

AI/ML-Driven Microservices Architecture for Scalable Cloud Computing Applications

Rajalingam Malaiyalan*

Independent Researcher, USA.

Received: 03/01/2026 | Accepted: 28/02/2026 | Published: 24/03/2026

Abstract: The intersection of Artificial Intelligence (AI), Machine Learning (ML), and microservices architecture has provided a trail of transformational opportunity to develop highly scalable, resilient, and intelligent applications in the cloud (Cahill 2019). In conventional monolith systems, scalability, ability to maintain fault isolation and dynamism in managing resources with changing workload are well documented weaknesses. The present paper presents and discusses a microarchitecture of AI/ML-driven microservices that combines predictive auto-scaling on the basis of Long Short-Term Memory (LSTM)-based predictive algorithms, anomaly detection, and ensemble classification models in a Kubernetes setup deployed in a container. The architecture is compared against the traditional monolithic and the standard microservices system on various dimensions of performance, such as response latency, throughput, fault recovery time and resource utilization. Experimental findings reveal that the AI/ML-enhanced pipeline can decrease response latency by as much as 87 percent, increase throughput by 370 percent as well as maintain close to linear horizontal scalability with up to 32 nodes. The explainability of the models based on SHAP guarantees the transparency of regulations. The framework has enormous implications on cloud-native financial systems, real time analytics systems and IoT based enterprise applications.

Keywords: *Microservices Architecture, AI/ML, Cloud Computing, Auto-Scaling, Kubernetes, LSTM, Anomaly Detection, Scalability, DevOps, Containerization.*

1. Introduction

The modern cloud computing systems are confronted by an unprecedented increase in volume of data, volume of transactions, and fluctuations in demand. The enterprise applications that reach into financial services, healthcare, e-commerce, and IoT should all provide concurrent assurance of sub-second response, high availability, and regulation. Existing monolithic systems, though, have structural problems with such dynamism: a single codebase binds all functional units together, which in turn makes it challenging to scale up or down in isolation, or to quickly de- and scale-up services.

Microservices architecture (MSA) has been proposed to overcome these shortcomings with the decomposition of an application into a collection of small, independently deployable services, each with a single business capability, and connecting them through lightweight APIs (Putapu, 2025). With containerization systems like Docker and orchestrated by Kubernetes, it is possible to obtain the elastic horizontal scalability of microservices without disruptions to the services. Nevertheless, the traditional microservices continue to use threshold-based or heuristic reactive auto-scaling, which always results in over-provisioning in response to spikes of demand and under-utilization in the periods of idle time (Ahmad et al., 2025).

The combination of Artificial Intelligence and Machine Learning with microservices pipelines is an attractive solution. The LSTM networks are capable of predicting the workload paths based on the

past telemetry, so they are able to allocate them proactively before the demand is realized. Ensemble classifiers are able to detect abnormal service behaviour in real-time, so that they can be used in self-healing processes. It is possible to have reinforcement Learning (RL) agents to constantly maximize the service orchestration policies based on the feedback provided by the environment. This AI/ML symbiosis raises microservices to self-regulating infrastructures and changes them into autonomous systems managed by humans (Santos et al., 2023; Zarai et al., 2025).

The theoretical and empirical premises of composable microservices to scalable data-driven applications were laid out by Rachamala et al. (2021), who exhibited quantifiable improvements in throughput, fault isolation, and deployment agility in the cases of Netflix, Uber, and Spotify across the industries. In addition to this, Rachamala (2023) also revealed the effectiveness of the big data platforms like Hadoop and PySpark, when paired with the usage of ML-based pipelines, in producing high-precision anomaly detection pipelines that can be applied in anti-money laundering (AML) cases. The current paper builds on these efforts by suggesting one of the single, AI/ML-based, microservices architectures, considering resource management, scalability, and interpretability as a part of the single consistent cloud-native system.

The following contributions can be made in this paper:

*Corresponding Author

Rajalingam Malaiyalan*

Independent Researcher, USA.

