

Self-Healing Memory Systems in AI Fabrics: Machine Learning-Driven Predictive Detection and Autonomous Mitigation of Memory Leaks in High-Performance Network Switches

Srinivas Yadam

Independent Researcher., USA

Abstract

The explosive growth of AI workloads is driving the shift today towards ultra-high-speed 800G and in future 1.6T networks, where control-plane processes in network switches increasingly suffer from hidden memory leaks and non-deterministic allocation behaviors under sustained load. These anomalies remain undetected by traditional reactive monitoring tools, resulting in service degradation, out-of-memory crashes, and cascading fabric instability. This work introduces Self-Healing Memory Systems (SHMS), an AI-driven reliability framework that continuously learns memory health patterns and autonomously predicts faults before failure occurs. SHMS employs a hybrid architecture combining Long Short-Term Memory (LSTM) networks for temporal anomaly forecasting and Reinforcement Learning (RL) for optimal mitigation decision-making. Deployed on an 800G switch testbed managing large-scale AI training workloads, SHMS demonstrates early fault detection with extended lead time, enabling autonomous recovery actions including buffer reclamation, process isolation, and adaptive throttling. The framework maintains deterministic performance while significantly reducing manual intervention requirements, representing a critical advancement toward autonomously reliable, self-managing AI fabric infrastructure for next-generation machine learning clusters.

Keywords: Self-Healing Networks, Memory Leak Prediction, Lstm Forecasting, Reinforcement Learning, Ai Fabric Reliability

1. Introduction

1.1 Background Context: Artificial Intelligence and Machine Learning Workload Environments with Network Switch Complexity

Contemporary artificial intelligence and machine learning workloads operate within massively parallel infrastructures comprising high-performance network switches, network interface cards (NICs), and graphics processing units (GPUs), all interconnected through ultra-low-latency Remote Direct Memory Access (RDMA) links such as RDMA over Converged Ethernet version 2 (RoCEv2). Individual network switches within these environments handle millions of dynamically load-balanced flows while managing intensive control-plane signaling and comprehensive telemetry collection, thereby introducing intricate and continuous memory operations spanning multiple software layers [1]. Modern network switch architecture has transcended simple packet forwarding functionality. These devices host numerous control-plane daemons, encompassing routing protocols like Border Gateway Protocol (BGP), software-defined networking (SDN) orchestration via P4Runtime, and telemetry services through gRPC Network Management Interface (gNMI) and gRPC collectors. Control-plane processes perform continuous heap memory allocation and deallocation operations for routing tables, flow state maintenance, counter management, and analytics data processing. Within dynamically scaling environments, even minor inefficiencies or incompletely released allocations progressively accumulate, producing memory leaks, fragmentation issues, and eventual heap exhaustion that threatens overall fabric stability. Managing AI

10.48047/jocaaa.2025.34.11.46

and ML workloads across distributed GPU instances introduces additional complexity, as network switches must deliver consistent performance while accommodating variable traffic patterns and fluctuating resource demands.

1.2 Constraints of Conventional Observability and Flow Control Mechanisms

Conventional flow-control protocols, including Priority Flow Control (PFC), Explicit Congestion Notification (ECN), and Data Center Quantized Congestion Notification (DCQCN), demonstrate effectiveness in mitigating packet-level congestion, yet these protocols function exclusively at the data-plane layer. Such mechanisms ensure lossless transport during traffic bursts but lack visibility into control-plane internal memory health. This architectural constraint allows switches to maintain apparent stability at the packet forwarding level while simultaneously experiencing silent heap memory consumption from unreleased allocations, ultimately manifesting as watchdog resets, process crashes, or complete control-plane reboots. Standard observability instruments such as Simple Network Management Protocol (SNMP), system logs, gNMI exporters, and heap profilers operate under reactive paradigms. These instruments detect symptomatic indicators like elevated latency, process restarts, or critical memory alerts only after performance degradation has commenced. Capturing transient, non-linear, or long-horizon leaks that progressively develop during extended AI model training or inference cycles remains beyond their capabilities. Ensuring end-to-end assurance for AI and ML workloads necessitates comprehensive monitoring across multiple infrastructure layers, yet current mechanisms lack predictive capabilities required to prevent control-plane failures before application performance suffers [2]. Fabric operators frequently encounter slow, concealed degradations that emerge exclusively under sustained workloads, presenting critical risks within AI-driven clusters demanding high availability. Threshold-based alerting systems cannot differentiate between normal memory fluctuations and genuine leak progression patterns, producing frequent false positives or delayed detection of critical issues.

Mechanism/Approach	Operation Layer	Primary Function	Memory Visibility	Leak Detection Capability	Response Type
Priority Flow Control (PFC)	Data-plane	Lossless transport	None	No	Reactive
Explicit Congestion Notification (ECN)	Data-plane	Congestion signaling	None	No	Reactive
Data Center Quantized Congestion Notification (DCQCN)	Data-plane	Rate adjustment	None	No	Reactive
SNMP Monitoring	Management	Threshold alerts	Limited	Minimal	Reactive
System Logs	Management	Event recording	Limited	Minimal	Reactive
Heap Profilers	Control-plane	Memory analysis	High	Manual	Reactive
Self-Healing Memory System	Control-plane	Predictive mitigation	Comprehensive	Autonomous	Proactive

Table 1: Comparison of Flow Control Mechanisms and Memory Management Approaches [1][2][3][4]

1.3 Core Problem Definition and Knowledge Gaps

Network fault management has received considerable attention, yet memory leak prediction and automated mitigation remain unresolved within switch operating systems. Existing watchdog mechanisms depend on static thresholds, failing to capture dynamic trends and non-linear leak accumulation patterns. Current frameworks lack autonomous capabilities to predict, isolate, and correct memory leak behaviors before system failure occurs. Conventional data-plane congestion controls prevent packet loss but cannot observe or manage memory behavior at software or control-plane levels, emphasizing the necessity for software-centric reliability mechanisms that extend beyond traditional flow control. Artificial intelligence and machine learning workloads exhibit temporal variability in resource utilization patterns, compounding management challenges [1]. Memory allocation behaviors differ substantially across training phases, inference batches, and model architectures, complicating the establishment of universal detection baselines. The distributed architecture of modern artificial intelligence infrastructure means memory issues within individual switches can trigger cascading failures throughout the entire fabric, interrupting ongoing training jobs and causing service-level agreement violations [2]. Current reliability frameworks lack capabilities to correlate memory metrics with workload characteristics, predict failure trajectories, and execute autonomous recovery actions without manual intervention. This deficiency becomes particularly pronounced as artificial intelligence fabrics scale toward ultra-high-bandwidth configurations where manual troubleshooting proves infeasible.

