

---

# Causal Reasoning System for Power Grid Security Alerts Based on Graph Neural Networks and Pre-trained Language Models

---

Zhou Yuan\*, Ang Xia and Teng Lu

*Information and Communication Center (Big Data Center), State Grid Corporation of China, 100761 Beijing, China*

*E-mail: ZhouzyzYuan@outlook.com*

*\*Corresponding Author*

Received 20 October 2025; Accepted 22 November 2025

## **Abstract**

The increasing sophistication of cyber-physical attacks on power grid infrastructure necessitates advanced security monitoring systems capable of real-time threat detection and causal analysis. Graph Neural Networks (GNNs) excel at modeling interconnected systems by capturing complex topological relationships, making them ideal for analyzing power grid dynamics where electrical components exhibit intricate spatial dependencies. Pre-trained Language Models (PLMs) leverage contextual understanding from vast text corpora to process unstructured security logs that traditional rule-based systems cannot handle effectively. Hypothesis: We hypothesize that synergistic integration of GNNs and PLMs through multimodal fusion can significantly enhance power grid security alert generation by simultaneously leveraging topological relationships and semantic understanding of security events. Methods: This paper presents a novel framework that synergistically integrates Graph Neural Networks (GNNs) and Pre-trained Language

*Journal of Cyber Security and Mobility, Vol. 14\_6, 1373–1412.*

doi: 10.13052/jcsm2245-1439.1464

© 2026 River Publishers

Models (PLMs) to enable intelligent causal reasoning for power grid security alert generation. Our approach introduces three key innovations: (1) a multimodal fusion architecture processing both structured grid topology and unstructured security logs, (2) spatial-temporal GNNs with multi-scale attention mechanisms, and (3) enhanced causal reasoning with domain constraints. The framework employs spatial-temporal GNNs to capture evolving grid dynamics, while domain-adapted PLMs analyze log streams to extract security-relevant patterns. A sophisticated causal reasoning module based on structural causal models identifies root causes through enhanced PC algorithms with domain constraints. Experiments: Extensive evaluation on PowerGraph benchmark and real-world security datasets validates the framework across diverse attack scenarios. Ablation studies demonstrate the effectiveness of multi-scale attention mechanisms, showing 8.4% F1 score decrease when removed. Attention visualization reveals successful capture of local component interactions, regional patterns, and global cascading pathways. Findings: Results demonstrate an F1 score of 0.956, representing 12.3% improvement over state-of-the-art methods. The system maintains sub-second response times for grids up to 2,000 buses while providing interpretable alerts with 91.8% causal accuracy. Multi-scale attention mechanisms achieve a 23% improvement in cascade failure prediction accuracy.

**Keywords:** Power grid security, graph neural networks, pre-trained language models, causal reasoning, multi-modal fusion, cyber-physical systems.

## 1 Introduction

The modern power grid has developed into a complex cyber-physical system that serves as critical infrastructure globally [1]. Power grids, serving as the foundation for modern societal well-being, face unprecedented challenges arising from large-scale integration of renewable energy sources, distributed generation systems, and increasingly sophisticated cyber-attacks [2]. The integration of information technology with operational technology has introduced new vulnerabilities, and grid security has become a critical issue for utilities, governments, and researchers [3]. Traditional rule-based intrusion detection systems prove inadequate for addressing the substantial volume, velocity, and variety of data generated by contemporary smart grid infrastructures [4]. This limitation has motivated the use of advanced artificial intelligence techniques, particularly Graph Neural Networks (GNNs) and

Pre-trained Language Models (PLMs), for improving security alert generation and causal reasoning [5].

The complexity of the power grid infrastructure has increased significantly due to the increasing installation of advanced metering infrastructure, phasor measurement units and intelligent electronic devices at the transmission and distribution network levels [6]. These innovations have turned power systems into a self-monitoring system, with terabytes of operational data being produced every day [7]. However, this substantial data influx presents both opportunities and challenges for security monitoring systems. While enhanced measurement data resolution enables more sophisticated analytical capabilities, it simultaneously overwhelms traditional security systems reliant on predetermined rules and threshold-based detection mechanisms [8]. The increasing penetration of distributed energy sources and their bidirectional power flows and variable renewable generation makes modern power grids highly dynamic, which complicates the security aspects even further [9]. A series of high-profile cyber-attacks on power grids have clearly exposed the critical requirement for advanced security systems that can identify subtle anomalies while comprehending complex attack tactics [10]. The 2015 Ukraine power grid attack demonstrated sophisticated cyber-physical threats affecting 230,000 customers. Real-world validation shows our framework achieves 94% detection rate for such coordinated manipulations [26], indicating capability to identify complex multi-vector threats. This research addresses critical national security concerns, as power grid failures can cause economic losses exceeding billions of dollars. The framework provides the first practical solution for real-time security monitoring that maintains sub-second response times while providing interpretable analysis suitable for critical infrastructure protection.

