



# Autonomous Kubernetes Cluster Healing using Machine Learning

Sai Bharath Sannareddy

Senior Cloud Infrastructure Engineer, Illinois, USA

[saibharathdevsecops@gmail.com](mailto:saibharathdevsecops@gmail.com)

**ABSTRACT:** Modern Kubernetes environments underpin mission-critical applications across healthcare, finance, and cloud-native enterprises. While Kubernetes provides robust primitives for container orchestration, it still relies heavily on manual intervention, static rules, and reactive alerts to recover from failures such as pod crashes, node instability, resource exhaustion, and cascading service degradation. As clusters scale in size and complexity, traditional monitoring and rule-based remediation mechanisms become insufficient to meet strict reliability objectives.

This paper presents an autonomous Kubernetes cluster healing framework driven by machine learning, designed to proactively detect anomalies, predict failure patterns, and execute self-healing actions without human intervention. The proposed system combines telemetry from Kubernetes control planes, observability platforms, and application-level signals with machine learning models that learn normal and abnormal operational behavior. By integrating predictive analytics with automated remediation workflows, the framework enables clusters to recover from failures faster, reduce mean time to detect (MTTD), and significantly lower mean time to recovery (MTTR).

Unlike conventional auto-scaling or threshold-based alerting, the proposed approach leverages historical incident patterns, resource utilization trends, and service-level indicators (SLIs) to make context-aware healing decisions. The architecture supports common remediation actions such as intelligent pod restarts, node cordoning, workload rescheduling, configuration rollbacks, and policy-driven scaling. The framework is cloud-agnostic and applicable to Kubernetes platforms deployed on Azure Kubernetes Service (AKS), Amazon EKS, and hybrid environments.

The study demonstrates how machine learning-driven autonomous healing improves cluster resilience, reduces operational toil, and enhances service reliability in regulated, production-grade environments. This work contributes a practical foundation for next-generation self-managing Kubernetes systems and establishes a pathway toward fully autonomous cloud-native infrastructure.

**KEYWORDS:** Kubernetes; autonomous systems; self-healing infrastructure; machine learning; site reliability engineering; cluster resilience; observability; cloud-native systems; AKS; DevOps automation.

## I. INTRODUCTION

Kubernetes has become the de facto orchestration platform for cloud-native applications, enabling organizations to deploy, scale, and manage distributed workloads across diverse infrastructure environments. Enterprises increasingly rely on Kubernetes to host latency-sensitive, compliance-critical, and revenue-generating systems. However, as Kubernetes clusters grow in size and heterogeneity, maintaining operational reliability becomes significantly more complex. Failures no longer occur as isolated events; instead, they propagate across microservices, nodes, and infrastructure layers, often requiring rapid and informed intervention to prevent service degradation.

Traditional Kubernetes resilience mechanisms—such as pod restarts, health probes, horizontal pod autoscaling, and static alerting rules—are primarily **reactive** and **rule-driven**. These mechanisms assume that failure conditions are known in advance and can be captured through predefined thresholds. In real-world production environments, this assumption frequently fails. Resource contention, noisy neighbors, cascading retries, misconfigurations, and partial outages generate failure patterns that are subtle, multi-dimensional, and context-dependent. As a result, human operators are often required to triage alerts, correlate metrics, identify root causes, and execute remediation actions under time pressure.



From a Site Reliability Engineering (SRE) perspective, this operational model introduces significant **toil**, increases **mean time to detect (MTTD)** and **mean time to recovery (MTTR)**, and places reliability at risk—particularly in environments with strict service-level objectives (SLOs). As clusters scale across regions, clouds, and compliance domains, the cost of manual intervention grows nonlinearly, making human-driven remediation unsustainable.

Machine learning offers an opportunity to fundamentally change how Kubernetes clusters respond to failure. Instead of relying solely on static thresholds and handcrafted rules, machine learning models can learn normal operational behavior, identify anomalous patterns, and anticipate failure conditions before they fully manifest. When combined with automated remediation workflows, these predictive capabilities enable **autonomous healing**, where clusters not only detect failures but also act to resolve them without human involvement.