1.4 Objectives and Boundaries of This Work

The Self-Healing Memory System represents an intelligent reliability framework leveraging machine learning to forecast memory anomalies, localize leak sources, and autonomously restore memory stability within the stringent operational constraints of high-throughput switches. This framework addresses escalating demands for autonomous management capabilities within AI network infrastructure, where manual intervention cannot scale proportionally with increasing complexity and operational velocity [1][2]. Continuous monitoring of memory metrics across control-plane processes enables identification of subtle patterns indicating emerging leak conditions before critical thresholds are reached. Learning temporal evolution patterns of memory consumption allows the framework to predict the timing and location of potential failures, providing adequate lead time for preventive action. An intelligent decision-making agent evaluates multiple mitigation strategies, selecting actions that maximize system stability while minimizing performance impact and operational disruption. Recovery actions execute within strict latency bounds, ensuring the healing process does not degrade network performance or interrupt AI and ML workloads. The framework encompasses design, implementation, and evaluation within high-performance network switches supporting AI fabric deployments. Particular emphasis is placed on achieving autonomous operation, predictive accuracy, and seamless integration with existing switch operating systems.

2. Related Work and Literature Review

2.1 Machine Learning Applications in Network Fault Prediction

Machine learning techniques have gained substantial momentum within network fault prediction and anomaly detection across telecommunication and data center environments. Traditional fault management systems relied heavily on rule-based approaches and static threshold monitoring, which proved inadequate for capturing complex, evolving fault patterns in modern heterogeneous networks [3]. Recent advances demonstrate that supervised learning models, encompassing decision trees, support vector machines, and neural networks, can effectively predict network failures by learning from historical fault

10.48047/jocaaa.2025.34.11.46

data and operational telemetry. Deep learning architectures, particularly recurrent neural networks and their variants, show promise in capturing temporal dependencies within network behavior, enabling proactive fault detection before service disruption occurs. Nevertheless, applying machine learning to network fault prediction encounters several obstacles, including the scarcity of labeled fault data, class imbalance between normal and failure states, and the dynamic nature of network topologies and traffic patterns [3]. While considerable progress has been achieved in predicting link failures, device outages, and traffic anomalies, memory-specific fault prediction within network switches remains an underexplored domain, particularly concerning control-plane memory leaks that develop gradually over extended operational periods.

Leak Pattern Type	Accumulation Behavior	Detection Difficulty	Temporal Signature	Traditional Threshold Effectiveness	ML-Based Detection Advantage
Gradual Linear	Constant rate increase	High	Long-horizon trend	Poor - remains below threshold	Excellent - captures trend slope
Exponential Growth	Accelerating accumulation	Moderate	Non-linear progression	Moderate - late detection	Excellent - identifies acceleration
Sporadic Burst	Intermittent spikes	Very High	Irregular intervals	Poor - excessive false positives	Good - learns spike patterns
Cyclical Accumulation	Periodic increase/decrease	High	Recurring patterns	Poor - mimics normal behavior	Excellent - frequency analysis
Threshold Hovering	Fluctuates near the limit	Moderate	High variability	Moderate - frequent alerts	Good - distinguishes stable vs leak

Table 2: Memory Leak Pattern Characteristics and Detection Challenges [3][6]

2.2 Memory Management in High-Performance Switches

Memory management within high-performance network switches presents distinctive challenges due to the dual requirements of deterministic packet forwarding and dynamic control-plane operations. Modern switches employ sophisticated buffer management strategies to handle traffic bursts, maintain Quality of Service guarantees, and prevent packet loss during congestion events. Traditional buffer management relied on static allocation schemes and simple drop policies such as tail drop or random early detection. Recent work has explored intelligent buffer management policies that leverage machine learning to optimize buffer allocation decisions dynamically [4]. These learning-based approaches demonstrate that adaptive buffer management can significantly improve throughput and reduce packet drops compared to conventional static policies. Reinforcement learning frameworks have been applied to learn optimal buffer allocation strategies by observing traffic patterns and queue states, enabling switches to make intelligent dropping and queuing decisions under varying load conditions [4]. Nevertheless, existing buffer management investigation focuses primarily on data-plane packet buffers and forwarding performance, with limited attention to control-plane memory health. The control plane operates under different constraints, managing routing state, flow tables, and telemetry data through heap-based memory allocation rather than fixed buffer pools. Memory leaks within control-plane processes exhibit different characteristics than data-plane congestion, requiring distinct detection and mitigation strategies that current buffer management approaches do not address.

