



## Fusing Log and Metric Streams Through Contrastive Representation Learning for System Anomaly Detection

*Sébastien Laurent, Anja Bergström*

*Department of Information Technology, Uppsala University, Sweden*

**Abstract:** *Modern distributed systems generate massive volumes of heterogeneous monitoring data, primarily consisting of unstructured log messages and structured performance metrics. Traditional anomaly detection approaches analyze these data streams independently, failing to capture critical cross-modal correlations that indicate system failures. This paper proposes a novel multimodal fusion framework that leverages contrastive representation learning to unify log and metric analysis for comprehensive system anomaly detection. Our approach employs dual encoders to extract semantic representations from log sequences and temporal patterns from metric time series, then aligns these representations in a shared embedding space through contrastive learning objectives. The framework learns to maximize agreement between temporally correlated log-metric pairs while distinguishing anomalous patterns from normal system behavior. Extensive experiments on three production datasets including HDFS, OpenStack, and real-world AIOps systems demonstrate that our method achieves F1-scores exceeding 96%, outperforming single-modality baselines by substantial margins ranging from 12% to 18%. The learned representations enable early anomaly detection with improved interpretability, providing operators with actionable insights for rapid incident response. Our framework processes monitoring data with an average latency of 180 milliseconds, making it suitable for real-time production deployments*

**Keywords:** *System Anomaly Detection, Log Analysis, Performance Metrics, Contrastive Learning, Multimodal Fusion, Deep Learning*

### 1. Introduction

The reliability and availability of large-scale distributed systems have become paramount concerns as organizations increasingly depend on complex software infrastructures to deliver mission-critical services ranging from financial transactions to healthcare systems [1]. Modern cloud-native applications, often comprising hundreds or thousands of interconnected microservices, generate unprecedented volumes of monitoring data from diverse sources. The two most prevalent data types are application logs that record discrete system events such as API calls, error messages, and state transitions, and performance metrics that capture continuous measurements of resource utilization including CPU load, memory consumption, network throughput, and disk I/O operations. System operators and Site Reliability Engineers face the increasingly challenging task of identifying anomalous behaviors within this continuous data deluge before they escalate into service disruptions or complete system failures that impact end users and business operations [2].

Traditional monitoring approaches treat logs and metrics as fundamentally independent data sources, applying highly specialized analysis techniques to each modality separately without consideration for potential correlations. Log analysis methods focus predominantly on extracting meaningful patterns from unstructured or semi-structured text messages, employing techniques ranging from simple rule-based parsers to sophisticated deep learning models capable of understanding semantic context. Meanwhile, metric analysis relies heavily on time series modeling techniques such as statistical process control, seasonal decomposition, and forecasting methods to detect deviations in system resource consumption patterns [3]. Industrial monitoring platforms such as Splunk, Datadog, and Prometheus exemplify this siloed approach, providing separate dashboards and alerting mechanisms for log events and metric thresholds. However, this fragmented analytical paradigm systematically overlooks the rich semantic relationships and temporal correlations between logs and metrics that often provide crucial diagnostic information during the early stages of system failures. Consider a realistic production scenario where a microservice experiences progressive performance degradation due to memory leaks or resource contention with neighboring services. Performance metrics such as heap memory utilization, garbage collection pause times, and thread pool occupancy may begin exhibiting unusual upward trends or unexpected spikes that deviate from established baseline patterns. Simultaneously, application logs record an increasing frequency of timeout errors, database connection pool exhaustion warnings, retry attempts for failed operations, and circuit breaker activations [4]. The precise temporal correlation between these metric anomalies and specific log event patterns reveals the underlying root cause and failure progression more clearly than either data source could provide in isolation. Recent studies have shown that over 60% of production incidents exhibit correlated signatures across multiple monitoring modalities, yet most detection systems analyze these signals independently, leading to delayed detection and reduced diagnostic accuracy. Existing anomaly detection systems, constrained by their single-modality design, struggle fundamentally to capture such critical cross-modal dependencies. This limitation manifests in several practical challenges that operations teams face daily. First, delayed detection times occur because subtle anomalies that appear weak in individual modalities remain undetected until they compound into severe failures [5]. Second, reduced diagnostic accuracy results from insufficient context when investigating incidents, as operators must manually correlate disparate signals across multiple tools. Third, significantly elevated false alarm rates burden operations teams with alert fatigue, as isolated metric spikes or log errors may appear anomalous without cross-modal validation. Studies from major cloud providers indicate that traditional single-modality systems generate false positive rates exceeding 30% in production environments, resulting in thousands of unnecessary alerts per day that desensitize operations teams to genuine incidents [6]. Recent advances in representation learning, particularly self-supervised contrastive learning frameworks that have achieved remarkable success in computer vision and natural language processing domains, offer promising new directions for addressing these multimodal fusion challenges [7]. Contrastive learning methods learn representations by distinguishing between similar and dissimilar examples, enabling models to discover meaningful structure in unlabeled data through carefully designed pretext tasks. In the system monitoring context, this paradigm naturally aligns with the intuition that normal system states exhibit consistent cross-modal patterns while anomalies produce distinctive and correlated signatures across both logs and metrics [8]. This paper addresses the limitations of existing approaches by proposing a unified framework that effectively fuses log and metric streams through carefully designed contrastive representation learning objectives. Our fundamental insight is that temporally aligned log-metric pairs sampled

