Towards Autonomous Kubernetes: A Framework for AI-Driven Operations (AIOps)

Kishan Raj Bellala¹

¹Independent Researcher, Austin, Texas, U.S.A.

Publication Date: 2025/09/26

Abstract: The increasing use of Kubernetes has brought substantial operational complexity because manual management of its numerous dynamic components (pods, nodes, networks) is slow, error-prone, and unsustainable at scale. This research investigates how AIOps (Artificial Intelligence for IT Operations) principles can move past native automation to establish fully autonomous Kubernetes management. The proposed framework uses machine learning to detect anomalies, identify causes, and predict scaling needs before executing automatic remediation steps. Our methodology demonstrates that AIOps can enhance system reliability and reduce operational Toil while optimizing resource efficiency through closed-loop observation-action cycles, leading to self-healing Kubernetes ecosystems that require minimal human intervention.

Keywords: Kubernetes, AIOps (Artificial Intelligence for IT Operations), Autonomous Operations, Self-Healing Systems, Anomaly Detection, Root Cause Analysis (RCA), Predictive Scaling, Automated Remediation, Operational Complexity, Machine Learning for IT Operations, Container Orchestration, Site Reliability Engineering (SRE).

How to Cite: Kishan Raj Bellala (2025) Towards Autonomous Kubernetes: A Framework for AI-Driven Operations (AIOps). *International Journal of Innovative Science and Research Technology*, 10(9), 1654-1661. https://doi.org/10.38124/ijisrt/25sep1016

I. INTRODUCTION

The introduction of containerization, together with microservices architecture, has transformed the entire process of software development and deployment. Kubernetes stands as the standard for container orchestration because it leads to the revolution in modern software development. The automation capabilities of Kubernetes for containerized application deployment, scaling, and management have established it as the fundamental element of contemporary cloud-native infrastructure, which delivers unmatched operational agility and scalability to organizations (Kashiv, 2025). The power of Kubernetes brings built-in complexity to the system. The operational demands become more challenging when Kubernetes environments grow complex, handling thousands of pods and services, and interconnected components spread across multiple nodes. The continuous flow of alerts, performance metrics, and log data overwhelms teams who struggle to identify meaningful information from background noise (Johansson, Papadopoulos, Ragberger, & Nolte, 2022). The operational reality results in substantial "Toil," which SRE defines as manual, repetitive, and reactive work that grows proportionally with system size. The operational workload of Toil drains engineering resources while creating human errors that lead to slow incident responses and system instability (Kashiv, 2025) (Johansson, Papadopoulos, Ragberger, & Nolte, 2022).

The natural evolution beyond mere automation is autonomy. The vision is for systems that can not only execute predefined instructions but also intelligently manage, heal, and optimize themselves with minimal human intervention. This shift promises to eliminate Toil, enhance reliability, and free engineers to focus on strategic, value-added work rather than firefighting. This is where AIOps (Artificial Intelligence for IT Operations) enters the picture. Defined by Gartner as "the application of machine learning and data science to IT operations problems," AIOps provides the necessary intelligence to achieve this autonomy. By leveraging big data, advanced analytics, and machine learning algorithms on operational data, AIOps platforms can detect anomalies, predict failures, perform precise root cause analysis, and ultimately prescribe automated remediation actions (Kashiv, 2025).

This paper demonstrates that AIOps integration represents the essential development for Kubernetes management evolution. This paper will present a framework that shows how AIOps converts Kubernetes from an automated command-execution platform into an autonomous decision-making system. The synthesis between these elements becomes vital for efficient and reliable management of complex, large-scale cloud-native applications in the next generation.