This paper is motivated by the need to bridge the gap between **observability** and **actionability** in Kubernetes operations. While modern platforms provide rich telemetry—metrics, logs, traces, and events—most systems stop at alert generation. The proposed framework extends beyond detection by introducing a closed-loop system in which machine learning models continuously analyze operational signals and trigger context-aware healing actions. These actions are selected based on historical effectiveness, current cluster state, and policy constraints, allowing the system to respond intelligently rather than blindly.

Unlike prior approaches that focus exclusively on auto-scaling or simple anomaly detection, this work emphasizes **end-to-end autonomy**. The framework integrates failure prediction, decision-making, and remediation execution into a unified control loop. It is designed to operate in real production Kubernetes environments, supporting managed services such as Azure Kubernetes Service (AKS) and Amazon Elastic Kubernetes Service (EKS), as well as hybrid deployments. The focus is not on replacing Kubernetes primitives, but on augmenting them with intelligence that adapts over time.

The primary contributions of this paper are threefold. First, it introduces a machine learning-driven architecture for autonomous Kubernetes cluster healing that aligns with real-world SRE practices. Second, it demonstrates how predictive models can reduce operational latency by anticipating failures rather than merely reacting to them. Third, it provides a practical blueprint for integrating machine learning into Kubernetes control planes in a safe, auditable, and policy-governed manner.

By combining machine learning with cloud-native automation, this work advances the state of Kubernetes reliability engineering and moves cluster operations closer to fully self-managing systems. The remainder of the paper elaborates on the technical background, architectural design, machine learning models, and operational impact of the proposed autonomous healing framework.

## II. BACKGROUND AND RELATED WORK

Kubernetes is built around a set of control loops that continuously reconcile the desired state of the system with its observed state. Core components such as the scheduler, kubelet, and controller manager enable basic fault tolerance through mechanisms like pod restarts, replica reconciliation, and node health checks. While these primitives provide baseline resilience, they are intentionally generic and rely on predefined rules that do not adapt to complex or evolving failure conditions.

In production environments, failures rarely occur in isolation. Resource saturation, network instability, misconfigured deployments, and cascading retries often interact across multiple layers of the system. Traditional monitoring and alerting frameworks, even when integrated with Kubernetes-native tools, are primarily reactive. They depend on static thresholds and predefined rules that trigger alerts after degradation has already occurred. This results in delayed detection, alert fatigue, and heavy reliance on human operators for diagnosis and remediation.

Rule-based self-healing approaches have been proposed to reduce manual intervention. These systems typically automate common responses such as restarting pods, scaling workloads, or rescheduling workloads when nodes become unhealthy. While effective for well-understood failure modes, rule-based healing lacks contextual awareness and cannot generalize to unseen or compound failures. As clusters evolve, maintaining an ever-growing set of remediation rules becomes operationally expensive and brittle.



Recent research has explored the use of anomaly detection techniques in cloud and containerized environments. Statistical methods and machine learning models have been applied to detect abnormal resource usage, performance deviations, and service-level violations. These approaches improve failure detection but often stop short of autonomous remediation. In many cases, machine learning is used only to enhance alerting accuracy, leaving decision-making and recovery actions to human operators.

More advanced work has investigated reinforcement learning for resource management and auto-scaling in cloud systems. These efforts focus primarily on optimizing performance or cost by adjusting resource allocations. However, they typically operate in constrained problem spaces and do not address broader cluster-level failures such as configuration drift, control-plane instability, or cascading service degradation.

This paper differs from prior work by proposing a **closed-loop autonomous healing framework** that integrates machine learning-based failure prediction with automated remediation actions. Instead of treating detection and recovery as separate concerns, the proposed approach unifies observability, prediction, decision-making, and execution into a single feedback-driven system. Machine learning models continuously learn normal and abnormal cluster behavior, while an autonomous decision engine selects and executes healing actions based on current context, historical effectiveness, and policy constraints.

By shifting from static rules to adaptive learning, the framework addresses a critical gap in existing Kubernetes reliability solutions. It moves beyond reactive alerting and isolated auto-scaling toward **self-managing clusters capable of anticipating and resolving failures autonomously**. This work builds on existing Kubernetes primitives but extends them with intelligence that evolves as the system and its workloads change.

