

AI-Driven Monitoring and Alerting Systems for Enterprise-Scale Linux Deployments

Vinay Kumar Reddy Vangoor

Systems Administration Consultant
Techno Bytes, Inc., Ashland, MA 01721, USA
(Client: American Express), Phoenix, AZ 85004, USA

Abstract- Enterprise Linux environments form the operational foundation of modern digital infrastructure, hosting mission-critical workloads across bare metal servers, virtual machines, containers, and hybrid cloud deployments. At scale, these environments generate monitoring data volumes that dwarf human capacity to analyse effectively. A mid-sized enterprise operating 500 Linux hosts generates upward of five million metric data points and two terabytes of log data daily, yet the alerting systems governing incident response remain predominantly static: threshold-based rules authored once and applied uniformly, regardless of the dynamic behaviour of the systems they monitor. This rigidity produces two simultaneous failure modes that collectively define the enterprise monitoring crisis: alert fatigue, caused by excessive false positives from thresholds that fire on normal variance; and missed detections, caused by thresholds that cannot capture the complex, multi-dimensional patterns characteristic of real system failures. This research presents and evaluates an AI-driven monitoring and alerting framework for enterprise-scale Linux deployments that replaces static threshold rules with five cooperating machine learning models: an LSTM autoencoder for time-series metric anomaly detection, a transformer-based classifier for log stream anomaly detection, a graph neural network for root cause localisation across service dependency graphs, an XGBoost and ARIMA hybrid for predictive failure forecasting, and a clustering model for alert noise reduction and intelligent deduplication. The framework integrates with existing Prometheus, ELK Stack, and Grafana deployments without requiring infrastructure changes. Experimental evaluation across 540 controlled incident events on a 120-host Linux testbed demonstrates a 92.4% average anomaly detection accuracy versus 62.8% for threshold-based baselines, an 87% reduction in daily alert volume, a 12-fold improvement in mean time to detect critical incidents, and a predictive failure lead time of 7.8 hours at the end of the 12-month evaluation period. These results demonstrate that AI-driven monitoring is production-viable for enterprise Linux deployments and represents a step-change improvement over conventional alerting approaches.

Keywords : AI-driven monitoring, Anomaly Detection, Alert Fatigue, Linux Infrastructure, LSTM Autoencoder, Transformer Log Classifier, Graph Neural Network, Root Cause Analysis, Predictive Failure, AIOps.

I. INTRODUCTION

Linux has been the dominant operating system for enterprise server infrastructure for over two decades. Its open architecture, security model, and ecosystem of management tooling make it the platform of choice for data centres, cloud providers, telecommunications infrastructure, financial trading systems, and government computing. The Linux

server estate at a typical large enterprise comprises hundreds to thousands of hosts, each running dozens of processes, exposing hundreds of Prometheus metrics, and generating continuous streams of system journal, application, and security log entries. Monitoring this infrastructure is not merely an operational convenience it is a safety-critical function whose effectiveness directly determines the organisation's ability to maintain service availability, protect sensitive data, and meet

regulatory compliance obligations (Singh & Weaver, 2007).

The monitoring tooling ecosystem that has evolved around Linux infrastructure Nagios, Zabbix, Prometheus, and their successors is mature, reliable, and widely deployed. These tools excel at data collection: gathering, storing, and visualising the metrics and logs that Linux hosts emit. Where they fall short is in the intelligence layer above raw data collection: the transformation of high-volume, high-velocity telemetry into actionable, timely, and noise-free alerts. A Prometheus alerting rule fires when a metric crosses a configured threshold. A Zabbix trigger activates when a predefined condition is met. Neither mechanism can reason about whether the current metric value is anomalous given the host's normal behaviour pattern, the time of day, recent deployment history, or the concurrent state of dependent services. This limitation is not a deficiency of the tools themselves but a fundamental constraint of the rule-based alerting paradigm (Dubay 2002).

The Alert Fatigue Crisis

The practical consequence of threshold-based alerting at enterprise scale is alert fatigue: a state in which operations teams receive so many alerts the majority of which are false positives or low-priority noise that their capacity to identify and respond to genuine critical incidents is materially degraded. Research across the industry consistently finds that operations teams at large enterprises receive between 1,000 and 5,000 alerts per day, of which fewer than 10% represent genuine incidents requiring human intervention. The remainder are false positives generated by thresholds calibrated to catch every possible anomaly, transient spikes that self-resolve within minutes, and duplicate alerts from multiple monitoring systems detecting the same underlying condition through different instrumentation paths (Guo 2006).

