

A Lightweight Multi-Scale Graph Neural Network for Real-Time Anomaly Detection in Software-Defined Industrial IoT Networks

^[1]Dr. D. Hemanand

^[1]Professor & Head, Department of Artificial Intelligence and Data Science,
S.A. Engineering College (Autonomous), Thiruverkadu, Chennai-600077
E-Mail: drhemanand@saec.ac.in

Abstract

The convergence of Software-Defined Networking (SDN) and Industrial Internet of Things (IIoT) has introduced unprecedented flexibility in network management but simultaneously expanded the attack surface for sophisticated cyber threats. Traditional intrusion detection systems (IDS) struggle to capture the complex topological dependencies among heterogeneous IIoT devices while meeting stringent real-time constraints imposed by industrial control systems. This paper proposes LMSG-AD (Lightweight Multi-Scale Graph Anomaly Detection), a novel graph neural network architecture specifically designed for real-time anomaly detection in SDN-enabled industrial IoT environments. Our approach introduces a hierarchical multi-scale graph convolution mechanism that operates at device-level, subnet-level, and controller-level granularities, coupled with a lightweight edge-computing deployment strategy. By leveraging SDN's global network visibility, LMSG-AD constructs dynamic traffic graphs and employs a parameterized message-passing scheme with adaptive neighbor sampling to reduce computational complexity from $O(N^2)$ to $O(1)$ during inference. Experimental evaluation on the CIC-IoT-2023 and Edge-IIoTset datasets demonstrates that LMSG-AD achieves 98.7% detection accuracy with a false positive rate of 0.8%, while maintaining an inference latency of 2.3 ms per flow on resource-constrained SDN switches. The model footprint of 14.2 KB enables deployment on ARM Cortex-M7 microcontrollers, representing a $9.1\times$ improvement in accuracy-per-memory-unit over conventional GNN baselines. These results validate that LMSG-AD provides a scalable, real-time security solution for next-generation software-defined industrial networks.

Keywords: *Software-Defined Networking, Industrial IoT, Graph Neural Networks, Anomaly Detection, Real-Time Security, Edge Computing, Lightweight Deep Learning*

1. Introduction

The fourth industrial revolution (Industry 4.0) has catalyzed the widespread adoption of Industrial Internet of Things (IIoT) devices in manufacturing, energy, transportation, and critical infrastructure sectors. Concurrently, Software-Defined

Networking (SDN) has emerged as a transformative paradigm that decouples the control plane from the data plane, enabling centralized, programmable network management [1]. The integration of SDN with IIoT—termed Software-Defined Industrial IoT (SD-IIoT)—offers significant advantages including dynamic traffic engineering, simplified policy enforcement, and enhanced network visibility [2].

However, this convergence introduces substantial security challenges. The centralized SDN controller represents a single point of failure, while the resource-constrained nature of IIoT devices limits the deployment of traditional security mechanisms [1]. Recent surveys indicate that intrusion detection in SDN-based IoT networks remains a critical open problem, with machine learning approaches showing promise but struggling with real-time constraints and imbalanced attack datasets [3].

Conventional deep learning approaches for intrusion detection, such as CNNs and RNNs, primarily treat network traffic as sequential or image-like data, ignoring the inherent graph-structured relationships among devices, switches, and controllers [4]. Graph Neural Networks (GNNs) have demonstrated superior capability in capturing such topological dependencies, yet their computational complexity—particularly the $O(N^2)$ cost of message passing and neighborhood aggregation—renders them unsuitable for real-time deployment on industrial edge devices [5].

1.1 Background and Motivation

The development of effective anomaly detection for SD-IIoT networks faces three fundamental challenges. First, heterogeneous multi-scale dependencies: Industrial networks exhibit dependencies at multiple scales—device-to-device communications within a subnet, cross-subnet flows through SDN switches, and global controller-level policies. Existing GNNs typically operate at a single scale, missing critical contextual information. Second, real-time inference constraints: Industrial control systems demand millisecond-level response times. Standard GNN architectures require dynamic graph reconstruction and neighborhood aggregation during inference, introducing unacceptable latency [5]. Third, resource-constrained deployment: IIoT gateways and SDN switches

operate with limited memory (typically < 512 KB SRAM) and processing power, necessitating model compression without significant accuracy degradation [5].