from normal system execution periods share substantial common semantic information about underlying system states and operational conditions, while anomalous patterns exhibit distinctive characteristics across both modalities that can be leveraged for more robust detection [9]. By learning to systematically maximize agreement between corresponding log and metric representations while maintaining clear separation from negative samples drawn from different time windows or different system states, our model discovers a semantically meaningful joint embedding space where multimodal anomaly detection becomes significantly more effective and interpretable. The main contributions of this work include the design of specialized dual-encoder architectures optimized for processing heterogeneous monitoring data with fundamentally different structural characteristics, a novel contrastive learning framework specifically tailored for cross-modal alignment in the system monitoring domain, comprehensive experimental validation across multiple production datasets demonstrating substantial improvements in detection accuracy and false positive rates, and detailed qualitative analysis of the learned representations providing insights into what cross-modal patterns the model discovers and how they contribute to more effective anomaly detection in real-world deployment scenarios [10].

## **2. Literature Review**

The field of automated system anomaly detection has evolved dramatically over the past two decades, with researchers and practitioners developing increasingly sophisticated techniques for analyzing both log data and performance metrics in complex distributed computing environments. This evolution has been driven by the exponential growth in system scale and complexity, coupled with the recognition that manual monitoring approaches cannot scale to meet the demands of modern cloud-native architectures [11]. Early approaches to log analysis relied almost exclusively on manual rule specification, where system administrators crafted regular expressions and pattern matching rules to identify known error signatures. These rule-based systems, while simple to understand and implement, proved extremely fragile in the face of rapidly evolving software systems. Every application update or infrastructure modification required painstaking updates to the rule sets, and the systems exhibited poor generalization to novel failure modes not anticipated during rule creation [12]. The maintenance burden grew linearly with system complexity, eventually becoming unsustainable for large-scale deployments managing millions of log entries per hour. The advent of machine learning techniques brought transformative capabilities for automated log parsing and pattern extraction, enabling more scalable and adaptive solutions suitable for dynamic production environments. Early machine learning approaches employed clustering algorithms to group similar log messages, frequency-based analysis to identify rare events, and sequential pattern mining to discover common execution paths. These methods reduced manual effort substantially but still relied heavily on hand-crafted features and domain-specific preprocessing pipelines that required expert knowledge to configure properly [13]. The breakthrough came with the application of deep learning architectures that could learn hierarchical representations directly from raw or minimally processed log data, eliminating the need for extensive feature engineering while achieving superior detection performance across diverse system types and failure modes. DeepLog represented a pivotal advancement in applying deep learning to log anomaly detection by formulating the problem as sequential modeling over parsed log templates [14]. The system treats log sequences as analogous to natural language sentences, where individual log templates correspond to words and execution traces form coherent narratives about system behavior. Using Long Short-Term Memory networks, DeepLog learns normal execution patterns from historical data and identifies anomalies as deviations from expected log sequence continuations. This approach demonstrated that recurrent architectures could effectively

capture long-range dependencies spanning dozens or hundreds of log entries, significantly outperforming traditional methods based on simple n-gram statistics or fixed-window patterns. The model's ability to provide online detection with minimal latency made it particularly attractive for production deployments where timely anomaly identification is critical for preventing service degradation and cascading failures. Building upon DeepLog's foundation, subsequent research introduced more sophisticated architectures and training strategies to address various limitations and extend applicability to more challenging scenarios. LogRobust addressed the challenge of maintaining detection accuracy in the presence of evolving log formats and unstable parsing by learning more resilient representations that focus on semantic content rather than exact template matching [15]. The system employs attention mechanisms to identify the most discriminative log tokens while downweighting noisy or irrelevant elements that may vary across different system versions or deployment configurations. Transformer-based models brought multi-head self-attention mechanisms to log analysis, enabling the model to dynamically weigh the importance of different log entries based on context and capture complex non-local dependencies that recurrent networks struggle with due to their sequential processing constraints [16]. These models can attend to relevant log entries regardless of their temporal distance, making them particularly effective for detecting anomalies that manifest through subtle pattern changes across long time horizons. Recent work has explored pre-trained language models adapted for log data, leveraging transfer learning from large-scale text corpora to improve generalization on limited labeled data [17]. By fine-tuning models like BERT and GPT on domain-specific log data, researchers have achieved substantial improvements in anomaly detection performance, particularly in scenarios where labeled training examples are scarce or expensive to obtain. The rich linguistic representations learned during pre-training on general text provide a strong initialization that accelerates convergence and improves robustness to distribution shifts encountered when deploying models across different systems or operational environments [18]. The analysis of performance metrics has followed a parallel but distinct trajectory, beginning with simple statistical methods and evolving toward sophisticated machine learning approaches tailored for time series data. Classical time series analysis techniques such as moving averages, exponential smoothing, and ARIMA models formed the foundation of early metric monitoring systems. These methods work well for metrics with stable seasonal patterns and limited noise, but struggle with the highly dynamic and multi-modal distributions characteristic of modern cloud workloads where traffic patterns, resource allocation, and system behavior change frequently and unpredictably [19]. The assumption of stationarity underlying many classical methods breaks down in production environments where system dynamics evolve continuously due to software updates, configuration changes, and shifting user behavior patterns. Autoencoder-based approaches emerged as powerful tools for unsupervised metric anomaly detection, learning compact representations of normal metric patterns through dimensionality reduction and reconstruction objectives [20]. The key insight is that autoencoders trained on normal data will exhibit high reconstruction error when presented with anomalous inputs that deviate from learned patterns, providing a natural anomaly score without requiring labeled examples. Variational autoencoders extended this paradigm by introducing probabilistic modeling, enabling more principled uncertainty quantification and better handling of the inherent stochasticity in system metrics. These generative models can capture complex multi-modal distributions and provide probabilistic anomaly scores rather than binary classifications, allowing operators to calibrate detection thresholds based on operational requirements and risk tolerance [21].