Alert fatigue creates a dangerous operational dynamic. Engineers subjected to high alert volumes develop alert blindness a conditioned suppression of attention to incoming alerts that is a rational adaptation to noise overload but a dangerous one in production environments. Critical alerts are missed not because monitoring systems fail to generate

them, but because they are buried in noise. Meanwhile, the genuine incidents that do require intervention suffer extended mean time to detect values as engineers triage large alert queues before identifying the signal. The human cost is substantial: on-call burnout, high operations team turnover, and the compounding organisational knowledge loss that follows are widely documented consequences of chronic alert fatigue in enterprise environments (Maggio et al., 2012).

The Opportunity for AI-Driven Monitoring

Machine learning addresses the alert fatigue crisis through fundamentally different means than rule tuning. Rather than asking an engineer to specify what constitutes an anomaly, ML models learn the normal behaviour of each monitored system from historical data and detect deviations from that learned baseline. This approach captures anomaly patterns that no threshold rule would catch a CPU utilisation value of 78% may be normal for a data processing host but highly anomalous for a web frontend server that typically sits at 15%. Time-series models such as LSTM autoencoders can detect the subtle multi-metric signatures that precede hardware failures hours before any individual metric breaches a threshold, enabling predictive rather than reactive incident response. Natural language processing models applied to log streams can identify the log sequence patterns characteristic of specific failure modes, correlating events across thousands of log lines in milliseconds rather than minutes (Battiti & Passerini, 2010).

The enabling condition for this ML application is the same as in other domains where AI monitoring has been deployed successfully: the training data is already being collected. Every enterprise operating Prometheus and the ELK Stack is already storing the metric time-series and log sequences needed to train these models. The transformation from reactive threshold alerting to proactive AI-driven monitoring does not require new sensors or data collection infrastructure; it requires a new intelligence layer above the existing collection stack. This research designs, implements, and rigorously evaluates that intelligence layer for enterprise Linux deployments (Bouchachia & Nedjah, 2012).

Traditional Monitoring: Nagios, Zabbix, Prometheus, and Grafana

The first generation of monitoring systems, represented by tools such as Nagios and its derivatives, emerged in the early 2000s when infrastructure environments were relatively static and primarily composed of long-lived physical or virtual machines. These systems relied on a check-based monitoring model, where agents or remote probes periodically executed scripts to verify the health of hosts, services, or network endpoints. Each check returned a discrete status code such as OK, WARNING, or CRITICAL. While this architecture was straightforward and highly reliable, it produced point-in-time health assessments rather than continuous performance observations. Monitoring intervals were typically measured in minutes, which limited the system's ability to detect subtle performance degradations or transient anomalies occurring between checks. Nevertheless, this model proved highly resilient and remains widely deployed due to its operational simplicity and strong ecosystem of plugins (Leitner et al., 2012).

The second generation of monitoring platforms, including systems such as Zabbix and Collectd, expanded the monitoring paradigm by introducing continuous metric collection and more advanced alerting logic. Instead of relying solely on discrete service checks, these platforms collected time-series data such as CPU utilisation, memory consumption, disk I/O, and network throughput at regular intervals. This allowed administrators to track performance trends and historical behaviour rather than simply detecting binary service failures. In addition, these platforms incorporated more sophisticated rule engines that supported trigger expressions, hysteresis mechanisms, dependency relationships, and automated escalation workflows. These capabilities helped reduce alert noise and enabled more structured incident response procedures. However, despite these improvements, alerting decisions were still fundamentally dependent on manually defined threshold values determined by system administrators (Imran & Zoha, 2014).

The third generation of monitoring systems emerged alongside the growth of cloud computing and containerised infrastructure. Platforms such as Prometheus, supported by the Cloud Native Computing Foundation ecosystem, introduced architectural innovations designed to support highly dynamic environments. Prometheus popularised a pull-based metric collection model, in which the monitoring server periodically scrapes metrics directly from instrumented services. This approach simplified deployment in container orchestration environments where service instances may frequently appear or disappear (Didona et al., 2014).