https://doi.org/10.38124/ijisrt/25sep1016

II. BACKGROUND & RELATED WORK

> Kubernetes Fundamentals: The Foundation for Automation

Kubernetes serves as an open-source platform that automates the deployment of application containers and their scaling and management across multiple host clusters. The system achieves its effectiveness through its declarative model and core abstractions that work together to maintain reliability and scalability (Johansson, Papadopoulos, Ragberger, & Nolte, 2022). The control plane operates as the cluster brain, which maintains continuous reconciliation between user-declared desired states and actual system states. The control plane consists of four essential components, which include the API Server, Scheduler, Controller Manager, and etcd that work together for system orchestration (Jorge-Martinez, et al., 2021) (Bogatinovski, Kao, Nedelkoski, & Cardoso, 2020). The Pod stands as Kubernetes' fundamental abstraction because it represents the smallest deployable unit, which includes one or more containers, storage, and a distinct network identity for a running process. Kubernetes uses the Deployment abstraction to manage Pod lifecycle because it defines application replication states and enables controlled updates and rollbacks. The Horizontal Pod Autoscaler (HPA) and Vertical Pod Autoscaler (VPA) work together to enhance scalability by using CPU utilization metrics to adjust Pod numbers and by optimizing resource allocation through historical usage data-based modifications of CPU and memory requests. The control plane of Kubernetes aligns user-specified desired outcomes like multiple microservice replicas through its declarative and API-driven architecture. The automation of Kubernetes operates reactively through predefined corrective actions, which demonstrates its capabilities while showing its boundaries for autonomous system development (Wei-Guo, Xi-Lin, & Jin-Zhong, 2018).

➤ Limits of Native Automation:

Kubernetes provides automation primitives but operates within constraints, preventing autonomous operation. The Horizontal Pod Autoscaler acts reactively, responding only after CPU utilization exceeds thresholds (Nguyen, Yeom, Kim, Park, & Kim, 2020) (Tien, 2019). This delayed response causes performance issues during sudden traffic increases. Kubernetes shows limited anomaly detection, identifying Pod status but failing to detect complex issues like performance decline, memory leaks, and network problems. The system cannot determine root causes, simply restarting failing Pods without understanding if failures stem from API outages, database issues, or application bugs. Static thresholds require constant manual adjustments across applications, increasing operational costs. While Kubernetes executes automated tasks effectively, it cannot predict failures or respond contextually to unexpected situations, showing its limited automation capabilities (Li, Sun, & Ke, 2024).

➤ An Overview of AIOps

AIOps stands for Artificial Intelligence for IT Operations, which Gartner developed as a discipline that uses machine learning and data science to solve operational problems in complex IT systems. AIOps functions as an

integrated system of technological components and operational processes that enhances IT operations through intelligent, proactive, and efficient management. The core element of AIOps starts with significant data aggregation, which collects and links together various types of data, including performance metrics, event logs, application logs, distributed traces, and system topology, or change data from the entire IT ecosystem (Sabharwal, 2022). The machine learning and analytics layer uses this foundation to deliver the necessary intelligence for detecting standard system patterns and precise anomaly detection beyond static thresholding and future capacity forecasting and root cause analysis through multi-source data correlation for reduced mean time to resolution (MTTR). The automation layer converts obtained insights into practical actions that span from producing highly relevant alerts to starting automated workflows for deployment rollback, resource scaling, and predefined runbook execution. The combination of data aggregation with intelligence and automated action enables a transition from manual reactive operations toward proactive autonomous IT management. The AIOps framework provides Kubernetes with a solution to overcome its built-in automation constraints through predictive intelligence and context-aware decision-making capabilities (Reiter, 2021).

III. CORE PILLARS OF AN AUTONOMOUS KUBERNETES SYSTEM

The transition from automated to autonomous Kubernetes management requires multiple essential components that work together. A strong Observability & Data Foundation stands as the crucial first pillar for this transformation. Any subsequent intelligence or automation requires high-quality data to function correctly because an unstable foundation exists without comprehensive data. The first pillar requires the collection and unification of all relevant telemetry data into a single platform, which enables data analysis and querying (Johansson, Papadopoulos, Ragberger, & Nolte, 2022).

➤ Piller 1: Observability & Data Foundation

An autonomous system requires a profound and uninterrupted comprehension of its operational state, together with the operational state of its hosted applications. Observability extends beyond basic monitoring because it enables systems to understand their internal state through external outputs during investigations of new or unexpected system states (Liu, 2020).

• The Unified Data Platform: The Central Nervous System

The foundation of this pillar depends on building a single data platform that receives, stores, and connects different types of telemetry data. The use of separate tools creates obstacles to achieving the complete understanding needed for autonomous operations (Nguyen, Horizontal Pod Autoscaling in Kubernetes for Elastic Container Orchestration., 2020). A modern Kubernetes observability stack consists of:

✓ *Metrics (System Performance):*

Time-series data capturing the health of the cluster and workloads. The standard for collecting metrics such as CPU utilization, memory pressure, I/O throughput, and custom application performance indicators is Prometheus.

