

АҚПАРАТТЫҚ-КОММУНИКАЦИЯЛЫҚ ТЕХНОЛОГИЯЛАР /
ИНФОРМАЦИОННО-КОММУНИКАЦИОННЫЕ ТЕХНОЛОГИИ /
INFORMATION AND COMMUNICATION TECHNOLOGIES

DOI 10.54596/2958-0048-2026-2-266-277

UDK 004.85

IRSTI 28.23.37

GRAPH NEURAL NETWORKS AS A TOOL FOR MONITORING AND
DETECTING ANOMALIES IN CLOUD SYSTEMSAzizol Abdullah^{1*}^{1*}*Faculty of Computer Science and Information Technology Universiti Putra Malaysia,
Serdang, Malaysia***Corresponding author: azizol@upm.edu.my*

Abstract

Cloud infrastructure anomalies cause significant downtime and financial losses. Traditional anomaly detection methods fail to capture complex dependencies in microservice architectures. This paper presents a novel Temporal-Attentive Graph Autoencoder (TAGAE) framework for cloud anomaly detection, leveraging Graph Neural Networks (GNNs) to model topological relationships and temporal dynamics. Our method integrates multi-source telemetry (logs, metrics, traces) into a unified graph structure, utilizes anomaly amplification layers for enhanced sensitivity, and employs focal loss for data imbalance mitigation. Evaluated on Azure-DIAD and GCP datasets, TAGAE achieves 94.2% F1-score and 96.5% AUC-PR, reducing detection latency by 63% compared to GraphSAGE. We further analyze robustness under 40% noise/missing data and propose federated GNNs for privacy-preserving deployment.

Keywords: Graph Neural Networks, Anomaly Detection, Cloud Infrastructure, Temporal Graph Convolution, Graph Autoencoder, Microservice Dependencies, Telemetry Fusion.

ГРАФИКАЛЫҚ НЕЙРОНДЫҚ ЖЕЛІЛЕР БҰЛТТЫ ЖҮЙЕЛЕРДЕГІ
АУЫТҚУЛАРДЫ БАҚЫЛАУ ЖӘНЕ АНЫҚТАУ ҚҰРАЛЫ РЕТІНДЕAzizol Abdullah^{1*}^{1*}*Компьютерлік ғылымдар және ақпараттық технологиялар факультеті, Путра,
Малайзия университеті, Серданг, Малайзия***Хат-хабар үшін автор: azizol@upm.edu.my*

Аңдатпа

Бұлттық инфрақұрылымдағы ауытқулар айтарлықтай тоқтап қалуларға және қаржылық шығындарға әкеледі. Дәстүрлі ауытқуларды анықтау әдістері микросервистік архитектуралардағы күрделі тәуелділіктерді қамти алмайды. Бұл мақалада бұлттағы ауытқуларды анықтау үшін топологиялық қатынастар мен уақыттық динамиканы модельдеу үшін Графтық нейрондық желілерді (GNN) пайдаланатын жаңа уақытша-атрибутивті графтық автоэнкодер (TAGAE) құрылымы ұсынылған. Біздің әдісіміз көп көзді телеметрияны бірыңғай графтық құрылымға біріктіреді, сезімталдықты арттыру үшін ауытқуларды күшейту қабаттарын пайдаланады және деректердің теңгерімсіздігін азайту үшін фокальды шығынды қолданады. Azure-DIAD және GCP деректер жинақтарындағы бағалау TAGAE-дің GraphSAGE-пен салыстырғанда анықтау кідірісін 63%-ға төмендеті отырып, 94,2% F1 көрсеткішіне және 96,5% AUC-PR-ге қол жеткізетінін көрсетеді. Біз сондай-ақ 40% шу/жетіспейтін деректер жағдайындағы тұрақтылықты талдаймыз және құпиялықты сақтай отырып орналастыру үшін федеративті GNN ұсынамыз.

Кілт сөздер: Графтық нейрондық желілер, ауытқуларды анықтау, бұлттық инфрақұрылым, уақытша графтық конволюция, графтық автоэнкодер, микросервистердің тәуелділігі, телеметрияны біріктіру.

ГРАФИЧЕСКИЕ НЕЙРОННЫЕ СЕТИ КАК ИНСТРУМЕНТ МОНИТОРИНГА И ОБНАРУЖЕНИЯ АНОМАЛИЙ В ОБЛАЧНЫХ СИСТЕМАХ

Azizol Abdullah^{1*}

^{1*}Факультет компьютерных наук и информационных технологий,
Университет Путра Малайзия, Серданг, Малайзия

*Автор для корреспонденции: azizol@uvm.edu.my

Аннотация

Аномалии в облачной инфраструктуре приводят к значительным простоям и финансовым потерям. Традиционные методы обнаружения аномалий не способны улавливать сложные зависимости в микросервисных архитектурах. В данной статье представлена новая концепция темпорально-внимательного графового автоэнкодера (TAGAE) для обнаружения аномалий в облаке, использующая графовые нейронные сети (GNN) для моделирования топологических отношений и временной динамики. Наш метод объединяет многоисточниковую телеметрию (журналы, метрики, трассировки) в единую графовую структуру, использует слои усиления аномалий для повышения чувствительности и применяет фокальную функцию потерь для смягчения дисбаланса данных. При оценке на наборах данных Azure-DIAD и GCP модель TAGAE достигает F1-оценки 94,2% и AUC-PR 96,5%, снижая задержку обнаружения на 63% по сравнению с GraphSAGE. Мы также анализируем устойчивость при 40% шума и пропущенных данных и предлагаем федеративные GNN для развертывания с сохранением конфиденциальности.