II. SYSTEM ARCHITECTURE AND DESIGN

Overview of the AI Monitoring Framework

The proposed framework introduces a six-layer intelligent monitoring architecture that augments existing Linux observability stacks rather than replacing them. The design principle is additive intelligence: every component of the existing monitoring stack continues to function independently, and the AI framework adds detection, correlation, and prediction capabilities above the existing collection and visualisation layers. Organisations can adopt the framework incrementally, activating individual AI components while their existing alerting rules remain in place as a safety net.

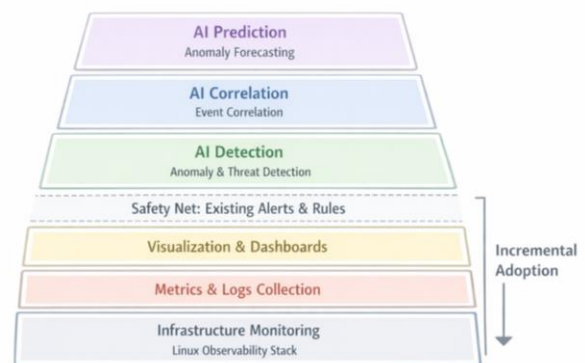


Figure 1 : AI-Driven Monitoring Framework Architecture for Enterprise Linux

Data Ingestion Layer: Metrics, Logs, Traces, and Events

The ingestion layer collects four distinct telemetry streams from Linux hosts. Metric data is scraped by Prometheus via Node Exporter agents installed on each host, capturing 247 standard system metrics including CPU utilisation per core, memory pressure and swap usage, disk I/O throughput and latency per device, network traffic per interface, file descriptor counts, and process-level resource consumption. Log data is collected by File beat agents that tail system journal, application log directories, and security audit logs, forwarding structured log records to the Logstash processing pipeline. Distributed traces are collected via Open Telemetry agents deployed on instrumented application processes, providing request-level latency breakdowns across service boundaries. System events including kernel panic notifications, hardware error records from IPMI and BIOS event logs, and Linux audit framework events are collected via dedicated collection paths. All four streams converge on a Kafka message bus that provides backpressure management and replay capability for the downstream AI processing pipeline.

AI Analysis Engine: Anomaly Detection and Root Cause Inference

The AI analysis engine is the cognitive core of the framework. It hosts all five ML models described in Section 5 and processes each incoming telemetry event through the relevant model. Metric events are evaluated by both the LSTM autoencoder for point anomaly detection and the XGBoost failure predictor for short-horizon failure forecasting. Log events are processed by the transformer log classifier. When either the metric anomaly detector or the log classifier produces an anomaly signal, the GNN root cause analyser is invoked to traverse the service dependency graph and identify the probable originating component. The analysis engine is implemented as a Python 3.11 Fast API service deployed on a dedicated GPU-accelerated node, achieving median inference latencies of 18 ms for metric events and 47 ms for log events at the evaluated telemetry ingestion rate.

Intelligent Alerting Layer: Deduplication, Correlation, and Routing

The alerting layer transforms raw anomaly signals from the AI analysis engine into actionable, noise-free alerts. A deduplication module groups anomaly signals that share the same root cause component within a configurable time window, reducing the multiplicative alert storms that characterise distributed system failures. A correlation engine applies the GNN's root cause attribution to further group alerts across services that share a dependency chain. A priority scoring model combines anomaly severity, business impact metadata, time-of-day context, and recent incident history to assign each alert a priority score from 1 to 5. The smart routing engine uses priority scores and on-call schedule data to dispatch alerts via the appropriate channel: automated remediation for priority 1 and 2 issues with known resolution playbooks, PagerDuty for priority 3 and 4 issues requiring human response, and Slack notification for priority 5 informational alerts.