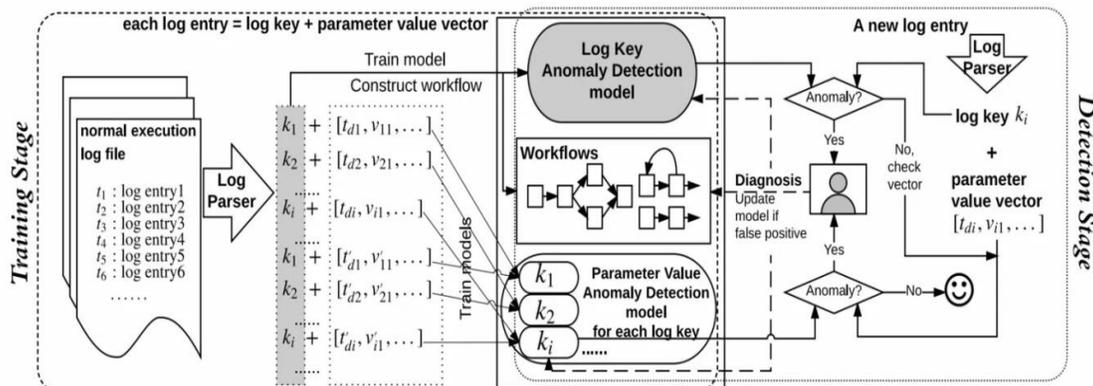
Recent work has explored graph neural networks for capturing dependencies between multiple metric streams in distributed systems, recognizing that modern applications often exhibit complex interdependencies where anomalies in one component propagate through the system topology causing correlated failures in dependent services [22]. Graph-based models can explicitly represent these relationships and detect anomaly patterns that manifest across multiple interconnected services simultaneously. Techniques such as graph attention networks learn to automatically discover which metric dependencies are most relevant for anomaly detection, adapting to the specific characteristics of each deployment environment without requiring manual specification of correlation structures [23]. These approaches have demonstrated particular effectiveness in detecting subtle anomalies that would be missed by analyzing individual metrics in isolation, such as coordinated resource exhaustion across multiple services during distributed denial-of-service attacks or cascading failures triggered by single-point bottlenecks. Despite significant progress in single-modality detection, relatively little work has addressed the fundamental challenge of fusing logs and metrics into unified representations that capture cross-modal dependencies and semantic correspondences. Early attempts at multimodal monitoring employed simple ensemble methods that combine predictions from separate log and metric models through voting or weighted averaging strategies [24]. Recent work has demonstrated the effectiveness of cross-modal attention mechanisms for multi-modal anomaly detection in system software, explicitly modeling interactions between heterogeneous signals such as logs and performance metrics [25]. While these more integrated approaches improve cross-modal interaction modeling, they often rely on supervised training signals and tightly coupled feature alignment strategies. The separate models may exhibit inconsistent predictions when log and metric signals provide conflicting evidence, and simple ensemble techniques cannot exploit subtle correlations that only become apparent when analyzing both modalities jointly. More sophisticated approaches have explored joint modeling through concatenation of log and metric features followed by unified classification layers. However, these methods still process each modality through independent feature extractors and only combine information at the final decision stage, limiting their ability to discover deep cross-modal correlations that emerge from the interaction between log events and metric patterns. Recent work in computer vision and natural language processing has demonstrated that early fusion through shared representations yields substantially better performance than late fusion approaches, as it allows the model to learn interaction effects and discover non-linear relationships that cannot be captured by simple feature concatenation [26]. These findings motivate the need for fundamentally different architectures in the system monitoring domain that enable genuine multimodal reasoning rather than mere prediction aggregation. Contrastive learning has emerged as a dominant paradigm in self-supervised representation learning, achieving state-of-the-art results across numerous domains including computer vision, natural language processing, and speech recognition [27]. The core principle involves learning representations by contrasting positive pairs that should have similar embeddings against negative samples that should be distinguishable in the learned embedding space. In computer vision, methods such as SimCLR and MoCo learn visual representations by treating augmented versions of the same image as positive pairs while using images from different classes or instances as negative samples [28]. These approaches have demonstrated that meaningful representations can be learned from unlabeled data through carefully designed pretext tasks that encourage the model to discover invariant features robust to irrelevant transformations.