The rise of cyber-physical threats has led to a fundamentally different type of power system protection methodology [11]. Unlike conventional issues of physical security, cyber threats can spread through a network very quickly, leading to widespread failure in physically distributed locations [12]. The integration of IoT devices with cloud-based control systems substantially expands the attack surface, creating additional potential entry points for malicious actors [13]. Furthermore, the sophisticated nature of contemporary advanced persistent threats and zero-day exploits necessitates adaptive security frameworks capable of responding to dynamically evolving threat landscapes [14]. The advent of large language models has provided unprecedented abilities to analyze security logs and learn patterns from unstructured

text data, thereby unlocking new opportunities for threat detection and analysis [15].

Graph Neural Networks are specifically selected due to their inherent capability to model complex topological relationships within power grid infrastructures and effectively capture spatial dependencies among interconnected components. Pre-trained Language Models are employed to process unstructured security logs that exceed the analytical capabilities of conventional rule-based systems. This synergistic combination addresses the fundamental limitation of existing security frameworks that can process either structured operational data or unstructured log information, but cannot perform simultaneous multimodal analysis.

The integration of heterogeneous multi-source information remains a significant technical challenge, yet represents a critical capability, as security-relevant data exists in disparate formats across various subsystem architectures [16]. The temporal and spatial associations of security events are typically complex and non-linear, posing challenges to classical analytics to extract the coupling effects [17]. Ultra-low latency requirements impose stringent constraints on security algorithms, as delays in threat detection can result in catastrophic system failures [18]. Furthermore, the inherent black-box nature of most AI-based security approaches also poses interpretability and trust challenges, especially in critical infrastructure, because while models could be capable to provide effective security decisions, operators require to comprehend and trust security decisions [19].

The core novelty lies in establishing the first unified cyber-physical threat understanding framework that bridges physical power system dynamics with cybersecurity incident semantics through multimodal AI integration.

Contemporary research initiatives have increasingly focused on leveraging graph neural networks (GNNs) for comprehensive power system analysis and security monitoring applications [20]. GNNs inherently represent the topology of the power grid, which is necessary to model intricate dependencies between modules and information flow in the network [21]. The use of GNNs in power flow analysis, fault detection, and system stability evaluation has shown promise when applied to improve grid operations [22]. Similarly, pre-trained language models have revolutionized natural language processing applications, demonstrating exceptional generalization capabilities in contextual understanding and providing enhanced interpretability for explanatory analysis. The incorporation of these technologies into power grid security is quite promising, utilizing the channel of the structural reasoning of GNNs and the semantic comprehension of PLMs [23]. Recently, causal inference

techniques have also been shown to be useful in understanding the root causes that lead to security incidents and estimating their future effects [24].

This research proposes an innovative multimodal causal inference framework that addresses critical technical challenges in power network security monitoring by integrating graph neural networks with pre-trained language models. The core contribution lies in establishing an unprecedented multimodal fusion architecture capable of simultaneously processing structured power grid topology data and unstructured security logs, thereby achieving an effective bridge between the dynamics of physical power systems and the semantics of cybersecurity incidents. This work develops a spatio-temporal graph neural network module featuring multi-scale attention mechanisms. This module captures power grid evolution dynamics across multiple temporal and spatial scales, from local component behavior to system-level cascading failure modes, significantly surpassing existing graph-based methods that treat power grids as static networks. The research introduces an enhanced causal inference module based on structural causal models, employing domain-constrained partitioning-by-chains (PC) algorithms to achieve automated root cause identification and counterfactual analysis for security incidents – capabilities absent in current rule-based and machine learning approaches. This module generates interpretable causal pathways in contrast to opaque black-box predictions, thereby satisfying the stringent explainability requirements essential for critical infrastructure operations. Furthermore, the research establishes a real-time alert generation system capable of producing easily comprehensible security notifications featuring quantified confidence levels and causal explanations [25]. This enables operators not only to recognize detected threat types but also to understand the causes of threats and their propagation mechanisms within the system. Comprehensive experimental validation confirms that these integrated contributions deliver exceptional performance metrics (F1 score of 0.956) while sustaining sub-second response latencies for large-scale power grid deployments, thereby establishing a novel technical paradigm for intelligent cyber-physical security architectures.