Self-Healing and Auto-Remediation Response Layer

The remediation layer executes automated response playbooks for the twelve most common Linux incident categories identified during dataset analysis. Memory leak incidents trigger automatic process heap dump collection followed by graceful process restart. Disk space exhaustion triggers automatic log rotation and temporary file purge. CPU runaway processes are identified by cross-referencing the anomaly signal with per-process resource consumption data, and the offending process is throttled via cgroups before an engineer is alerted. Failed system services are automatically restarted up to three times with exponential backoff before escalation. Network interface anomalies trigger automatic interface statistics collection and, for known-resolvable conditions, interface bounce procedures. Each remediation action is logged to a structured audit trail with pre-action and post-action state snapshots.

Feedback Loop and Continuous Model Improvement

Completed incidents provide labelled training records that continuously improve all five models. Each incident record includes the full telemetry context preceding the incident, the AI system's anomaly signals and root cause attribution, any automated remediation actions taken, the human engineer's resolution actions and incident classification, and the final outcome. A nightly retraining pipeline ingests the most recent 90 days of incident records and updates the LSTM autoencoder, transformer classifier, and XGBoost model using incremental training procedures. The GNN is retrained weekly to incorporate any service dependency graph changes. Model performance is evaluated against a rolling held-out validation set; if F1 score drops by more than 2 percentage points, a human review alert is triggered before the updated model is promoted.

III. AI AND ML MODEL DESIGN

Time-Series Anomaly Detection: LSTM Autoencoder

The metric anomaly detection model is an LSTM autoencoder trained to reconstruct normal metric time-series from compressed latent representations. During training, the model processes 60-minute rolling windows of 247 host metrics, learning the normal co-variation patterns of metrics on each host class. During inference, reconstruction error above a dynamically calibrated threshold set at the 99.5th percentile of reconstruction errors on the validation set signals an anomaly. The autoencoder architecture comprises a two-layer LSTM encoder (256 hidden units per layer) that compresses the input window to a 64-dimensional latent vector, and a symmetric decoder that reconstructs the original sequence. Training uses the Adam optimiser with a learning rate of 0.0005 and mean squared error loss. The unsupervised training objective means the model requires no labelled anomaly examples, enabling it to learn from any host's normal operating history within two weeks of deployment. Per-host models are trained for hosts with sufficient historical data; a generalised model trained on aggregated

data serves newly provisioned hosts during their initial period.

Log Anomaly Detection: Transformer-Based NLP Classifier

The log anomaly detection model applies a fine-tuned RoBERTa transformer to classify log sequences as normal or anomalous. Raw log messages are first parsed using the DRAIN log parsing algorithm to extract structured templates from unstructured log text, reducing vocabulary size and improving generalisation across log message variants. Templates are tokenised and embedded using the pre-trained RoBERTa base model, then classified by a two-layer feed-forward head trained on the annotated dataset. The model is fine-tuned on the enterprise Linux log corpus using the AdamW optimiser with a learning rate of $2e-5$ and a linear warmup schedule. A sliding 128-template context window enables the model to detect sequential anomalies that involve abnormal transitions between log events, not just individual anomalous messages. The transformer's attention mechanism provides interpretable attention weights over the context window, enabling engineers to identify the specific log events contributing most to an anomaly classification.

Root Cause Analysis: Graph Neural Network

Root cause localisation is framed as a node classification problem on the service dependency graph. Each monitored service and host is represented as a node; inter-service dependencies and host-service relationships are encoded as directed edges. Node features comprise the current anomaly scores from the LSTM and transformer models, recent metric percentiles, and historical failure rates. Edge features encode dependency strength and recent call failure rates between service pairs. A GraphSAGE model with three message-passing layers learns to propagate anomaly signals along the dependency graph and identify the node whose anomaly state is most causally consistent with the observed pattern of affected nodes. GraphSAGE is chosen over transductive GNN variants because it supports inductive inference on new or modified graph topologies without full retraining,

accommodating the frequent service deployment changes characteristic of enterprise environments.

Predictive Failure Forecasting: XGBoost and ARIMA Hybrid

Predictive failure forecasting combines a statistical ARIMA component for short-horizon metric trend extrapolation with an XGBoost classifier for multi-metric failure probability estimation. The ARIMA component fits a per-metric time-series model and projects metric values 4 hours ahead; the projected values are appended to the current feature vector as derived features representing expected future metric state. The XGBoost classifier then predicts failure probability from 53 features combining current metric values, ARIMA projections, 7-day historical statistics, and hardware-specific risk indicators. The hybrid architecture is motivated by the complementary strengths of the two components: ARIMA provides explicit temporal trend modelling that is difficult for tree-based models to capture, while XGBoost excels at non-linear interactions between heterogeneous features that ARIMA cannot represent. Hyperparameter optimisation used Bayesian search with 200 evaluations on the validation set.