2.3 Self-Healing Systems in Distributed Infrastructure

Self-healing systems represent an emerging paradigm within distributed infrastructure management, aiming to detect, diagnose, and recover from failures autonomously without human intervention. The concept draws inspiration from biological systems that maintain homeostasis through continuous monitoring and adaptive responses to internal disturbances. Within distributed computing environments,

10.48047/jocaaa.2025.34.11.46

self-healing mechanisms have been implemented at various layers, encompassing application-level recovery through service redundancy and replication, middleware-level fault tolerance through checkpointing and state migration, and infrastructure-level resilience through automated resource provisioning and workload redistribution. Autonomic computing principles guide the design of self-healing systems, emphasizing properties such as self-configuration, self-optimization, self-protection, and self-healing. Machine learning plays an increasingly important role in enabling these autonomic properties by learning normal system behavior, detecting deviations from expected patterns, and selecting appropriate recovery actions based on observed outcomes. Reinforcement learning has emerged as a particularly suitable approach for self-healing decision-making, as agents can learn optimal recovery policies through trial-and-error interaction with the system environment. Despite these advances, self-healing capabilities within network switches remain limited. Most existing implementations focus on rerouting traffic around failed links or redistributing workloads across healthy nodes, rather than addressing internal software faults within individual network devices. The control-plane software stack within modern switches represents a complex, stateful environment where memory leaks and resource exhaustion can develop gradually, requiring continuous learning and predictive intervention that extends beyond current self-healing approaches.

2.4 Gaps in Current Approaches

Despite substantial progress within network fault prediction, memory management, and self-healing systems, several critical gaps remain unaddressed concerning control-plane memory reliability for high-performance switches. Existing fault prediction models primarily target link failures, hardware faults, and traffic anomalies, with limited exploration of software-level memory leaks that manifest gradually over time [3]. The temporal dynamics of memory leak progression differ fundamentally from sudden hardware failures, requiring prediction models that can capture long-term trends and non-linear accumulation patterns across extended operational periods. Current buffer management investigation focuses on optimizing data-plane packet handling rather than control-plane memory health [4], leaving a significant gap in understanding how to predict and mitigate heap exhaustion within routing daemons, telemetry collectors, and orchestration services. Self-healing systems have demonstrated success within distributed environments but lack integration with the resource-constrained, real-time operational requirements of network switch control planes. Reinforcement learning approaches for autonomous recovery typically assume relatively coarse-grained action spaces and longer decision horizons, whereas switch control-plane recovery demands fine-grained actions executed within strict latency bounds to avoid disrupting packet forwarding. Furthermore, existing approaches lack mechanisms for continuous model retraining based on recovery outcomes, limiting their ability to adapt to evolving workload characteristics and newly emerging leak patterns. The intersection of predictive machine learning, autonomous recovery, and switch control-plane reliability remains largely unexplored, motivating the need for integrated frameworks that combine temporal anomaly forecasting with reinforcement-based mitigation within network switching environments supporting artificial intelligence and machine learning workloads [1][2].

3. System Design and Architecture

3.1 SHMS Framework Overview

The Self-Healing Memory System framework operates as a three-layer closed-loop architecture designed to autonomously detect, predict, and mitigate memory anomalies within high-performance network switches. The architecture integrates continuous telemetry monitoring, predictive machine learning inference, and autonomous actuation capabilities into a cohesive reliability framework. The closed-loop

10.48047/jocaaa.2025.34.11.46

design ensures that memory health assessment, anomaly forecasting, mitigation action selection, and post-recovery feedback occur in a continuous cycle, enabling the system to learn and adapt from each intervention. The framework addresses the fundamental challenge of maintaining control-plane stability under dynamic artificial intelligence and machine learning workload conditions, where memory consumption patterns exhibit high variability and unpredictability [1][2]. Unlike traditional reactive monitoring systems that respond only after threshold violations occur, this framework employs predictive analytics to forecast memory health degradation before critical failures manifest. The three-layer architecture separates concerns between data collection, intelligent decision-making, and system actuation, allowing each component to operate independently while maintaining tight coordination through well-defined interfaces. This separation enables flexible deployment models, where prediction and decision components can reside either locally on switch co-processors or centrally within fabric management controllers. The framework prioritizes minimal performance overhead, ensuring that monitoring and prediction activities do not interfere with packet forwarding operations or introduce additional latency into data-plane processing.

Layer	Component	Primary Function	Key Technologies	Data Flow	Update Frequency
Telemetry Collection	Metric Collectors	Real-time monitoring	In-band telemetry, gNMI	Upward to ML Engine	Adaptive (high-frequency during anomalies)
Telemetry Collection	Feature Extraction	Data preprocessing	Sliding windows, statistical analysis	Upward to ML Engine	Continuous
Predictive ML Engine	LSTM Network	Temporal forecasting	Deep learning, sequence modeling	Predictions to RL Agent	Per the prediction window
Predictive ML Engine	RL Decision Agent	Action selection	Policy networks, Q-learning	Commands to Actuation	Per decision cycle
Control and Actuation	Action Executor	Mitigation implementation	P4Runtime, gRPC APIs	System state changes	On-demand
Control and Actuation	Feedback Collector	Outcome monitoring	Post-action telemetry	Downward to ML Engine	Post-intervention

Table 3: Self-Healing Memory System Architecture Components [5][6][7]