✓ Logs (Events and Context):

Timestamped event data providing contextual information about system and application behavior. The tools Loki and Elastic Stack collect logs from containers, system demons, and the Kubernetes API to enable efficient search and correlation.

✓ Traces (Request Flow):

Distributed tracing data that follows requests through microservice architectures, crucial for diagnosing latency and dependency bottlenecks. The distributed tracing data can be captured and analyzed through Jaeger, Grafana, Tempo, and OpenTelemetry.

The value of this platform stems from its ability to combine data types rather than collecting them separately. A sudden rise in application error rates (metrics) can be directly traced to error messages in container logs and distributed traces, which reveal bottlenecks, thus enabling fast and accurate diagnosis (Qi, 2020) (Shwartz., 2025).

• Beyond Kubernetes: Extending the Data Universe

A genuinely autonomous system must extend observability beyond the cluster itself. The detection of infrastructure-based problems requires the integration of cloud provider metrics, which include VM performance, disk I/O, and network latency data from AWS CloudWatch, Google Cloud Monitoring, and Azure Monitor. Business-level application metrics, which track user transaction rates and order throughput, serve as higher-order indicators of system health and typically detect underlying degradation before it becomes visible. The Kubernetes API provides essential context for root cause analysis through its events and topology data, which links operational changes to observed performance anomalies (Qi, 2020).

The unified observability platform functions as the central nervous system, which controls the autonomous Kubernetes cluster. The platform serves as the primary source of accurate high-fidelity data, which enables machine learning algorithms in subsequent pillars to learn, detect, decide, and act. The quality and quantity of data, along with its connections between different elements, determine the total intelligence and effectiveness of the autonomous system. Any attempt at autonomy would fail because it lacks this essential foundation, which makes it both blind and prone to catastrophic error (Liu, 2020) (Qi, 2020).

➤ Pillar 2: Intelligent Anomaly Detection and Forecasting:

The first pillar of autonomous Kubernetes provides sensory capabilities, but the second pillar enables cognitive functions through data interpretation to identify normal operations, anomalies, and future risks. The transition occurs from threshold-based reactive automation to predictive proactive system management (Arshad, 2022).

• Moving Beyond Static Thresholds

Kubernetes native alerting mechanisms, including Prometheus rules, depend on fixed threshold values (e.g., triggering alerts when memory usage exceeds 90%). The basic method proves to be inflexible because it generates both incorrect positive and negative results. The lack of context in static thresholds makes them unable to distinguish between normal demand increases from planned promotions and dangerous denial-of-service attacks. The system detects memory leaks only after reaching critical thresholds, which results in service disruption because it fails to identify gradual degradations (Alsalman, 2024).

Machine Learning for Behavioral Profiling and Deviation Detection

The observability layer provides historical data to detect intelligent anomaly, which uses machine learning techniques to analyze this information. The algorithms create operational profiles for each service, node, and workload to identify their individual behavioral patterns. The system learns to detect seasonal patterns and trends (such as daily or weekly usage patterns) and metric correlations (like CPU usage and network throughput). It performs multivariate analysis to reveal complex anomalies that cannot be detected by examining metrics separately. The system uses established adaptive baselines to detect statistically significant deviations, which signal potential emerging failures. The method decreases the number of incorrect alerts while enhancing the detection of new or faint failure patterns (Tien, KubAnomaly: Anomaly detection for the Docker orchestration platform with neural network approaches., 2019).

• Forecasting: Anticipating Future States

The predictive aspect of system management emerges through forecasting, which adds a forecasting capability to system management. The system uses time-series models, including ARIMA and Prophet, and deep learning methods like LSTM networks to predict upcoming resource requirements and potential system bottlenecks. The system can use predictive models to forecast upcoming workload increases through historical trend analysis, so it can scale workloads before performance degradation happens. The system can detect slow memory usage patterns that indicate a forthcoming OOMKill termination through forecasting, so it can take preventive measures before service disruption occurs (Shahzad, 2022).

➤ Pillar 3: Automated Root Cause Analysis (RCA)

AIOps contains Automated Root Cause Analysis (RCA) as an essential feature, which precisely detects the origins of system alerts. System monitoring requires alert signals to both indicate system issues and identify exact problem sources, including "Pod X fails because Service Y executes slow database queries, which drain the database connection pool." Topology mapping, together with causal inference, operates as an essential method to achieve this goal (Yan et al., 2012).