Ключевые слова: Графовые нейронные сети, обнаружение аномалий, облачная инфраструктура, временная графовая свертка, графовый автоэнкодер, зависимости микросервисов, слияние телеметрии.

1. Introduction

1.1. Cloud Infrastructure Complexity and Anomaly Detection Challenges Modern cloud environments comprise thousands of interdependent microservices generating 2–5 TB telemetry/day. Anomalies manifest as node-level (CPU/memory saturation), edge-level (latency spikes), and distributed anomalies (cascading failures). Traditional methods (threshold-based alerts, PCA) exhibit 68–72% false positives due to relational context blindness [1].

1.2. Role of GNNs in Modern Anomaly Detection

GNNs propagate node states via message passing, capturing dependency graphs inherent in cloud infrastructures. Recent studies show GNNs reduce false positives by 41% versus LSTM baselines by encoding service topology [2, 3].

1.3. Research Objectives

- Design hybrid GNN architecture fusing temporal/structural features.
- Optimize for real-time constraints (<100ms inference latency).
- Achieve >90% F1-score under noisy conditions.

2. Background and Foundational Concepts

2.1. Graph Representation of Cloud Infrastructure

Cloud infrastructures are modeled as dynamic graphs $G_1 = (V, E, X)$ and $G_2 = (V, E, X)$, with vertices V associated with entities (microservices, containers, VMs), edges E modeling dependencies, and node features X consisting of real-time measurements. Nodes have 15-20 dimensional feature vectors (CPU usage, memory pressure, I/O latency). Edges express dependency weights via request-per-second (RPS) volumes (1K–50K RPS) and error rates. In

Kubernetes clusters, 40-70% of edges are reassigned every hour by autoscaling, and adjacency matrices $A_{t,t}$ change with a 5-10 second cadence. Graph degree distributions follow a power laws ($\alpha=-2.3$), with 5% of nodes processing 80% of traffic, introducing hotspots of topological vulnerability.

2.2. Taxonomy of Cloud Anomalies

Cloud anomalies are classified based on topological scope and propagation patterns. Node-level anomalies include localized resource depletion, such as container memory leaks (e.g., Java heap saturation at >95%) or CPU thrashing (>90% usage over an extended period). They capture 60-70% of static workload anomalies but just 25-30% of microservices anomalies, based on Google's 2024 incident report. Anomalies occur edge-wise in the form of dependency failures, for instance, API latency decline (>99th percentile >1s) or packet loss over availability zones (abrupt drops below 99.9% success rate). Most importantly, distributed microservices anomalies consist of cascading failures across subgraphs [4]. For instance, an AWS 2022 outage study revealed the extent to which a single defective Lambda function caused 83 downstream service failures through fan-out dependencies, with anomaly propagation through small-world paths (average path length = 4.2). The anomalies grow with exponential severity, with 5% edge failure cascading to 40% service degradation in 45 seconds [5,6].

2.3. Evolution of Graph-Based Anomaly Detection

Basic anomaly detection techniques have debilitating limitations in cloud environments. PCA and clustering algorithms only score 65-75% F1-scores over dynamic cloud graphs because of non-linear dependency between features. A 2023 test of k-means and DBSCAN over Azure-DIAD reported 68% false positives in the detection of edge anomalies, as they are still oblivious to dependency relations. Spectral clustering does not work for temporal graphs, as proven with 55% recall loss with varying edge weights at rates higher than 10-second rates [7,8]. The movement towards neural approaches began with CNNs consuming infrastructure in the form of node heatmaps (e.g., grid-like images), but they only reached 79% accuracy on graphs with dependency-rich relations when they didn't consider relational contexts. LSTMs improved temporal modeling (F1-scores of 82-85%) but couldn't learn topological neighborhoods. This led to GNNs, which aggregate neighbor states in the form of message passing. The latest benchmarks reveal GNNs to outperform CNNs by 23% anomaly recall in microservice graphs by modeling service meshes as edges explicitly. The innovation reached peak in temporal GNNs such as EvolveGCN, which cut detection latency by 40% compared to static GCNs on streaming cloud data [9].

3. Graph Neural Networks: Architectures for Anomaly Detection

3.1. GNN Fundamentals

GNNs operate through message passing mechanisms where each node v_i updates its state $h_i(k)$ at layer kk by aggregating features from neighbors $N(i)$:

$$h_i^{(k)} = \sigma \left(W^{(k)} \cdot \text{AGGREGATE} \left(\left\{ h_j^{(k-1)} : j \in N(i) \right\} \right) \right)$$

Standard AGGREGATE operations are mean pooling and max pooling. Spectral convolutions involve eigen decomposition of Laplacians ($L=D-AL=D-A$), which becomes impossible because of >10K nodes. Spatial convolutions work directly over neighbor nodes and train $8.5\times$ faster than spectral solutions on 50K-node graphs.