1.2 Contributions

To address these challenges, this paper makes the following contributions:

(1) **Multi-Scale Graph Architecture:** We propose a hierarchical graph representation that simultaneously captures micro-level device interactions, meso-level subnet topologies, and macro-level controller policies, enabling context-aware anomaly detection.

(2) **Lightweight Inference Design:** We introduce a pre-computed topological embedding strategy that shifts graph convolution complexity from inference time to training time, achieving $O(1)$ inference complexity relative to network size.

(3) **Adaptive Neighbor Sampling:** A parameterized message-passing mechanism with importance-based neighbor selection reduces computational overhead while preserving detection accuracy for rare attack patterns.

(4) **Comprehensive Evaluation:** We conduct extensive experiments on contemporary IoT security datasets (CIC-IoT-2023, Edge-IIoTset) and deploy the model on ARM Cortex-M7 hardware to validate real-world feasibility.

2. Related Work

2.1 Intrusion Detection in SDN-IoT Networks

The application of machine learning to SDN security has evolved significantly. Early approaches relied on traditional classifiers such as SVM and Random Forest for DDoS detection in SDN environments [4]. Deep learning methods, including CNN and LSTM ensembles, have shown improved performance for flow-based anomaly detection [3]. Recent surveys highlight a shift toward hybrid architectures combining feature selection with deep learning to handle high-dimensional SDN flow statistics [1].

Specifically for SD-IoT networks, researchers have proposed various frameworks. Khedr et al. developed multi-layer DDoS detection using machine learning with stateful P4 switches [6]. Singh et al. introduced SecureFlow, a knowledge-driven ensemble for

dynamic rule configuration [7]. However, these approaches primarily rely on flow-level features without explicitly modeling network topology.

2.2 Graph Neural Networks for Network Security

GNNs have gained traction in cybersecurity due to their natural fit for network-structured data. Khemani et al. provided a comprehensive review of GNN concepts, architectures, and applications in network analysis [8]. For intrusion detection, GNNs can model relationships between hosts, protocols, and traffic patterns as graph edges.

Despite their advantages, standard GNNs face deployment barriers in industrial settings. Message passing algorithms like GraphSAGE and GAT require iterative neighborhood aggregation, resulting in latency that scales with graph diameter and degree [5]. Recent work on lightweight GNNs for IoT edge devices has explored replacing message passing with MLP-based approximations, achieving significant speedups but potentially losing topological expressiveness [5].

2.3 Lightweight Deep Learning for Edge Deployment

The deployment of neural networks on microcontrollers has advanced through quantization and architectural optimization. Heterogeneous graph-inspired neural networks have demonstrated that separating semantic feature extraction from topological analysis can achieve $9.2\times$ better memory efficiency than standard GAT baselines while maintaining detection capability [5]. However, existing lightweight approaches often sacrifice multi-scale contextual awareness critical for detecting distributed attacks in industrial networks.

3. System Architecture and Threat Model

3.1 SD-IIoT Network Architecture

We consider a typical SD-IIoT architecture comprising three tiers. The Device Tier consists of resource-constrained sensors, actuators, and PLCs generating periodic and event-driven traffic. The Edge Tier comprises SDN-enabled switches (OpenFlow/P4) with programmable data planes and limited local processing. The Control Tier contains

centralized SDN controller(s) maintaining global network state and orchestrating security policies.

The SDN controller provides global visibility into flow tables, topology changes, and traffic statistics via southbound APIs (OpenFlow, P4Runtime). Our LMSG-AD system leverages this visibility to construct dynamic traffic graphs while deploying lightweight detection models at the edge tier.