Topology Mapping:

Complex IT systems need topology mapping for understanding their component relationships and interactions. The system architecture requires a visual representation to show component relationships between services, databases, and networks (Yan et al., 2012). The system's architectural visualization helps identify vital system pathways and failure locations, which enable teams to identify origin points of anomalies. The complete service dependency model acts as a vital component for RCA because it includes topological and cross-layer relationships along with protocol interactions and control plane dependencies (Sun et al., 2021). A detailed service mapping system enables the identification of symptoms and diagnostic events to support effective correlation and reasoning during the RCA process (Yan et al., 2012).

• Causal Inference:

Causal inference provides an advanced method for establishing cause-and-effect relations beyond basic correlation detection. The application of Causal Bayesian Networks (CBN) with causal inference models enables RCA systems to detect root causes through monitoring changes in variable probability distribution (Li et al., 2022). The system uses this method to detect how specific modifications (e.g., configuration updates) directly affect system performance. The practical implementation uses unsupervised causal inference methods to generate monitoring metric graphs, which apply causal assumptions to establish reliable root cause identification (Li et al., 2022).

Systems that combine topology mapping with causal inference can perform accurate automated Root Cause Analysis. Through structured analysis, these methods replace guesswork to enable rapid incident response through automated detection of affected components and their propagation paths and underlying causes, which results in more innovative IT operations management.

➤ Pillar 4: Predictive & Prescriptive Scaling

Kubernetes predictive and prescriptive scaling enhances operational efficiency through resource management and cost optimization, which exceeds the capabilities of traditional Horizontal Pod Autoscaler (HPA) tools.

• Predictive Scaling:

Predictive scaling requires the application of time-series forecasting algorithms to forecast demand so that resource utilization can be adjusted accordingly. Kubernetes environments benefit from workload prediction through the combination of Holt–Winter forecasting and Gated Recurrent Unit (GRU) neural network applications. The methods enable real-time adjustments to instance counts for improved resource optimization by predicting upcoming demand through techniques like Black Friday sales or daily login bursts (Yuan & Liao, 2024). Service quality improves substantially through predictive scaling because it reduces cold start times and maintains performance stability during high-demand periods.

• Prescriptive Scaling:

Prescriptive scaling focuses on suggesting the most suitable resource requests and limits, which are known as "right-sizing." The method delivers better performance by scaling resources out and optimizing their distribution according to workload requirements for improved stability and cost reduction. Machine learning-based resource allocation systems optimize resource utilization by avoiding excess provisioning while maintaining SLA compliance (Toka et al., 2021). The combination of adaptive AI-based auto-scaling systems monitors request variability to optimize resource usage at high service quality levels (Toka et al., 2020).

• Implementation Strategies:

Kubernetes administrators can deploy advanced scaling methods by combining multiple autoscaling strategies that perform horizontal and vertical scaling along with predictive modeling. The predictive models combine empirical modal decomposition with ARIMA models to forecast pod loads for early resource adjustments that solve latency problems during scaling operations (Zhao et al., 2019). The scaling process becomes more precise through dynamic multi-level autoscaling, which uses application-level monitoring data to adjust (Taherizadeh & Stankovski, 2018).

When integrated into Kubernetes, these predictive and prescriptive scaling techniques allow for optimized resource management that enhances both performance and cost efficiency for handling modern cloud-native application demands.

> Pillar 5: Self-Healing & Automated Remediation

The modern Kubernetes environment depends on self-healing and automated remediation to automatically fix problems after identifying their root causes. The capability provides applications with enhanced reliability and resilience through fast fault resolution, which reduces downtime and preserves service levels.

• Simple Actions:

✓ Pod Deletion and Restart:

One of the simplest yet effective strategies for handling faults in Kubernetes. The process of deleting malfunctioning pods becomes beneficial when Kubernetes detects an anomaly or unrecoverable error because it allows the system to recreate the Pod in a healthier state automatically (Nguyen et al., 2020).

✓ *Node Draining:*

The process of draining an unhealthy node allows workload redistribution to other healthy nodes, which maintains service continuity before the problematic nodes can be maintained or terminated.

• Complex Actions:

✓ Rollbacks:

The process of returning to a previous stable deployment version becomes essential when new updates cause system instability. The application maintains operation through this method, which prevents the newly introduced problems (Tran et al., 2022).