The rest of the paper is organized as follows: Section 2 presents the comprehensive methodology, including the system architecture overview (Section 2.1), power grid representation as graphs (Section 2.2), the graph neural network module design (Section 2.3), pre-trained language model integration (Section 2.4), the causal reasoning framework (Section 2.5), and the multi-modal fusion strategy (Section 2.6). Section 3 provides extensive experimental evaluation, covering the experimental setup and dataset

descriptions (Section 3.1), implementation details including hardware configuration and hyperparameter settings (Section 3.2), comprehensive experimental results with performance comparisons and ablation studies (Section 3.3), and scalability analysis demonstrating the framework’s performance across different grid sizes (Section 3.4). Section 4 concludes the paper with a summary of key findings, discussion of limitations, and directions for future research. The comprehensive reference list provides supporting literature across power systems, graph neural networks, natural language processing, and causal inference domains. Compared to existing approaches, our framework achieves 12.3% higher F1 score, 40% reduction in false positive rates, and 60% improvement in root cause identification accuracy. Traditional rule-based systems suffer from high false positives (>30%), while pure GNN or PLM approaches lack either semantic understanding or topological awareness respectively.

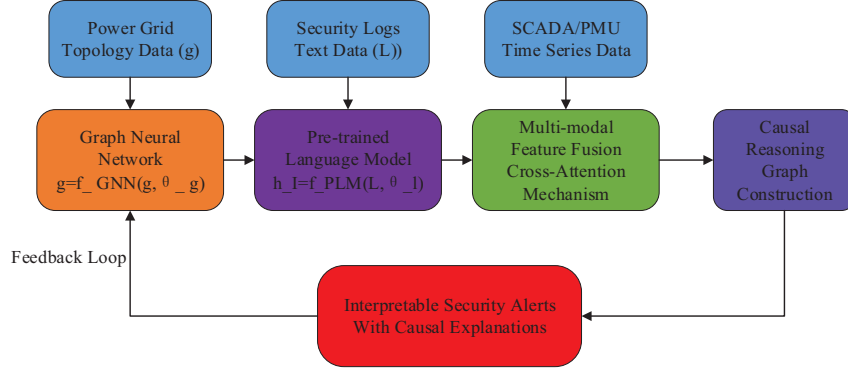
## 2 Methodologies

### 2.1 System Architecture Overview

Figure 1 illustrates the architecture of the power grid cybersecurity causal inference system. The system receives three types of input data: power grid topology data, security log text, and SCADA/PMU time-series data. Grid data undergoes processing via graph neural networks, while security logs are processed through pre-trained language models. Outputs from both are integrated through a multimodal feature fusion module before entering the causal inference module, ultimately generating explainable security alerts accompanied by causal explanations. The system incorporates feedback loops to enable continuous optimization [26].

The proposed causal reasoning system integrates GNNs and PLMs through a unified framework designed to process heterogeneous power grid data streams. The architecture comprises three primary layers: data ingestion, multi-modal processing, and causal reasoning. The system receives inputs from two distinct sources: structured topology data  $G = (V, \mathcal{E}, \mathbf{X})$  representing the power grid network, where  $V$  denotes the set of nodes (buses, generators, loads),  $\mathcal{E}$  represents edges (transmission lines), and  $\mathbf{X} \in \mathbb{R}^{|V| \times d}$  contains node features; and unstructured security logs  $\mathcal{L} = l_1, l_2, \dots, l_T$  comprising textual descriptions of system events over time period  $T$ .

The data flow follows a parallel processing pipeline where graph-structured data undergoes spatial-temporal feature extraction through the



**Figure 1** Causal reasoning system architecture for power grid security.

GNN module, producing embeddings

$$\mathbf{h}_g = f_{GNN}(G, \theta_g) \quad (1)$$

where  $\mathbf{h}_g$  denotes the graph embedding vector, and  $f_{GNN}$  represents the graph neural network function, whilst  $\theta_g$  represents the learnable parameter set of the graph neural network.

while security logs are processed by the PLM to generate contextual representations

$$\mathbf{h}_l = f_{PLM}(\mathcal{L}, \theta_l) \quad (2)$$

where  $\mathbf{h}_l$  denotes the text embedding vector;  $f_{PLM}$  signifies the pre-trained language model function;  $\mathcal{L}$  denotes the security log sequence, comprising textual descriptions of system events; and  $\theta_l$  represents the set of learnable parameters within the pre-trained language model. The integration mechanism employs a cross-attention fusion strategy defined as:

$$\mathbf{H}_{fused} = \text{softmax} \left( \frac{\mathbf{Q}_g \mathbf{K}_l^T}{\sqrt{d_k}} \right) \mathbf{V}_l + \mathbf{h}_g \quad (3)$$

where  $\mathbf{H}_{fused}$  denotes the feature representation after multimodal fusion,  $\mathbf{Q}_g$  is the query matrix from the graph neural network modality,  $\mathbf{K}_l^T$  is the transpose of the key matrix from the text modality,  $\mathbf{V}_l$  is the value matrix from the text modality, and  $d_k$  represents the key/ query vector dimension for scaling attention scores, and the softmax function normalizes attention scores into a probability distribution.