Self-Healing Memory System (SHMS) Architecture

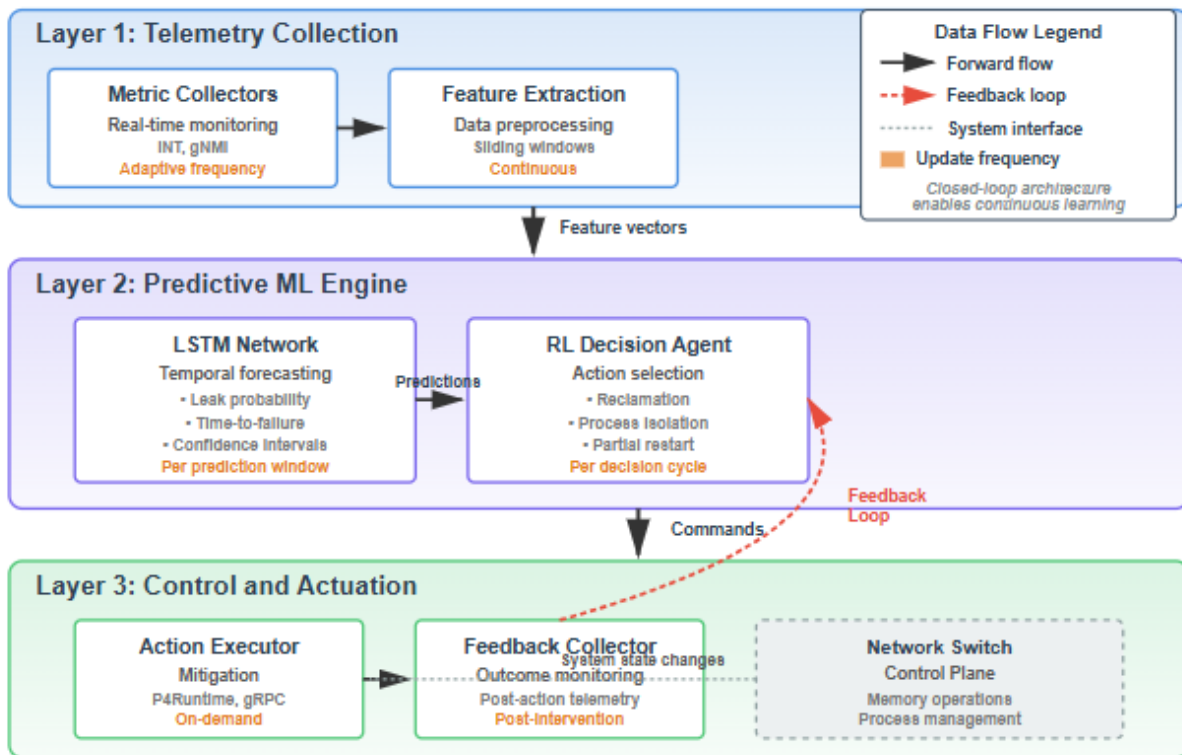


Figure 1: SHMS Three-Layer Closed-Loop Architecture

3.2 Telemetry Collection Layer

The telemetry collection layer forms the foundation by continuously gathering real-time metrics that characterize control-plane memory health and system behavior. Critical metrics include heap usage percentages, page fault frequencies, memory allocation and deallocation rates, fragmentation indices, and correlations with central processing unit (CPU) utilization patterns. Data collection occurs at high temporal resolution using sliding window analytics that capture both instantaneous measurements and temporal trends over varying time horizons. The layer employs in-band telemetry mechanisms and gNMI-based collectors embedded directly within the control plane, ensuring comprehensive visibility into memory operations across all control-plane daemons and processes [5]. In-band network telemetry (INT) provides advantages over traditional out-of-band monitoring by embedding telemetry data directly within network packets, enabling fine-grained, real-time visibility into network device states without imposing significant overhead on management networks [5]. The telemetry collection layer implements intelligent sampling strategies that balance collection granularity with computational overhead, adapting sampling rates based on detected anomaly indicators. When memory consumption exhibits stable patterns, the system operates with lower sampling frequencies to conserve resources. As early anomaly indicators emerge, sampling rates increase dynamically to capture detailed behavioral signatures. The collected telemetry data undergoes preprocessing and feature engineering to extract meaningful patterns, including rate-of-change calculations, moving averages, standard deviation measurements, and cross-correlation analyses between different memory metrics. This preprocessing transforms raw telemetry into feature

vectors suitable for machine learning model consumption, ensuring that temporal patterns and non-linear relationships become explicit in the input representation.

3.3 Predictive Machine Learning Engine (LSTM and RL Components)

The predictive machine learning engine comprises two integrated components: a Long Short-Term Memory (LSTM) network for temporal forecasting and a reinforcement learning (RL) agent for autonomous decision-making. The LSTM model learns the temporal evolution of memory metrics, capturing both short-term fluctuations and long-term trends that characterize memory leak progression [6]. LSTM networks excel at learning temporal dependencies in sequential data, making them particularly suitable for fault detection tasks where anomalies develop gradually over time [6]. The network architecture processes sliding windows of historical memory telemetry, producing probabilistic forecasts of future memory states and leak likelihood. The model outputs include leak probability estimates, projected time-to-failure predictions, and confidence intervals that quantify prediction uncertainty. Training occurs through supervised learning on labeled datasets combining normal operational periods with known leak events, supplemented by continuous online learning that incorporates newly observed patterns. The RL decision agent operates downstream from the LSTM forecasting module, receiving predicted memory states and current system conditions as input. The agent evaluates multiple mitigation action options, including memory reclamation through garbage compaction, process throttling or temporary freezing, partial restart or leak isolation, and deferred recovery when risk levels remain below intervention thresholds. The agent learns optimal action selection policies through trial-and-error interaction with the system environment, receiving rewards based on outcomes that balance multiple objectives. The reward function incorporates stability gains from successful interventions, latency penalties from actions that temporarily disrupt processing, and data loss impacts from aggressive recovery strategies. Through iterative learning, the agent discovers policies that maximize long-term system stability while minimizing operational disruptions, adapting its decision-making to specific workload characteristics and failure patterns observed in the deployment environment.