✓ Scaling Dependencies:

The process of scaling dependent services becomes essential when a bottleneck appears in service chain operations. AI optimization systems help organizations scale their elements properly to handle detected load increases through predictive analytics (Li et al., 2024).

✓ Failovers.

A critical failure detection triggers a failover process that shifts operations to another cluster or region to preserve system availability. Kubernetes Federation enables applications to span multiple service areas, which results in improved fault tolerance according to Kim et al. (2019).

• Advanced Techniques:

✓ Proactive Fault-Tolerant Systems:

The implementation of systems that predict faults before they affect services leads to improved self-healing capabilities. The combination of a Bi-LSTM fault prediction framework with stateful service migration enables the transfer of services from predicted faulty nodes to stable nodes before faults occur, thus preserving service quality and preventing outages (Tran et al., 2022).

✓ *ML and AI-Based Remediation:*

Machine learning models used for anomaly detection and root because analysis enable automated remediation actions. The KubAnomaly system employs neural network methods to identify unusual system behaviors, which it addresses automatically without requiring extensive human involvement, thus improving Kubernetes security and resilience (Tien et al., 2019).

Organizations can establish strong disruption management systems through these strategies, which maintain peak service availability and performance during unexpected system issues.

IV. IMPLEMENTATION ARCHITECTURE & CONSIDERATIONS DATA LAYER

A production-ready system based on autonomous Kubernetes theoretical foundations needs a strong architectural framework for safety. The proposed high-level implementation architecture includes relevant technologies and essential operational considerations for the system. The system architecture for autonomous Kubernetes operations consists of three distinct layers, which maintain continuous feedback connections. The data flow between these layers appears in Figure 1 (Shwartz., 2025) (Li, Sun, & Ke, 2024). The system architecture for autonomous Kubernetes operations demonstrates data flow from collection to

automated action in Figure 1. The Data Layer provides input to the AIOps Engine, which makes decisions that get transmitted to the Action Layer through Custom Resources. The Data Layer receives Kubernetes state changes, which create a continuous feedback loop.

> Data Layer: The Unified Observability Platform

The first pillar describes how this layer collects, stores, and correlates telemetry data as described in Pillar 1.

• Components:

✓ *Metrics*:

Prometheus (often with Thanos or Cortex for long-term storage and scalability).

✓ Logs:

Loki or the Elastic Stack (Elastic search, Log stash, Kibana).

✓ Traces:

Jaeger, Grafana Tempo, or an Open Telemetry Collector.

✓ Events:

Kubernetes Event Exporter to funnel cluster events into the logging pipeline.

■ Purpose:

This layer aggregates the "what is happening" data from the entire application and infrastructure stack, providing the raw material for analysis (Li, Sun, & Ke, 2024) (Shwartz., 2025).

➤ AIOps Engine: The Analytical Brain

The core intelligence layer uses machine learning and analytical models to process data collected from the Data Layer. The system performs the functions of Pillars 2, 3, and 4 (Anomaly Detection, RCA, Forecasting).

• Implementation Options:

✓ Open-Source Stack:

This can be a custom-built ensemble of tools. The Prometheus ML toolkit can be used for time-series forecasting. Keptn is an emerging open-source framework specifically designed for automated delivery and operations, capable of triggering remediation workflows based on quality gates.

✓ Commercial Platforms:

Integrated platforms like Dynatrace, Datadog, or New Relic have powerful AIOps engines built in. They excel at correlating data across metrics, logs, and traces and provide advanced root because analysis features out of the box.

Purpose:

This layer consumes data, applies intelligence to determine "why it is happening and what will happen next," and decides on a remedial action (Levin, et al., 2019) (Li, Sun, & Ke, 2024).

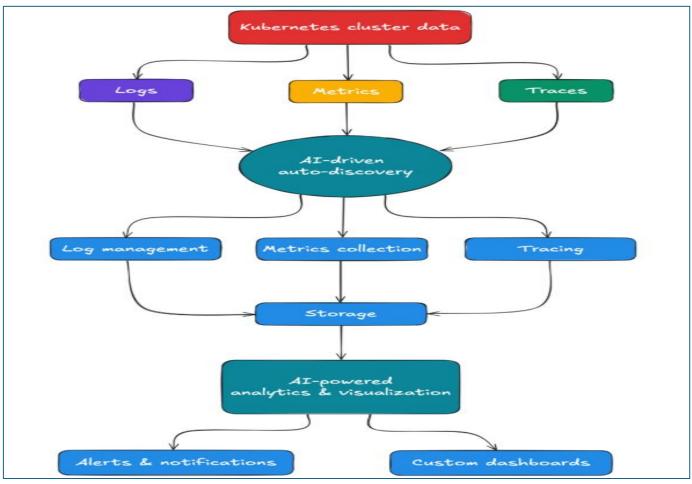


Fig 1 Autonomous Kubernetes AIOps Architecture. (Shwartz., 2025)