### III. PROBLEM STATEMENT AND DESIGN GOALS

#### 3.1. Problem Statement

Despite widespread adoption of Kubernetes, operating clusters at scale remains a fundamentally reactive process. Existing resilience mechanisms depend on static health checks, predefined thresholds, and manual intervention to diagnose and recover from failures. These approaches assume that failure modes are known in advance and can be addressed through fixed remediation logic. In real-world production environments, this assumption does not hold.

Modern Kubernetes clusters experience complex failure scenarios that arise from interactions across infrastructure, platform, and application layers. Examples include cascading pod restarts caused by transient resource exhaustion, control-plane instability triggered by misconfigurations, and partial service degradation resulting from network latency or dependency failures. Such issues often manifest gradually, evade simple threshold-based detection, and require contextual understanding to resolve correctly.

Current operational workflows place a heavy burden on human operators to correlate metrics, logs, and events, determine root causes, and execute recovery actions under time pressure. This manual process increases mean time to detect (MTTD) and mean time to recovery (MTTR), introduces operational risk, and does not scale with cluster size or workload diversity. While rule-based self-healing mechanisms automate known responses, they lack adaptability and fail when confronted with novel or compound failure patterns.

The core problem addressed in this work is the **absence of an autonomous, learning-driven mechanism capable of predicting failures and executing context-aware healing actions in Kubernetes clusters without human intervention**.

#### 3.2. Design Goals

To address this problem, the proposed system is designed around the following goals:

##### 3.2.1 Predictive Failure Awareness

The system must identify early signals of degradation and predict impending failures before service-level objectives are violated, rather than reacting after outages occur.

##### 3.2.2 Autonomous Decision-Making

Healing actions should be selected and executed automatically based on learned patterns, current cluster state, and historical effectiveness, minimizing reliance on human operators.



### 3.2.3 Context-Aware Remediation

Recovery actions must account for workload characteristics, resource constraints, and dependency relationships to avoid disruptive or counterproductive interventions.

### 3.2.4 Policy-Governed Safety

Autonomous actions must operate within predefined safety boundaries, ensuring compliance with operational, security, and reliability policies.

### 3.2.5 Platform Compatibility

The framework must integrate with existing Kubernetes control planes and be deployable across managed services such as AKS and EKS without requiring invasive changes.

### 3.2.6 Operational Scalability

The system should scale with cluster size and workload complexity while maintaining low-latency decision-making suitable for production environments.

## 3.3 Scope of the Solution

This work focuses on cluster-level and workload-level failures observable through Kubernetes telemetry and application-level signals. It does not attempt to replace Kubernetes core components but rather augments existing mechanisms with machine learning-driven intelligence that enables proactive and autonomous healing.

## IV. SYSTEM ARCHITECTURE

The proposed autonomous Kubernetes cluster healing framework is designed as a **closed-loop control system** that continuously observes cluster behavior, predicts failures, selects remediation actions, and learns from outcomes. The architecture emphasizes modularity, safety, and compatibility with existing Kubernetes control planes, allowing deployment in managed environments such as AKS and EKS without invasive modifications.

At a high level, the system consists of five tightly integrated layers: **(1) Telemetry and Signal Collection**, **(2) Feature Aggregation and State Modeling**, **(3) Machine Learning-Based Failure Prediction**, **(4) Autonomous Healing Decision Engine**, and **(5) Execution and Feedback Loop**.

### 4.1 Telemetry and Signal Collection Layer

This layer continuously ingests operational signals from the Kubernetes cluster and hosted workloads. It aggregates telemetry from multiple sources to form a unified view of cluster health.

Primary inputs include:

- **Kubernetes metrics:** CPU, memory, disk, and network utilization at pod, node, and namespace levels
- **Control-plane events:** pod restarts, scheduling delays, node readiness transitions, eviction events
- **Application-level indicators:** latency, error rates, request throughput, custom health signals
- **Infrastructure signals:** node pressure conditions, storage I/O saturation, network latency

Data is collected using standard Kubernetes observability mechanisms (e.g., metrics APIs, event streams, log collectors), ensuring platform portability and minimal operational overhead.

### 4.2 Feature Aggregation and State Modeling Layer

Raw telemetry is noisy and high-dimensional. This layer transforms incoming signals into structured representations suitable for machine learning.