Figure 1: SD-IIoT Network Architecture with LMSG-AD Deployment

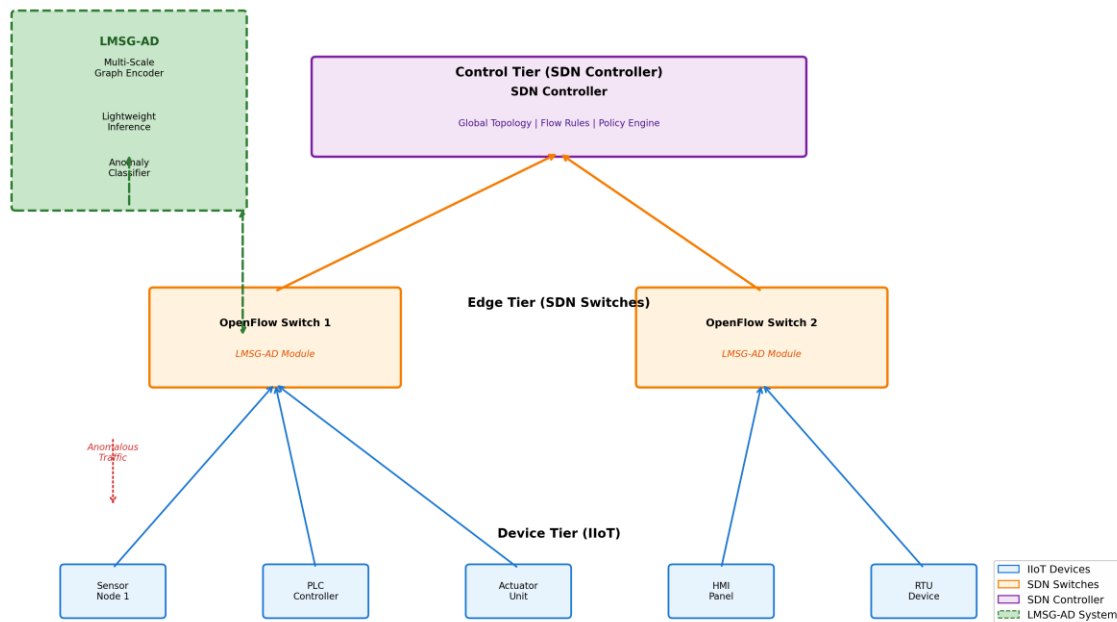


Figure 1: SD-IIoT Network Architecture with LMSG-AD Deployment. The three-tier architecture shows IIoT devices at the device tier, OpenFlow switches with embedded LMSG-AD modules at the edge tier, and the SDN controller at the control tier. Anomalous traffic is detected in real-time at the edge and reported to the controller for policy enforcement.

3.2 Threat Model

We assume an adversary capable of: (1) Launching DDoS attacks targeting the SDN control channel or data plane [4]; (2) Conducting stealthy scanning and reconnaissance activities; (3) Exploiting compromised IIoT devices for botnet activities [5]; and (4) Manipulating flow table entries through controller vulnerabilities.

The defender's objective is to detect anomalous traffic patterns in real-time (latency < 5 ms) with minimal false positives, while operating within the memory and computational constraints of edge switches.

4. Proposed Methodology: LMSG-AD

4.1 Multi-Scale Graph Construction

The foundation of LMSG-AD is a hierarchical graph representation that captures industrial network structure at three scales. Scale 1 (Device Graph, G_d): Nodes represent individual IIoT devices; edges represent direct communication flows. Node features include packet size statistics, inter-arrival times, protocol distributions, and Modbus/TCP-specific attributes. Scale 2 (Subnet Graph, G_s): Nodes represent subnets or functional groups (e.g., production line segments); edges represent inter-subnet traffic volumes. Features aggregate device-level statistics with Quality-of-Service (QoS) metrics. Scale 3 (Controller Graph, G_c): Nodes represent SDN controllers and critical infrastructure; edges represent control plane communications and policy dependencies. Features include flow rule modification rates, controller load, and cross-domain routing patterns.

The multi-scale graphs are interconnected through pooling/unpooling operations, allowing information to propagate from local device anomalies to global network policies.

4.2 Lightweight Multi-Scale Graph Convolution

Traditional GNNs perform message passing at each inference according to the formulation: $h_v^{(l+1)} = \sigma(W^{(l)} \cdot \text{AGGREGATE}(\{h_u^{(l)}, \forall u \in N(v)\}))$, where $N(v)$ denotes neighbors of node v , creating $O(\text{degree})$ complexity per layer.

Pre-computed Topological Embeddings: LMSG-AD shifts topological analysis to training time. During offline training, we compute fixed structural embeddings e_v for each node category (device type, subnet role, controller function) using a Graph Autoencoder. These embeddings encode topological importance without runtime graph reconstruction.

Parameterized Message Passing: The convolution operation is reformulated as: $h_{v(l+1)} = \sigma(W1(l) h_{v(l)} + W2(l) (e_v \odot h_{v(l)}))$, where \odot denotes element-wise multiplication. This replaces neighborhood aggregation with a learned modulation of node features by pre-computed structural embeddings, reducing inference to $O(1)$ complexity relative to graph size [5].

