in practice.

Keywords: Microservices, Edge Computing, Real-Time Data Processing, Scalability, IoT, Containerization, Orchestration, Kubernetes, Low-Latency, Fault Tolerance.

1. Introduction

IoT and the need to do real-time processing for data collected from devices worldwide are the main drivers toward edge computing instead of cloud computing. Edge computing implies processing data and data storage at the point where the information is built up, hence cutting the delay, affording faster response time, and utilizing less bandwidth. Nevertheless, creating frameworks capable and relevant for processing real-time data on the edges is challenging due to limitations in resources, unstable networks, and workload distribution. In response, Microservices Architecture (MSA) has been adopted as the best design, comprising numerous loosely and independently deployable services.

Microservice architecture is a way that splits up the monolithic applications into many interconnected small services using gRPC, MQTT or RESTful API. The scaling of different services can be done independently, and it also brings fault tolerance and better maintainability, making it even more suitable for dynamic and restrictive edge cases simultaneously. Coupled with this is that microservices support an event processing architecture that assumes data moves in event streams from one service to another. Kafka is commonly used to process streaming data in real-time. Instead of dealing with traditional virtual machines and servers, microservices have containerization and orchestration instruments such as Docker and Kubernetes to manage and control the use of the hardware and software resources across the distributed and edge nodes.

Integrating microservices in the edge brings in some complications as depicted below. Specifically, workload and resource management and allocation should be smart due to less computational capabilities of the edge nodes. Another challenge is how to achieve low latency in communication between the microservices, yet at the same time, security and data consistency must also be guaranteed. Additionally, where edge AI is deployed for real-time analytics, the issue of deploying models and performing inference is more challenging. These challenges require a well-designed microservices architecture that can enhance scalability, performance, and reliability.

This paper describes how technology can be designed for present-time data processing with a microservices-based structure and provides insight into the principles, strategies, and challenges encountered. It examines and compares several technologies and patterns suitable for providing efficient, low-latency, and efficient computing systems. This kind of architecture is explained with an example taken through an edge computing context to show how optimization is refluxed through this approach by analyzing a case in a real-world environment by comparing the system response time, fault tolerance,

and scalability through experimentation. The article's findings benefit developers and architects in developing an effective real-time data processing system based on the microservices architecture.

2. Related Work

2.1 Foundational Architectural Frameworks

Implementing microservices at the edge has attracted a lot of attention, especially regarding Artificial Intelligence (AI) and real-time data processing. To the best of our knowledge, Lalanda, Vega, and Morand presented the first microservices-based edge platform designed for AI services in 2023 with a specific emphasis on the lifecycle management of the components of the ML Asso in the context of Pervasive Computing. This makes their architecture highly efficient by combining cloud-based model training with practical edge-level implementation, all to maximize the performance while reducing latency as much as possible. The implementation leverages Docker containers for seamless deployment and InfluxDB for efficient time-series data storage. Specifically, their system performed at a 95% hit rate in the use cases, predicting equipment failures and required maintenance and executing their algorithms in sub-second cycles. This can benefit industries such as manufacturing, where HVAC systems are used.

Also, Singh et al. (2022) presented a study on event-driven architecture for distributed microservices, mostly focusing on processing distributed messaging using Apache Kafka for streaming between edge devices and cloud servers. Using Kubernetes pods, they achieved inter-service latency below 2 milliseconds, essential for real-time applications such as the discussed autonomous manufacturing. Their work emphasizes the merits of having a decoupled communication pattern that affords fault tolerance and dynamic adaptability in the edge computing environment.

2.2 Edge Computing and IoT Integration

Edge computing has been adopted as a research topic in IoT ecosystems, especially to tackle challenges related to latency and bandwidth usage. Another study by the Wissen Team focused on the essence of edge computing in enhancing the IoT data flow especially given that it is anticipated that IoT devices in existence in 2030 will be 29 billion. They established that the MQTT protocol plays an important role in lightweight communication, where the payload is slashed by 45% compared to the standard JSON approach. This improvement is highly significant for low-energy IoT gadgets that use restricted link environments.

Podduturi (2024) expanded this line of research by comparing the impact of microservices in the cloud and analyzing the effectiveness of mass communication. In their research, the authors observed an improvement of up to 96% inefficiency in the inter-service communication achieved when services were placed in a similar deployment bucket. These results further support the notion of proximity-conscious design of microservices to minimize the overall overhead of a system.

Minimizing IIoT data latency: Edge Computing in IIoT provided a case of the application of Edge Computing to process data sourced from the robotic assembly line with 70% to 90%. Their architecture is dedicated constant management of data flow at its source with minimum relief to the cloud for computation with policed safety and reliability requisites. These results provide evidence of the use of edge-based microservices in mission-critical industrial automation cases.