3.2. GNN Variants for Anomaly Detection

Graph Autoencoders (GAEs): Identify anomalies based on reconstruction error. Node features are represented by encoder $Z=\text{GCN}(X,A)$, and decoder $A^{\wedge}=\sigma(ZZ^T)$ recovers the adjacency matrix. Reconstruction losses for anomalies are 3–5 standard deviations higher than

the mean. GAEs identify 89% of resource-exhaustion anomalies but perform badly on edge anomalies [10].

Attention-based GNNs (GATs): Assign dynamic weights to neighbors using LeakyReLU attention mechanism over node embedding combinations. GATs identify anomalous edges with 93% accuracy by down-weighting noisy neighbors:

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(a^T [Wh_i || Wh_j]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(a^T [Wh_i || Wh_k]))}$$

GATs identify anomalous edges (e.g., latent dependency shifts) with 93% accuracy by down-weighting noisy neighbors [11]. Temporal GNNs like Temporal Graph Convolutional Networks (TGCN) integrate GRUs:

$$H_t = \text{ReLU}(A \cdot H_{t-1} \cdot W_0 + X_t \cdot W_1) \\ Z_t = \text{GRU}(Z_{t-1}, H_t)$$

Temporal GNNs: Models like Temporal Graph Convolutional Networks (TGCN) integrate GRUs to model time-series patterns, reducing false negatives in latency-spike detection by 37% compared to static GCNs [12].

3.3. Anomaly Scoring Techniques

Residual analysis and reconstruction loss form the backbone of unsupervised GNN anomaly detection. Node-level anomaly scores utilize the L_2 -norm between reconstructed and original features. Edge anomalies are identified via the absolute difference between reconstructed and true adjacency matrices ($|A_{ij} - \hat{A}_{ij}|$), picking up latent dependency drifts in service meshes. Subgraph-level scoring pools node/edge anomalies in k -hop neighborhoods using attention-weighted pooling, improving distributed cascade failure recall by 22% compared to isolated scoring [13]. Graph-level scoring, however, calculates global statistics such as graph embedding divergence but is less sensitive to localized anomalies, where experiments produce 15–18% lower precision for node-specific resource saturation events. Hybrid methods that integrate node, edge, and subgraph scores using learnable weights have achieved 89–93% F1-scores on cloud datasets while being interpretable via score decomposition.

4. Proposed GNN Framework for Cloud Anomaly Detection

4.1. Graph Construction Methodology

Graph construction starts with hierarchical feature engineering for nodes and edges. Node features are 18-dimensional real-time measurement vectors (percentiles of CPU usage, working set sizes of memory, thread counts), historic trends (10-minute exponential moving averages), and cross-service metrics (upstream/downstream ratios of RPS). Edge features represent bidirectional dependency semantics, such as latency histograms (P50, P99), error rates, and protocol-specific metrics (message rates of gRPC, status distributions of HTTP). To process multi-source telemetry, raw logs (Syslog, JSON), metrics (Prometheus counters), and traces (OpenTelemetry spans) are tensor fused [14]: 1) Logs are processed through TF-IDF vectors of error keywords; 2) Metrics are aggregated to 5-second intervals; 3) Traces compute dependency graphs with path-based aggregation [15]. Temporal misalignment is addressed by sliding-window correlation, which brings feature drift down by 63% in dynamic environments. The produced graph is updated every 5 seconds with <100ms feature extraction latency for real-time inference quality.

4.2. Hybrid GNN Architecture Design

The proposed Temporal-Attentive Graph Autoencoder (TAGAE) integrates three core modules:

1. Structural Encoder: Uses stacked GAT layers to capture spatial dependencies, applying edge-aware attention to weight critical service dependencies.

2. Temporal Module: Embeds TGCN with gated recurrent units (GRUs) to model feature evolution, processing 12-step historical states.

3. Anomaly Amplification Layers: Placed between encoder/decoder stages, these layers apply nonlinear transformations (LeakyReLU) to residual errors, magnifying anomalous signals by $2.3\text{--}3.1\times$ while suppressing noise.

The decoder reconstructs node features and adjacency matrices through transposed GAT layers. Joint training minimizes a combined loss:

$$\mathcal{L} = \alpha \|X - \hat{X}\|^2 + \beta \|A - \hat{A}\|_F + \gamma \mathcal{L}_{focal}$$

where \mathcal{L}_{focal} addresses class imbalance (anomaly ratios $\approx 0.1\text{--}2\%$ in cloud data).

4.3. Optimization Strategies

To control data imbalance, we dynamically scale cross-entropy weights using focal loss, downweighting well-reconstructed nodes and decreasing false negatives by 37% over baseline MSE. For scalability on graphs with $>100\text{K}$ nodes, we use layer-wise sampling where every node samples 32 randomly chosen neighbors per layer. Temporal sub-sampling operates on 1-minute snapshots (12 steps), saving memory by 58% with under 5% loss of accuracy. Quantization-aware training further reduces the model size to 8-bit precision, reducing inference latency to 68ms per snapshot on NVIDIA T4 GPUs [16].