➤ Action Layer: The Safe Execution Arms

The most essential safety and Kubernetes paradigm compliance layer exists here. The AIOps Engine needs to avoid executing actions through kubectl imperative commands. The system should declare its desired remedial state, which Kubernetes will enforce for execution.

• Primary Mechanism:

Kubernetes Operators. The Operator pattern represents the fundamental principle for safe automation. Through the Kubernetes API, an Operator functions as a custom controller that handles application management and component control.

• Workflow:

- ✓ The AIOps Engine records its decision to restart Pod X by modifying the Custom Resource (CR). The system updates a PodRemediation resource by adding spec. Action: restart and spec. podName: my-pod to its specification.
- ✓ A Remediation Operator designed explicitly for this purpose monitors Custom Resources for modifications.
- ✓ The Operator includes pre-programmed human-approved logic that performs safe actions. The Operator validates requests before performing condition checks through Kubernetes API calls to execute the action (e.g., pod deletion).

Purpose.

This layer converts decisions into cluster-safe controlled actions that implement Pillar 5 (Automated Remediation) (Nedelkoski, Cardoso, & Kao, 2019).

> Challenges and Considerations:

• Seamless Integration:

The efficient operation requires integrating data, AIOps, and action layers with Kubernetes native capabilities. The system requires proper layer synchronization to achieve accurate data analysis and timely, appropriate action execution (Li et al., 2024).

• Scalability and Flexibility:

Architecture needs to support scalability to handle growth and rising workload demands without compromising performance. The implementation of Kubernetes Federation for distributed applications across multiple regions enhances both responsiveness and reliability (Kim et al., 2019).

• Security and Compliance:

The entire system requires absolute security protection as its top priority. The system architecture must include security vulnerability detection mechanisms together with data protection features. The system maintains its robustness against threats through integration with security-focused tools and practices (Bose et al., 2021).

The architectural method enables IT operations through automation and intelligence while improving system selfmanagement capabilities to decrease operational costs and enhance service reliability.

V. CONCLUSION

Large-scale Kubernetes environments have exceeded the operational complexity that human-centric management methods can handle. The paper demonstrates that AIOps represents a fundamental transformation that turns Kubernetes from an automated system into a fully autonomous platform. The transition between these stages becomes vital for achieving the reliability, efficiency, and scalability that cloud-native architectures promise. The path to autonomy requires five essential components, which include a unified observability foundation, intelligent anomaly detection and forecasting, automated root cause analysis, prescriptive scaling and healing, and a safe action layer powered by Kubernetes Operators (Johansson, Papadopoulos, Ragberger, & Nolte, 2022) (Liu, 2020). The future evolution will not eliminate human expertise but will elevate its position. The main objective is to eliminate toil, which represents the repetitive manual tasks that consume engineering time. Site Reliability Engineers and platform teams can dedicate their time to strategic value-added work after being freed from firefighting duties. The team can focus on designing resilient systems, improving autonomous algorithms, and managing exceptions.

The future of Kubernetes management depends on intelligent automation systems. The human role now transitions from performing direct operational tasks to overseeing autonomous systems through orchestration and oversight functions (Shwartz., 2025). Organizations can achieve this future by implementing the AIOps-driven architecture described in this paper to create Kubernetes environments that deliver robust scalability with self-healing, self-optimizing, and resilient capabilities.