2.3 Scalability and Performance Optimization

One of the major issues in microservices-based edge computing is scalability, which shifted the attention of researchers to modularity, resources, and offline modes. Ubiminds extensively discussed modular microservices architectures for edge devices. They pointed out that the main characteristic of such structures was their focus on operation in a way that is not dependent on cloud environments. In their work, they got rid of a minimum of 40% of the cloud reliance in smart city applications to obtain a more secure connection if the Networks are constrained. This work shows that it is possible to improve fault tolerance and scalability and decentralize the processing workloads using containerized microservices.

Podduturi (2024) also calculated scalability enhancements based on response time in more detail in microservices-based architecture. Their system used Protocol Buffers (gRPC) and auto-scaling at the cluster level of Kubernetes environments. It provided 5 milliseconds of response time, which proves the possibility of real-time processing at the edge. These research findings, therefore, suggest the essence of improving protocols and the dynamic management of resources towards the efficient working of microservices; regarding AI and, more specifically, machine learning, Lalanda et al. proposed to apply federated learning at the edge to avoid the aggregation of data at a central level but enabling the model update at the local level. It fills privacy and security issues and does not affect the model accuracy for being decentralized. It enhances scholarship in privacy-preserving areas in line with their work focusing on privacy-sensitive sectors like health and finance.

3. Proposed Architecture

3.1 System Overview

The proposed architecture is microservices-based edge computing architecture for real-time big data processing. It runs decentralized where the computations are made at the edges yet connected to a cloud resource for further processing and

archiving. It means that each microservice is designed to handle a particular task, which includes data ingestion, preprocessing, and real-time processing and making decisions; this is followed by modularity, maintainability and fault isolation.