Key functions include:

- Temporal aggregation (sliding windows, rolling statistics)
- Correlation of metrics across pods, services, and nodes
- Derivation of higher-level indicators such as saturation trends, anomaly scores, and dependency stress signals

The output of this layer is a **cluster state vector** representing the current operational context. This state captures both instantaneous conditions and short-term trends, enabling the system to distinguish transient fluctuations from meaningful degradation patterns.



### 4.3 Machine Learning–Based Failure Prediction Layer

The prediction layer applies machine learning models to estimate the likelihood and type of impending failures. Rather than predicting single metrics, the models learn multi-dimensional failure signatures associated with historical incidents and degraded states.

Typical prediction targets include:

- Probability of pod or service instability
- Risk of node-level failure or eviction cascades
- Likelihood of resource exhaustion leading to throttling or restarts
- Early indicators of cascading service degradation

Models operate continuously on updated state vectors and generate probabilistic risk assessments. This predictive capability enables the system to act **before** failures propagate into visible outages or SLO violations.

### 4.4 Execution and Feedback Loop

Selected actions are executed through Kubernetes-native APIs and automation mechanisms, ensuring consistency with existing operational practices. Crucially, every action is observed and evaluated.

The feedback loop:

- Monitors post-action cluster behavior
- Measures impact on key indicators (stability, latency, error rates)
- Records action outcomes and effectiveness
- Feeds results back into the learning pipeline

This continuous feedback enables the system to refine both its predictions and decision-making strategies over time, gradually improving autonomy and reliability.

### Architectural Summary

The proposed architecture transforms Kubernetes cluster management from a reactive, human-driven process into an adaptive, self-healing system. By tightly coupling machine learning–based prediction with autonomous remediation and continuous feedback, the framework enables clusters to anticipate failures, respond intelligently, and evolve with changing workloads and environments.

## V. MACHINE LEARNING MODELS FOR FAILURE PREDICTION

Autonomous healing in Kubernetes depends on the system’s ability to accurately identify early signals of failure amid noisy, high-volume telemetry. Unlike traditional anomaly detection tasks, cluster failures often emerge from **temporal correlations across multiple dimensions**, rather than abrupt metric spikes. This section describes the machine learning models used for failure prediction and explains how they are tailored to Kubernetes operational characteristics.

The framework adopts a **hybrid modeling strategy**, combining temporal learning, statistical anomaly detection, and probabilistic risk estimation to balance accuracy, robustness, and interpretability.

### 5.1 Failure Characteristics in Kubernetes Environments

Kubernetes failures exhibit several defining properties that influence model selection:

- **Temporal dependency:** Failures often develop over time (e.g., gradual memory pressure leading to eviction cascades).
- **Multi-dimensional coupling:** Metrics across pods, nodes, and services interact in non-linear ways.
- **Context sensitivity:** Similar metric patterns may be benign or critical depending on workload type and deployment state.
- **Class imbalance:** Failure events are rare relative to normal operation.
- **Operational noise:** Short-lived spikes and transient events are common and should not trigger remediation.

These characteristics make simple thresholding and static classifiers ineffective.



## 5.2 Feature Representation and Learning Signals

The prediction models operate on **cluster state vectors** derived from the architecture described in Section V. Features are grouped into three categories:

### 1. Resource dynamics

CPU saturation trends, memory pressure growth rates, disk I/O contention, network latency variance.

### 2. Control-plane behavior

Pod restart frequency, scheduling delays, eviction events, readiness transitions.

### 3. Service health indicators

Latency distributions, error-rate gradients, retry amplification signals.

Temporal aggregation is applied using sliding windows to preserve short-term and medium-term trends without overwhelming the models with raw data.

## 5.3 Temporal Sequence Models

To capture evolving degradation patterns, the framework employs **temporal sequence models**, such as recurrent neural networks (RNNs) or lightweight LSTM variants.

These models learn mappings of the form:

$$P(F_{t+k} | X_{t-L:t})$$

where  $X_{t-L:t}$  represents historical cluster states and  $F_{t+k}$  denotes the probability of failure within a future horizon.

Temporal models are particularly effective for:

- Predicting cascading pod failures
- Detecting slow resource leaks
- Identifying instability preceding node failures

Their primary advantage lies in recognizing failure *trajectories*, not just static anomalies.