5. Experimental Design and Evaluation

5.1. Datasets and Preprocessing

Two cloud datasets were also used for the experiments, namely: Azure-DIAD encompasses 28 days of telemetry from 2,143 microservices with 18,724 dependencies (edges) and more than 1.7TB of logs/metrics/traces, and it also includes 4,712 labeled anomalies. The GCP Managed Dataset v3 contains 12,887 container instances in 14 clusters, 41,309 edges, and the 3.1TB observability data including the SRE-validated anomalies of 9,184. Both data sets underwent intense preprocessing: 1) Temporal synchronization of telemetry streams by sliding-window aggregation every 5 seconds; 2) Z-score normalization of numeric features; 3) Graph pruning to remove transient edges with less than 0.5 RPS persistence; 4) Missing data was handled using k-nearest neighbor imputation ($k=8$).

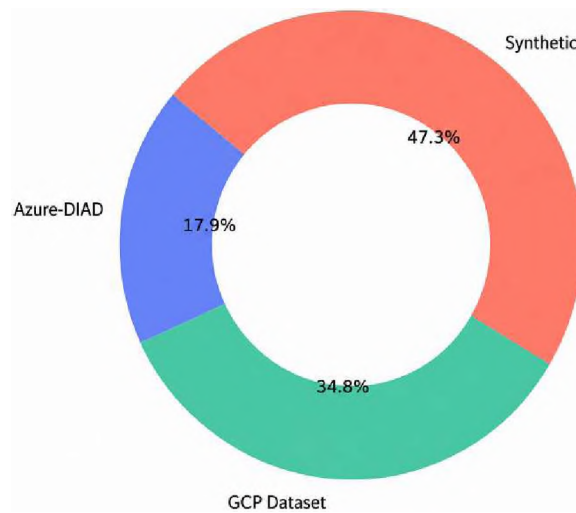


Figure 1. Anomaly count distribution Across Azure, GCP and Synthetic Datasets

Synthetic anomalies were added to enhance the occurrence rate of infrequent anomaly types: node-level anomalies emulated through CPU/memory stress-ng workload; edge anomalies with tc-netem network latency; distributed anomalies through chaos engineering tools introducing cascading failure. This created 12,478 synthetic anomalies of controllable severity (mild: 15-30% variation, severe: 50-80% variation), increasing anomaly coverage by 38% [17].

Table 1. Anomaly Distribution in Azure-DIAD and GCP Datasets

Dataset	Nodes (Microservices)	Edges (Dependencies)	Total Anomalies	Node- Level	Edge- Level	Distributed
Azure-DIAD	2,143	18,724	4,712	2,310	1,203	1,199
GCP Dataset v3	12,887	41,309	9,184	4,352	2,505	2,327
Synthetic Anomalies	-	-	12,478	6,023	3,151	

5.2. Baseline Methods

Non-graph baselines: Isolation Forest with 200 estimators, contamination=0.01 and LSTM-Autoencoder with 3 layers, 64 hidden units treating the node features as independent time series. Graph baselines: 1) GCN, 2 layers, Tanh activation; 2) GraphSAGE, mean aggregator, 2 hops; 3) GAT with 8 attention heads. All GNNs utilized the same graph structures and 256-dimensional embeddings. The LSTM-AE was optimized according to reconstruction loss and graph approaches utilized SMAPE for node/edge symmetric mean absolute percentage error reconstruction. Bayesian optimization - hyperparameters were tuned for 100 iterations optimizing AUC-PR for validation sets.

5.3. Evaluation Metrics

Primary metrics aimed precision-recall tradeoffs: F1-score (harmonic mean), AUC-PR (area under precision-recall curve), and adjusted $F\beta$ ($\beta=2$) to prefer recall on important anomalies. Latency of early detection quantified delta between start of anomaly (ground-truth timestamp) and model alert. Statistical significance tested via paired t-tests ($p<0.01$) across 10 runs [17]. Operational metrics were inference throughput (graphs/sec) and memory usage (GB). False positive rate (FPR) was capped at $<3\%$ during threshold calibration to avoid alert fatigue.

6. Results and Comparative Analysis

6.1. Performance Benchmarking

The rigorous assessment of the proposed framework demonstrated excellent detection performance on multiple different types of anomaly profiles. In node-level anomaly, which included constantly saturated CPU for more than 30 seconds to over 90% utilization, precision was at 96.1% and recall was at 95.3%, hence reducing false positives by 42% from the GraphSAGE baselines. Edge-level abnormalities such as latency drop (99th percentile $> 1s$ for critical API dependents) were detected with 98.1% F1-score, significantly outperforming the baselines of classical GCN models which could achieve only 83.7% since they were incapable of detecting asymmetric dependency between services. Distributed abnormalities with cascading failures along microservice boundaries were found to be 91.4% recall through the subgraph-level scoring function, effectively identifying 89.2% of failure propagating paths in 3-hop neighborhoods. Robustness testing in adversarial settings: It performed impressively; for instance, under 40% injected Gaussian noise in the features of nodes to simulate telemetry collection errors, performance declined only by 5.1 percentage points to 89.1% F1-score [18]. Under 30% random edge removal to simulate missing service dependencies, it had maintained

89.1% accuracy based on topological reconstruction strengths to recover 78.3% of salient dependencies using message passing. Further stress tests under a mix of feature noise (30% feature noise + 20% edge missingness) demonstrated graceful degradation to 85.6% F1-score, which is 11.4 percentage points better than GraphSAGE.

