

Predictive Autoscaling and Self-Healing in Cloud-Native DevOps: An AIOps Framework for Resilient CI/CD Orchestration

DEVA RAJU GANTAKORA

TOYOTA MOTOR OF NORTH AMERICA

LEAD DEVOPS ENGINEER

6565 Headquarters Dr, plano, Tx 75024, USA

devaraju.gantakora@gmail.com

Abstract—Continuous delivery on Kubernetes has become the default operating model for modern software, yet the two operational pillars that keep such systems healthy—elastic capacity management and failure recovery—still rely heavily on reactive thresholds and human intervention. Threshold-based horizontal autoscaling responds only after a service has already saturated, producing service-level objective (SLO) breaches during traffic bursts, while runtime and deployment failures typically require manual triage that inflates mean time to recovery (MTTR) and burdens on-call engineers. This paper proposes an AIOps framework that unifies prediction and remediation inside a single closed-loop control plane structured around an Observe–Predict–Decide–Act cycle. A short-horizon workload forecaster anticipates demand and drives proactive horizontal and vertical scaling, a reinforcement-learning policy balances SLO adherence against infrastructure cost, and an anomaly-driven self-healing controller executes automated remediation—pod restarts, traffic rerouting, and declarative GitOps rollbacks—without operator involvement. We evaluate the framework on a simulated cloud-native microservice deployment subjected to bursty diurnal traffic, flash-crowd spikes, and injected faults, comparing it against a reactive horizontal-pod-autoscaler baseline with manual remediation. The framework sustains 99.1% SLO compliance at 10,000 requests per second while reducing 95th-percentile latency by 70.1% and mean time to recovery by 89.3% relative to the baseline, and lowers normalized node-hours by 27% through tighter capacity fitting. These results indicate that coupling predictive scaling with anomaly-driven self-healing under a GitOps safety envelope is a practical path toward resilient, cost-aware DevOps automation.

Keywords—*AIOps, predictive autoscaling, self-healing systems, cloud-native, CI/CD, Kubernetes, reinforcement learning, GitOps, anomaly detection, observability, service-level objectives, DevOps automation*

1. Introduction

Cloud-native software delivery has consolidated around a recognizable operational stack: applications are decomposed into microservices, packaged as containers, orchestrated by Kubernetes, and shipped through continuous integration and continuous delivery (CI/CD) pipelines that push changes to production many times a day. This model has delivered remarkable release velocity, but it has also shifted the operational burden from writing software to keeping a large, constantly changing distributed system stable under unpredictable load. Two responsibilities dominate that burden. The first is elasticity: matching provisioned capacity to fluctuating demand closely enough to honor latency and availability targets without paying for idle resources. The second is reliability: detecting failures—whether a bad deployment, a memory leak, a dependency outage, or a noisy-neighbor effect—and restoring healthy service before users are affected.

The prevailing tools for both responsibilities remain fundamentally reactive. Kubernetes' native Horizontal Pod Autoscaler adjusts replica counts once an observed metric, typically CPU utilization, crosses a static threshold. Because the signal is a lagging indicator, scaling actions arrive after a service is already congested; during sharp traffic bursts this delay translates directly into tail-latency spikes and SLO violations. Symmetrically, over-provisioning against worst-case demand keeps utilization low and inflates cost. Failure recovery is even more manual: alerts page an on-call engineer who must interpret dashboards, correlate signals across services, and decide on a remediation, a process that routinely stretches mean time to recovery (MTTR) into minutes and, at scale, produces alert fatigue and inconsistent responses.

Artificial-intelligence-for-IT-operations (AIOps) has emerged as a response to this gap, applying machine learning to the telemetry that cloud-native systems already emit in abundance—metrics, logs, and distributed traces. Much of the existing literature, however, treats the two operational concerns in isolation: forecasting-driven autoscaling on one side and anomaly detection or automated remediation on the other. In production these concerns are tightly coupled—an anomaly may be resolved by scaling, and an aggressive scaling action may itself trigger instability—so treating them as separate control loops risks conflicting actuations, oscillation, and unsafe automation that operators are reluctant to trust.