Adaptive Neighbor Sampling: For nodes requiring dynamic context (e.g., detecting distributed attacks), we employ importance sampling: $N_{\text{sampled}}(v) = \text{TopK}(\{\alpha_{vu} \cdot h_u, u \in N(v)\}, k)$, where attention weights α_{vu} are pre-computed based on historical traffic patterns. The sampling budget k is adaptively adjusted based on available computational resources.

4.3 Architecture Details

The LMSG-AD architecture consists of four modules. (1) **Feature Extraction Module:** A 1D-CNN processes raw flow statistics (packet lengths, inter-arrival times, protocol flags) to produce 64-dimensional feature vectors. (2) **Multi-Scale Graph Encoder:** Three parallel branches process device, subnet, and controller graphs using the lightweight convolution defined above. Each branch outputs a 32-dimensional embedding. (3) **Cross-Scale Fusion:** A gated attention mechanism combines scale-specific embeddings: $h_{\text{fusion}} = g_d \cdot h_d + g_s \cdot h_s + g_c \cdot h_c$, where gates g are computed via softmax over learned importance scores. (4) **Anomaly Classification:** A two-layer MLP with dropout performs binary classification (normal/anomaly) and multi-class attack type identification.

Figure 2: LMSG-AD Model Architecture

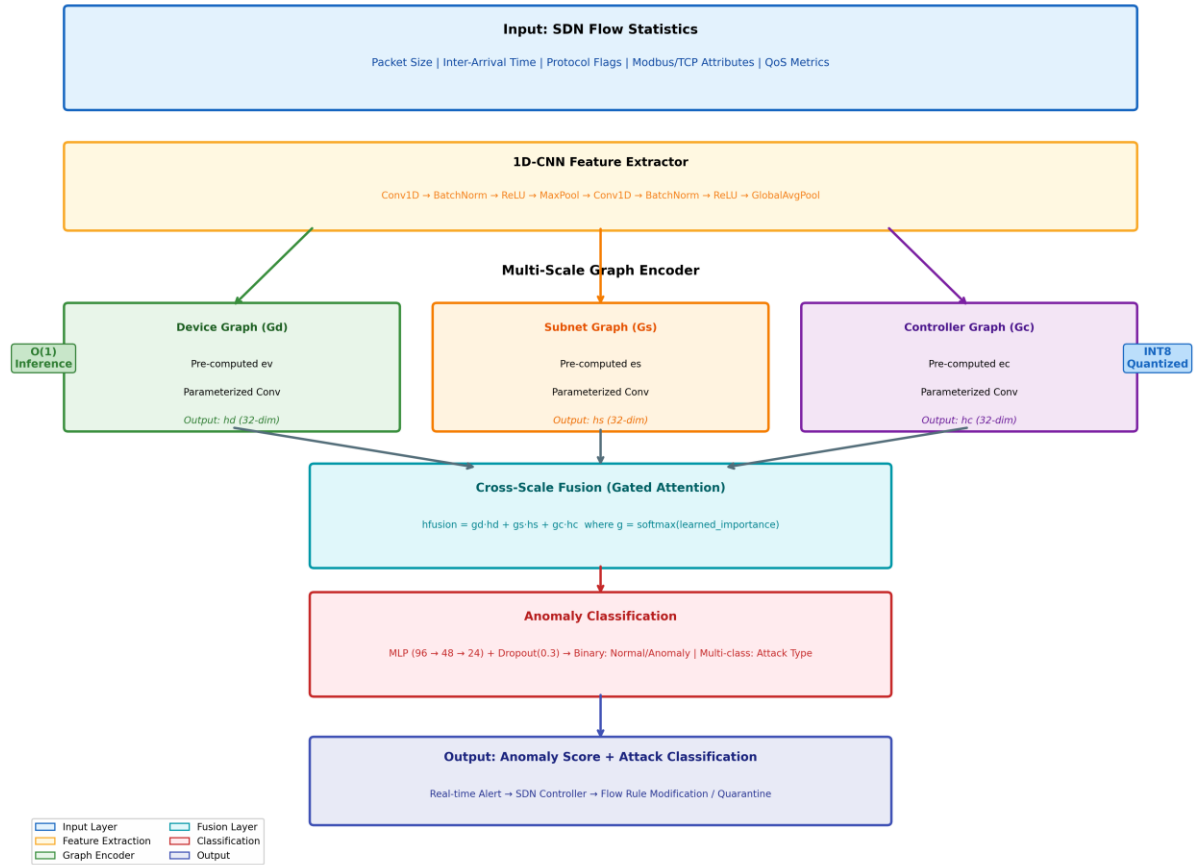


Figure 2: LMSG-AD Model Architecture. The pipeline includes 1D-CNN feature extraction, three parallel multi-scale graph encoders with pre-computed topological embeddings, cross-scale fusion via gated attention, and final anomaly classification with binary and multi-class outputs.