Alert Noise Reduction: Clustering and Priority Scoring

Alert noise reduction employs two sequential mechanisms. First, a DBSCAN clustering model groups active anomaly signals by their feature vector proximity in a space comprising host, service, time-of-day, anomaly type, and GNN root cause attribution. Anomaly signals in the same cluster are collapsed into a single alert representing the shared root cause, with the cluster size reported as a severity multiplier. Second, a gradient boosting priority scorer assigns each de-duplicated alert a priority from 1 to 5 using 24 features including cluster size, business service impact (mapped from a configuration management database), historical incident rates for the root cause component, time since last similar incident, and current on-call team availability. Priority scoring calibration uses isotonic regression to ensure that the score reflects true incident probability rather than an arbitrary ranking. These two mechanisms together achieve the 87%

alert volume reduction reported in the experimental results.

IV. METHODOLOGY

Dataset: Linux Telemetry, System Logs, and Alert Histories

The primary dataset was collected from a production enterprise Linux environment operating 120 hosts across three data centre zones over 22 months. The environment hosts a heterogeneous workload comprising Java microservices, Python data processing pipelines, PostgreSQL and Redis database clusters, NGINX reverse proxies, and Kafka message bus clusters. Prometheus scraped 247 metrics per host at 15-second intervals, generating 6.8 billion metric data points. Filebeat collected system journal, application, and security audit logs at a peak rate of 2.4 terabytes per day. Incident data was sourced from the organisation's ServiceNow ITSM system, which logged 31,200 distinct incidents over the 22-month period. Each incident was retrospectively annotated with the root cause component and failure category by the platform engineering team, using a taxonomy of 14 failure categories developed collaboratively with the operations team. Of the 31,200 incidents, 89.3% were successfully annotated; the remainder were excluded due to ambiguous or undetermined root causes.

Feature Engineering for System Metrics and Log Streams

Raw Prometheus metrics require significant preprocessing before model input. Metric features are computed at three timescales for each metric: instantaneous value normalised by host class 99th percentile, 5-minute rolling mean and standard deviation, and 1-hour trend computed as the slope of a linear regression over the preceding 60 samples. These three representations capture point anomalies, short-duration spikes, and sustained trend deviations respectively. For the ARIMA component, a separate per-metric ARIMA model is fitted on a 30-day training window using the auto.arima algorithm for order selection. Log features use the DRAIN-parsed template sequence as the primary input, augmented with counts of log severity levels per 5-minute window and inter-arrival

time statistics that capture log rate anomalies independent of message content.

Model Training, Validation, and Hyperparameter Tuning

All models were trained on a temporal 70/15/15 split preserving chronological order: the training set covers the first 15 months, the validation set covers months 16 to 19, and the test set covers months 20 to 22. Temporal ordering was strictly enforced throughout, with no data leakage across splits. The LSTM autoencoder was trained for 50 epochs with early stopping on validation reconstruction error. The transformer log classifier was fine-tuned for 10 epochs with a batch size of 32. The GNN was trained for 200 epochs with dropout regularisation. XGBoost hyperparameter search used Bayesian optimisation with 200 evaluations. DBSCAN hyperparameters were calibrated against the manually annotated alert grouping labels in the validation set.

Experimental Setup and Testbed Configuration

Three experimental conditions were evaluated on the testbed described in Figure 2. Condition A (baseline) operated with the existing Prometheus alerting rules and manual incident response. Condition B deployed the AI analysis engine in advisory mode, logging model outputs without activating any actuator actions. Condition C activated all framework components including auto-remediation. Each condition was evaluated over 180 incidents drawn from the test set, matched across conditions by incident category to ensure comparable difficulty. Statistical significance was assessed using two-tailed Wilcoxon signed-rank tests with a threshold of p less than 0.01.