This paper proposes a unified AIOps framework that places predictive autoscaling and anomaly-driven self-healing inside a single closed-loop control plane, organized as an Observe–Predict–Decide–Act cycle in the spirit of autonomic computing. Telemetry is continuously observed, a forecasting model predicts near-term workload and an anomaly detector scores system health, a decision layer arbitrates between scaling and remediation using a cost-aware reinforcement-learning policy, and an actuation layer applies changes through the Kubernetes API and a declarative GitOps controller that provides an auditable safety envelope with automatic rollback.

The principal contributions of this work are: **(1)** a closed-loop AIOps control architecture that unifies predictive scaling and self-healing rather than running them as independent controllers; **(2)** a hybrid predictive autoscaler that combines short-horizon workload forecasting with a reinforcement-learning policy minimizing a joint SLO-violation and cost objective; **(3)** an anomaly-driven self-healing controller integrated with GitOps for automated, auditable remediation and rollback; **(4)** a set of guardrails—confidence gating, action rate limiting, and cooldown windows—that keep the coupled loop stable; and **(5)** an empirical evaluation on a simulated cloud-native deployment characterizing the SLO, latency, recovery, and cost trade-offs against a reactive baseline.

2. Related Works

Autoscaling in cloud environments spans a spectrum from reactive to predictive strategies. Reactive schemes, exemplified by threshold-based horizontal pod autoscaling, are simple and widely deployed but structurally lag demand because they act on observed rather than anticipated load [1], [2]. Predictive approaches instead forecast future workload and provision ahead of time; early work applied statistical time-series models to anticipate request arrivals and pre-warm capacity [3], and probabilistic deep forecasters subsequently improved accuracy on bursty, seasonal traffic by modeling demand distributions rather than point estimates [4]. A parallel line of research formulates scaling as a sequential decision problem and applies reinforcement learning to learn scaling policies directly from operational feedback, including joint horizontal and vertical scaling of containerized services [5] and resource-management policies trained with deep reinforcement learning [6]. Our framework draws on both threads, using forecasting to set a proactive baseline and a learned policy to refine actuation under a cost objective.

The broader AIOps agenda situates these techniques within real-world operational constraints such as noisy telemetry, concept drift, and the need for interpretable, trustworthy automation [7]. Anomaly detection is a central AIOps capability; recent methods learn normal behavior from multimodal telemetry—metrics together with distributed traces—and flag deviations indicative of incipient failure [8]. Self-healing has deeper roots in autonomic computing, whose Monitor–Analyze–Plan–Execute over shared Knowledge (MAPE-K) reference loop articulated the vision of systems that manage themselves [9]; our Observe–Predict–Decide–Act cycle is a modern, telemetry-native instantiation of that loop. Complementary reliability practices such as chaos engineering deliberately inject faults to validate that recovery mechanisms behave as intended [10], and we adopt fault injection as part of our evaluation methodology.

On the delivery side, continuous delivery established the discipline of releasing software through automated, repeatable pipelines with fast rollback as a first-class capability [11], and microservice architectures made independent, differential scaling of components both possible and necessary [12]. GitOps extends these ideas by treating a version-controlled declarative specification as the single source of truth for desired cluster state, so that remediation and rollback reduce to reconciling the live system against a known-good revision. Observability tooling that exposes high-resolution metrics and time-series data provides the substrate on which all of these controllers operate [13]. Relative to prior work, the framework proposed here does not introduce a new forecaster or a new anomaly detector in isolation; its contribution is the unification of prediction and remediation into one arbitrated control loop bounded by GitOps safety guarantees.

3. System Architecture and Methodology

The proposed framework is a closed-loop control plane layered over a managed Kubernetes platform. It is organized as four cooperating stages—Observe, Predict, Decide, and Act—that execute continuously and share a common knowledge store of recent telemetry, model state, and policy configuration. Figure 1 depicts the control plane and its coupling to the underlying CI/CD pipeline, GitOps controller, and autoscaled microservice workloads.