## VI. AUTONOMOUS HEALING DECISION ENGINE

Predicting failures alone does not improve system reliability unless predictions are translated into timely and appropriate remediation actions. The autonomous healing decision engine serves as the **control brain** of the proposed framework, transforming probabilistic failure predictions into safe, context-aware corrective actions. Unlike static remediation scripts or rule-based playbooks, this engine operates dynamically, learning which actions are most effective under specific cluster conditions.

### 6.1 Decision-Making Objectives

The primary objective of the decision engine is to **restore or preserve system stability with minimal disruption**. Each decision balances multiple competing factors:

- Failure risk severity and confidence
- Scope of impact (single pod, service, node, or cluster)
- Operational cost of intervention
- Risk of unintended side effects
- Compliance with safety and policy constraints

Rather than selecting the most aggressive response, the engine prioritizes **graduated remediation**, escalating actions only when simpler interventions fail.

### 6.2 Action Space Definition

The decision engine operates over a predefined but extensible action space that reflects common Kubernetes remediation mechanisms. Actions include:

- Targeted pod restarts or rescheduling
- Node cordoning and workload migration
- Horizontal or vertical scaling adjustments
- Configuration rollback or redeployment
- Traffic throttling or isolation of unstable dependencies

Each action is associated with metadata describing its blast radius, execution cost, and historical success rate.



### 6.3 Explainability and Operator Visibility

To maintain trust, every autonomous decision is accompanied by an explanation describing:

- The predicted failure and associated risk score
- The selected action and rationale
- Alternative actions considered but rejected

These explanations are surfaced through dashboards and logs, allowing operators to audit decisions and intervene when necessary.

## VII. EVALUATION AND OPERATIONAL IMPACT

Evaluating autonomous healing systems in production Kubernetes environments requires a focus on **operational outcomes** rather than isolated model accuracy metrics. The effectiveness of the proposed framework is assessed through its impact on reliability indicators, operational efficiency, and system stability. This section presents a conceptual evaluation grounded in real-world SRE practices and production observability signals.

### 7.1 Evaluation Methodology

The framework is evaluated using historical incident patterns, simulated failure injections, and live telemetry replay in controlled environments. Rather than benchmarking against synthetic datasets, the evaluation focuses on **before-and-after operational metrics** commonly used by SRE teams.

Key evaluation dimensions include:

- Mean Time to Detect (MTTD)
- Mean Time to Recovery (MTTR)
- Service-Level Objective (SLO) adherence
- Alert volume and operator toil
- Stability of remediation actions

This approach aligns evaluation with how Kubernetes reliability is measured in practice.

### 7.2 Failure Detection and Prediction Effectiveness

The machine learning-based prediction layer identifies early degradation signals that are typically missed by static threshold alerts. By correlating temporal trends across multiple signals, the system detects failure trajectories before they manifest as visible outages.

Observed benefits include:

- Earlier detection of gradual resource exhaustion
- Improved identification of cascading failure patterns
- Reduction in false-positive alerts caused by transient spikes

This predictive capability shifts operations from reactive response to proactive intervention.

## VIII. CASE STUDIES

To illustrate the practical applicability of the proposed autonomous healing framework, this section presents representative case studies based on common failure scenarios observed in production Kubernetes environments. These cases are drawn from managed Kubernetes deployments and reflect realistic operational conditions encountered in AKS and EKS-style clusters.

### Case Study 1: Cascading Pod Failures Due to Memory Pressure Scenario

A stateful microservice experiences gradual memory growth due to increased request volume. As node-level memory pressure rises, Kubernetes begins evicting pods across multiple namespaces, triggering cascading restarts and intermittent service disruption.

#### Autonomous System Behavior

- The prediction layer detects a sustained upward trend in memory utilization combined with increased pod restart frequency.
- Failure risk for node-level instability crosses the predefined threshold.



- The decision engine selects a low-blast-radius action: targeted rescheduling of affected pods and controlled scaling of memory-intensive workloads.
- Node cordoning is deferred unless stabilization fails.

## Outcome

The system stabilizes memory usage before eviction cascades propagate. Service availability is preserved, and no manual intervention is required. MTTR is reduced by eliminating the need for human diagnosis.

## Case Study 2: Control-Plane Stress During Deployment Rollout Scenario

A rolling deployment introduces a misconfiguration that increases pod startup latency. The scheduler experiences delays, and readiness probes fail intermittently, creating partial service degradation.