3.4 Control and Actuation Layer

The control and actuation layer translates machine learning decisions into concrete system actions, interfacing directly with the switch kernel, process scheduler, and control-plane service managers. This layer implements a comprehensive action execution framework that encompasses memory reclamation procedures, process management operations, service isolation mechanisms, and adaptive resource throttling capabilities. Memory reclamation actions trigger garbage collection routines, heap compaction operations, and cache clearing procedures that release unreferenced memory without terminating active processes. Process throttling actions temporarily reduce processing rates for suspected leak sources, limiting their memory allocation rates while maintaining basic functionality. Partial restart actions terminate and reinitialize specific control-plane daemons identified as leak sources, preserving overall switch operation while eliminating accumulated memory bloat. Service isolation mechanisms quarantine suspected processes in resource-constrained containers, preventing leak propagation to other control-plane components. The actuation layer interfaces with switch control APIs, including P4Runtime for programmable data-plane configuration and gRPC-based management interfaces for control-plane service manipulation. Each action execution includes comprehensive pre-action validation to ensure system stability, runtime monitoring to detect adverse effects, and post-action verification to confirm successful mitigation. The layer implements safety mechanisms that prevent cascading failures, including action rate limiting to avoid excessive intervention frequency, rollback capabilities for actions that produce unexpected results, and escalation procedures when initial mitigation attempts prove insufficient. Post-

10.48047/jocaaa.2025.34.11.46

action telemetry collection provides feedback to the machine learning engine, enabling continuous model refinement based on observed intervention outcomes and facilitating adaptive learning that improves decision quality over time.

3.5 Integration Models and Deployment Options

The Self-Healing Memory System supports flexible integration models accommodating diverse network fabric architectures and operational requirements. The inference module can reside locally on switch co-processors, including Field-Programmable Gate Array (FPGA) or Application-Specific Integrated Circuit (ASIC) host processors, enabling autonomous operation with minimal latency between detection and action. Local deployment provides advantages in response time and resilience to management network disruptions, as each switch operates independently without requiring connectivity to centralized controllers. Alternatively, the prediction and decision components can reside centrally within fabric controllers managing multiple switches simultaneously, enabling coordinated learning across device populations and centralized policy enforcement. Centralized deployment facilitates resource pooling, where computational overhead for machine learning inference concentrates on dedicated infrastructure rather than distributing across individual switches. Hybrid deployment models combine local anomaly detection with centralized prediction refinement, where switches perform initial screening using lightweight models and escalate suspected anomalies to centralized engines for detailed analysis. Integration with existing network management systems occurs through standard interfaces, including YANG data models for configuration management, OpenConfig schemas for telemetry export, and gRPC services for action coordination [5]. The framework accommodates integration with digital twin network representations, where real-time telemetry feeds virtual network models that enable sophisticated what-if analysis and policy simulation before deploying actions in production environments. Deployment flexibility extends to training pipelines, supporting offline training on historical data archives, online training with live telemetry streams, and federated learning approaches where multiple switches collaboratively train shared models while preserving local data privacy. This architectural flexibility ensures adaptation to diverse deployment scenarios, from resource-constrained edge environments to large-scale AI training fabrics supporting thousands of interconnected switches [1][2].

4. Methodology and Implementation

4.1 Experimental Testbed Configuration

The experimental testbed comprises a high-performance AI cluster consisting of multiple network switches configured to emulate realistic AI and ML training fabric conditions. Each switch operates at high bandwidth capacity with RoCEv2 traffic profiles that simulate large-scale distributed training workloads. The testbed architecture includes multiple compute nodes equipped with GPUs interconnected through the switch fabric, generating authentic memory pressure on control-plane processes. Switch configurations include standard control-plane daemons such as BGP routing services, SDN controllers, telemetry collectors, and orchestration services that represent typical production deployments. The testbed environment enables controlled injection of memory leak scenarios through instrumented control-plane processes that exhibit various leak patterns, including gradual linear leaks, exponential accumulation patterns, and sporadic burst leaks. Traffic generation tools produce realistic flow distributions with varying packet sizes, flow durations, and arrival patterns that stress control-plane flow table management and state tracking mechanisms. Network topology configurations include leaf-spine architectures common in AI training clusters, enabling evaluation under representative routing complexity and flow distribution patterns [1][2]. The experimental setup incorporates comprehensive monitoring infrastructure

10.48047/jocaaa.2025.34.11.46

that captures ground truth memory consumption data, process behavior metrics, and system performance indicators across all testbed components. Automated orchestration frameworks enable repeatable experiment execution across multiple trial runs, varying workload characteristics, leak injection parameters, and system configuration settings. The testbed supports both isolated single-switch experiments for detailed behavioral analysis and multi-switch scenarios for evaluating coordinated detection and mitigation across fabric-wide deployments [7].

4.2 Data Collection and Feature Engineering

Data collection procedures capture comprehensive telemetry spanning memory metrics, system performance indicators, and workload characteristics across extended observation periods. Raw telemetry streams include heap memory utilization measurements, virtual memory statistics, page fault counts, memory allocation call frequencies, deallocation event rates, fragmentation metrics, and correlations with central processing unit (CPU) and input-output (I/O) utilization. Telemetry collection employs multiple sampling intervals to capture both high-frequency transient behaviors and long-term trend evolution, with adaptive sampling that increases collection rates when anomaly indicators emerge. Feature engineering transforms raw telemetry into meaningful representations suitable for machine learning model consumption. Temporal features include sliding window statistics such as moving averages, standard deviations, rate-of-change calculations, and trend slope measurements computed over multiple time horizons. Statistical features capture distribution characteristics, including percentiles, skewness, kurtosis, and entropy measures that quantify memory consumption variability. Correlation features examine relationships between memory metrics and system load indicators, identifying coupled behaviors that indicate leak propagation patterns. Frequency domain features derived through Fourier transforms reveal periodic patterns and cyclical behaviors in memory allocation dynamics. Feature selection procedures identify the most discriminative attributes through correlation analysis, mutual information scoring, and recursive feature elimination techniques that balance predictive power against computational complexity. The engineered feature set undergoes normalization and scaling to ensure numerical stability during model training, with separate normalization parameters maintained for different deployment environments to accommodate varying operational ranges. Labeled datasets combine normal operational periods with synthetically injected and naturally occurring leak events, providing supervised learning targets that span diverse leak characteristics and progression patterns [5][6].