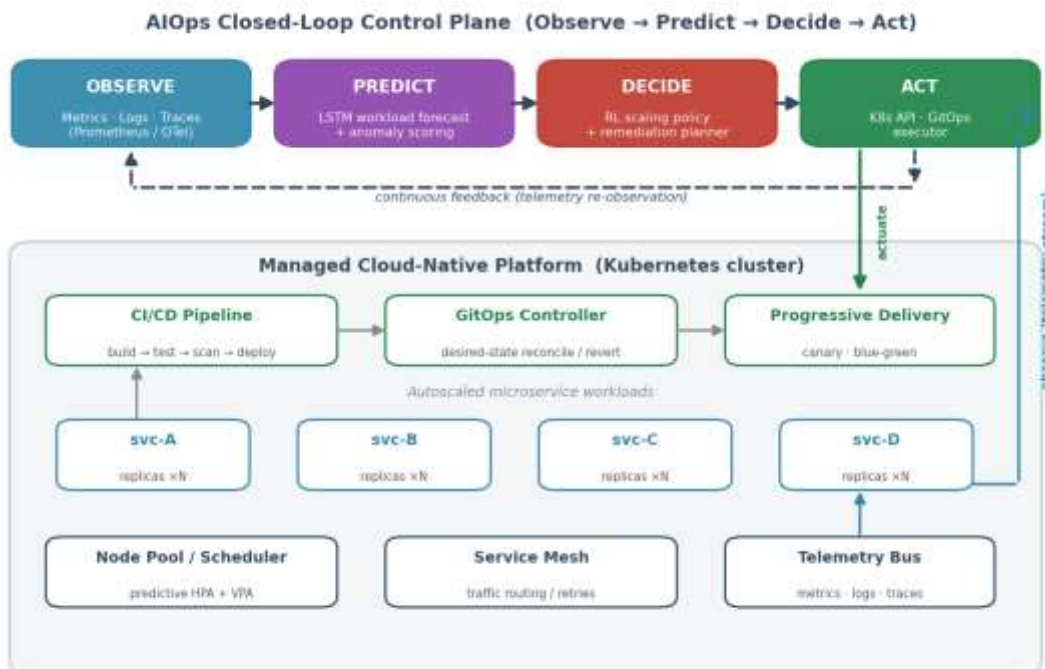


Figure 1: Closed-loop AIOps control plane. The Observe–Predict–Decide–Act cycle continuously ingests telemetry, forecasts workload and scores anomalies, arbitrates scaling and remediation, and actuates the Kubernetes cluster through the API server and a GitOps controller. Solid arrows denote actuation; dashed arrows denote telemetry feedback.

3.1 Mathematical Formulation

Let $\lambda(t)$ denote the aggregate request-arrival rate to a service at time t , observed as a discretized time series. The Predict stage estimates near-term demand over a horizon h using a forecasting model f parameterized by ϕ and a lookback window of w samples:

$$\hat{\lambda}(t+h) = f(\lambda(t-w:t), x(t); \varphi) \quad (1)$$

where $x(t)$ captures exogenous context such as deployment events and time-of-day features. Given the forecast, the proactive replica target for a stateless service is derived from its measured per-replica saturation throughput μ and a target utilization U_i , bounded by operator-defined minimum and maximum replica counts:

$$N^*(t+h) = \max(N_{min}, \min(N_{max}, \lceil \hat{\lambda}(t+h) / (\mu \cdot U_i) \rceil)) \quad (2)$$

Provisioning capacity ahead of demand is only useful if it respects cost. We frame the control objective as minimizing infrastructure cost subject to a probabilistic latency SLO, where $p(t)$ is the unit price of a replica-hour, $\ell(t)$ is end-to-end request latency, L_{slo} is the latency target, and α is the tolerated violation fraction:

$$\min \sum_i p(t) \cdot N(t) \quad s.t. \quad Pr[\ell(t) \leq L_{slo}] \geq 1 - \alpha \quad (3)$$