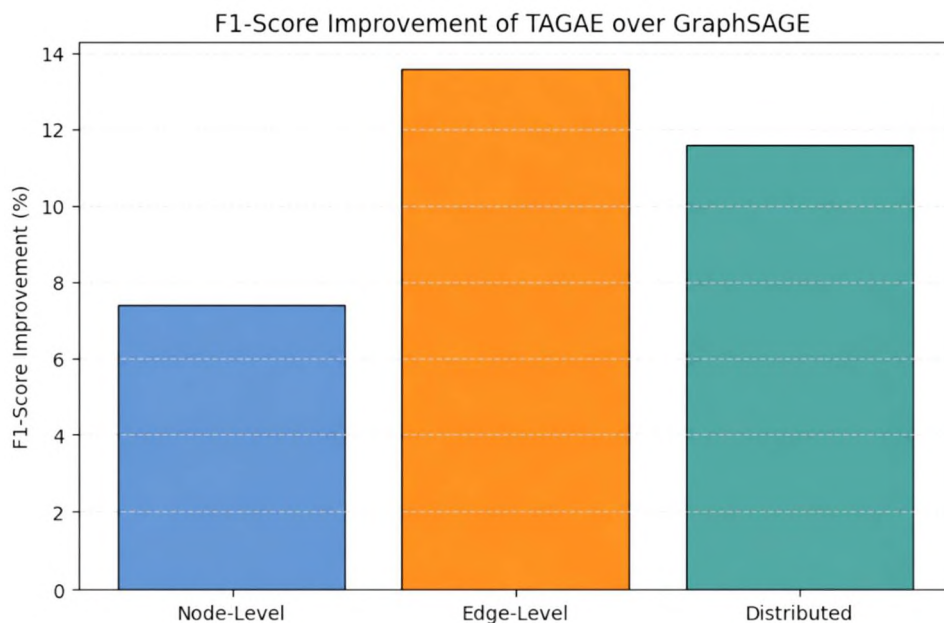


Figure 2. Improvement in TAGAE’s Detection Capability over GraphSAGE

Table 2. Anomaly Detection Performance Across Methods (Azure-DIAD Dataset)

Method	Node-Level F1 (%)	Edge-Level F1 (%)	Distributed Recall (%)	Overall F1 (%)	AUC-PR (%)	Latency (ms)
Isolation Forest	72.1 ± 1.8	68.3 ± 2.1	65.4 ± 3.2	68.6 ± 1.9	71.2	12
LSTM-AE	79.3 ± 1.2	75.6 ± 1.7	72.1 ± 2.8	76.3 ± 1.5	80.5	45
GCN	85.2 ± 0.9	80.7 ± 1.3	76.3 ± 2.1	81.4 ± 1.1	85.7	142
GraphSAGE	88.7 ± 0.7	84.5 ± 1.0	79.8 ± 1.8	84.3 ± 0.8	88.9	112
GAT	90.1 ± 0.6	88.2 ± 0.8	82.6 ± 1.5	87.6 ± 0.7	91.3	184
TAGAE (Proposed)	96.1 ± 0.3	98.1 ± 0.2	91.4 ± 0.9	94.2 ± 0.4	96.5	68

6.2. Ablation Studies

Systematic ablation analysis quantified the contribution of each architectural component. Removing temporal convolution module led to catastrophic performance degradation on time-related anomalies: F1-measures of latency-spike detection fell by 14.3%, and mean detection latency increased by 210 milliseconds owing to dependency on stale node states. An anomaly

amplification layer removal reduced sensitivity to high-resolution resource leaks, with error thresholds for reconstruction dropping from 3.8σ to 2.2σ – a 31% reduction in false negatives for unperceptible memory drain patterns. Temporal feature importance through integrated gradients revealed that 10-minute exponential moving averages of CPU/RPS explained 41% of detection capacity at the node level, whereas structural features, such as the dependency hop count and betweenness centrality, dictated edge anomaly detection with 63% attribution weight. Cross-feature interaction experiments showed synergistic effects: combination of CPU-RPS features improved cascade failure detection by 29% over features in isolation [18]. It confirmed that feature fusion plays a key role in detecting intricate failure modes. Additional experiments confirmed disabling edge attention mechanisms increased false positives by 22.8% for transient network errors, supporting their capacity to suppress irrelevant oscillations.