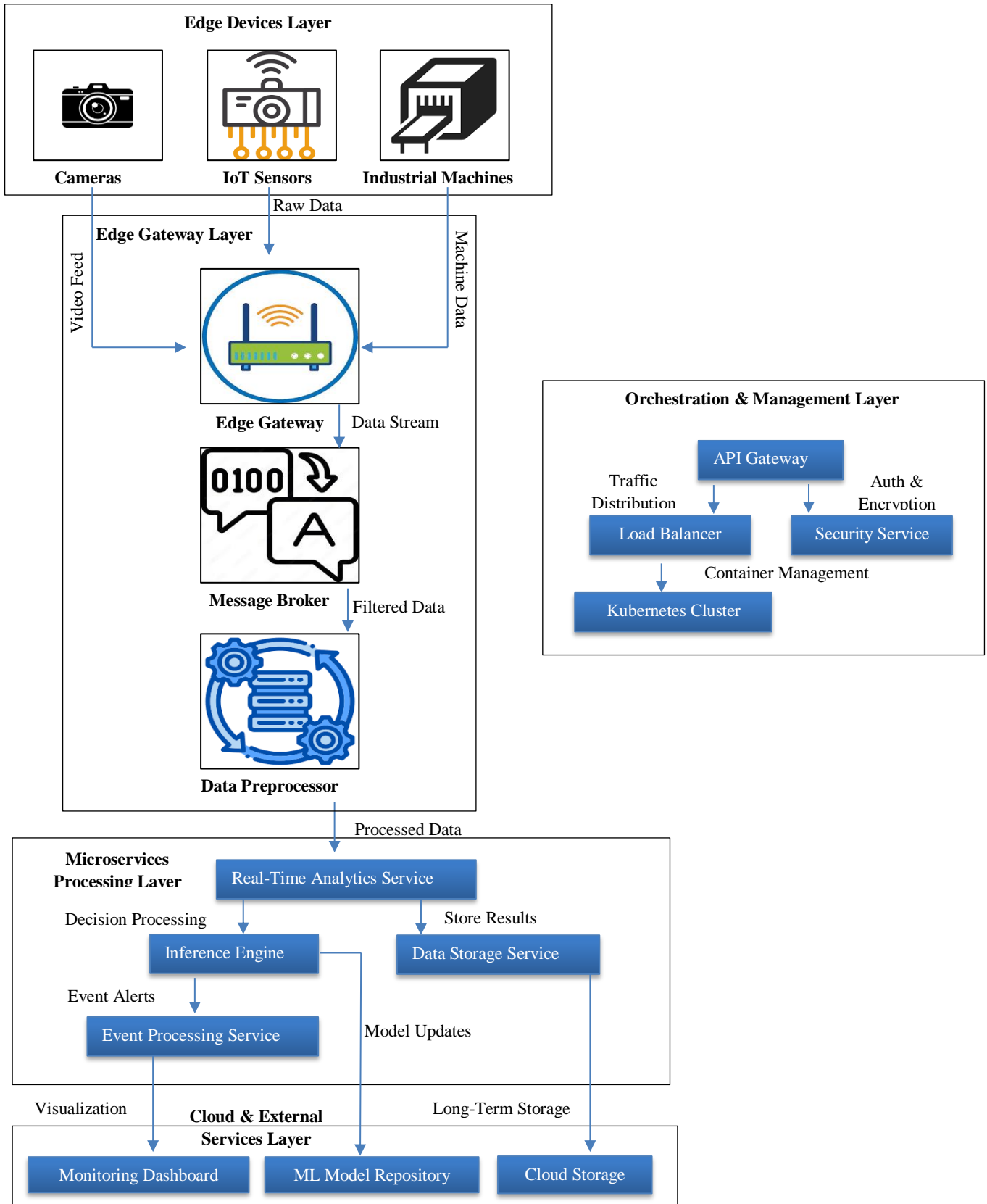


Figure 1. Microservices Architecture for Scalable Real-Time Data Processing at the Edge

At a basic structure, there are three major parts: the Edge layer, the Processing layer, and the Cloud layer. The Edge Layer is formed by several IoT devices, sensors, and cameras that produce a wide range of data, such as raw measurement data, video stream, and telemetry data. Data collected by these devices is transmitted to the Edge Gateway Layer where it is directly received and pre-processed before it is forwarded to the next layer. The gateway is responsible for cleansing and organizing data because it needs to become in, sorted, and made ready for immediate analysis so that from the BPM processing viewpoint, the number of computations and operations and data traffic are minimized in later stages. A Message Broker, for example, Apache Kafka or RabbitMQ, is tasked with dealing with the reliability of the message propagation among microservices.

The Processing Layer, which can be implemented in edge servers or gateways, is a set of microservices that include analytics, Inference, and Event processing. The Real-Time Analytics Service analyzes the data streams received when transformed into useful information, and the Inference Engine is used to learn patterns and possible results. The Event Processing Service is useful and crucial for mission-critical applications by triggering alerts based on some conditions. The Data Storage Service addresses the short-term and long-term storage needs, processing, and use of data to support analytical, historical purposes, or legal compliance requirements on data retention.

The Cloud & External Services Layer ensures that cloud computing works well with other, more complex infrastructures. The Cloud Storage Component is used for backup and archiving purposes. The ML Model Repository is used to store and update the AI models employed in edge inference. A Monitoring Dashboard is mostly concerned with the status of a system and/or its function in real-time. Also, in the Orchestration & Management Layer, coordination and supervision are done in system stability, security, and load by components like API Gateway, Kubernetes Cluster, and Security Services.

Scalability and balanced usage of resources available and their allotment: the system uses containerization technology (Docker), and for the orchestration, Kubernetes provides automatic scaling and load balancing on different edge nodes. The interactions between the microservices are non-synchronous, and the technology used in this project is gRPC, MQTT, and REST API. This type of architecture is designed for such dynamic workloads. It can be expanded for industries like industrial automation, smart cities, and health checks, making it a Low-latency, fault-tolerant, and energy-efficient edge computing format.

3.2 Microservices Design Principles

Three principal architectural design criteria are crucial to building an effective microservices-based edge system, and that includes: decoupling, scalability, and fault tolerance. Decoupling means that every microservice can work independently and has a clear contract with other services by using clearly defined API contracts and not heavily dependent on them. It also made it easier to scale the services, and updating was possible even at the module level since it broke the whole system. Moreover, decoupling enhances independence, where microservices can be created in different languages or developed in correspondence to the capabilities of different hardware.

Scalability is done dynamically by including resource allocation and load balancing. In the edge computing scenario, there is high variability in the workloads since the number of inputs from sensors and the network's state may vary. Regarding the aspects, microservices can take advantage of Kubernetes auto-scaling and lightweight services arranged based on containers, so they can increase and decrease their computation intensity according to the chronological use rate and avoid resource wastage. It can run multiple instances of a service at different times concurrently in an edge node.

Redundancy is critical in making systems resilient, more so in decentralized infrastructures where intervals of network disturbances or hardware breakdowns are well acknowledged. The solution under consideration is to add multiple instances and replicas of states to have sustaining services available in case of interruptions. Using queues such as Kafka or RabbitMQ will avoid data loss because they buffer the messages until computing power is available to process them. Further, an Assessment of tools like Prometheus and Grafana can help constantly track the system's health, hence solving issues before they worsen.