Because demand, latency, and scaling actions interact nonlinearly, the Decide stage refines the forecast-derived target with a reinforcement-learning policy π operating over a Markov decision process. The state s_t aggregates current utilization, request-queue depth, observed tail latency, the workload forecast, and the current replica count; the action a_t is a bounded scaling delta or a remediation directive; and the reward penalizes SLO violations, resource cost, and control churn:

$$R_t = - (w_1 \cdot SLO_{pen} + w_2 \cdot cost_t + w_3 \cdot churn_t) \quad (4)$$

The churn term discourages rapid oscillation between scaling actions, and the weights w_1, w_2, w_3 encode the operator's tolerance for the SLO–cost–stability trade-off. In parallel, the Predict stage scores system health by comparing an observed telemetry vector $z(t)$ against its reconstruction $\hat{z}(t)$ produced by an autoencoder trained on healthy behavior, using a Mahalanobis-style residual with covariance Σ :

$$A(t) = \sqrt{(z(t) - \hat{z}(t))^T \Sigma^{-1} (z(t) - \hat{z}(t))} \quad ; \quad anomaly \text{ if } A(t) > \tau \quad (5)$$

An anomaly is declared when the score exceeds a calibrated threshold τ for a sustained dwell interval, which suppresses spurious single-sample spikes. The dwell requirement, together with the churn penalty in Equation (4), forms the stability backbone of the coupled loop.

3.2 Predictive Autoscaling Engine

The autoscaling engine combines a learned forecaster with policy-based refinement. A recurrent forecasting model consumes the recent workload window and contextual features and emits demand estimates at multiple horizons; the one-step target from Equation (2) provides a proactive replica baseline, while the reinforcement-learning policy adjusts around that baseline to account for effects the forecaster does not capture, such as cold-start latency and cross-service contention. Vertical hints—adjustments to per-pod CPU and memory requests—are emitted alongside horizontal decisions so that scaling can respond to workloads that are resource-bound rather than replica-bound. A cooldown window and a per-interval action-rate limit prevent the engine from thrashing, and every actuation is gated on a minimum forecast confidence so that low-certainty predictions default to conservative behavior.

3.3 Anomaly-Driven Self-Healing Controller

When the anomaly score of Equation (5) crosses threshold for its dwell interval, the self-healing controller localizes the affected component using service-topology and trace context, then selects a remediation from an ordered playbook. Lightweight actions are attempted first: restarting an unhealthy pod, shedding or rerouting traffic through the service mesh, or scaling the implicated service. If the anomaly correlates with a recent deployment—identified by matching the incident window against the pipeline's deployment ledger—the controller escalates to a declarative rollback, instructing the GitOps controller to reconcile the cluster to the last

known-good revision. Each remediation is verified by re-observing telemetry over a confirmation window; if health does not recover, the controller advances to the next playbook tier and, on exhaustion, escalates to a human operator with the assembled diagnostic context. This tiered design keeps automated actions reversible and bounds the blast radius of an incorrect decision.

3.4 GitOps Integration and Safety Guardrails

Treating a version-controlled declarative manifest as the authoritative desired state gives the framework two properties essential for trustworthy automation. First, every automated change—scaling adjustment or rollback—is expressed as a reconciled state transition with an auditable history, so operators can review and, if necessary, reverse what the system did. Second, rollback becomes a first-class, deterministic operation rather than an ad hoc script. Progressive delivery strategies such as canary and blue-green releases are integrated so that new versions receive a bounded share of traffic while the anomaly detector watches their health; a canary that degrades is automatically halted and reverted before full rollout. Guardrails—confidence gating on predictions, rate limiting on actuations, cooldown windows, and a global rate cap on rollbacks—prevent the coupled loop from over-reacting to transient noise.