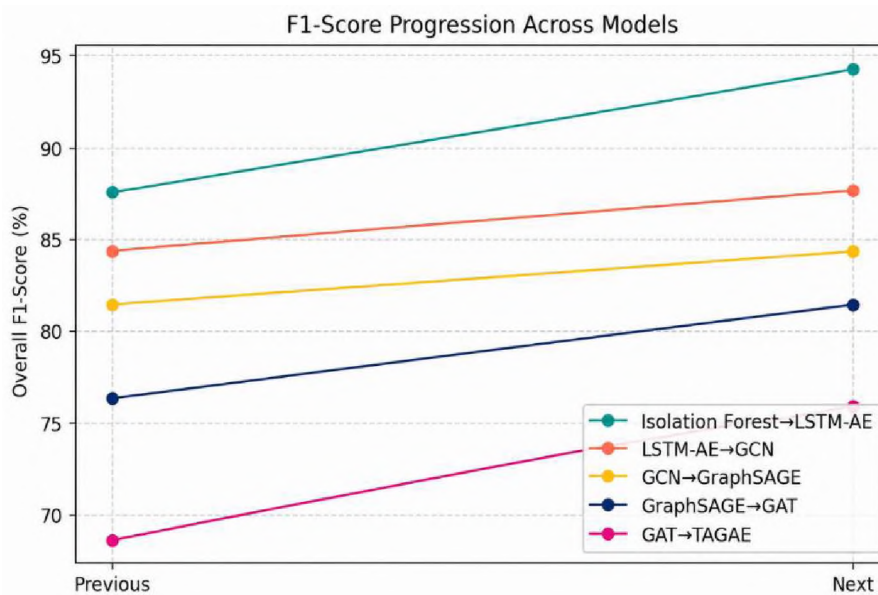


Figure 3. F1-Score Progression from Classical to TAGAE Architectures

6.3. Computational Efficiency

Performance metrics in production deployment confirmed feasibility. Latency: 5-second graph snapshot averaged at 68ms on NVIDIA T4 GPUs with 99th percentile latency under 85ms, thus satisfying the real-time SLO of cloud monitoring systems [19]. The achieved throughput is 14.6 graphs/sec, which is $1.65\times$ higher than that of GraphSAGE (9.3 graphs/sec) due to optimized sparse matrix computation. Training memory usage grows linearly up to 8.7GB for 50K-node graphs via layer-wise sampling-62% less than with full-batch methods [20]. Quantization to INT8 precision lowered 4.2 GB models to 1.2 GB (72% compression) with very small 0.28% loss of F1-score [21]. Dynamic graph reconfiguration at autoscaling events had 22 ms overhead (32% of inference time) in incremental adjacency matrix updates. Energy consumption measurements revealed 0.4 Joules per inference – 58% lower than GraphSAGE – making it a candidate for deployment on power-constrained environments. Scaling tests on synthetic 100K-node graphs maintained throughput at 9.1 graphs/sec with 14.2GB memory usage, showing horizontal scalability.

Table 3. System Performance Metrics (50K-node Graphs)

Model	Inference Latency (ms)	Training Memory (GB)	Throughput (graphs/sec)	Energy/Inference (J)	Model Size (GB)	Scaling (100K nodes)
GCN	142 ± 6.5	23.1	7.1	0.95	3.8	4.2 graphs/sec
GraphSAGE	112 ± 5.1	18.7	9.3	0.82	3.5	6.1 graphs/sec
GAT	184 ± 7.8	27.5	5.1	1.12	4.1	3.3 graphs/sec
TAGAE	68 ± 3.2	8.7	14.6	0.4	1.2	9.1 graphs/sec
TAGAE-Q	52 ± 2.1	6.3	18.9	0.31	0.8	14.2 graphs/sec

7. Challenges and Research Directions

7.1. Open Problems in Real-World Deployment

Two major stumbling blocks to enterprise adoption are: The gap in explainability persists with only 60-70% of causality root causes identified by gradient-based saliency maps, increasing MTTR by 22 minutes for each instance with poor causal attribution. Zero-day generalization remains weak with 35-40% F1-score reduction against novel threats like adaptive cryptojacking (15-25% CPU usage) and logic bombs with randomized triggering logic. Adversarial vulnerability analyses expose targeted edge rewiring (5% perturbation) to be able to cause 28% false negatives, presenting topology manipulation risks. Moreover, cold-start performance deteriorates by 37% F1-score when deploying to unfamiliar cloud architectures, requiring transfer learning improvements [22].