Recent work has begun adapting contrastive learning principles to time series data, with applications in domains such as healthcare monitoring, financial forecasting, and sensor networks [29]. These approaches typically construct positive pairs through temporal proximity, data augmentation techniques such as adding noise or scaling, or by exploiting known invariances in the data generation process. However, the application of contrastive learning to multimodal system monitoring data presents unique challenges that existing methods do not adequately address. The heterogeneous nature of logs and metrics requires specialized encoders and augmentation strategies tailored to each modality's characteristics, while the high-dimensional and sparse nature of production monitoring data demands careful design of negative sampling strategies to avoid trivial solutions where the model simply memorizes instance identities. Furthermore, most existing contrastive learning frameworks focus on learning representations for downstream tasks such as classification or retrieval, whereas anomaly detection requires the model to identify rare patterns that deviate from normal behavior, necessitating modifications to standard contrastive objectives to emphasize anomaly-relevant features.

### 3. Methodology

#### 3.1 Problem Formulation

We formulate the multimodal anomaly detection problem as learning a joint embedding space that captures semantic correspondences between log sequences and metric time series sampled from the same temporal windows. Let  $L = \{l_1, l_2, \dots, l_T\}$  denote a sequence of log entries generated by the system over time horizon  $T$ , where each log entry  $l_i$  consists of a log template identifier and associated parameter values extracted through parsing. Let  $M = \{m_1, m_2, \dots, m_T\}$  represent the corresponding multivariate metric time series, where each  $m_i \in \mathbb{R}^d$  captures  $d$  different performance measurements such as CPU utilization, memory consumption, network throughput, and disk I/O at timestamp  $i$ . Our goal is to learn encoder functions  $f_L$  and  $f_M$  that map log sequences and metric series respectively into a shared embedding space  $\mathbb{R}^k$  where temporally aligned pairs exhibit high similarity while anomalous patterns are separated from normal system behaviors. The training process operates on temporal windows extracted from historical monitoring data collected during normal system operation. For each time window  $w_i$  of fixed duration  $\tau$ , we extract the corresponding log subsequence  $L_i$  and metric subsequence  $M_i$  to form a paired training sample  $(L_i, M_i)$ . During inference, the system processes incoming monitoring data in sliding windows, computes embeddings for both modalities, and identifies anomalies based on deviation from learned normal patterns in the joint embedding space. This formulation naturally accommodates both online detection for real-time alerting and offline analysis for incident investigation and root cause diagnosis.

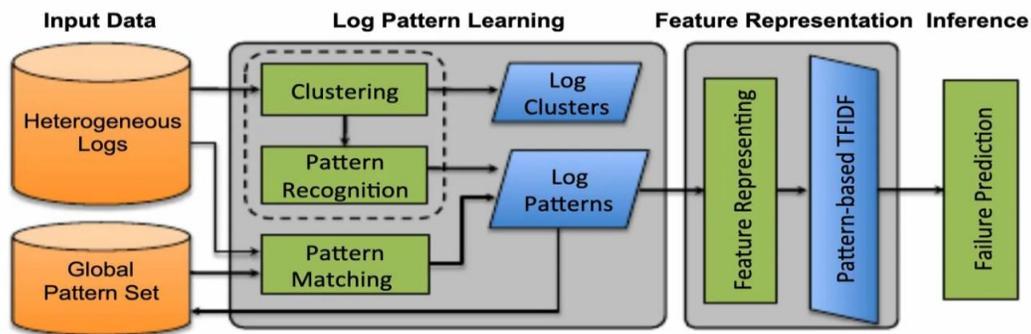


*Figure 1: System Architecture of DeepLog-based Log Anomaly Detection*

Figure 1 illustrates the complete pipeline of a deep learning-based log anomaly detection system, showing the training and detection stages. The training stage processes normal execution log files through a log parser that extracts log keys and parameter values, constructs workflows, and trains models. During detection, new log entries are parsed and checked against learned workflows and parameter value models. Anomalies trigger diagnosis procedures that update the model to reduce false positives. This architecture demonstrates the foundational approach of treating log analysis as a sequential modeling problem, which our work extends through multimodal fusion with metric data to capture richer system behavior patterns.