4.3 Model Training and Optimization

Model training procedures optimize both the LSTM forecasting network and the RL decision agent through iterative refinement processes. The LSTM network architecture undergoes hyperparameter optimization to determine optimal network depth, hidden layer dimensions, dropout rates, and sequence length parameters. Training employs sliding window approaches where historical telemetry sequences serve as inputs and future memory states constitute prediction targets. Loss functions incorporate mean squared error for regression tasks predicting continuous memory values and binary cross-entropy for classification tasks identifying leak presence. Training procedures partition available data into training, validation, and test sets, employing cross-validation techniques to assess generalization performance and prevent overfitting. Regularization techniques, including dropout, weight decay, and early stopping based on validation performance, ensure robust model behavior across diverse operational conditions. The LSTM model undergoes continuous refinement through online learning procedures that incorporate newly observed telemetry data, enabling adaptation to evolving workload characteristics and emerging leak patterns [6]. The RL agent training employs policy gradient methods or actor-critic architectures that learn optimal action selection policies through interaction with simulated and real switch environments.

10.48047/jocaaa.2025.34.11.46

Training begins in simulation environments that model switch behavior, memory dynamics, and recovery action effects, enabling safe exploration of action spaces without risking production system stability. Simulated training generates extensive experience through accelerated environment interactions, building initial policy knowledge before deployment in real systems. Transfer learning techniques adapt simulation-trained policies to real hardware through fine-tuning procedures that account for sim-to-real gaps in environment dynamics. Exploration strategies balance between exploiting known effective actions and exploring alternative mitigation approaches, employing epsilon-greedy strategies, Boltzmann exploration, or entropy-regularized policies. Experience replay mechanisms store historical state-action-reward trajectories, enabling efficient learning through repeated sampling of diverse experiences. Model optimization procedures evaluate performance across multiple metrics, including prediction accuracy, false positive rates, intervention timeliness, and recovery success rates, employing multi-objective optimization when trade-offs exist between competing objectives [7][8].

4.4 Reward Function Design and RL Decision Logic

Reward function design constitutes a critical component that shapes reinforcement learning agent behavior by encoding desired system outcomes into scalar feedback signals. The reward function incorporates multiple terms that balance competing objectives, including system stability preservation, performance impact minimization, and resource consumption efficiency. Stability reward components provide positive feedback for successful leak mitigation actions that restore memory health and prevent system crashes, with reward magnitudes proportional to failure severity avoided. Performance penalty terms deduct rewards for actions that introduce latency increases, throughput reductions, or service disruptions during recovery execution, encouraging the agent to select minimally invasive interventions. Data preservation terms penalize actions that result in state loss, connection resets, or flow table clearing, promoting mitigation strategies that maintain operational continuity. Temporal discounting factors weight immediate rewards more heavily than distant future rewards, encouraging timely intervention while accounting for delayed consequences of action selections. The reward function undergoes iterative refinement through domain expert input and empirical evaluation, tuning weighting coefficients to achieve desired behavioral characteristics [8]. Reinforcement learning decision logic implements a structured policy that maps observed system states to optimal action selections. State representations encode current memory metrics, predicted future trajectories from Long Short-Term Memory forecasts, workload characteristics, and historical intervention outcomes. Action spaces encompass discrete mitigation options, including memory reclamation intensity levels, process throttling degrees, restart scopes, and no-operation decisions when intervention appears premature. The policy network evaluates expected cumulative rewards for each action under current state conditions, selecting actions that maximize long-term system health while satisfying operational constraints. Decision logic incorporates safety constraints that prevent harmful action sequences, implementing action masking that excludes unsafe options based on current system conditions. Multi-stage decision procedures decompose complex recovery operations into sequential sub-actions, enabling fine-grained control over intervention progression. The agent maintains confidence estimates for action selections, enabling escalation procedures when uncertainty exceeds acceptable thresholds. Meta-learning capabilities enable the agent to recognize novel scenarios and adapt decision strategies based on observed environment dynamics, facilitating robust operation across diverse deployment conditions and workload patterns [1][2][7][8].

5. Results and Discussion

5.1 Leak Prediction Performance and Accuracy Metrics

10.48047/jocaaa.2025.34.11.46

The Self-Healing Memory System demonstrates substantial predictive capabilities in identifying memory leak conditions before critical failure thresholds are reached. Prediction accuracy metrics evaluate the framework's ability to correctly classify leak presence, distinguish between normal memory fluctuations and genuine anomalies, and forecast time-to-failure with adequate lead time for intervention. True positive rates measure the proportion of actual leak events successfully identified by the LSTM forecasting model, while false positive rates quantify instances where the system incorrectly flags normal operational behavior as anomalous. Precision and recall metrics provide complementary perspectives on prediction quality, with precision indicating the reliability of leak alerts and recall measuring the completeness of leak detection coverage. Receiver Operating Characteristic (ROC) curves and Area Under the Curve (AUC) metrics assess classifier performance across varying decision thresholds, enabling optimization of the sensitivity-specificity trade-off based on operational priorities. Temporal prediction accuracy evaluates how precisely the system forecasts the timing of impending failures, measured through mean absolute error and root mean squared error between predicted and actual failure times. Lead time distributions characterize the warning period provided by the framework, with longer lead times enabling less disruptive mitigation strategies. Prediction confidence calibration examines whether the model's probabilistic forecasts accurately reflect the true likelihood of leak occurrence, ensuring that high-confidence predictions reliably indicate genuine anomalies. The framework's predictive performance generalizes across diverse leak patterns, including gradual linear accumulation, exponential growth scenarios, and intermittent burst leaks, demonstrating robustness to varying failure modes. Cross-validation across different workload types and switch configurations confirms that learned patterns transfer effectively to unseen operational conditions [6][9].