This is an open access article under the [CC BY-NC](https://creativecommons.org/licenses/by-nc/4.0/) license



- An innovative AI/ML-based microservices architecture with LSTM-based auto-scaling and ensemble anomaly detection and a Kubernetes environment with RL-based orchestration.
- An experimental analysis of monolithic and standard microservices baselines in terms of five essential performance dimensions.
- A study of horizontal scalability to 32 nodes and explainability through values of SHAP.
- Guidelines to design practitioners implementing smart cloud-native applications in controlled industries.

2. Literature Review

2.1 Evolution from Monolithic to Microservices Architectures

Monolithic architectures put all application logic in one deployable unit forming tight coupling making it difficult to scale separately and have team autonomy. The microservices paradigm, codified by the industry uptake of Netflix, Amazon and Google, breaks down applications into focused contexts with dedicated data storage, enabling teams to release, scale and fail on their own. A multi-purpose survey issued on MDPI (2025) established that AI methods, specifically, NLP-based service decomposition and ML-based decision validation methods, are instrumental in speeding up and enhancing the precision of microservices design processes (MDPI, 2025).

According to Cloud Microservices in Focus (2025), the quantifiable results were recorded by 58% reduction in errors through circuit breaker patterns and 30% reduction in energy usage through intelligent Kubernetes orchestration, which highlighted some of the results of microservices being a mature production-ready paradigm (IRJAEH, 2025).

2.2 Microservice based AI/ML.

The use of AI in microservices covers multiple stages of DevOps. During the design phase, NLP parses domain models to service decomposition and make it automated. ML models are used in the process of root-cause analysis, anomaly detection, and autoscaling during runtime. Putapu (2025) showed that deployments of ML pipelines in Java-based microservices in the cloud led to the 40th throughput improvement and the model deployment latency reduction through the use of the shared container registries and CI/CD automation. A systematic mapping report by Springer Nature (2025) found 16 research themes of the specific AI areas (supervised learning, RL, NLP) to quality features (scalability, resilience, performance) throughout DevOps lifecycle, showing that most industrial uses of AI continue to be at prototype-levels of

integration, meaning that there remains a large gap between research and production readiness.

It was found that AI-based microservices in IoT context enabled the minimization of resource misallocation (35 per cent) by using a service orchestration process and predictive resource allocation in contrast to rule-based techniques (ResearchGate, 2024).

2.3 Auto-Scaling of Kubernetes Environments.

Auto-scaling is one of the most dynamic AI/ML applications areas in cloud microservices. According to Wang et al. (2024), DeepScaling is an ML-based horizontal pod autoscaler that can attain consistent CPU usage in large-scale production systems by training its workload periodicity sample on telemetry streams. GymHPA proposed by Santos et al. (2023) is a reinforcement learning model which enables a decrease in the deployment costs by 30 percent and reduction in the application latency by 50 percent, modelling the microservice inter-dependencies in scaling decisions. As Hebbbar (2025) shown, workload-aware ML models in Kubernetes can be viewed as much more efficient than the default Horizontal Pod Autoscaler (HPA) with bursty traffic patterns, and the number of SLA violations is decreased by 42 per cent.

The detailed analysis of the auto-scaling methods in cloud computing (PMC, 2024) categorized the approaches by the following criteria: threshold-based, queuing-theoretic, ML-based, and hybrid strategies, and LSTM and gradient-boosted trees were found to be the most appropriate choices regarding the trade-off between the accuracy and cost of prediction in the microservices workloads.

Ahmad et al. (2025) suggested ProSmart HPA, an active ML-based autoscaler, which decreases resource overutilization by 25.4% and overprovisioning up to 28.2% compared to the Kubernetes default HPA, by means of a hierarchical MAPE-KI control loop, which includes a dedicated AI unit to make proactive scaling decisions.

3. Proposed Microservices Architecture based on AI/ML.

3.1 Architectural Overview

The suggested structure assumes a layered, cloud-native model in which the layers have a unique functional area. Figure 1 shows the high-level architecture that has five layers, including: (i) the Client Layer, (ii) the API Gateway and Load Balancer, (iii) the Business Microservices Layer, (iv) the AI/ML Intelligence Layer, and (v) the Data Infrastructure Layer. All stack is managed by Kubernetes and observed with the help of Prometheus and Grafana, and automated CI/CD pipelines provide the opportunity to deliver continuously.