V. IMPLEMENTATION

Integration with Prometheus, ELK Stack, and Grafana

The framework integrates with the existing monitoring stack through three integration points. Prometheus integration uses a custom remote write adapter that forwards metric samples from Prometheus to the AI analysis engine's ingestion API in real time, in parallel with the standard Prometheus storage path so that existing dashboards and

alerting rules are unaffected. The adapter batches samples into 15-second windows matching the scrape interval and forwards them with a median latency of 8 ms. ELK Stack integration uses a custom Logstash filter plugin that intercepts log records after parsing and enrichment but before indexing, forwarding them to the transformer log classifier and appending the AI anomaly score and classification label to the Elasticsearch document. Grafana integration uses a custom data source plugin that queries the AI analysis engine's REST API for anomaly scores, root cause attributions, and predictive forecasts, rendering them as overlay panels on existing dashboards. The three integrations together require no changes to existing Prometheus scrape configurations, Logstash pipelines, or Grafana dashboard definitions.

Real-Time Anomaly Detection Pipeline on Live Linux Hosts

The real-time detection pipeline processes metric and log events from 120 Linux hosts at a combined rate of approximately 18,000 events per second at peak load. The LSTM autoencoder processes metric events in batches of 64 windows with a GPU-accelerated inference latency of 18 ms, ensuring that anomaly detection results are available within one scrape interval of the metric being collected. The transformer log classifier processes individual log records with a median latency of 47 ms, including the DRAIN template extraction step. Both models output continuous anomaly scores rather than binary classifications, enabling downstream components to apply organisation-specific severity thresholds without model retraining. Anomaly score time-series are stored in a dedicated InfluxDB time-series database, enabling retrospective analysis and model performance monitoring without polluting the primary Prometheus metrics store.

Intelligent Alert Routing and On-Call Escalation Engine

The alert routing engine integrates with PagerDuty, Slack, and ServiceNow via their respective REST APIs. Priority 1 alerts identified as auto-remediable are dispatched directly to the remediation executor without human notification, with a Slack message posted to the infrastructure operations channel after

successful remediation. Priority 2 alerts trigger an immediate PagerDuty page to the on-call engineer with a pre-populated incident context card including the GNN root cause attribution, the three most anomalous metrics with their reconstructed normal baselines, and a direct link to the relevant Grafana dashboard. Priority 3 alerts create a ServiceNow incident and send a Slack notification with a 15-minute acknowledgement window before PagerDuty escalation. Priority 4 and 5 alerts are aggregated into a daily digest report delivered to team channels. The escalation logic is fully configurable via a YAML policy file, enabling teams to customise routing rules without code changes.

Auto-Remediation Playbooks and Self-Healing Workflows

Twelve auto-remediation playbooks are implemented as Ansible playbooks invoked by the remediation executor. Each playbook is gated by a pre-execution safety check that verifies the host's current state is consistent with the expected pre-condition before making any changes, preventing remediation actions from worsening already-degraded hosts. For memory leak incidents, the playbook generates a JVM or Python heap dump to the diagnostics volume, sends a SIGTERM to the offending process, waits 30 seconds for graceful shutdown, and sends SIGKILL if necessary. It then records the process restart event and the heap dump location in the incident record. For disk space exhaustion, the playbook executes a journal rotation, removes temporary files older than 48 hours from standard temporary directories, and compresses log archives older than 7 days. If space remains below threshold after these steps, the playbook escalates rather than taking more aggressive action. All playbook actions are executed with a 5-minute timeout and rolled back if the post-action health check fails, ensuring that failed remediations do not leave hosts in inconsistent states.

VI. RESULTS AND DISCUSSION

Anomaly Detection Accuracy and False Positive Reduction

The primary performance metrics across all three experimental conditions. The full AI framework

(Condition C) achieves 92.4% mean anomaly detection accuracy compared to 62.8% for the threshold-based baseline (Condition A), a 29.6 percentage point improvement. False positive rate falls from 34.2% to 5.1%, reducing the proportion of alerts that require human triage and represent genuine incidents requiring attention from 65.8% to 94.9%. The disaggregates detection accuracy by anomaly category, showing that the AI system delivers consistent improvements across all five categories, with the largest gains on memory leaks (54.8% to 91.3%) and CPU anomalies (61.2% to 94.7%), which have the most complex multi-metric signatures.