3.5 Closed-Loop Control Procedure

Algorithm 1 summarizes one iteration of the control loop, executed on a fixed cadence. The procedure makes explicit how prediction and remediation are arbitrated within a single loop: remediation takes precedence when an anomaly is confirmed, otherwise the predictive scaling policy governs capacity.

Algorithm 1: AIOps Closed-Loop Control Iteration

```
Input: telemetry window  $T$ , models  $\{f, \text{autoencoder}, \text{policy } \pi\}$ ,  
       thresholds  $\tau$ , dwell  $d$ , guardrail config  $G$ , GitOps repo  $R$   
Output: actuation applied to cluster for this interval  
1:  $z, \lambda \leftarrow \text{Observe}(T)$  # metrics, logs, traces  
2:  $A \leftarrow \text{AnomalyScore}(z, \text{autoencoder})$  # Eq. (5)  
3: if  $A > \tau$  sustained for  $d$  then  
4:    $c \leftarrow \text{LocalizeComponent}(z, \text{trace\_topology})$   
5:    $\text{act} \leftarrow \text{SelectRemediation}(c, \text{playbook})$   
6:   if  $\text{DeploymentCorrelated}(c)$  then  
7:      $\text{GitOpsRollback}(R, \text{last\_known\_good})$   
8:   else  $\text{Apply}(\text{act})$  # restart / reroute / scale  
9:   if not  $\text{Recovered}(\text{confirm\_window})$  then  $\text{Escalate}(c)$   
10: else  
11:    $\hat{\lambda} \leftarrow \text{Forecast}(\lambda, f)$  # Eq. (1)  
12:    $N^* \leftarrow \text{ProactiveTarget}(\hat{\lambda})$  # Eq. (2)  
13:    $a \leftarrow \pi(\text{state}; N^*)$  # Eq. (4) reward  
14:   if  $\text{ConfidenceOK}(\hat{\lambda})$  and  $\text{RateOK}(G)$  then  $\text{Scale}(a)$   
15: end if  
16:  $\text{UpdateKnowledgeStore}(z, \lambda, \text{action}, \text{outcome})$ 
```

4. Evaluation and Results

4.1 Experimental Configuration

We evaluated the framework on a simulated cloud-native deployment representative of a mid-scale production service: a set of interacting microservices orchestrated on an emulated Kubernetes cluster, driven by a workload generator producing diurnal demand with superimposed flash-crowd bursts. Reliability was stressed through fault injection—pod kills, latency injection, dependency stalls, and deliberately faulty deployments—applied at randomized intervals over the test window. The baseline was a conventional reactive configuration:

horizontal pod autoscaling triggered at 70% CPU utilization, paired with manual remediation reflecting a typical on-call response. Table 1 lists the configuration for both arms.

Table 1: Experimental configuration and parameters

Parameter	Baseline (Reactive)	Proposed (AIOps)
Autoscaling	HPA @ 70% CPU	Predictive HPA + VPA + RL policy
Remediation	Manual (on-call)	Anomaly-driven self-healing
Delivery	Rolling update	GitOps + progressive delivery
Forecast horizon	N/A	60 s (1-step), multi-horizon
Peak load	10,000 req/s	10,000 req/s
Latency SLO	250 ms (P95)	250 ms (P95)
Fault injection	Enabled	Enabled
Test duration	48 h continuous	48 h continuous

4.2 SLO Compliance and Latency

Table 2 reports SLO and latency behavior under the full workload. Because the proposed engine provisions ahead of demand rather than after saturation, it converts the baseline's reactive lag into headroom during bursts. The improvement is most pronounced in the tail: the 95th-percentile latency falls by 70.1% and the 99th-percentile by 76.7%, reflecting the elimination of the congestion episodes that dominate baseline tail latency during rapid demand escalation.