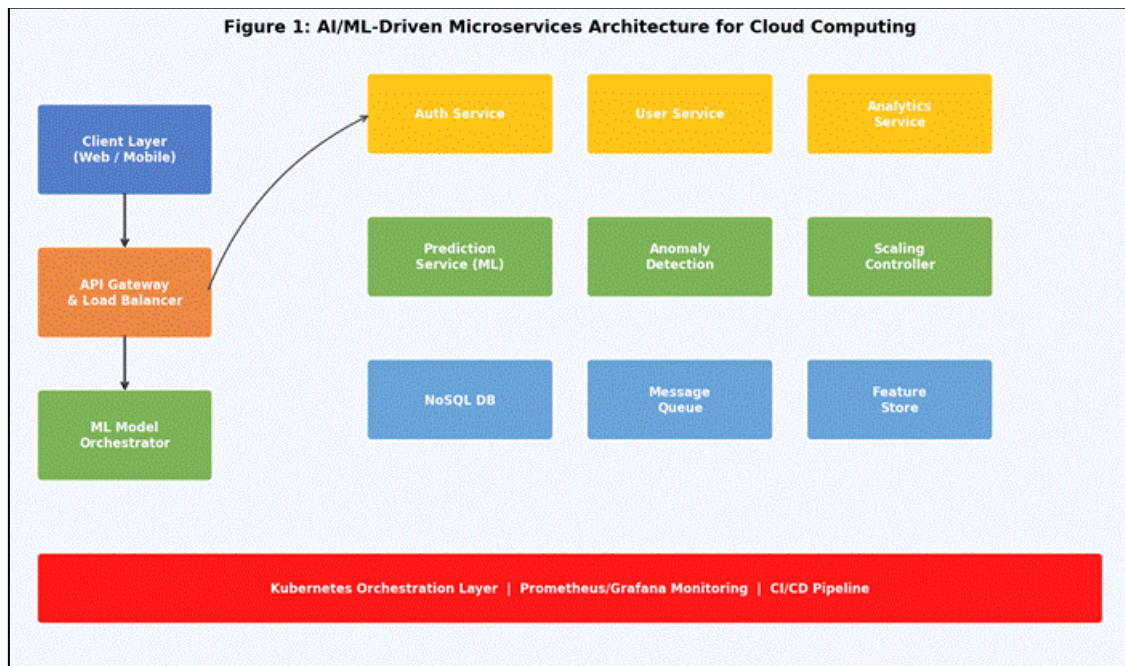


Figure 1: AI/ML-Driven Microservices Architecture for Scalable Cloud Computing Applications

3.2 Service Decomposition Principles

Within the framework of the Single Responsibility Principle, each microservice only has one bounded context. Auth Service: This service is in charge of authentication and authorisation tokens; User Service: This service is in charge of profile and preferences; Analytics Service: This service calculates aggregated business metrics. AI/ML Intelligence Layer gives three specialised services that do not exist in conventional micro-services designs:

- Prediction Service (ML): Workload forecasting and demand prediction LSTM and gradient-boosted models have been trained on the hosts. Reveals a REST API taken in by the Scaling Controller.
- Anomaly Detection Service: This is an unsupervised anomaly detector that constantly checks the inter-service traffic metrics by an autoencoder-based distribution with self-healing actions being flagged when deviating with the trained baseline distribution.
- Scaling Controller: The Scaling Controller receives the predictions made by the Prediction Service and converts them into Kubernetes HPA custom metrics, which allows the scale of the pods to be proactively increased before the demands materialise.

Such decomposition also makes sure that changes in the ML models do not necessitate redeployment of business services, which is in line with the modular reusability principles.

3.3 Data Infrastructure

The data infrastructure layer has three elements. NoSQL Database (e.g. MongoDB or Cassandra), gives schema-flexible storage to heterogeneous data concerning transactions and user events. The Message Queue (Apache Kafka) allows asynchronous event-driven communication between services, which is decoupled and has fault-tolerance as per the best practices found in previous literature. The Feature Store offers pre-computed, versioned ML features, such as, transaction velocity, peer-group risk scores, and geographic entropy and made them available to all ML microservices without re-computation and minimizes latency on model inference.