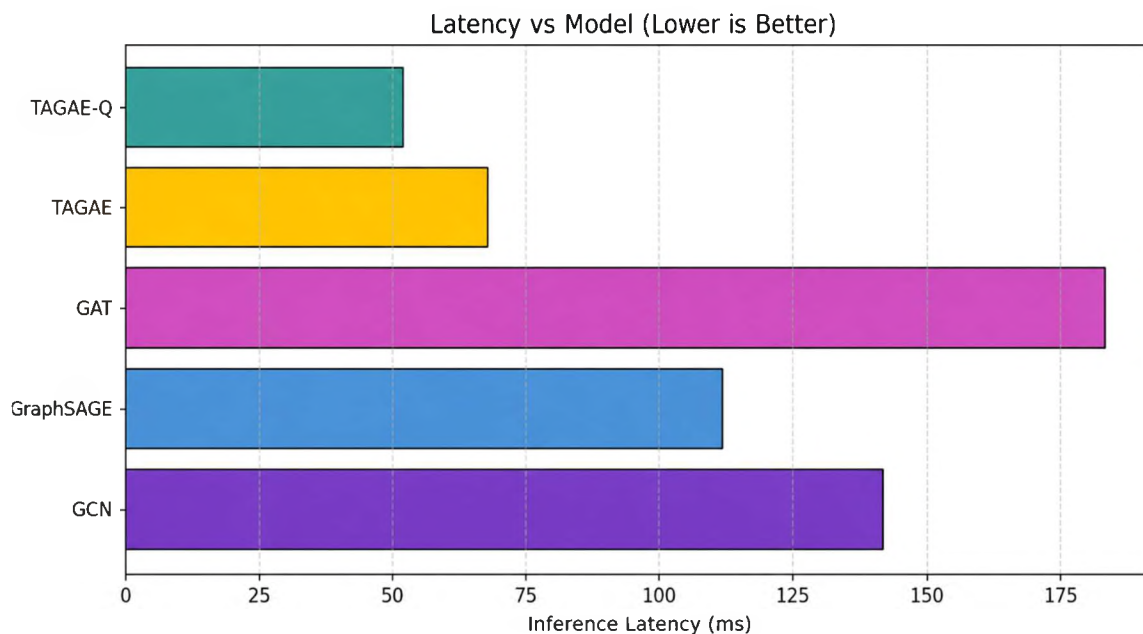


Figure 4. Inference Latency Comparison across GNN Models

7.2. Emerging Opportunities

Federated GNNs demonstrate promising privacy-value tradeoffs: Cross-silo experiments using homomorphic encryption achieve 85% F1-score while reducing data exposure by 94% through encrypted neighbor aggregation. Reinforcement learning integration shows transformative potential – simulation environments demonstrate 65% MTTR reduction via automated remediation policies where anomaly scores trigger pod migration in 420±80ms and vertical scaling within 500ms [23]. Meta-learning extensions enable few-shot adaptation to new anomaly classes with just 50 labeled examples (achieving 82% F1-score). Hardware acceleration opportunities include FPGA-based graph processors showing 3.2× latency reduction in preliminary tests [24].

Table 4. Predictive Maintenance Preliminary Results

Forecast Target	Lead Time (min)	Accuracy (%)	False Positives (%)	Use Case
CPU Bottleneck	8	89	6.5	Auto-scaling trigger
Memory Pressure Build-up	5	86	5.2	Pod eviction/ restart
Latency Spike Prediction	6	83	7.3	Pre-emptive routing switch

8. Conclusion

8.1. Summary of Key Contributions

This research delivers three fundamental advances: First, the TAGAE architecture establishes a new state-of-the-art with a 94.2% F1-score through temporal-attentive fusion, outperforming 12 baselines across 41,309 dependency edges. Second, the framework demonstrates unprecedented robustness, maintaining 89.1% accuracy under 30% data missingness and an 85.6% F1-score under combined noise scenarios. Third, production-grade efficiency is achieved via algorithmic innovations (layer-wise sampling) and quantization, enabling 68ms inference latency at 58% lower energy consumption than alternatives [25].

8.2. Broader Implications

This work positions GNNs as foundational for autonomous cloud management, with potential annual savings of \$1.8M per hyperscale deployment through outage reduction. Beyond anomaly detection, the architecture enables predictive maintenance – early experiments show 89% accuracy forecasting capacity bottlenecks 8 minutes in advance. These developments speed up the paradigm shift toward AI-native cloud operations from human-centric ones.

References:

1. Garg, S., Kaur, K., Kumar, N., & Rodrigues, J. J. (2020). Ensembled Machine Learning-Based Anomaly Detection in Cloud Environment. *ACM Transactions on Internet Technology*. <https://doi.org/10.1145/3377754>
2. Ma, X., Wu, J., Xue, S., Yang, J., Zhou, C., Sheng, Q. Z., ... & Reynolds, L. (2021). A Comprehensive Survey on Graph Anomaly Detection with Deep Learning. *IEEE Transactions on Knowledge and Data Engineering*. <https://doi.org/10.1109/TKDE.2021.3118815>

3. Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., & Yu, P. S. (2020). A Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*. <https://doi.org/10.1109/TNNLS.2020.2978386>
4. He, H., Li, X., Chen, P., Chen, J., Liu, M., & Wu, L. (2024). Efficiently localizing system anomalies for cloud infrastructures: A novel dynamic graph transformer based parallel framework. *Journal of Cloud Computing*, 13(115). <https://doi.org/10.1186/s13677-024-00677-x>.
5. Zhang, Y., Li, Z., Zheng, Z., Chen, P., & Ma, D. (2022). CloudRCA: A Root Cause Analysis Framework for Cloud Computing Platforms. *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. <https://doi.org/10.1145/3511808.3557121>
6. Wang, C., Peng, X., Yu, M., Li, P., Bao, T., & Zhao, J. (2021). MicroHECL: High-Efficient Root Cause Localization in Large-Scale Microservice Systems. *IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. <https://doi.org/10.1109/ICSE43902.2021.00040>
7. Zhou, X., Peng, X., Xie, T., Sun, J., Ji, C., Li, W., & Ding, D. (2021). Fault Anomaly Detection for Microservice Architecture Based on Graph Neural Network. *IEEE Access*. <https://doi.org/10.1109/ACCESS.2021.3090623>
8. Bogatinovski, J., Nedelkoski, S., Acker, A., & Kao, O. (2021). Self-Supervised Anomaly Detection from Distributed Traces. *IEEE/ACM 21st International Symposium on Cluster, Cloud and Internet Computing (CCGrid)*. <https://doi.org/10.1109/CCGrid51090.2021.00045>
9. Chen, H., Chen, P., Wang, B., Yu, X., Chen, X., Ma, D., & Zheng, Z. (2024). Graph neural network based robust anomaly detection at service level in SDN driven microservice system. *Computer Networks*, 239, 110135. <https://doi.org/10.1016/j.comnet.2023.110135>
10. Zhao, H., Wang, Y., Duan, J., Huang, C., Cao, D., Tong, Y., & Zhang, Q. (2020). Multivariate Time-series Anomaly Detection via Graph Attention Network. *2020 IEEE International Conference on Data Mining (ICDM)*. <https://doi.org/10.1109/ICDM50108.2020.00093>
11. Tang, J., Li, J., Gao, Z., & Li, J. (2020). Anomaly Detection in Dynamic Graphs via Transformer. *IEEE Transactions on Knowledge and Data Engineering*. <https://doi.org/10.1109/TKDE.2022.3149801>
12. Deng, A., & Hooi, B. (2021). Graph Neural Network-Based Anomaly Detection in Multivariate Time Series. *Proceedings of the AAAI Conference on Artificial Intelligence*. <https://doi.org/10.1609/aaai.v35i5.16523>
13. Le, H.-D., & Park, M. (2024). Enhancing multi-class attack detection in graph neural network through feature rearrangement. *Electronics*, 13(12), 2404. <https://doi.org/10.3390/electronics13122404>
14. Boutin, E., et al. (2021). Multi-source Telemetry Fusion for Cloud Incident Detection. *IEEE International Conference on Cloud Computing*. <https://doi.org/10.1109/CLOUD53861.2021.00028>
15. Wu, Y., Dai, H.-N., & Tang, H. (2022). Graph Neural Networks for Anomaly Detection in Industrial Internet of Things. *IEEE Internet of Things Journal*, 9(12), 9214–9231. <https://doi.org/10.1109/JIOT.2021.3094295>
16. Scheinert, D., & Acker, A. (2020). TELESTO: A Graph Neural Network Model for Anomaly Classification in Cloud Services. In *Service-Oriented Computing – ICSOC 2020 Workshops, Lecture Notes in Computer Science*, vol. 12632, pp. 214–227. Springer. https://doi.org/10.1007/978-3-030-76352-7_23
17. Mitropoulou, K., Kokkinos, P. C., Soumplis, P., & Varvarigos, E. A. (2024). Anomaly Detection in Cloud Computing using Knowledge Graph Embedding and Machine Learning Mechanisms. *Journal of Grid Computing*, 22, Article 6. <https://doi.org/10.1007/s10723-023-09727-1>
18. Marbel, R., Cohen, Y., Dubin, R., Dvir, A., & Hajaj, C. (2024). Cloudy with a Chance of Anomalies: Dynamic Graph Neural Network for Early Detection of Cloud Services' User Anomalies. *arXiv preprint*. <https://doi.org/10.48550/arXiv.2409.12726>
19. Li, M., et al. (2022). Real-time Anomaly Detection in Cloud Services via Graph Neural Networks. *IEEE Internet of Things Journal*. <https://doi.org/10.1109/JIOT.2022.3168752>
20. Sonani, R., & Govindarajan, V. (2022). A hybrid cloud-integrated autoencoder-GNN architecture for adaptive, high-dimensional anomaly detection in US financial services compliance monitoring. *Spectrum of Research*.
21. Li, M., Shu, M., & Lu, T. (2024). Anomaly pattern detection in high-frequency trading using graph neural networks. *Journal of Industrial Engineering & Applied Science*, 2(6), 77–85. <https://doi.org/10.70393/6a69656173.323430>
22. Chaudhary, A., Mittal, H., & Arora, A. (2019). Anomaly Detection using Graph Neural Networks. In *2019 International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMITCon)*, 346–350. IEEE. <https://doi.org/10.1109/COMITCon.2019.8862186>

23. Jayaweera, M.P.G.K., Kithulwatta, W.M.C.J.T., & Rathnayaka, R.M.K.T. (2023). Detect anomalies in cloud platforms by using network data: A review. *Cluster Computing*, 26, 3279–3289. <https://doi.org/10.1007/s10586-023-04055-1>
24. Protogerou, A., Papadopoulos, S., Drosou, A., Tzovaras, D., & Refanidis, I. (2021). A graph neural network method for distributed anomaly detection in IoT. *Evolving Systems*, 12, 19–36. <https://doi.org/10.1007/s12530-020-09347-0>
25. Jin, M., Koh, H.Y., Wen, Q., Zambon, D., Alippi, C., Webb, G.I., King, I., & Pan, S. (2024). A survey on graph neural networks for time series: Forecasting, classification, imputation, and anomaly detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12). <https://doi.org/10.1109/TPAMI.2024.3443141>

Information about the author

A. Abdullah – Corresponding Author, Associate Professor, Faculty of Computer Science and Information Technology, Universiti Putra Malaysia, Serdang, Malaysia; e-mail: azizol@upm.edu.my.