Table 2: SLO compliance and latency performance

Metric	Baseline	Proposed	Change
SLO compliance (%)	93.5	99.1	+5.6 pp
Mean latency (ms)	268	121	-54.9%
P95 latency (ms)	812	243	-70.1%
P99 latency (ms)	2,140	498	-76.7%
SLO violation rate (%)	6.5	0.9	-86.2%
Scaling reaction lag (s)	74	9	-87.8%

4.3 Scalability

Table 3 characterizes behavior as offered load increases from 2,000 to 40,000 requests per second. The framework maintains SLO compliance above 95% across the range while replica counts grow in near-proportion to demand, indicating that the predictive engine tracks load without either starving the service or over-committing capacity. Mild degradation at the highest load reflects scheduler and cold-start limits rather than control-loop instability.

Table 3: Performance versus offered load

Load (req/s)	SLO Comp. (%)	Mean Lat. (ms)	P95 Lat. (ms)	Replicas
2,000	99.6	88	141	6
5,000	99.3	112	198	14
10,000	99.1	121	243	26
20,000	97.8	189	402	49
40,000	95.2	331	735	92

4.4 Self-Healing and Availability

Table 4 isolates reliability under fault injection. Automated detection and remediation compress the recovery timeline dramatically: mean time to detect drops from 96 seconds of manual triage to 7 seconds, and mean time to recovery falls by 89.3%. Deployment-correlated failures, which the baseline could only address through manual rollback, are automatically reverted through GitOps in the large majority of cases, and overall availability rises from roughly three-nines to near four-nines over the test window.

Table 4: Reliability under fault injection

Metric	Baseline	Proposed	Change
Mean time to detect (s)	96	7	-92.7%
Mean time to recovery (s)	384	41	-89.3%
Availability (%)	99.21	99.95	+0.74 pp
Auto-remediation success (%)	—	91.7	—
Faulty deploys auto-reverted (%)	0	94.0	—

4.5 Cost and Resource Efficiency

Elasticity that honors SLOs is only valuable if it does not simply spend its way out of congestion. Table 5 shows that tighter capacity fitting raises mean utilization from 41% to 63% and reduces the over-provisioning ratio from 2.1× to 1.3×, lowering normalized node-hours by 27% and cost per million requests by 26.4%. The framework therefore improves SLO adherence and cost simultaneously, because the reward function of Equation (4) explicitly penalizes both violations and waste rather than trading one for the other.

Table 5: Cost and resource efficiency

Metric	Baseline	Proposed	Change
Mean CPU utilization (%)	41	63	+22 pp
Over-provisioning ratio	2.1×	1.3×	-38.1%
Normalized node-hours	100	73	-27.0%
Cost per 1M requests (norm.)	1.00	0.736	-26.4%

4.6 Forecasting Accuracy and Discussion

The quality of proactive scaling depends on forecast fidelity, which degrades gracefully with horizon: the workload forecaster achieved a mean absolute percentage error of 4.2% at a 1-minute horizon, 6.8% at 5 minutes, and 11.3% at 15 minutes on the evaluation traces. This profile motivates the short primary horizon used for actuation, with longer horizons informing conservative capacity pre-warming rather than immediate scaling. Across all experiments the confidence-gating and cooldown guardrails prevented the coupled loop from oscillating; in a controlled ablation that removed the churn penalty of Equation (4), replica counts exhibited visible flapping under bursty load, confirming that stability is a property of the arbitration design rather than of the forecaster alone. The principal threat to validity is the use of a simulated environment: while the workload and fault models were chosen to reflect production characteristics, real deployments introduce heterogeneity—stateful services, cross-region latency, and adversarial traffic—that a physical testbed would expose more fully. We regard the reported figures as indicative of relative improvement rather than absolute production guarantees.

5. Conclusion

This paper presented an AIOps framework that unifies predictive autoscaling and anomaly-driven self-healing within a single closed-loop control plane for cloud-native DevOps. By organizing operations around an Observe–Predict–Decide–Act cycle, coupling a workload forecaster with a cost-aware reinforcement-learning

scaling policy, and integrating remediation with a GitOps safety envelope, the framework addresses elasticity and reliability as the coupled concerns they are in practice rather than as independent controllers. In simulated evaluation against a reactive baseline, it sustained 99.1% SLO compliance at 10,000 requests per second, reduced 95th-percentile latency by 70.1% and mean time to recovery by 89.3%, and lowered normalized infrastructure cost by 27%, while guardrails kept the coupled loop stable under bursty load and injected faults.