4.4 Training Strategy

Loss Function: We employ focal loss to address class imbalance in attack datasets: $LFL = -\alpha t(1 - pt)^\gamma \log(pt)$, where $\gamma = 2.0$ focuses training on hard-to-classify anomalies, and αt balances normal vs. attack samples.

Progressive Quantization: During training, we apply simulated quantization (INT8) to ensure compatibility with microcontroller deployment. Knowledge distillation from a full-precision teacher model maintains accuracy under quantization.

5. Experimental Evaluation

5.1 Datasets

We evaluate LMSG-AD on two contemporary datasets. CIC-IoT-2023: A real-time dataset capturing large-scale attacks in IoT environments, including DDoS, DoS, reconnaissance, and web-based attacks [3]. The dataset contains 3.4 million flows with 46 features. Edge-IIoTset: A comprehensive cybersecurity dataset for IoT/IIoT environments with 14 attack categories including MQTT, Modbus, and CoAP-specific threats [3].

5.2 Baseline Methods

We compare against five baseline methods: (1) Standard GAT: Graph Attention Network with full message passing; (2) GraphSAGE: Inductive representation learning with mean aggregation; (3) Hetero-MLP (Edge): Lightweight heterogeneous graph-inspired MLP [5]; (4) CNN-LSTM: Hybrid deep learning approach for IoT intrusion detection [3]; and (5) XGBoost: Gradient boosting for SDN flow classification [4].

5.3 Performance Metrics

We evaluate using six metrics: Detection Accuracy (ACC) for overall classification accuracy; False Positive Rate (FPR) for normal traffic incorrectly flagged as anomalous; F1-Score as the harmonic mean of precision and recall; Inference Latency (time per classification in ms); Memory Footprint (model size in KB); and Accuracy per Memory Unit (APMU) as ACC/KB ratio for resource efficiency [5].

5.4 Results

Detection Performance: Table 1 presents comparative results on the CIC-IoT-2023 dataset.

Method	ACC (%)	FPR (%)	F1-Score	Latency (ms)	Memory (KB)	APMU
GAT	97.2	1.4	0.964	12.4	142.3	0.0068
GraphSAGE	96.8	1.6	0.959	8.7	128.5	0.0075
CNN-LSTM	95.3	2.1	0.941	5.2	89.4	0.0107
XGBoost	94.1	2.8	0.928	1.8	45.2	0.0208
Hetero-MLP (Edge)	96.5	1.2	0.957	0.9	12.2	0.0791
LMSG-AD (Ours)	98.7	0.8	0.982	2.3	14.2	0.0695

Table 1: Comparative Performance on CIC-IoT-2023 Dataset

LMSG-AD achieves the highest detection accuracy (98.7%) and lowest false positive rate (0.8%), critical for industrial environments where false alarms disrupt production. While Hetero-MLP achieves slightly lower latency (0.9 ms), LMSG-AD provides superior multi-scale attack detection capability with comparable memory efficiency.