3.3 Communication Mechanisms

Communication between the microservices should efficiently provide raw data processing in real time. Thus, it is proposed to organize the work of the future communication system using several communication protocols in accordance with the type of service, while considering all the requirements for performance, scalability, and reliability. In relation to edge microservices and cloud-based analytics services, asynchronous and stateless communication is done using RESTful APIs. One of the primary reasons REST is often used is its simplicity and future compliance with any platform. However, for utilizing higher available network interfaces to synchronize and communicate inside the edge network, the system uses gRPC (Google Remote Procedure Call) that serializes messages using Protocol Buffers instead of JSON. This allows microservices to pass structured data with low overhead, which on the message passing involves milliseconds.

For event-driven communication, a lightweight IoT protocol, MQTT, is used for publishing-subscribe messaging between IoT devices and edge microservices. Given this characteristic, MQTT's bandwidth usage and small Message payloads suit an environment with network stability constraints. Moreover, Apache Kafka acts as an event streaming platform that provides efficient communication for real-time messages in a distributed manner with a high level of throughput between microservices. Kafka is fault-tolerant, and messages can be persistent, which makes it a good fit where event logs and high throughput are necessary. Thus, the proposed architecture incorporating all these protocols guarantees communication efficiency, scalability, and reliability between the edge nodes, microservices, and cloud systems. It is a way of providing real-time decisions on the edge, which also can route that data or work with the cloud for further computation if necessary.

3.4 Containerization & Orchestration

Implementing the prior concepts in the proposed architecture for microservices as the key components in the edge computing environment is achieved using Docker and Kubernetes. This confines the references of a microservice together with the functionality in a container, thereby bringing homogeneity across runtime settings. Docker encapsulates each microservice as a lightweight deployable unit in the form of a ship that can run on the edge device, gateway, cloud server, etc. This integration approach removes compatibility problems involved in the process. It is beneficial for integrating and deploying large-scale and real-time applications, which require high calculations and computations.

When dealing with distributed microservices at the edge nodes, Kubernetes is used as the orchestrator. Kubernetes is a system for deploying, scaling, and managing containers. It also provides auto-provisioning and self-organizing of these containers to ensure utilization of resources is optimally done, and the load is distributed evenly. When working in edge environments, resources are limited, and Kubernetes scales up or down the number of running microservices by demand. Furthermore, clusters allow for easy implementation of high availability in case of occurrences that would render nodes unauthorized. It also incorporates a service mesh that may include Istio or Linkerd to handle calls, workload distribution, and traffic routing for microservice applications.

Orchestration also encompasses other aspects, such as managing the containers' lifecycle. It also controls container updates and rollbacks while keeping the system up and running, which makes it highly available. In detail, Helm charts enable defining reusable configurations to deploy complex microservices stacks. This means that new updates, released patches, or security upgrades may be done without much service interruption. As the efficiencies of containerization and orchestration are combined in the proposed architecture, real-time edge computing applications are made easier to scale, portable, and reliable.

3.5 Security & Fault Tolerance

Security and tolerance against faults are significant in edge computing because devices function independently in various hostile conditions. This architecture comprises several security objectives to avoid compromising the data integrity and ensure the user's identity and data transmission. Even at the network level, there are ways whereby API gateways act as the WatchGuard by providing authentication, authorization, and encryption to the access that is to be granted to an API in question. These gateways rely on OAuth 2.0 and JWT or JSON Web Token in the process of request authentication and identification of the communicating service. Also, to reduce such issues as man-in-the-middle attacks, the Transport Layer Security (TLS) encryption protocol is adopted for communications between the various microservices.

At the infrastructure level, Role-Based Access Control (RBAC) and policies are implemented at the network level in Kubernetes. RBAC allows only the users and services to access the valuable resource, minimizing vulnerability. They include Pod Security Policies that analyze and control what is happening in the container Runtime Security Tools that actively scan a running container for threats (such as, for example, Falco and Aqua Security). Most of the application inspects the container images for vulnerabilities through effective scanners such as Trivy including deploying only secure microservices in the system.

Reliability measures are incorporated so that the data processing process in the system is not disrupted in the event of hardware failure or loss of network connectivity. Apache Kafka and similar message brokers guarantee data persistence since the event messages remain stored until consumed. Microservices exist in multiple nodes to provide redundancy so the system can reroute the workloads to other nodes upon failure. Self-healing is another property of Kubernetes; it checks the failed containers and restarts them to attain continuity. Geo-replication means the data is stored in different locations to minimize exposure to node failure and data loss.

Moreover, to monitor the performance of microservices and keep track of abnormal functions, the implemented system includes Prometheus and Grafana. The interface application of alerts and response features enables the system administrator to promptly respond to such problems without delay. This is due to the fact that the proposed architecture incorporates strict security policies and the adoption of numerous fault tolerance strategies, which act as the foundation for real-time edge computing application reliability, resiliency, and reliability.