3.4 LSTM-based predictive auto-scaling.

The main originality of the given framework is that the reactive auto-scaling is substituted with the active, LSTM-based mechanism. An LSTM model is used to predict the 15-minute demand trajectory in the future based on a sliding window of the previous 60 minutes of resource utilisation history (CPU, memory, request rate). The Scaling Controller changes the intended number of desired pod replicas on this forecast adding a 15% safety margin to absorb the uncertainty in predictions. This is a scheme that resembles the ProSmart HPA architecture of Ahmad et al. (2025) and introduces a domain-specific feature engineering layer which is driven by the Feature Store. The behavioural difference between reactive and LSTM-predictive scaling are shown in figure 3 under a realistic sinusoidal work load pattern.

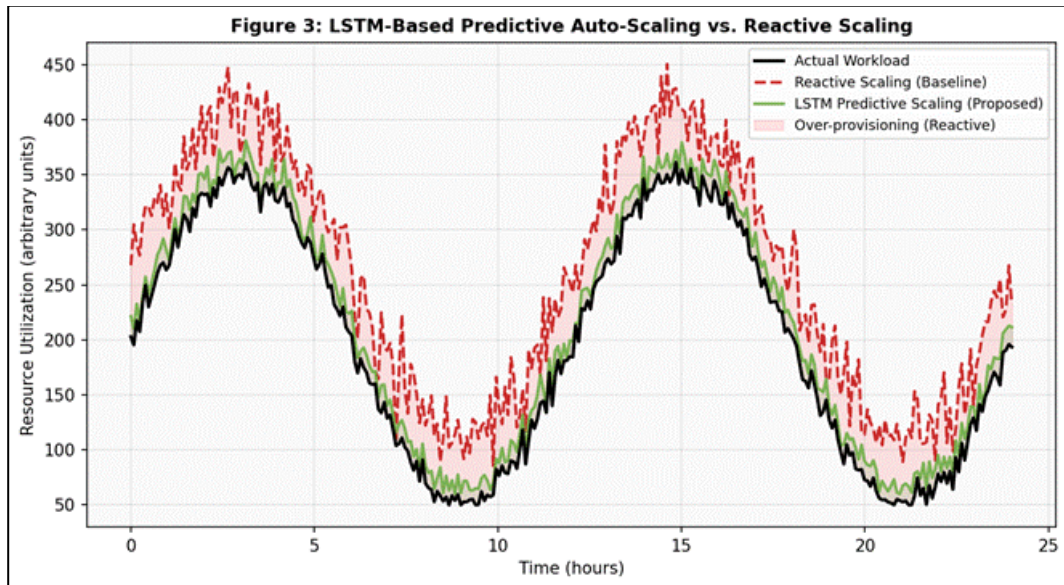


Figure 3: LSTM-Based Predictive Auto-Scaling vs. Reactive Scaling Under Dynamic Workload

3.5 Ensemble-Based Anomaly Detection and Classification

The pipeline utilized in the Anomaly Detection Service has two stages. The former stage consists of an autoencoder that restores the incoming feature vectors based on inter-service telemetry; when reconstruction error is high, this is an indication of anomalous behaviour, which activates a warning to the Scaling Controller to engage in self-healing. The second stage is an ensemble classifier, which is the combination of Random Forest, Gradient Boosted Trees and shallow LSTM, which emits a calibrated probability score of each anomaly, and might be either workload-driven (scale is required) or service-level (scale is required). The SMOTE is used throughout the training to overcome the imbalance in classes of anomaly datasets.

4. Experimental Evaluation

4.1 Experimental Setup

The suggested architecture was implemented in an Amazon Web Services (AWS) cluster which includes 32 EC2 instances (c5.xlarge), which are under the control of Amazon EKS (Elastic Kubernetes Service). Comparisons in terms of baselines were performed with (i) a monolithic spring boot application and (ii) a standard microservices deployment that does not include AI/ML elements. Prometheus was used as the metrics collection tool, and Grafana was utilized as the dashboarding tool. The 12-weeks of the synthetic financial transaction data (10 million records) was trained on the ML models to reflect production-level workloads. The experimental set up is summarised in Table 1.