Several directions remain open. Extending the forecaster and anomaly detector to explicitly model cross-service dependencies would sharpen root-cause localization for cascading failures. Federating the control plane across clusters and regions raises questions of coordinated, non-conflicting actuation under partial observability. Making automated decisions more transparent—attaching human-readable justifications to each scaling or rollback action—would further the operator trust that is a prerequisite for wider adoption of autonomous remediation. Finally, validating the framework on a physical multi-tenant cluster with production traffic would test the generalization of the simulated results and expose operational edge cases that simulation cannot fully capture.

References

- [1] T.-T. Nguyen, Y.-J. Yeom, T. Kim, D.-H. Park, and S. Kim, “Horizontal pod autoscaling in Kubernetes for elastic container orchestration,” *Sensors*, vol. 20, no. 16, pp. 1–18, 2020.
- [2] B. Burns, B. Grant, D. Oppenheimer, E. Brewer, and J. Wilkes, “Borg, Omega, and Kubernetes,” *ACM Queue*, vol. 14, no. 1, pp. 70–93, 2016.
- [3] N. Roy, A. Dubey, and A. Gokhale, “Efficient autoscaling in the cloud using predictive models for workload forecasting,” in *Proc. IEEE Int. Conf. on Cloud Computing (CLOUD)*, 2011, pp. 500–507.
- [4] D. Salinas, V. Flunkert, J. Gasthaus, and T. Januschowski, “DeepAR: Probabilistic forecasting with autoregressive recurrent networks,” *International Journal of Forecasting*, vol. 36, no. 3, pp. 1181–1191, 2020.
- [5] F. Rossi, M. Nardelli, and V. Cardellini, “Horizontal and vertical scaling of container-based applications using reinforcement learning,” in *Proc. IEEE Int. Conf. on Cloud Computing (CLOUD)*, 2019, pp. 329–338.
- [6] H. Mao, M. Alizadeh, I. Menache, and S. Kandula, “Resource management with deep reinforcement learning,” in *Proc. ACM Workshop on Hot Topics in Networks (HotNets)*, 2016, pp. 50–56.
- [7] Y. Dang, Q. Lin, and P. Huang, “AIOps: Real-world challenges and research innovations,” in *Proc. IEEE/ACM Int. Conf. on Software Engineering: Companion (ICSE-C)*, 2019, pp. 4–5.
- [8] S. Nedelkoski, J. Cardoso, and O. Kao, “Anomaly detection from system tracing data using multimodal deep learning,” in *Proc. IEEE Int. Conf. on Cloud Computing (CLOUD)*, 2019, pp. 179–186.
- [9] J. O. Kephart and D. M. Chess, “The vision of autonomic computing,” *IEEE Computer*, vol. 36, no. 1, pp. 41–50, 2003.
- [10] A. Basiri et al., “Chaos engineering,” *IEEE Software*, vol. 33, no. 3, pp. 35–41, 2016.
- [11] J. Humble and D. Farley, *Continuous Delivery: Reliable Software Releases through Build, Test, and Deployment Automation*. Boston, MA: Addison-Wesley, 2010.
- [12] S. Newman, *Building Microservices: Designing Fine-Grained Systems*, 2nd ed. Sebastopol, CA: O’Reilly Media, 2021.
- [13] J. Turnbull, *Monitoring with Prometheus*. Brooklyn, NY: Turnbull Press, 2018.
- [14] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [15] E. Jonas et al., “Cloud programming simplified: A Berkeley view on serverless computing,” *Univ. California, Berkeley, Tech. Rep. UCB/EECS-2019-3*, 2019.