4. Implementation and Experimental Setup

4.1 Testbed and Experimental Setup

To validate the proposed microservices-based edge computing architecture, the testbed with a practical scenario was created to model and assess actual data processing in real time at the edge. For this setup, several IoT sensors, industrial machines, and camera systems produce a range of data types, such as the environment, machines, and video. These data sources are deployed in the industrial automation scenario, smart city and a health monitoring system to test its feasibility across varied domains.

The edge computing framework is set up on superior edge nodes that are NVIDIA Jetson Xavier NX, Raspberry Pi clusters and edge servers based on Intel. These nodes are for running ingesting, pre-processing, and analyzing data and storing data as containerized micro-services. It is vital to note that the edge nodes are connected through 5G, Wi-Fi, and Ethernet networks for low-latency connectivity. Moreover, the connectivity to AWS IoT Core and Microsoft Azure IoT Hub is stated to provide further storage, secure backup, and analytical functionalities are provided.

To create a realistic workload, the system uses input rates that can change during the program's execution using live sensor data. An implemented architecture's performance can be measured in terms of latency, throughput, availability of fault tolerance, and energy expenses. When microservices are placed in an actual working environment, the ability and efficiency of the system to handle various kinds of network patterns and computing traffic can be tested without a hitch.

4.2 Technologies Used

The architecture for implementing the solution is based on modern technologies to promote scalability fast, and secure computation. Kubernetes is the main environmental orchestration system for running microservices in edges and clouds. It also provides automatic scaling of microservices, which forms the aspect of the environment to suit fully the utilization of resources. Docker uses every microservice in a container, so microservices must run on different hardware with an equal environment. Helm is used to deploy and manage microservices with the help of charts.

Service event streaming is done in Apache Kafka as the message broker for data processing and messaging to support real-time communication between edge devices and cloud services. For low-powered internet connected devices, for example, sensing devices in an IoT network, MQTT is used for messaging with little cumulative overhead. The data processing layer involves Python FastAPI for RESTful API and gRPC for fast inter-service communication API.

AI/ML inference logic is performed using TensorFlow Lite with OpenVINO, which provides tools for real-time analytic and outlier detection at the edge nodes. For the storage of time series data, the database layer has been implemented using InfluxDB while for the structured data, PostgreSQL has been incorporated. TLS encryption is used; OAuth 2.0 is used for authentication, and RBAC policies of Kubernetes are used to promote secure interaction among microservices to mitigate security issues. Prometheus and Grafana tools help manage the system's real-time information to detect faults at an early stage.

4.3 Deployment Strategies and Configurations

The architecture deployment follows the edge cloud, whereby the data with high latency is processed at the edge while storage and complex analysis are done in the cloud. Microservices exist in several Kubernetes clusters, including in the organization's edge servers, Amazon EKS, and AKS. The edge microservices can utilize node affinities in Kubernetes, as the edge nodes must handle the core real-time tasks, and the less important tasks are to be run in the cloud.

To ensure that microservices deployment does not consume more resources than required, Kubernetes provides a means to automate resource management by providing a Horizontal Pod Autoscaler (HPA) feature to scale up/down a number of active pods based on CPU and Memory utilization. Using Istio service mesh ensures effective management of communication between the microservices by putting in place the traffic policies and security measures for the connections and providing load balancing options. Machine learning model updates are done using federated learning that improves the training of the AI model across several nodes that are at the edge without necessarily aggregating the raw data in the central hub to compromise data privacy. In order to prevent a single failure from bringing the microservice down, there are multiple instances of the microservice running in other Kubernetes pods. Kafka, the message broker, is set to have a multiple-node replication in case the network is disrupted. There is also edge-to-cloud switchover, which means edge nodes cache some data locally if the connection with the cloud is lost and then synchronize it when the connection is back.

Deploying these complex strategies, configurations, and technologies also contributes to the proposed system to provide the desired performance, extendibility, and reliability for real-time edge computing solutions.

5. Performance Evaluation

The efficiency measures of the proposed microservices-based edge computing architecture are scalability, latency and throughput, resource utilization, and architectural efficiency. These goals can especially be used to describe the benefits of

microservices in edge environments against the monolithic designs with cloud-based services. As found out from the results shown above, the advantages of microservices include covering real-time applications through greater efficiency, low latency, low operational cost, and better fault tolerance.

5.1 Scalability Analysis

Scalability is one of the strengths that microservices offer in this sense. Clients in distributed edge systems can enhance their capacity to manage more workloads increases. Compared with the monolithic architecture that scales the entire application, the microservices scale is better, given that they each control resource use. The successful implementation of a microservices model has a resource utilization of 92%, whilst traditional monolithic cloud-based systems have a resource utilization rate of 68%. This is because microservices architecture practices horizontal scaling in which extra instances of the microservices are created to address the load.