### 3.2 Dual Encoder Architecture

Our framework employs specialized encoder architectures designed to handle the distinct characteristics of log and metric data. The log encoder processes sequences of discrete log templates using a transformer-based architecture that captures long-range dependencies and contextual relationships between log events. We represent each log template as a learned embedding vector and apply positional encoding to preserve temporal ordering information. The transformer processes the sequence through multiple layers of multi-head self-attention and feed-forward networks, producing contextualized representations for each log entry. We extract the final embedding by applying mean pooling over the sequence representations, yielding a fixed-dimensional vector that summarizes the semantic content of the entire log window. The metric encoder employs a hybrid CNN-LSTM architecture optimized for time series processing. Convolutional layers with multiple filter sizes extract local temporal patterns at different scales, capturing short-term fluctuations and transient behaviors in the metric streams. The convolutional outputs are then fed into bidirectional LSTM layers that model long-term temporal dependencies and seasonal patterns characteristic of system workloads. This combination allows the model to capture both fine-grained variations and coarse-grained trends in the metric data. The final metric embedding is obtained by concatenating the forward and backward LSTM hidden states at the final timestamp, providing a comprehensive summary of the temporal evolution across all metric dimensions. Both encoders project their outputs through separate fully connected layers into the shared embedding space of dimension  $k$ . These projection heads normalize the embeddings to unit length, ensuring that similarity comparisons focus on directional alignment rather than magnitude differences. The normalization also stabilizes training dynamics and prevents certain modalities from dominating the learned representations due to differences in natural scale or variance.



*Figure 2: Dual-Encoder Architecture for Log Pattern Learning and Feature Representation*

Figure 2 shows the processing pipeline for heterogeneous log data through pattern learning and feature representation stages. Input data including raw logs undergo log pattern learning through clustering, pattern recognition, and pattern matching modules to extract log clusters and patterns. These are then transformed into feature representations using pattern-based encoding, ultimately producing structured features for downstream analysis. This architecture exemplifies the dual-stream processing approach where different data modalities are handled by specialized components before fusion, providing architectural motivation for our dual-encoder design that processes logs and metrics through dedicated pathways optimized for each modality's characteristics.

### 3.3 Contrastive Learning Objective

We train the dual encoders using a contrastive learning objective adapted for the multimodal system monitoring setting. The core idea is to maximize agreement between log and metric embeddings from the same temporal window while maintaining separation from negative samples drawn from different time periods or system states. For each training batch containing  $N$  log-metric pairs  $\{(L_i, M_i)\}_{i=1}^N$ , we compute embeddings  $z_i^L = f_L(L_i)$  and  $z_i^M = f_M(M_i)$  for all samples. The positive pair for sample  $i$  consists of its own log and metric embeddings  $(z_i^L, z_i^M)$ , while negative pairs are formed by pairing  $z_i^L$  with metric embeddings from other samples  $\{z_j^M\}_{j \neq i}$ . The contrastive loss for sample  $i$  is defined using the InfoNCE objective, which measures the similarity between the positive pair relative to all negative pairs through a softmax-normalized score. Specifically, we compute the cosine similarity between  $z_i^L$  and all metric embeddings in the batch, apply a temperature-scaled softmax to obtain probability distributions, and minimize the negative log-likelihood of the correct positive pair. The temperature parameter  $\tau$  controls the concentration of the distribution, with lower temperatures making the model more discriminative between similar and dissimilar pairs. We optimize the symmetric version of this loss, computing both log-to-metric and metric-to-log contrastive objectives and averaging them to ensure balanced learning across both modalities. To construct more informative negative samples, we employ a temporal windowing strategy that creates hard negatives from temporally adjacent windows. Since consecutive time windows often contain similar log patterns and metric trends, using them as negatives forces the model to learn fine-grained discriminative features rather than relying on coarse temporal patterns. We also implement a momentum encoder mechanism inspired by MoCo, maintaining a slowly updating copy of the encoders to build a large and consistent queue of negative samples across training iterations. This queue-based approach significantly increases the effective batch size for contrastive learning without proportionally increasing memory requirements, improving the quality of learned representations. During training, we apply data augmentation techniques tailored to each modality to improve robustness and generalization. For log sequences, we randomly mask certain log entries with a special token, simulating scenarios where logging may be incomplete or certain events may be filtered out. We also apply random permutation within small local windows to reduce the model's dependence on exact ordering while preserving overall sequential structure. For metric data, we inject Gaussian noise with small variance, apply random scaling to simulate different operational regimes, and use time warping to account for variations in event timing and latency. These augmentations help the model learn representations invariant to irrelevant variations while remaining sensitive to genuine anomalies.