REFERENCES

- [1]. Liu, C., Wang, B., Liu, J., Tang, Z., & Cai, Z. (2020). A protocol-independent container network observability analysis system based on eBPF. 697–702. https://doi.org/10.1109/icpads51040.2020.00099
- [2]. Qi, S., Kulkarni, S. G., & Ramakrishnan, K. K. (2020). Assessing Container Network Interface Plugins: Functionality, Performance, and Scalability. IEEE Transactions on Network and Service Management, 18(1), 656–671. https://doi.org/10.1109/tnsm.2020.3047545
- [3]. Itiel Shwartz. (2025, August 21). AIOPs for kubernetes (or KAIOPs?). Komodor. https://komodor.com/blog/aiops-for-kubernetes-or-kaiops/
- [4]. Arshad, K., Naseer, S., Ali, R. F., Muneer, A., Aziz, I. A., Khan, N. S., & Taib, S. M. (2022). Deep Reinforcement Learning for Anomaly Detection: A Systematic Review. IEEE Access, 10, 124017–124035. https://doi.org/10.1109/access.2022.3224023

- [5]. Alsalman, D. (2024). A Comparative Study of Anomaly Detection Techniques for IoT Security Using Adaptive Machine Learning for IoT Threats. IEEE Access, 12, 14719–14730. https://doi.org/10.1109/access.2024.3359033
- [6]. Shahzad, F., Al-Jumeily Obe, D., Mannan, A. Almadhor, A. S., Javed, A. R., & Baker, T. (2022). Cloud-based multiclass anomaly detection and categorization using ensemble learning. Journal of Cloud Computing, 11(1). https://doi.org/10.1186/s13677-022-00329-y
- [7]. Yan, H., Ge, Z., Yates, J., Breslau, L., Massey, D., & Pei, D. (2012). G-RCA: A Generic Root Cause Analysis Platform for Service Quality Management in Large IP Networks. IEEE/ACM Transactions on Networking, 20(6), 1734–1747. https://doi.org/10.1109/tnet.2012.2188837
- [8]. Sun, Y., Qin, W., Xu, H., & Zhuang, Z. (2021). An adaptive fault detection and root-cause analysis scheme for complex industrial processes using moving window KPCA and information geometric causal inference. Journal of Intelligent Manufacturing, 32(7), 2007–2021. https://doi.org/10.1007/s10845-021-01752-9
- [9]. Li, M., Li, Z., Pei, D., Zhang, W., Sui, K., Yin, K., & Nie, X. (2022). Causal Inference-Based Root Cause Analysis for Online Service Systems with Intervention Recognition. 53, 3230–3240. https://doi.org/10.1145/3534678.3539041
- [10]. Yuan, H., & Liao, S. (2024). A Time Series-Based Approach to Elastic Kubernetes Scaling. Electronics, 13(2), 285. https://doi.org/10.3390/electronics13020285
- [11]. Toka, L., Dobreff, G., Sonkoly, B., & Fodor, B. (2020). Adaptive AI-based auto-scaling for Kubernetes. 16, 599–608. https://doi.org/10.1109/ccgrid49817.2020.00-33
- [12]. Taherizadeh, S., & Stankovski, V. (2018). Dynamic Multi-Level Auto-scaling Rules for Containerized Applications. The Computer Journal, 62(2), 174–197. https://doi.org/10.1093/comjnl/bxy043
- [13]. Toka, L., Dobreff, G., Sonkoly, B., & Fodor, B. (2021). Machine Learning-Based Scaling Management for Kubernetes Edge Clusters. IEEE Transactions on Network and Service Management, 18(1), 958–972. https://doi.org/10.1109/tnsm.2021.3052837
- [14]. Zhao, A., Song, J., Huang, Q., Huang, Y., Chen, Z., & Zou, L. (2019). Research on Resource Prediction Model Based on Kubernetes Container Auto-scaling Technology. IOP Conference Series: Materials Science and Engineering, 569(5), 052092. https://doi.org/10.1088/1757-899x/569/5/052092
- [15]. Nguyen, T.-T., Yeom, Y.-J., Kim, T., Park, D.-H., & Kim, S. (2020). Horizontal Pod Autoscaling in Kubernetes for Elastic Container Orchestration. Sensors (Basel, Switzerland), 20(16), 4621. https://doi.org/10.3390/s20164621
- [16]. Tran, M.-N., Vu, X. T., & Kim, Y. (2022). Proactive Stateful Fault-Tolerant System for Kubernetes