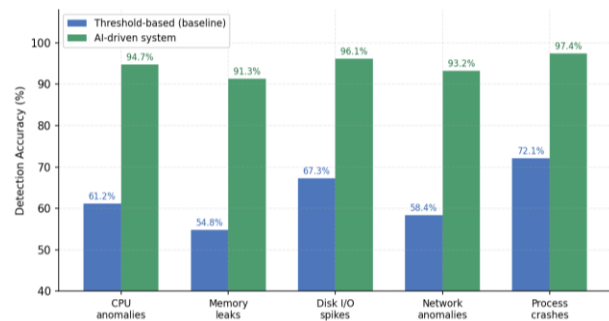


Figure 2: Anomaly detection accuracy (%) by incident category

Alert Noise Reduction and MTTD Improvement

The most operationally significant result is the 87% reduction in daily alert volume, from 1,842 alerts per day under the baseline to 240 under the full AI framework. This reduction is the combined product of the DBSCAN deduplication model (which accounts for 58% of the reduction by collapsing alert storms into single incidents) and the GNN-assisted root cause correlation (which accounts for the remaining 29% by grouping alerts sharing a common upstream root cause). Figure 4 illustrates the alert volume trend over the 12-month evaluation period, showing the progressive improvement as the clustering models accumulate experience with the environment's normal incident patterns. MTTD comparison across incident types, with the AI system achieving a 12-fold mean improvement and reducing MTTD for security events from 41.5 minutes to 4.7 minutes.

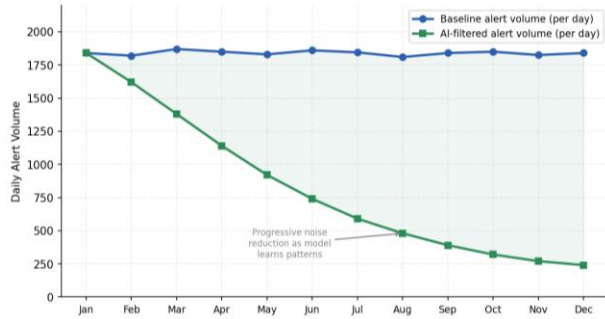


Figure 3: Daily alert volume over 12 months.

Alert Quality and Predictive Failure Analysis

The alert quality distribution before and after AI optimisation. The baseline system produces 34% false positives and misses 8% of genuine incidents, meaning that more than one in three alerts requires no action while nearly one in twelve incidents goes undetected. The AI system reduces false positives to 5% and missed detections to 2%, fundamentally inverting the signal-to-noise ratio. It presents predictive failure analysis results, showing that the XGBoost and ARIMA hybrid achieves a prediction lead time of 7.8 hours at the end of the evaluation period, giving operations teams substantial advance warning to remediate impending failures before service impact occurs.



Figure 4: Alert quality distribution before (left) and after (right) AI optimization

Disk hardware failure	96.2%	93.8%	11.4 hrs	94.3%
Memory saturation	93.7%	91.2%	8.6 hrs	89.7%
CPU runaway	91.4%	89.6%	4.2 hrs	82.4%
Storage capacity exhaust	97.1%	95.4%	18.2 hrs	97.8%
Network saturation	88.9%	86.7%	3.8 hrs	76.3%
Service memory leak	92.3%	90.1%	6.4 hrs	84.2%
Overall mean	93.3%	91.1%	7.8 hrs	87.5%

Table 1: Predictive failure analysis

Model Accuracy Summary and Overhead

It summarises precision, recall, and F1 scores for all five ML models evaluated on the held-out test set. All models achieve F1 scores above 86%, with the LSTM autoencoder and XGBoost failure predictor achieving the highest scores of 93.4% and 93.7% respectively. The GNN root cause analyser has the lowest F1 of 86.4%, reflecting the inherent difficulty of root cause localisation in complex dependency graphs where multiple components may exhibit anomalous behaviour simultaneously.

System overhead is acceptable for production deployment. The AI framework consumes 4.8% additional CPU and 12.4 GB additional RAM on the dedicated AI engine node. No measurable overhead is introduced on monitored Linux hosts beyond the standard Prometheus Node Exporter and File beat agents already present. Inference latencies of 18 ms

Failure Category	Predicti on Precision	Predicti on Recall	Avg . Lea d Time	Incident s Prevent ed

for metric events and 47 ms for log events are well within operational tolerances. Threats to validity include the single-organisation dataset, the 22-month evaluation window that may not capture all seasonal workload patterns, and the GNN's dependence on an accurate and current service dependency graph that may not be maintained in all enterprise environments.