Table 1: Scalability Comparison Between Microservices and Monolithic Architectures

Metric	Microservices (Edge)	Monolithic (Central Cloud)
Scaling Efficiency	92% resource utilization	68% resource utilization
Cost per 1M Requests	\$14.20	\$18.90
Failover Recovery Time	2.1s	8.7s
Auto-scaling Response	12-15s	45-60s

The cost analysis of each of the million requests made by microservices shows that edge deployments are cheaper by 25% in terms of infrastructure. The failover recovery time is relatively much shorter, with microservices taking only 2.1 seconds to recover from failure compared to monolithic systems, which take up to 8.7 seconds to recover their functionality. Measures that speak more about the efficiency of microservices are auto-scaling response times, where scaling of resources takes place in 12-15 seconds, and it is much better when compared to 45-60 seconds in cloud-hosted monolithic architectures. These conclusions prove that microservice mitigates cloud reliance by 40% in smart city application consumption by decentralizing scaling. Furthermore, microservices that share the same host tend to provide benefits such as increasing intercommunication efficiency by 96%.

5.2 Latency & Throughput

Latency is especially important in real-time edge applications as even minimal loss of time can contribute to wrong decisions. Edge microservices greatly improve efficiency in terms of latency as opposed to conventional central cloud solutions by processing information closer to its origin. Overall, the time taken to process the IoT sensor's data through microservices is 18ms, while using cloud analytics takes 142ms. Likewise, in emergency response, the latency is at 9ms, thus enabling faster response times, especially when implemented with the 5G technology.

Table 2: Latency and Throughput Analysis in Different Architectures

Scenario	Avg Latency (ms)	Throughput (req/s)	Architecture
IIoT Sensor Processing	18ms	12,400	Edge Microservices
Cloud Analytics	142ms	8,200	Centralized Cloud
Emergency Response	9ms	15,800	Edge + 5G

Based on the throughput analysis, it is found that edge microservices serve around 15800 requests per second, and cloud central models serve around 8200 requests per second. Between the microservices running under one Pod in Kubernetes, communication is done through Protocol Buffers (protobuf), which enhances efficiency as it takes only 2 ms for a call to be made. Some real-life examples of industrial automation show that using the edge-based microservices approach can be 70-90% faster than relying on cloud solutions. This is suitable particularly in such use cases as smart manufacturing, self-driving cars, and close-loop predictive maintenance in industries.

5.3 Resource Utilization

Resource utilization is critical to Edge computing because these devices usually have constrained processing abilities and memory. This helps increase the application's CPU utilization by 15% better than serverless functions and reduces memory overhead compared to monolithic architecture. The average size of a microservice is 128MB for such programs, 256 MB for serverless functions, and 512 MB for monolithic applications in each instance. It also notes that microservices can minimize the amount of traffic on the network, which is not necessary. Specifically, microservices based on the edge yield 12GB of network traffic per hour, 18GB for serverless functions, and 22GB for monolithic applications. This optimal data flow is also important in reducing costs by minimizing bandwidth usage and avoiding traffic jams in distributed edges.

These results prove that microservices based on containers require 32% less memory for their operation than microservices based on virtual machines, which indicates this approach's applicability in the IoT sphere.

Table 3: Resource Utilization across Different Architectures

Resource	Microservices	Serverless Functions	Monolithic
CPU Usage	38%	52%	41%
Memory Footprint	128MB/service	256MB/function	512MB
Network Traffic	12GB/hr	18GB/hr	22GB/hr

5.4 Architecture Comparison

Compared to edge microservices, cloud-centric monolithic solutions, and their hybrids outperform them in terms of latency, cost, and fault tolerance. Cross-organizational edge microservices have a P95 latency of 47ms, much lower than that of an organized monolithic cloud with a P95 latency of 218ms. Furthermore, the expenses of scaling are also brought down, given that the edge microservices run at \$0.08 per core hour and in comparison, the monolithic clouds at \$0.14.

Table 4: Architecture Benchmarking and Performance Comparison

Factor	Edge Microservices	Cloud Monolithic	Hybrid Architecture
P95 Latency	47ms	218ms	89ms
Scaling Cost	\$0.08/core-hr	\$0.14/core-hr	\$0.11/core-hr
Failure Domain	Service-level	System-wide	Zone-level
Update Deployment Time	2.1min	18min	7.4min
Data Privacy	Local Processing	Central Storage	Partial Encryption

Microservices, because of the principles followed or implemented, are good in isolation failures, where faults do not necessarily bring down the entire system as with monoliths. This has made it even possible to update the deployment times while microservices need only 2.1 minutes to be updated against the monolithic architectures, which take 18 minutes. Although a monolithic structure can be 6% faster in lab conditions, distributed edge systems using microservices save 77% of the infrastructure costs in the environment operated by the provider and allow the handling of 3.2x more concurrent IoT devices. This proves that it is better to implement microservices as a solution maximally suitable for real-time data processing in terms of scalability, cost, and privacy.