## Autonomous System Behavior

- Temporal models detect correlated increases in scheduling latency and readiness failures.
- The decision engine identifies the failure as configuration-induced rather than resource-related.
- A controlled rollback of the deployment is initiated, constrained by policy to limit impact scope.
- Post-action monitoring confirms stabilization.

## Outcome

The rollback is executed automatically within minutes, preventing escalation into a full outage. The system avoids unnecessary scaling or restarts that could worsen control-plane stress.

## Cross-Case Observations

Across all cases, the autonomous healing framework demonstrates the ability to:

- Differentiate between transient and persistent failures
- Select minimally disruptive remediation actions
- Prevent escalation into larger outages
- Reduce dependency on human intervention

These scenarios highlight how machine learning-driven autonomy enhances Kubernetes reliability in practical, production-grade environments.

## IX. CHALLENGES AND LIMITATIONS

While the proposed autonomous Kubernetes cluster healing framework demonstrates significant potential, several challenges and limitations must be acknowledged. These considerations are important for realistic deployment and guide future refinement of autonomous reliability systems.

### 9.1 Data Quality and Signal Noise

Machine learning models rely on high-quality telemetry to produce accurate predictions. In practice, observability data may be incomplete, delayed, or noisy due to sampling, instrumentation gaps, or transient infrastructure issues. Poor data quality can degrade prediction accuracy and increase uncertainty, particularly during rare or complex failure scenarios.

### 9.2 Limited Failure Label Availability

Accurate labeling of historical failure events is often incomplete or inconsistent in production environments. While the framework mitigates this through semi-supervised learning and anomaly detection, the absence of precise labels can limit model specificity and slow convergence for rare failure modes.

### 9.3 Concept Drift and Evolving Workloads

Kubernetes clusters evolve continuously as workloads, configurations, and infrastructure change. Machine learning models trained on historical data may lose accuracy over time due to concept drift. Continuous retraining and drift detection are necessary but introduce operational overhead and require careful governance to avoid instability.



## X. CONCLUSION

This paper presented an autonomous Kubernetes cluster healing framework that integrates machine learning-based failure prediction with context-aware, policy-governed remediation. As Kubernetes environments continue to grow in scale and complexity, traditional rule-based monitoring and manual recovery approaches are increasingly insufficient to meet modern reliability expectations. The proposed system addresses this gap by transforming cluster operations from a reactive, human-driven process into a proactive, self-healing capability.

By combining temporal learning, anomaly detection, and probabilistic risk estimation, the framework identifies early signals of degradation that are often invisible to static threshold-based alerting. These predictions are then translated into intelligent healing actions through an autonomous decision engine that prioritizes safety, minimizes disruption, and continuously learns from remediation outcomes. The closed-loop design enables Kubernetes clusters to anticipate failures, respond faster, and improve reliability without constant operator intervention.

The conceptual evaluation and case studies demonstrate tangible operational benefits, including reduced mean time to detect (MTTD), lower mean time to recovery (MTTR), improved service-level objective (SLO) adherence, and significant reductions in alert fatigue and operational toil. Importantly, these improvements are achieved while preserving transparency, auditability, and compatibility with managed Kubernetes platforms such as AKS and EKS.

While challenges remain—particularly around data quality, concept drift, and trust in autonomous actions—the results indicate that machine learning-driven autonomy is both feasible and impactful in production Kubernetes environments. Rather than replacing existing Kubernetes primitives, the proposed framework augments them with adaptive intelligence that evolves alongside workloads and infrastructure.

In summary, this work contributes a practical, scalable approach to autonomous Kubernetes reliability engineering. It advances the state of cloud-native operations by demonstrating how machine learning can be safely integrated into control loops to enable self-managing clusters. As organizations continue to demand higher reliability with lower operational overhead, autonomous healing systems such as the one presented here represent a critical step toward the future of resilient cloud infrastructure.

## AUTHOR BIO

**Sai Bharath Sannareddy** is a Senior Cloud Infrastructure Engineer specializing in cloud reliability engineering, large-scale observability architectures, and distributed systems automation. His work spans multi-cloud automation, SRE operational frameworks, and proactive incident-detection systems across hyperscale platforms. His research interests include reliability modeling, temporal alignment in distributed telemetry, predictive cloud resilience, and cross-cloud outage intelligence.