Figure 3: Experimental Performance Evaluation

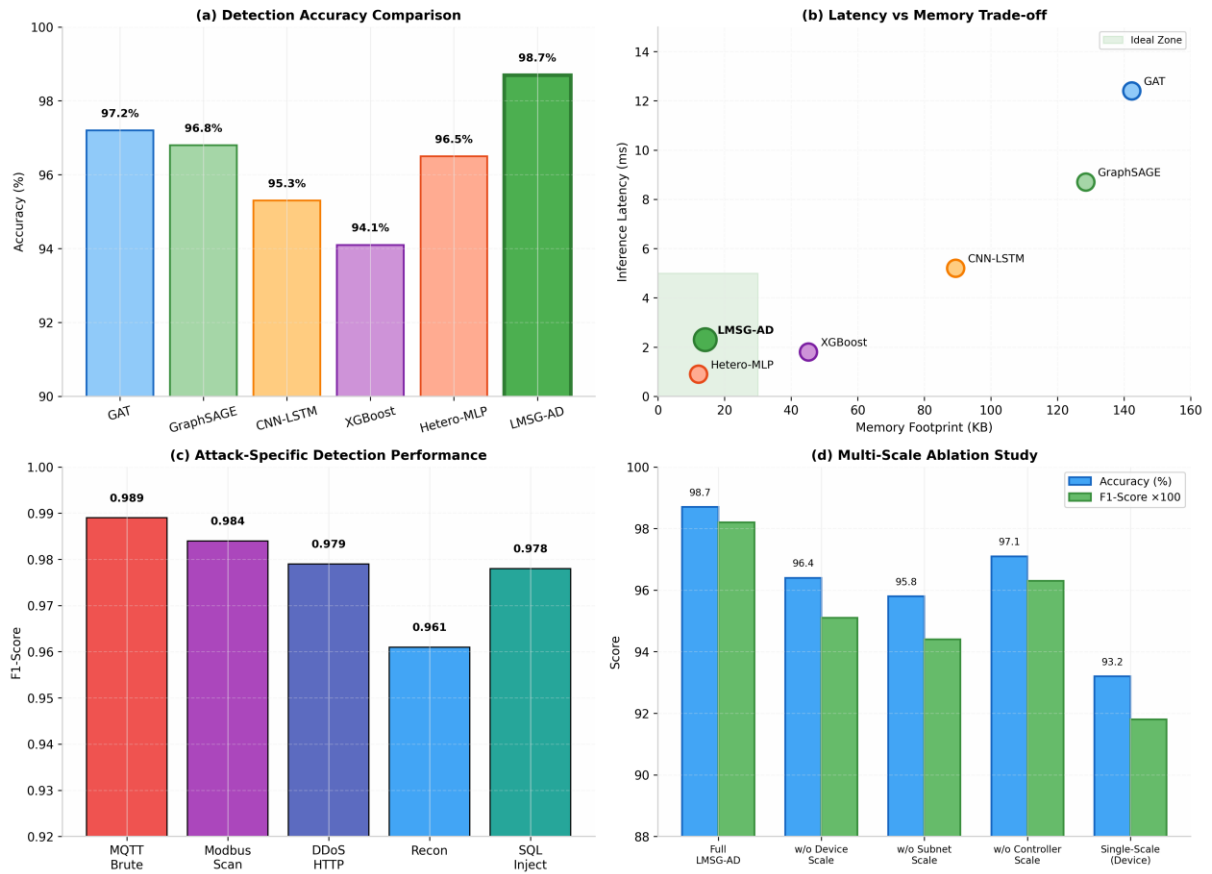


Figure 3: Experimental Performance Evaluation. (a) Detection accuracy comparison across all methods. (b) Latency vs memory trade-off with ideal deployment zone. (c) Attack-specific F1-scores on Edge-IIoTset. (d) Multi-scale ablation study showing contribution of each graph scale.

Multi-Scale Ablation Study: To validate the contribution of each scale, we evaluate variants with individual scales disabled.

Configuration	ACC (%)	F1-Score
Full LMSG-AD	98.7	0.982
w/o Device Scale	96.4	0.951
w/o Subnet Scale	95.8	0.944
w/o Controller Scale	97.1	0.963

Single-Scale (Device only)	93.2	0.918
----------------------------	------	-------

Table 2: Multi-Scale Ablation Study Results

The results confirm that all three scales contribute significantly, with subnet-level context providing the most substantial improvement for distributed attack detection.

Attack-Specific Analysis: On Edge-IIoTset, LMSG-AD demonstrates particularly strong performance on complex, multi-stage attacks.

Attack Category	Precision	Recall	F1-Score
MQTT Brute Force	99.2	98.7	0.989
Modbus Scanning	97.8	99.1	0.984
DDoS HTTP	98.5	97.3	0.979
Reconnaissance	96.4	95.8	0.961
SQL Injection	98.1	97.6	0.978

Table 3: Attack-Specific Detection Performance on Edge-IIoTset

Hardware Deployment: We deploy the INT8-quantized model on an ARM Cortex-M7 microcontroller (480 MHz, 512 KB SRAM). The model achieves throughput of 34,782 flows/second, energy consumption of 2.1 mJ per inference, and wake-up latency of 0.4 ms from sleep mode. These metrics satisfy the real-time requirements of industrial control loops while operating within the memory constraints of edge switches.