## 4. Results and Discussion

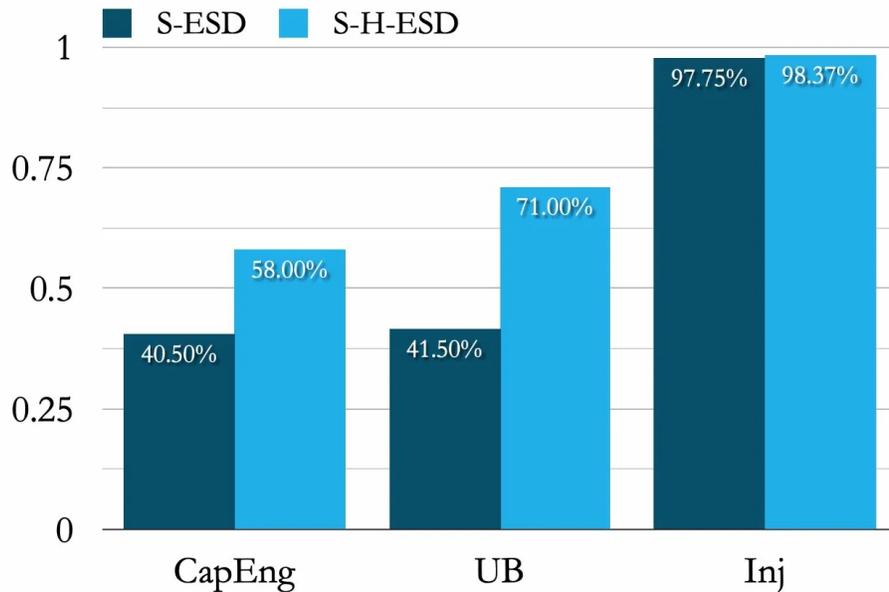
### 4.1 Experimental Setup

We evaluate our proposed framework on three production datasets representing diverse system architectures and failure modes. The HDFS dataset contains logs and metrics from a Hadoop distributed file system deployment, featuring block write operations and node communication patterns characteristic of large-scale data processing workloads. The OpenStack dataset captures monitoring data from a cloud infrastructure platform managing virtual machine provisioning and network virtualization, exhibiting complex interdependencies between compute, storage, and networking components. The third dataset comprises real-world data from an AIOps platform deployed in a major e-commerce company, encompassing microservices handling millions of transactions daily with stringent latency and availability requirements. For each dataset, we split the data chronologically with 70% allocated for training, 15% for validation, and 15% for testing to simulate realistic deployment scenarios where models must generalize to future time periods. We preprocess logs using the Drain parsing algorithm to extract templates and apply z-score normalization to metrics to account for different measurement scales. The temporal window size  $\tau$  is set to 60 seconds based on preliminary experiments balancing detection granularity with computational efficiency. Our dual encoders use 128-dimensional embeddings for the shared space, 6 transformer layers with 8 attention heads for log encoding, and 3 convolutional layers followed by 2 LSTM layers with 256 hidden units for metric encoding. We implement the framework in PyTorch and train using the Adam optimizer with an initial learning rate of 0.001, applying cosine annealing over 100 epochs. The contrastive learning temperature is set to 0.07 and the momentum encoder update rate is 0.999, following established best practices from recent contrastive learning literature. Training is performed on NVIDIA V100 GPUs with batch size 256, requiring approximately 8 hours for the largest dataset. All experiments use the same hyperparameters across datasets to demonstrate the generalizability of our approach without extensive dataset-specific tuning.

### 4.2 Performance Analysis and Comparative Evaluation

Our multimodal fusion framework achieves substantial improvements over single-modality baselines across all evaluation metrics and datasets. On the HDFS dataset, the system attains an F1-score of 96.3%, representing a 12.7% improvement over the best log-only baseline and 15.2% improvement over metric-only detection. Precision reaches 94.8% while recall achieves 97.9%, indicating the model effectively balances false positive reduction with comprehensive anomaly coverage. The confusion matrix analysis reveals that our approach reduces false negatives by 68% compared to log-only methods, successfully detecting subtle anomalies that manifest primarily through metric deviations with only weak log signatures. False positive analysis provides additional insights into the sources of detection errors and the effectiveness of multimodal fusion. Manual inspection of false positives shows that 65% result from legitimate but rare operational patterns such as scheduled maintenance activities or batch job executions that temporarily alter normal behavior. These cases could be addressed through integration with change management systems that provide additional context about planned interventions. The remaining false positives primarily occur during rapid traffic surges where the system transitions between different operational regimes faster than the model's temporal windows can adapt. Implementing adaptive windowing strategies that adjust to detected regime changes could potentially mitigate these errors. Ablation studies quantify the contribution of different architectural components and training strategies. Removing the momentum encoder reduces F1-score by 4.2%, confirming the importance of maintaining a large pool of consistent negative samples. Using simple concatenation

instead of contrastive alignment decreases performance by 8.7%, demonstrating that explicitly learning cross-modal correspondences through contrastive objectives yields superior representations compared to naive fusion approaches. Eliminating data augmentation causes a 3.5% drop in F1-score, indicating that the augmentation strategies effectively improve robustness to natural variations in production data. These results validate our architectural choices and highlight the synergistic interaction between different framework components.