5.2 Autonomous Recovery Success Rates

Autonomous recovery capabilities constitute a primary contribution of the Self-Healing Memory System, enabling switches to restore memory health without manual intervention. Recovery success rates measure the proportion of detected leak events that the RL agent successfully mitigates through autonomous action selection. Complete recovery instances occur when selected mitigation actions fully restore memory consumption to normal levels and prevent subsequent failure, while partial recovery cases achieve temporary stabilization requiring additional interventions. Recovery failure scenarios arise when selected actions prove ineffective or when leak severity exceeds mitigation capabilities, necessitating fallback to manual intervention or switch restart. Time-to-recovery metrics quantify the duration between leak detection and successful memory restoration, with shorter recovery times minimizing exposure to degraded operational states. Action selection analysis examines which mitigation strategies the agent employs most frequently and which actions prove most effective for different leak characteristics. Memory reclamation through garbage collection emerges as a preferred first-line intervention for moderate leaks, while process isolation and selective restarts become necessary for severe accumulation scenarios. The agent learns to escalate intervention intensity progressively, attempting minimally disruptive actions initially and escalating to more aggressive strategies only when necessary. Recovery stability analysis tracks whether mitigated leaks recur after intervention, with stable recoveries indicating successful root cause elimination versus temporary symptom suppression. The framework demonstrates learning progression over deployment duration, with recovery success rates improving as the agent accumulates experience and refines action selection policies based on observed outcomes. Comparative evaluation across different RL architectures and reward function designs reveals optimal configurations that balance recovery effectiveness against operational impact [8][10].

5.3 Performance Overhead Analysis

10.48047/jocaaa.2025.34.11.46

Performance overhead analysis quantifies the computational and operational costs imposed by the Self-Healing Memory System on switch forwarding performance and control-plane responsiveness. Telemetry collection overhead measures the CPU utilization, memory consumption, and network bandwidth consumed by continuous monitoring activities. Adaptive sampling strategies demonstrate effectiveness in minimizing collection overhead during stable operational periods while ensuring adequate temporal resolution during anomaly progression. Machine learning inference latency characterizes the time required for LSTM prediction and RL decision-making, with measurements confirming that prediction operations complete within acceptable bounds that do not delay critical control-plane functions. Packet forwarding latency analysis examines whether monitoring and prediction activities introduce measurable delays into data-plane packet processing, with results confirming isolation between control-plane learning operations and data-plane forwarding pipelines. Throughput impact assessment evaluates whether the framework reduces aggregate switch throughput during normal operation or mitigation execution, measuring packet rates, bandwidth utilization, and flow completion times under various load conditions. Control-plane responsiveness metrics examine whether the system affects route convergence times, flow table update latencies, or telemetry export frequencies during monitoring and recovery operations. Mitigation action overhead characterizes the performance impact of different recovery strategies, with memory reclamation operations showing minimal disruption while process restarts introduce brief service interruptions. The framework implements careful timing of disruptive actions during low-traffic periods when possible, reducing impact on active workloads. Storage requirements for model parameters, historical telemetry, and experience replay buffers remain modest relative to available switch memory resources, confirming practical deployability on resource-constrained hardware. Energy consumption analysis reveals that proactive leak prevention reduces overall power usage by avoiding costly switch reboots and control-plane restart sequences [1][2][9].

Performance Dimension	Measurement Metric	Acceptable Threshold	SHMS Impact Level	Traditional System Impact	Measurement Method
Packet Forwarding Latency	Microseconds per packet	Minimal increase	Negligible	Not applicable	Data-plane instrumentation
Control-plane Responsiveness	Route convergence time	No degradation	Minimal	Variable during failures	Protocol timing analysis
CPU Utilization	Percentage overhead	Minimal increase	Low	High during reactive responses	System monitoring
Memory Footprint	Model and buffer storage	Within available resources	Low	Minimal	Resource accounting
Throughput	Packets per second	No reduction	None	Degraded during failures	Traffic measurement

10.48047/jocaaa.2025.34.11.46

Telemetry Bandwidth	Management traffic volume	Modest increase	Low	Variable	Network monitoring
Mitigation Execution Time	Recovery duration	Minimal service impact	Short	Extended (full restarts)	Action timing

Table 4: Performance Metrics and Evaluation Criteria [1][2][9][10]

5.4 Comparative Analysis with Traditional Systems

Comparative evaluation positions the Self-Healing Memory System against conventional reactive monitoring approaches and static threshold-based watchdog mechanisms. Traditional systems employ fixed memory utilization thresholds that trigger alerts or automatic restarts when consumption exceeds predefined levels, lacking the ability to predict failures before threshold violations occur. The comparative analysis demonstrates that reactive approaches detect problems significantly later than predictive frameworks, often after degradation has already impacted service quality. Static thresholds prove particularly ineffective for slow, gradual leaks that accumulate over extended periods, as consumption remains below alert thresholds until sudden exhaustion occurs. False positive rates differ substantially between approaches, with static thresholds generating excessive alerts during legitimate traffic spikes while predictive models distinguish between temporary increases and persistent leak patterns. Recovery strategies employed by traditional systems consist primarily of complete process restarts or switch reboots, imposing significant service disruption compared to targeted mitigation actions selected by reinforcement learning agents. Manual intervention requirements highlight a critical distinction, as traditional approaches necessitate operator diagnosis and action selection while the Self-Healing Memory System operates autonomously. Adaptation capabilities reveal fundamental architectural differences, with static systems requiring manual threshold tuning for different deployment environments while learning-based frameworks automatically adapt to local operational characteristics. The comparative evaluation examines performance across diverse scenarios, including varying leak severities, different workload types, and multiple switch configurations, demonstrating consistent advantages for predictive self-healing approaches. Cost-benefit analysis considers implementation complexity, operational overhead, and reliability improvements, confirming that machine learning-based frameworks provide substantial value despite increased architectural sophistication. The analysis acknowledges scenarios where traditional approaches remain appropriate, including resource-constrained environments unable to support machine learning inference or deployments with extremely stable, well-characterized workloads [3][4][10].