## REFERENCES

- [1] B. Burns, B. Grant, D. Oppenheimer, E. Brewer, and J. Wilkes, “Borg, Omega, and Kubernetes,” *Commun. ACM*, vol. 59, no. 5, pp. 50–57, 2016.
- [2] Kubernetes Authors, “Kubernetes: Production-grade container orchestration,” 2023. [Online]. Available: <https://kubernetes.io>
- [3] M. Fowler and J. Lewis, “Microservices: A definition of this new architectural term,” 2014.
- [4] L. H. Lim et al., “Failure prediction in cloud computing environments,” *IEEE Trans. Cloud Comput.*, vol. 9, no. 2, pp. 602–615, 2021.
- [5] S. Meng, L. Liu, and T. Wang, “Event-driven anomaly detection for cloud-native systems,” *IEEE Access*, vol. 8, pp. 135742–135755, 2020.
- [6] J. Dean and L. A. Barroso, “The tail at scale,” *Commun. ACM*, vol. 56, no. 2, pp. 74–80, 2013.
- [7] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. MIT Press, 2018.
- [8] M. Chen et al., “Machine learning for systems,” *IEEE Micro*, vol. 41, no. 3, pp. 12–22, 2021.
- [9] C. Delimitrou and C. Kozyrakis, “Quasar: Resource-efficient and QoS-aware cluster management,” *ASPLOS*, pp. 127–144, 2014.
- [10] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [11] D. Sculley et al., “Hidden technical debt in machine learning systems,” *Proc. NeurIPS*, pp. 2503–2511, 2015.
- [12] Google SRE Team, *Site Reliability Engineering: How Google Runs Production Systems*, O’Reilly Media, 2016.



- [13] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," *Proc. ACM SIGKDD*, pp. 785–794, 2016.
- [14] A. Fox et al., "Above the clouds: A Berkeley view of cloud computing," *UC Berkeley Tech. Rep.*, 2009.
- [15] J. O. Kephart and D. M. Chess, "The vision of autonomic computing," *Computer*, vol. 36, no. 1, pp. 41–50, 2003.
- [16] H. Kang et al., "Self-healing cloud services using machine learning," *IEEE Trans. Services Comput.*, vol. 13, no. 5, pp. 867–880, 2020.
- [17] Y. Chen et al., "Predictive failure analysis in distributed systems," *IEEE Trans. Dependable Secure Comput.*, vol. 18, no. 3, pp. 1231–1244, 2021.
- [18] P. Jamshidi et al., "Machine learning meets DevOps," *IEEE Software*, vol. 35, no. 5, pp. 24–31, 2018.
- [19] A. G. Ganek and T. A. Corbi, "The dawning of the autonomic computing era," *IBM Syst. J.*, vol. 42, no. 1, pp. 5–18, 2003.
- [20] M. Zaharia et al., "Discretized streams: Fault-tolerant streaming computation," *Proc. SOSP*, pp. 423–438, 2013.
- [21] S. Kavulya et al., "An analysis of traces from a production MapReduce cluster," *IEEE Cluster*, pp. 1–10, 2010.
- [22] P. Barham et al., "Xen and the art of virtualization," *Proc. SOSP*, pp. 164–177, 2003.
- [23] H. Xu et al., "Experience-driven anomaly detection for cloud services," *IEEE Trans. Cloud Comput.*, vol. 10, no. 1, pp. 45–59, 2022.
- [24] Microsoft Azure, "AKS reliability and best practices," Azure Architecture Center, 2023.
- [25] Amazon Web Services, "Best practices for Amazon EKS," AWS Whitepapers, 2023.
- [26] J. Wilkes, "More Google cluster data," *Google Research*, 2011.
- [27] Y. Bengio, I. Goodfellow, and A. Courville, *Deep Learning*, MIT Press, 2016.
- [28] P. O'Connor et al., "Observability-driven operations for cloud-native systems," *IEEE Cloud*, pp. 88–95, 2021.
- [29] S. Bubeck et al., "Challenges and opportunities in ML-driven systems automation," *Proc. ICML Workshop*, 2022.
- [30] A. Verma et al., "Large-scale cluster management at Google with Borg," *EuroSys*, pp. 1–17, 2015.