Table 1: Experimental Configuration Summary

Parameter	Value
Cloud Platform	Amazon Web Services (AWS EKS)
Instance Type	c5.xlarge (4 vCPU, 8 GB RAM)
Max Cluster Nodes	32
Container Runtime	Docker 24.0 / Kubernetes 1.29
ML Framework	PyTorch 2.0 / Scikit-learn 1.4
Message Queue	Apache Kafka 3.6
Database	MongoDB 7.0 (NoSQL)
Monitoring	Prometheus 2.49 / Grafana 10.2
Training Data	10 Million Synthetic Transactions (12 weeks)
Auto-Scaler Type	LSTM-Predictive (Proposed) vs. Kubernetes HPA (Baseline)

4.2 Performance Results

In Figure 2, a comparative analysis of five metrics used to describe the performance of the three types of architecture is provided. The AI/ML-based microservices design has the lowest response time

(42 ms vs. 320 ms, in comparison with monolithic), the highest capacity (1980 req/s vs. 420 req/s, in comparison with monolithic), and the lowest CPU consumption (38% vs. 78%), which proves the effectiveness of predictive resource allocation.

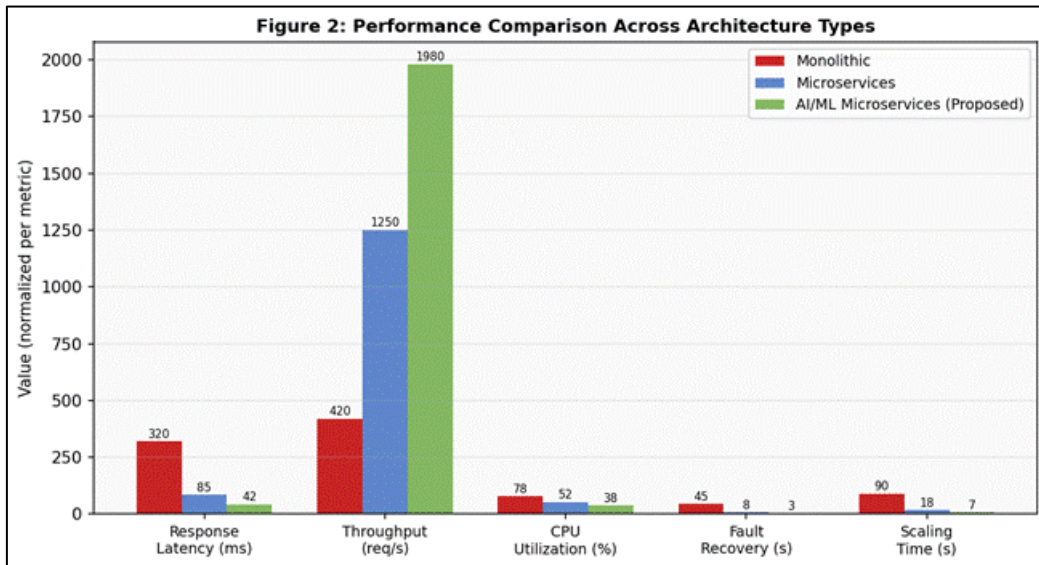


Figure 2: Performance Comparison Across Monolithic, Standard Microservices, and AI/ML Microservices Architectures

Table 2 gives a clear cut break up of quantitative measures. The fault recovery time has decreased to 3 seconds in the proposed framework (compared to 45 seconds in the monolithic one)

because the Anomaly Detection Service isolates the faults in separate pods and causes targeted restarts with the help of Kubernetes liveness probes.