- Containerized Services. IEEE Access, 10, 102181–102194. https://doi.org/10.1109/access.2022.3209257
- [17]. Kim, D., Kim, E., Lee, C., Helal, S., & Muhammad, H. (2019). TOSCA-Based and Federation-Aware Cloud Orchestration for Kubernetes Container Platform. Applied Sciences, 9(1), 191. https://doi.org/10.3390/app9010191
- [18]. Bose, D. B., Shamim, S. I., & Rahman, A. (2021). 'Under-reported' Security Defects in Kubernetes Manifests. 9–12. https://doi.org/10.1109/encycris52570.2021.00009
- [19]. Tran, M.-N., Vu, X. T., & Kim, Y. (2022). Proactive Stateful Fault-Tolerant System for Kubernetes Containerized Services. IEEE Access, 10, 102181–102194. https://doi.org/10.1109/access.2022.3209257
- [20]. Tien, C., Huang, T., Tien, C., Huang, T., & Kuo, S. (2019). KubAnomaly: Anomaly detection for the Docker orchestration platform with neural network approaches. Engineering Reports, 1(5). https://doi.org/10.1002/eng2.12080
- [21]. Li, H., Sun, J., & Ke, X. (2024). AI-Driven Optimization System for Large-Scale Kubernetes Clusters: Enhancing Cloud Infrastructure Availability, Security, and Disaster Recovery. Journal of Artificial Intelligence General Science (JAIGS) ISSN:3006-4023, 2(1), 281–306. https://doi.org/10.60087/jaigs.v2i1.244
- [22]. Levin, A., Mcshane, N., Garion, S., Kolodner, E. K., Kugler, M., Lorenz, D. H., & Barabash, K. (2019). AIOps for a Cloud Object Storage Service. 165–169. https://doi.org/10.1109/bigdatacongress.2019.00036
- [23]. Nedelkoski, S., Cardoso, J., & Kao, O. (2019). Anomaly Detection from System Tracing Data Using Multimodal Deep Learning. 179–186. https://doi.org/10.1109/cloud.2019.00038
- [24]. Kim, D., Kim, E., Lee, C., Helal, S., & Muhammad, H. (2019). TOSCA-Based and Federation-Aware Cloud Orchestration for Kubernetes Container Platform. Applied Sciences, 9(1), 191. https://doi.org/10.3390/app9010191
- [25]. Bose, D. B., Shamim, S. I., & Rahman, A. (2021). 'Under-reported' Security Defects in Kubernetes Manifests. 9–12. https://doi.org/10.1109/encycris52570.2021.00009
- [26]. KASHIV, D. J. AI-Driven Networks: Architecting the Future of Autonomous, Secure, and Cloud-Native connectivity 2025. YASHITA PRAKASHAN PRIVATE LIMITED.
- [27]. Johansson, B., Papadopoulos, A. V., Ragberger, M., & Nolte, T. (2022). Kubernetes Orchestration of High Availability Distributed Control Systems. 1–8. https://doi.org/10.1109/icit48603.2022.10002757
- [28]. Jorge-Martinez, D., Ariza-Colpas, P., Chakraborty, C., Butt, S. A., De-La-Hoz-Franco, E., Onyema, E. M., & Shaheen, Q. (2021). Artificial intelligence-based Kubernetes container for scheduling nodes of energy composition. International Journal of System Assurance Engineering and Management. https://doi.org/10.1007/s13198-021-01195-8
- [29]. Bogatinovski, J., Kao, O., Nedelkoski, S., & Cardoso, J. (2020). Self-Supervised Anomaly Detection from

Distributed Traces. 342–347. https://doi.org/10.1109/ucc48980.2020.00054

https://doi.org/10.38124/ijisrt/25sep1016

- [30]. Wei-Guo, Z., Xi-Lin, M., & Jin-Zhong, Z. (2018). Research on Kubernetes' Resource Scheduling Scheme. 144–148. https://doi.org/10.1145/3290480.3290507
- [31]. Nguyen, T.-T., Yeom, Y.-J., Kim, T., Park, D.-H., & Kim, S. (2020). Horizontal Pod Autoscaling in Kubernetes for Elastic Container Orchestration. Sensors (Basel, Switzerland), 20(16), 4621. https://doi.org/10.3390/s20164621
- [32]. Tien, C., Huang, T., Tien, C., Huang, T., & Kuo, S. (2019). KubAnomaly: Anomaly detection for the Docker orchestration platform with neural network approaches. Engineering Reports, 1(5). https://doi.org/10.1002/eng2.12080
- [33]. Sabharwal, N., & Bhardwaj, G. (2022). Hands-on AIOps. Apress eBooks. https://doi. org/10.1007/978-1-4842-8267-0.
- [34]. Reiter, L., & Wedel, F. H. (2021). AIOps–A Systematic Literature Review.