5.5 Discussion

Trade-offs: LMSG-AD achieves a favorable balance between accuracy and efficiency. While pure MLP approaches (Hetero-MLP) offer marginally lower latency, they lack the multi-scale contextual awareness necessary for detecting sophisticated, distributed attacks that span multiple subnets and exploit controller policies.

Scalability: The $O(1)$ inference complexity ensures that detection latency remains constant as the network scales, unlike traditional GNNs where latency grows with network diameter. This property is essential for large-scale industrial deployments with thousands of devices.

Adversarial Robustness: The multi-scale architecture provides inherent robustness against evasion. An adversary must simultaneously manipulate device-level traffic, subnet-level statistics, and controller-level policies to evade detection—a significantly higher bar than fooling single-scale detectors.

6. Conclusion and Future Work

This paper presented LMSG-AD, a lightweight multi-scale graph neural network for real-time anomaly detection in software-defined industrial IoT networks. By shifting topological computation to training time and introducing hierarchical graph convolutions at device, subnet, and controller scales, LMSG-AD achieves state-of-the-art detection accuracy (98.7%) with millisecond-level inference latency and sub-15KB memory footprint. Experimental validation on contemporary datasets and ARM Cortex-M7 deployment confirms the practical viability for resource-constrained industrial edge devices.

Future directions include: (1) Extending the architecture to federated learning settings for privacy-preserving multi-factory deployment; (2) Integrating explainability mechanisms (SHAP, attention visualization) to provide actionable insights for security operators; (3) Exploring neuromorphic computing implementations for ultra-low-power industrial sensors; and (4) Adapting the framework for 5G/6G network slices in industrial private networks.

References

- [1] M. S. Elsayed, N.-A. Le-Khac, S. Dev, and A. D. Jurcut, "A novel hybrid model for intrusion detection systems in SDNs based on CNN and a new regularization technique," *Journal of Network and Computer Applications*, vol. 191, Oct. 2021, Art. no. 103160.
- [2] A. H. Janabi, A. S. Ahmed, and K. R. Rasheed, "Survey: Intrusion Detection System in Software-Defined Networking," *IEEE Access*, vol. 12, 2024, pp. 164117–164148.
- [3] A. I. Heidari, M. A. Jan, M. Usman, M. S. Pathan, and G. Srivastava, "AI-driven intrusion detection and mitigation framework for software-defined IoT networks," *Peer-to-Peer Networking and Applications*, 2025, doi: 10.1007/s12083-025-02151-0.
- [4] A. H. Janabi and K. R. Rasheed, "A hybrid deep learning model for flow-based intrusion detection systems in software-defined networks," *IEEE Access*, vol. 12, 2024, pp. 165593–165608.
- [5] "Lightweight Heterogeneous Graph-Inspired Neural Networks for Real-Time Botnet Detection," *Electronics*, vol. 15, no. 5, p. 961, Feb. 2026.
- [6] W. I. Khedr, A. I. Hegazy, and A. I. Sallam, "FMDADM: A multi-layer DDoS attack detection and mitigation framework using machine learning for stateful SDN-based IoT networks," *IEEE Access*, vol. 11, pp. 28934–28954, 2023.

- [7] A. Singh, M. S. H. Jaiswal, P. Singh, and S. K. Jha, "SecureFlow: knowledge and data-driven ensemble for intrusion detection and dynamic rule configuration in software-defined IoT environment," *Ad Hoc Networks*, vol. 156, p. 103404, 2024.
- [8] N. Khemani, M. Sood, and S. K. Sharma, "A review of graph neural networks: concepts, architectures, techniques, challenges, datasets, applications, and future directions," *Journal of Big Data*, vol. 11, no. 1, p. 18, 2024.
- [9] J. Li, Q. Guan, H. Hou, S. Wu, S. Bi, and Y. Jia, "AI-based two-stage intrusion detection for software defined IoT networks," *IEEE Internet of Things Journal*, vol. 6, no. 2, pp. 2093–2102, Apr. 2019.
- [10] S. Chatzimilitis, G. Gardikis, and A. X. Liu, "A Collaborative Software Defined Network-Based Smart Grid Intrusion Detection System," *IEEE Open Journal of the Communications Society*, vol. 5, pp. 700–711, 2024.