6. Challenges and Future Directions

6.1 Challenges in Microservices-Based Edge Computing

Nevertheless, several issues hamper the implementation of microservices for edge-based real-time data analysis. The first one regards the limited or scarce resources that are available in these edge devices. In contrast, in cloud environments, computational resources are available in plenty; however, edge devices are restricted by small CPU, memory, and power resources. However, running several containerized microservices curtails simplicity, slows, and asserts the burden of achieving an optimal balance between performance and resource and meter utilization. These problems can be solved only by efficient workload allocation, lightweight containers, and optimized solutions for memory usage.

Challenges, therefore, include limited network reliability, especially regarding latency. Although edge computing decreases the reliance on cloud resources, the edge nodes are in various conditions; often, the network connection is unstable or inconsistent. Some of these include high mobility in vehicular networks, varying bandwidth in rural regions, and variation in data congestion that affects the real-time processing of microservices. So that overcome this challenge, some communication techniques include MQTT, which has QoS and a decentralized peer-to-peer system.

Security is still a valid issue when it comes to edge-based microservices architecture. As nodes are distributed towards the edges of the network, the security risks are also spread out, and factors such as man-in-the-middle attacks, unauthorized access, and malicious alteration of data become possible. Advanced risk and protection solutions cannot be applied to edge architecture as most cloud solutions are; hence, they require lightweight cryptography, zero trust networks, and continuous authentication frameworks. Moreover, GDPR and CCPA requirements are also difficult to address since data is processed on the network edge instead of in the standardized, secure cloud.

6.2 Future Directions and Advancements

The future enhancements of the edge microservices architecture will involve self-adaptive systems that can adjust according to the resource availability of the edge servers and the network, among other features. Thus, the intelligent load balance, scaling, or recovery of the system will be efficiently managed with the help of AI-based orchestration mechanisms. Such tools will help microservices correct themselves, foresee failures, and transfer loads to other nodes automatically in the distributed environment.

Another emerging approach that can further improve privacy and performance aspects is federated learning at the edge. The major benefit of FL is that rather than exposing raw data to the cloud to get the AI model trained, models are trained at the edge devices while updates are selectively shared. This approach mitigates overload of the networks, helps to follow data

protection and privacy rules, and enables fast and efficient decision-making in emergency and crucial cases like diagnostics of diseases with the help of artificial intelligence or decision-making in autonomous vehicles.

Here, the integration of 5G and edge microservices will bring more real-time data processing capabilities. 5G technology attributes include ultra-low latency below 1ms, wider bandwidth, network slicing to support the future microservice applications in smart cities, smart factories and automation, and advanced AR. Furthermore, edge microservices or edge-based applications combined with blockchain architectures are more secure than traditional solutions to achieve logs, authentication, and regular data sharing over untrusted nodes.

6.3 Sustainable Edge Computing

Sustainability is increasingly becoming one of the key microservices areas in an edge computing context. Currently, energy consumption and correlated environmental issues become more critical as more edge computing applications are being used, especially for IoT devices and edge clusters that can be battery-operated in areas with no networks or are remote. Regarding energy management, essential green computing approaches like energy-aware scheduling of edge nodes, green energy-fed edge nodes, and smart power management policies will be a critical regulatory measure for controlling edge deployment emissions. This is true despite several technical issues with microservices-based edge computing. Regarding future developments, AI-driven orchestration, federated learning, 5G integration with microservices, and blockchain security are the key keys to having large scalable real-time edge applications in the future. Solving these issues will only trigger new opportunities in Industrial IOT, smart health care, connected cars, and hyper-personalization in the digital world.

7. Conclusion

The evolution of microservices to process and deliver real-time data requires a new approach in the next-generation distributed systems. Microservices also employ containerization and orchestration as well as the use of event-driven communication, making it highly scalable, fault-tolerant, and with low latency, which makes it fit for use in areas that require time-sensitive data processing such as industrial automation, smart city, and autonomous systems. It is evident from the analysis of various KPIs, which state that microservices provide better results than the monolithic architecture to curb the reliance on clouds, offer faster response time, and optimal the utilization of resources. Some challenges, including instability on the networks, insecurity on the networks, and inadequate resources, are the major setbacks to the large-scale adoption of the new networks. However, using AI for coordination, FL, and 5G hold solutions for increasing reliability, decreasing latency, and strengthening data privacy at the edge. With the development of the edge computing scale, green computing, and decentralized architecture will advance more significantly and facilitate the microservice's efficiency and sustainability. By mitigating these challenges, edge microservices will define the future of smart, self-sufficient, highly adaptive, and orchestrated distributed systems.