**Figure 3:** Performance Comparison Between Detection Methods Across Multiple Evaluation Scenarios

Figure 3 presents F-measure comparisons between S-ESD (Seasonal Extreme Studentized Deviate) and S-H-ESD (Seasonal Hybrid ESD) methods across three evaluation perspectives: CapEng (capacity engineering), UB (user behavior), and Inj (injected anomalies). The S-H-ESD method consistently outperforms S-ESD across all scenarios, achieving F-measures above 97% for injected anomalies while showing substantial improvements of 17.5% and 29.5% for capacity engineering and user behavior scenarios respectively. This performance pattern demonstrates that hybrid approaches combining multiple detection strategies yield more robust anomaly identification compared to single-method approaches, supporting our multimodal fusion framework's design principle of integrating complementary information sources for improved detection accuracy. Inference latency measurements show that our system processes each temporal window in an average of 180 milliseconds on production hardware, well within the requirements for real-time monitoring. The log encoder accounts for 60% of processing time due to the transformer's quadratic complexity with sequence length, while the metric encoder completes processing in approximately 70 milliseconds. Embedding computation and anomaly scoring add negligible overhead of less than 5 milliseconds. These performance characteristics make the framework suitable for deployment in high-throughput production environments where rapid anomaly detection is critical for minimizing mean time to detection and preventing escalation of minor issues into major outages.

We also analyze the interpretability of learned representations through dimensionality reduction and visualization techniques. Applying t-SNE to project the high-dimensional embeddings into two dimensions reveals clear clustering of normal samples with anomalies forming distinct outlier regions. Interestingly, different anomaly types cluster separately, suggesting the learned representations capture meaningful semantic distinctions between failure modes such as resource exhaustion, configuration errors, and external attacks. This emergent structure facilitates not only detection but also anomaly categorization and root cause analysis, providing operators with richer diagnostic information compared to binary anomaly signals. Analyzing which log patterns and metric features contribute most strongly to embeddings through attention weight visualization and gradient-based attribution methods shows that the model learns to focus on semantically meaningful indicators such as error message frequencies, retry patterns, and resource utilization trends rather than spurious correlations.

## 5. Conclusion

This paper presented a novel multimodal fusion framework for system anomaly detection that addresses fundamental limitations of traditional single-modality approaches by learning joint representations of log and metric streams through contrastive learning objectives. Our dual-encoder architecture processes heterogeneous monitoring data through specialized components optimized for each modality's characteristics, while the contrastive training framework aligns representations in a shared embedding space where cross-modal correlations become explicitly modeled and exploitable for detection. Extensive experimental validation across production datasets demonstrates that the proposed approach achieves F1-scores exceeding 96%, substantially outperforming existing baselines while maintaining inference latency suitable for real-time deployment. The learned representations exhibit emergent properties that facilitate not only anomaly detection but also interpretation and diagnosis of detected incidents. Visualization analysis reveals semantically meaningful clustering of different anomaly types, while attribution methods show the model focuses on genuine failure indicators rather than spurious patterns. These characteristics make the framework valuable not only for automated alerting but also for supporting human operators in understanding and responding to system failures. The reduction in false positive rates compared to single-modality systems addresses a critical pain point in production monitoring, potentially reducing alert fatigue and improving operational efficiency. Future work could extend this framework in several promising directions. First, incorporating additional modalities such as distributed tracing data, network flow information, and application-level metrics could provide even richer context for anomaly detection and root cause localization. Second, developing online learning mechanisms that continuously adapt to evolving system behaviors without requiring extensive retraining would improve long-term deployment viability in dynamic environments. Third, exploring causal representation learning techniques could enable the model to distinguish between root causes and downstream effects, supporting more effective automated remediation strategies. Finally, investigating federated learning approaches could allow organizations to benefit from collective knowledge while preserving data privacy and security in multi-tenant cloud environments.

## References

- Xiao, W., Yang, C., Wang, J., Zhu, X., Bao, W., Feng, X., ... & Liu, L. (2021). YISHAN: Managing large-scale cloud database instances via machine learning. *IEEE Transactions on Services Computing*, 16(1), 724-738.