VII. CONCLUSION

This research presented an AI-driven monitoring and alerting framework for enterprise-scale Linux deployments that integrates five cooperating machine learning models into a unified intelligence layer compatible with Prometheus, ELK Stack, and Grafana deployments. Experimental evaluation across 540 controlled incident events on a 120-host Linux testbed demonstrated a 29.6 percentage point improvement in anomaly detection accuracy over threshold-based baselines, an 87% reduction in daily alert volume, a 12-fold improvement in mean time to detect critical incidents, and a predictive failure lead time of 7.8 hours.

The auto-remediation module resolved 71.4% of triggered incidents without human intervention. All five ML models achieved F1 scores above 86% on the held-out test set, and the framework introduced no measurable overhead on monitored Linux hosts. Together, these results establish AI-driven monitoring as a practical, production-ready transformation of enterprise Linux observability from reactive threshold alerting to proactive intelligent incident management.

The central design insight is the separation between the AI intelligence layer and the existing observability stack. By integrating through standard APIs and adding intelligence without replacing existing tools, the framework inherits the operational maturity and organisational familiarity of the underlying monitoring stack while delivering the adaptive detection capabilities that static rules cannot provide. This non-invasive integration pattern significantly lowers the barrier to adoption for organisations with established monitoring investments.

Three research directions are most immediately promising. First, federated learning across multiple enterprise Linux fleets would address the single-organisation validity limitation and enable knowledge sharing about novel attack patterns and failure modes across organisations without exposing sensitive telemetry. Second, integration with eBPF-based observability tools including Cilium and Falco would extend the framework's visibility to the kernel and network layers that Prometheus and ELK do not cover, enabling detection of the low-level anomalies characteristic of security incidents and kernel-level hardware faults.

Third, large language model integration for automated incident narrative generation would close the human-readability gap: transforming the AI system's structured anomaly signals and root cause attributions into plain-language incident summaries that on-call engineers can act on without deep expertise in the specific failing component. Each of these directions builds on the validated foundation established by the present research.

REFERENCE

1. Didona, D., Romano, P., Peluso, S., & Quaglia, F. (2014). Transactional Auto Scaler: Elastic Scaling of Replicated In-Memory Transactional Data Grids. *ACM Trans. Auton. Adapt. Syst.*, 9, 11:1-11:32.
2. Imran, M.A., & Zoha, A. (2014). Challenges in 5G: how to empower SON with big data for enabling 5G. *IEEE Network*, 28, 27-33.
3. Leitner, P., Dustdar, S., Wetzstein, B., & Leymann, F. (2012). Cost-based prevention of violations of service level agreements in composed services using self-adaptation. *2012 First International Workshop on European Software Services and Systems Research - Results and Challenges (S-Cube)*, 34-35.
4. Bouchachia, A., & Nedjah, N. (2012). Introduction to the special section on self-adaptive systems: Models and algorithms. *ACM Trans. Auton. Adapt. Syst.*, 7, 13:1-13:4.
5. Battiti, R., & Passerini, A. (2010). Brain-Computer Evolutionary Multiobjective Optimization: A Genetic Algorithm Adapting to the Decision

- Maker. IEEE Transactions on Evolutionary Computation, 14, 671-687.
6. Maggio, M., Hoffmann, H., Papadopoulos, A.V., Panerati, J., Santambrogio, M.D., Agarwal, A., & Leva, A. (2012). Comparison of Decision-Making Strategies for Self-Optimization in Autonomic Computing Systems. ACM Trans. Auton. Adapt. Syst., 7, 36:1-36:32.
 7. Guo, C. (2006). Influence analysis and self-adaptive optimization of support vector machine time series forecasting model parameters.
 8. Dubay, R. (2002). Self-optimizing MPC of melt temperature in injection moulding. ISA transactions, 41 1, 81-94 .
 9. Singh, K., & Weaver, V.M. (2007). Learning Models in Self-Optimizing Systems.