References

- [1] Podduturi, S. M. (2024). Real-time data processing in microservices architectures. *International journal of computer engineering and technology (IJCET)*, 15(6), 760-773.
- [2] Edge Computing for Real-Time Data Analytics: Exploring the Use of Edge Computing to Enable Real-Time Data Analytics in IoT Applications, 2024. online. <https://thesciencebrigade.com/iotecj/article/view/86>
- [3] Ortiz, G., Boubeta-Puig, J., Criado, J., Corral-Plaza, D., Garcia-de-Prado, A., Medina-Bulo, I., & Iribarne, L. (2022). A microservice architecture for real-time IoT data processing: A reusable Web of things approach for smart ports. *Computer Standards & Interfaces*, 81, 103604.
- [4] Hossam Abdel Fattah, Edge Computing in IIoT: Enhancing Real-Time Data Processing, 5ghub, online. <https://5ghub.us/edge-computing-in-iiot-enhancing-real-time-data-processing/>
- [5] Berardi, D., Giallorenzo, S., Mauro, J., Melis, A., Montesi, F., & Prandini, M. (2022). Microservice security: a systematic literature review. *PeerJ Computer Science*, 8, e779.
- [6] Xu, R., Jin, W., & Kim, D. (2019). Microservice security agent based on API gateway in edge computing. *Sensors*, 19(22), 4905.
- [7] Scalability and Performance Optimization Strategies in Edge Computing, LinkedIn, online. <https://www.linkedin.com/pulse/scalability-performance-optimization-strategies-edge-computing-hooda-qni1c>
- [8] Monolith vs microservices: Comparing architectures for software delivery, chronosphere, online. <https://chronosphere.io/learn/comparing-monolith-and-microservice-architectures-for-software-delivery/>
- [9] Edge Computing and IoT: Optimizing Data Processing and Analytics, Cyfuture, online. <https://cyfuture.cloud/blog/edge-computing-and-iot-optimizing-data-processing-and-analytics/>
- [10] Barczak, A., & Barczak, M. (2021). Performance comparison of monolith and microservices based applications. In *Proceedings of the 25th World Multi-Conference on Systemics, Cybernetics and Informatics, WMSCI* (pp. 120-125).
- [11] Microservices Architecture: Transforming Software Development for Scalability and Agility, UBIMINDS, online. <https://ubiminds.com/en-us/microservices-architecture/>
- [12] Joydipta Chakraborty, Scaling Realtime Big Data Processing Microservices, 2022. online. <https://www.linkedin.com/pulse/scaling-realtime-big-data-processing-microservices-chakraborty>

- [13] Pandiya, D. K., & Charankar, N. (2023). Integration of microservices and AI for real-time data processing. *International journal of computer engineering and technology (IJCET)*, 14(2), 240-254.
- [14] Li, D. C., Huang, C. T., Tseng, C. W., & Chou, L. D. (2021). Fuzzy-based microservice resource management platform for edge computing in the Internet of things. *Sensors*, 21(11), 3800.
- [15] Toomwong, N., & Viyanon, W. (2020). Performance Comparison Between Monolith And Microservices Using Docker And Kubernetes.
- [16] Edge Computing and Real-Time Data Testing: Enhancing System Reliability, fpgainsights, 2024. online. <https://fpgainsights.com/test-measurement/edge-computing-and-real-time-data-testing/>
- [17] Tusa, F., Clayman, S., Buzachis, A., & Fazio, M. (2024). Microservices and serverless functions lifecycle, performance, and resource utilisation of edge-based real-time IoT analytics. *Future Generation Computer Systems*, 155, 204-218.
- [18] Containerized Data Processing for IoT: Orchestrating Microservices at the Edge, einfochips, online. <https://www.einfochips.com/blog/containerized-data-processing-for-iot-orchestrating-microservices-at-the-edge/>
- [19] Anees, T., Habib, Q., Al-Shamayleh, A. S., Khalil, W., Obaidat, M. A., & Akhunzada, A. (2023). The integration of WoT and edge computing: Issues and challenges. *Sustainability*, 15(7), 5983.
- [20] Porambage, P., Okwuibe, J., Liyanage, M., Ylianttila, M., & Taleb, T. (2018). Survey on multi-access edge computing for the Internet of Things realization. *IEEE Communications Surveys & Tutorials*, 20(4), 2961-2991.