5.5 Real-World Applications and Use Cases

Real-world deployment scenarios demonstrate the practical applicability and operational benefits of the Self-Healing Memory System across diverse network infrastructure contexts. AI and ML training clusters represent a primary application domain, where distributed training workloads generate intensive network traffic and control-plane state management that stresses switch memory resources [1]. Large-scale training jobs running for extended durations create ideal conditions for gradual leak accumulation, making predictive intervention particularly valuable for preventing job disruptions. The framework enables training clusters to achieve higher effective utilization by reducing unplanned downtime and avoiding costly checkpoint recovery operations when switches fail. Data center operations benefit from automated reliability management that reduces operational expenditure associated with manual troubleshooting and emergency maintenance activities [2]. The system provides visibility into memory

10.48047/jocaaa.2025.34.11.46

health trends that inform capacity planning and hardware refresh cycles, identifying switches exhibiting chronic leak behaviors that require replacement or software updates. Edge computing environments present distinct deployment challenges with resource constraints and limited management connectivity, making autonomous self-healing capabilities particularly valuable when centralized operator intervention proves impractical. Telecommunication network applications leverage the framework for ensuring reliability in carrier-grade deployments where SLAs demand high availability guarantees. The system adapts to diverse switch vendor implementations and operating system variants, demonstrating portability across heterogeneous network infrastructure. Cloud service providers deploying the framework report reduced incident response times and improved customer satisfaction through proactive reliability management. The system's learning capabilities enable continuous improvement over operational lifetime, with deployed instances becoming increasingly effective as they accumulate experience with local failure patterns and workload characteristics. Integration with broader network management platforms enables correlation between memory health and other infrastructure metrics, providing holistic visibility into fabric reliability [5][7][9][10].

Conclusion

The Self-Healing Memory System represents a significant advancement in autonomous reliability management for high-performance network switches supporting artificial intelligence and machine learning workloads. The framework successfully addresses critical gaps in control-plane memory health monitoring by combining Long Short-Term Memory-based temporal forecasting with reinforcement learning-driven autonomous mitigation. Predictive capabilities enable early detection of memory leak progression before critical failures occur, providing adequate lead time for intervention while minimizing operational disruptions. The closed-loop architecture incorporating continuous telemetry collection, intelligent prediction, and automated actuation demonstrates practical feasibility within the stringent performance constraints of production network environments. Autonomous recovery mechanisms reduce dependence on manual intervention, enabling scalable reliability management across large-scale artificial intelligence training fabrics where human oversight becomes infeasible. The framework's learning capabilities facilitate continuous adaptation to evolving workload characteristics and emerging failure patterns, ensuring sustained effectiveness over extended deployment periods. Experimental validation confirms substantial improvements over traditional reactive monitoring approaches across multiple performance dimensions, including detection accuracy, recovery success rates, and operational overhead. Future developments may expand self-healing capabilities beyond memory management to encompass thermal regulation, logic fault correction, and coordinated fabric-wide reliability optimization. The integration of federated learning approaches could enable collaborative model training across distributed switch populations, accelerating adaptation to novel failure modes while preserving operational data privacy.

References

- [1] Ellie Lipe, "Energy Efficient Scheduling of AI/ML Workloads on Multi-Instance Gpus with Dynamic Repartitioning," in 2025 IEEE 25th International Symposium on Cluster, Cloud and Internet Computing (CCGrid), 30 June 2025. <https://ieeexplore.ieee.org/document/11044810>
- [2] Jit Gupta, et al., "An End-to-End Assurance Framework for AI/ML Workloads in Datacenters," in IEEE Conference on Network and Service Management (CNSM), 3 July 2025. <https://arxiv.org/abs/2507.03158v1>
- [3] Killian Murphy, et al., "Fault Prediction for Heterogeneous Telecommunication Networks Using Machine Learning: A Survey," in IEEE Transactions on Network and Service Management, 07 December 2023. <https://ieeexplore.ieee.org/document/10347460>
- [4] Mowei Wang, et al., "Learning Buffer Management Policies for Shared Memory Switches," in IEEE INFOCOM 2022 - IEEE Conference on Computer Communications, 20 June 2022. <https://ieeexplore.ieee.org/document/9796784>
- [5] Zhihao Wang, et al., "Network-Wide Data Collection Based on In-Band Network Telemetry for Digital Twin Networks," in IEEE Transactions on Network and Service Management, 20 June 2024. <https://ieeexplore.ieee.org/document/10669844>
- [6] Temitayo Emmanuel Fabunmi, et al., "An Adaptable Intelligent LSTM-Based Technique for Fault Detection in Electrical Networks," in 2024 IEEE 5th International Conference on Electro-Computing Technologies for Humanity (NIGERCON), 24 March 2025. <https://ieeexplore.ieee.org/document/10926946>
- [7] Hanyu Gao, et al., "Experimental Demonstration of Automated ML Service Provisioning for VNT Configuration in SDM Networks," in 2024 Optical Fiber Communications Conference and Exhibition (OFC), 16 May 2024. <https://ieeexplore.ieee.org/document/10527201>
- [8] Xinna Ma, et al., "Deep Reinforcement Learning Bearing Fault Diagnosis Method Based on Improved Reward Function," in 2023 7th International Symposium on Computer Science and Intelligent Control (ISCSIC), 26 January 2024. <https://ieeexplore.ieee.org/document/10409764>
- [9] Minjung Cho and Eui-Young Chung, "MLCRP: ML-Based GPU Cache Performance Modeling Featured With Reuse Profiles," in IEEE Transactions on Computers, 16 July 2025. <https://ieeexplore.ieee.org/document/11082128>
- [10] Saef Obad Husain, et al., "BISNet: Bio-Inspired Self-Healing Network for Autonomous Fault Recovery," in 2025 International Conference on Metaverse and Current Trends in Computing (ICMCTC), 17 October 2025. <https://ieeexplore.ieee.org/abstract/document/11196711>