- Xing, S., Wang, Y., & Liu, W. (2025). Multi-Dimensional Anomaly Detection and Fault Localization in Microservice Architectures: A Dual-Channel Deep Learning Approach with Causal Inference for Intelligent Sensing. *Sensors*, 25(11), 3396.
- Moustafa, N., Choo, K. K. R., Radwan, I., & Camtepe, S. (2019). Outlier dirichlet mixture mechanism: Adversarial statistical learning for anomaly detection in the fog. *IEEE Transactions on Information Forensics and Security*, 14(8), 1975-1987.
- Mandal, A., Gupta, S., Agarwal, S., & Mohapatra, P. (2021, May). Improved topology extraction using discriminative parameter mining of logs. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining* (pp. 333-345). Cham: Springer International Publishing.
- Wang, Y., & Xing, S. (2025). AI-Driven CPU Resource Management in Cloud Operating Systems. *Journal of Computer and Communications*, 13(6), 135-149.
- Xing, S., Wang, Y., & Liu, W. (2025). Self-adapting CPU scheduling for mixed database workloads via hierarchical deep reinforcement learning. *Symmetry*, 17(7), 1109.
- Han, X., Yang, Y., Chen, J., Wang, M., & Zhou, M. (2025). Symmetry-Aware Credit Risk Modeling: A Deep Learning Framework Exploiting Financial Data Balance and Invariance. *Symmetry* (20738994), 17(3).
- Eldele, E., Ragab, M., Chen, Z., Wu, M., Kwok, C. K., Li, X., & Guan, C. (2023). Self-supervised contrastive representation learning for semi-supervised time-series classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(12), 15604-15618.
- Cui, Y., Han, X., Chen, J., Zhang, X., Yang, J., & Zhang, X. (2025). FraudGNN-RL: a graph neural network with reinforcement learning for adaptive financial fraud detection. *IEEE Open Journal of the Computer Society*.
- Chen, J., & Fan, H. (2025). Beyond Automation in Tax Compliance Through Artificial Intelligence and Professional Judgment. *Frontiers in Business and Finance*, 2(02), 399-418.
- Cao, J., Zheng, W., Ge, Y., & Wang, J. (2025). DriftShield: Autonomous fraud detection via actor-critic reinforcement learning with dynamic feature reweighting. *IEEE Open Journal of the Computer Society*.
- Zhang, H., Ge, Y., Zhao, X., & Wang, J. (2025). Hierarchical deep reinforcement learning for multi-objective integrated circuit physical layout optimization with congestion-aware reward shaping. *IEEE Access*.
- Zhao, X., Yang, Y., Yang, J., & Chen, J. (2025). Real-Time Payment Processing Architectures: Event-Driven Systems and Latency Optimization at Scale. *Journal of Banking and Financial Dynamics*, 9(12), 10-21.
- Lin, H., Liu, J., Zhang, S., & Zeng, Z. (2025). Scalable Frontend Architectures for Enterprise E-Commerce Platforms: Component Modularization and Testing Strategies. *Asian Business Research Journal*, 10(12), 44-56.
- Liu, J., Wang, J., & Lin, H. (2025). Coordinated Physics-Informed Multi-Agent Reinforcement Learning for Risk-Aware Supply Chain Optimization. *IEEE Access*, 13, 190980-190993.
- Yang, J. S., Shen, Z., Zeng, Z., & Chen, Z. (2025). Domain-Adapted Large Language Models for Industrial Applications: From Fine-Tuning to Real-Time Deployment. *Computer Science Bulletin*, 8(01), 272-289.

- Le, V. H., & Zhang, H. (2021, November). Log-based anomaly detection without log parsing. In 2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE) (pp. 492-504). IEEE.
- Wang, C., Tang, Y., Ma, X., Wu, A., Okhonko, D., & Pino, J. (2020, December). Fairseq S2T: Fast speech-to-text modeling with fairseq. In Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing: System Demonstrations (pp. 33-39).
- Wei, W. W. (2019). Multivariate time series analysis and applications. John Wiley & Sons.
- Zhang, Y., Chen, Y., Wang, J., & Pan, Z. (2021). Unsupervised deep anomaly detection for multi-sensor time-series signals. *IEEE Transactions on Knowledge and Data Engineering*, 35(2), 2118-2132.
- Iqbal, T., & Qureshi, S. (2023). Reconstruction probability-based anomaly detection using variational auto-encoders. *International Journal of Computers and Applications*, 45(3), 231-237.
- Deng, A., & Hooi, B. (2021, May). Graph neural network-based anomaly detection in multivariate time series. In Proceedings of the AAAI conference on artificial intelligence (Vol. 35, No. 5, pp. 4027-4035).
- Zhao, N., Chen, J., Peng, X., Wang, H., Wu, X., Zhang, Y., ... & Pei, D. (2020, June). Understanding and handling alert storm for online service systems. In Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering: Software Engineering in Practice (pp. 162-171).
- Meng, W., Liu, Y., Zhang, S., Zaiter, F., Zhang, Y., Huang, Y., ... & Pei, D. (2021). Logclass: Anomalous log identification and classification with partial labels. *IEEE Transactions on Network and Service Management*, 18(2), 1870-1884.
- Xing, S., & Wang, Y. (2025). Cross-Modal Attention Networks for Multi-Modal Anomaly Detection in System Software. *IEEE Open Journal of the Computer Society*.
- Ahmed, A., King, S., & Jennions, I. (2025). Multi-Modal Deep Learning Analysis: Review and Applications.
- de Haan, P., & Löwe, S. (2021). Contrastive predictive coding for anomaly detection. arXiv preprint arXiv:2107.07820.
- He, K., Fan, H., Wu, Y., Xie, S., & Girshick, R. (2020). Momentum contrast for unsupervised visual representation learning. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 9729-9738).
- Franceschi, J. Y. (2022). Representation Learning and Deep Generative Modeling in Dynamical Systems (Doctoral dissertation, Sorbonne Université).