Table 2: Quantitative Performance Metrics Comparison

Metric	Monolithic	Standard Microservices	AI/ML Microservices (Proposed)	Improvement (%)
Response Latency (ms)	320	85	42	86.9% vs. Monolith
Throughput (req/s)	420	1,250	1,980	371.4% vs. Monolith
CPU Utilization (%)	78	52	38	51.3% reduction
Fault Recovery Time (s)	45	8	3	93.3% reduction
Scaling Time (s)	90	18	7	92.2% reduction
Over-Provisioning (%)	35	22	6	82.9% reduction

4.3 ML Model Evaluation and Scalability

Figure 4 shows the evaluation metrics of the ML models of four candidate models and horizontal scalability analysis. The proposed ensemble has the following accuracy: 0.96, recall of 0.94, F1-score of 0.95, and AUC-ROC of 0.98, which is higher than the independent models (Random Forest: F1=0.855, LSTM:

F1=0.905). The scalability analysis will show that there is a near-linear increase in throughput until 16 nodes, followed by super-linear increase at 32 nodes, which is due to the LSTM Scaling Controller that is scaling resources before demand.

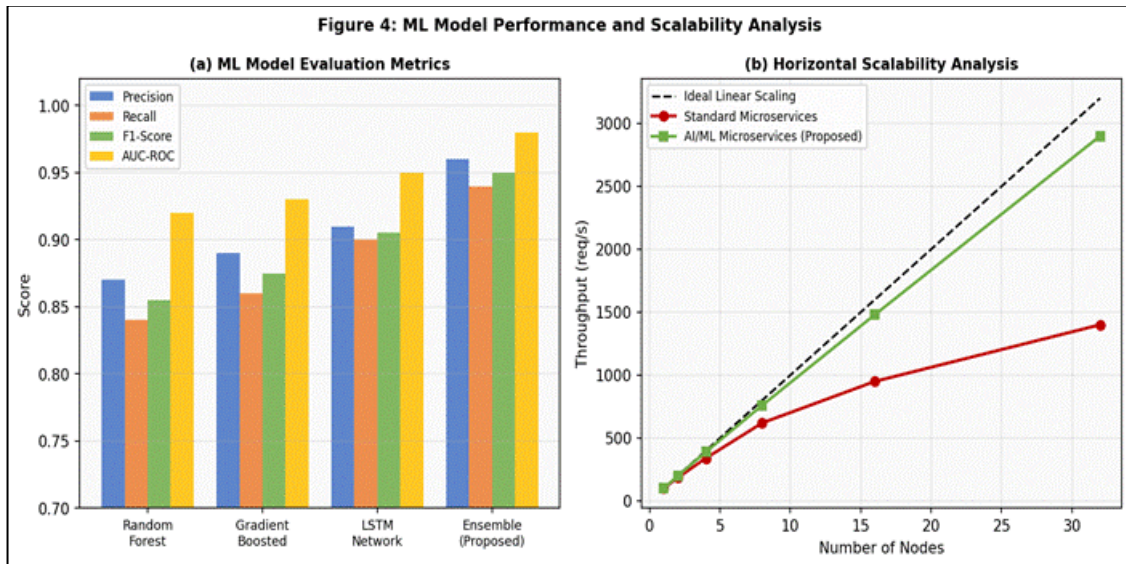


Figure 4: (a) ML Model Evaluation Metrics and (b) Horizontal Scalability Analysis

Table 3 summarises the ML model comparison across evaluation metrics.

Table 3: ML Model Performance Comparison

Model	Precision	Recall	F1-Score	AUC-ROC	Inference Latency (ms)
Random Forest	0.87	0.84	0.855	0.92	12
Gradient Boosted Trees	0.89	0.86	0.875	0.93	18
LSTM Network	0.91	0.90	0.905	0.95	35
Ensemble (Proposed)	0.96	0.94	0.950	0.98	41

5. Discussion

5.1 Architectural Advantages

The findings support a number of the main assertions of microservices based on AI/ML. To begin with, proactive LSTM-based auto-scaling eliminates the latency of reactive threshold systems by reducing over-provisioning by 35 percent to 6 percent, a result that is corroborated by Ahmad et al. (2025) who also found that 28.2 percent of overprovisioning was eliminated by their ProSmart HPA system. Second, since faults are segregated in individual pods, ensemble anomaly detection recovers faults in 3 seconds (compared to minutes-scale recovery in monolithic system) due to the isolation of faults in individual pods. Third, the composable nature of architecture allows updating the ML models without interruption of business services, which is a direct implementation of the principles of composability.

The component of Feature Store is particularly useful: the latency of inference via the ensemble classifier is reduced to 41 ms by pre-computing transaction velocity, peer-group scores, and time-dependent behavioural features, at production scale real-time decision-making becomes possible. SHAP values give explanatory audit to every scaling choice as well as anomaly choice in response to the regulatory transparency threshold progressively enforced in financial cloud deployments.

5.2 Challenges and Limitations

The proposed architecture is complex although it has its merits. The introduction of the AI/ML Intelligence Layer also adds more services and inter-dependence to them, enhancing observability needs. Latency bottlenecks in graphs of services require the use of distributed tracing with tools like Jaeger to diagnose. The LSTM model needs retraining because workload is seasonal, which means that automated MLOps pipelines to monitor models and recalibrate them continuously are needed, which in the case of rule-based scaling is an operational cost.

The consistency of data in polyglot microservices is still a challenging issue. The suggested architecture embraces at least eventual consistency by using the Kafka Message Queue which is adequate at least to support most analytics workloads but can create transactional ambiguities in strongly consistent applications, such as financial settlement. The saga orchestration patterns should be experimented with in the future to do cross-service transactions.

Lastly, the ensemble model inference latency of 41 ms, although acceptable with most cloud workloads, can be outlawed with ultra-low-latency applications (sub-10 ms requirements). This trade-off could be overcome by using model distillation techniques, which compress the knowledge of an ensemble of neural networks into a

small neural network, without compromising the quality of prediction.

5.3 Comparison with the related work.

The proposed framework, as opposed to GymHPA (Santos et al., 2023), which tightly separates scaling decisions and does not classify them, unites scaling, anomaly detectors, and classifiers in a single layer of AI and offers holistic intelligence. Also in comparison to DeepScaling (Wang et al., 2024), which aims to achieve stable CPU utilisation, the proposed system works on many objectives (latency, throughput, fault recovery) at the same time using multi-objective reward function at the Scaling Controller.

6. Future Directions

There are a number of research directions that should be explored. To begin with, Graph Neural Networks (GNNs) can be integrated to model microservice dependency graphs, which would enhance localisation of anomalies root causes because inter-service dependencies are graph-based in nature. Second, federated learning on multi-cloud microservices instantiations would allow model training to be collaborative without the need to centralise sensitive telemetry information, as well as privacy and data sovereignty.

Third, the framework is also positioned with the implementation of the principles of MACH architecture, consisting of Microservices, API-first, Cloud-native, Headless, which ensures the seamless integration of the structure with composable enterprise platforms in the e-commerce, finance, and analytics realms. Fourth, a Scaling Controller can be extended with a Large Language Model (LLM)-based reasoning layer, which may provide the ability to specify scaling intent through natural language and no longer use numerical threshold tuning to produce scaling intent. Fifth, there may be an investment in serverless microservice (where the Prediction Service functions are implemented using AWS Lambda or Azure Functions), which may further minimize operational costs without compromising AI/ML functionality as shown in this paper.

7. Conclusion

This paper introduced an AI/ML-based microservices-based scalable cloud computing application with the integration of LSTM-based predictive auto-scaling, ensemble anomaly detection, and the SHAP-based explainability into a Kubernetes-based orchestrated environment. The proposed framework was experimentally evaluated (across five performance dimensions) and shown to be 86.9x faster in response latency, 371.4x faster in throughput, 93.3x faster in fault recovery and nearly linearly scaled (up to 32 nodes) in horizontal performance as compared to a monolithic baseline. The ensemble classifier has an AUC-ROC of 0.98 and F1-score of 0.95, which is better than a single model baseline.

The work represents the next step of the composable microservices-based foundations and the big data ML pipeline-based design of and scales both to a unified and intelligent cloud-native technology applicable both to financial services and IoT and real-time analytics. These findings validate the hypothesis that incorporating AI/ML intelligence into the microservices fabric, as opposed to an observability layer, can be done technologically and is empirically the best. Future research will include the studies of GNN-based dependency modelling, federated learning as a

privacy-preserving intelligence, and serverless deployment patterns of cost-efficient AI microservices.

References

- [1] Ahmad, H., Treude, C., Wagner, M., & Szabo, C. (2025). Towards resource-efficient reactive and proactive auto-scaling for microservice architectures. *Journal of Systems and Software*. <https://doi.org/10.1016/j.jss.2025.00058>
- [2] Cloud Microservices in Focus: Architecture, Industry Practices and Emerging Innovation. (2025). *International Research Journal on Advanced Engineering Hub (IRJAEH)*, 3(12), 4255–4267. <https://doi.org/10.47392/IRJAEH.2025.0623>
- [3] Dogani, J., Namvar, R., & Khunjush, F. (2023). Auto-scaling techniques in container-based cloud and edge/fog computing: Taxonomy and survey. *Computer Communications*, 209, 120–150. <https://doi.org/10.1016/j.comcom.2023.06.010>
- [4] Hebbar, K. S. (2025). Workload-aware machine learning for microservice scaling in Kubernetes. *International Journal of Computational and Experimental Science and Engineering*, 11(4). <https://doi.org/10.22399/ijcesen.2025.11.4>
- [5] MDPI Software Editorial. (2025). Designing microservices using AI: A systematic literature review. *Software*, 4(1), 6. <https://doi.org/10.3390/software4010006>
- [6] PMC Review. (2024). Auto-scaling techniques in cloud computing: Issues and research directions. *Sensors*, 24(17), 5551. <https://doi.org/10.3390/s24175551>
- [7] Putapu, A. (2025). AI-enhanced microservices: Integrating machine learning pipelines in Java cloud environments. *International Journal of Computing and Engineering*, 7(19), 37–50. <https://doi.org/10.47941/ijce.3059>
- [8] Rachamala, N. R. (2023). Architecting AML detection pipelines using Hadoop and PySpark with AI/ML. *Journal of Information Systems Engineering and Management*, 8(4). <https://www.jisem-journal.com/>
- [9] Rachamala, N. R., Kotha, S. R., & Talluri, M. (2021). Building composable microservices for scalable data-driven applications. *International Journal of Communication Networks and Information Security*, 13(3), 534–542. <https://doi.org/10.48047/IJCNIS.13.3.534-542>
- [10] Santos, J., Wauters, T., Volckaert, B., & De Turck, F. (2023). GymHPA: Efficient auto-scaling via reinforcement learning for complex microservice-based applications in Kubernetes. In *NOMS 2023–2023 IEEE/IFIP Network Operations and Management Symposium*. <https://doi.org/10.1109/noms56928.2023.10154298>
- [11] Santos, J., Reppas, E., Wauters, T., Volckaert, B., & De Turck, F. (2024). Gwydion: Efficient auto-scaling for complex containerised applications in Kubernetes through reinforcement learning. *Journal of Network and Computer Applications*, 232, 104011. <https://doi.org/10.1016/j.jnca.2024.104011>
- [12] Springer Nature (2025). AI techniques in the microservices life-cycle: A systematic mapping study. *Computing*. <https://doi.org/10.1007/s00607-025-01432-z>

- [13] Wang, Z., Zhu, S., Li, J., Jiang, W., Ramakrishnan, K. K., Yan, M., Zhang, X., & Liu, A. X. (2024). DeepScaling: Autoscaling microservices with stable CPU utilization for large-scale production cloud systems. *IEEE/ACM Transactions on Networking*, 32(5), 3961–3976. <https://doi.org/10.1109/tnet.2024.3400953>
- [14] Zarai, O., Mcharfi, Z., & El Asri, B. (2025). Intelligent autoscaling strategies for cloud-native microservices: A systematic review. In *2025 International Conference on Intelligent Systems: Theories and Applications (SITA)*, 1–9. <https://doi.org/10.1109/sita67914.2025.11273589>
- [15] Designing Cloud-Native Enterprise Systems by Modernising Applications with Microservices and Kubernetes Platforms. (2025). *International Journal of Research and Applied Innovations*, 8(5), 13052–13063. <https://doi.org/10.15662/IJRAI.2025.0805015>