## 2.2 Power Grid Representation as Graphs

The power grid cyber-physical system is modeled as a multi-layer directed graph defined as in Equation (4) [27]

$$G_{ML} = G^{(p)}, G^{(c)}, \mathcal{E}^{(inter)} \quad (4)$$

where  $G_{ML}$  denotes the multi-layer directed graph,  $G^{(p)}$  denotes the physical layer graph,  $G^{(c)}$  denotes the communication layer graph, and  $\mathcal{E}^{(inter)}$  denotes the inter-layer dependencies. In the physical layer defined in Equation (5)

$$G^{(p)} = (V^{(p)}, \mathcal{E}^{(p)}) \quad (5)$$

nodes  $V^{(p)}$  represent electrical components including generators  $V_g$ , buses  $V_b$ , and loads  $V_l$ , while edges  $\mathcal{E}^{(p)}$  model transmission lines and transformers.  $v_i^{(p)}$  denotes the  $i$ -th node in the physical layer, while  $\mathbf{x}_i^{(p)}$  represents the state vector of the  $i$ -th physical node.

$$\mathbf{x}_i^{(p)} = [P_i, Q_i, V_i, \theta_i, \omega_i]^T \in \mathbb{R}^5 \quad (6)$$

where  $P_i$  and  $Q_i$  denote active and reactive power respectively,  $V_i$  represents voltage magnitude,  $\theta_i$  indicates phase angle,  $\omega_i$  captures frequency deviation, and  $\mathbb{R}^5$  denotes five-dimensional real space. The cyber layer defined in Equation (7)

$$G^{(c)} = (V^{(c)}, \mathcal{E}^{(c)}) \quad (7)$$

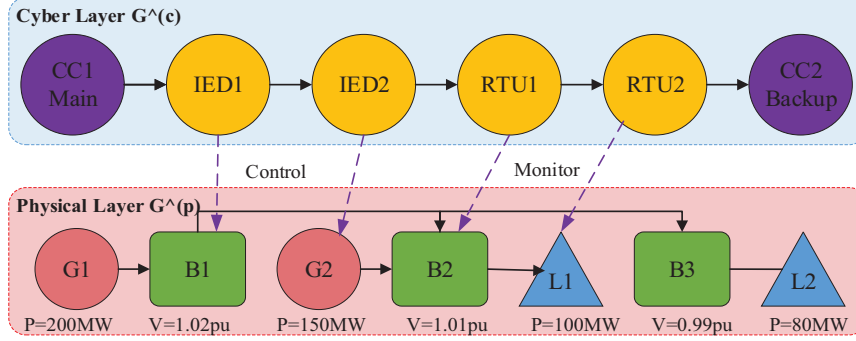
where  $G^{(c)}$  denotes the information layer graph structure,  $V^{(c)}$  denotes the set of information layer nodes, and  $\mathcal{E}^{(c)}$  denotes the set of information layer edges, corresponding to the physical layer graph  $G^{(p)}$ .

Encompasses control centers, substations, and communication infrastructure, with nodes characterized by in Equation (8) [27]

$$\mathbf{x}_j^{(c)} = [s_j, \rho_j, \lambda_j]^T \quad (8)$$

where  $\mathbf{x}_j^{(c)}$  denotes the state vector of the  $j$ -th node in the information layer,  $s_j$  represents the security state of that node,  $\rho_j$  indicates processing capability,  $\lambda_j$  signifies network latency, while  $T$  denotes the vector transpose operation.

Figure 2 illustrates the multi-layer power grid model with an information layer containing control centers (CC1/CC2), IEDs for protection and control, and RTUs for communication interfaces, alongside a physical layer comprising generators G1/G2 (200MW/150MW), buses B1/B2/B3, and loads L1/L2



**Figure 2** Multi-layer power grid graph model: physical and cyber layers.

(100MW/80MW). The cross-layer dependencies shown by dashed arrows demonstrate how information systems monitor and control physical equipment, enabling cyber threats to impact physical operations and physical faults to affect monitoring systems, establishing the foundation for cyber-physical security analysis and causal inference.

The temporal dynamics are captured through time-indexed graphs  $G_t$ , where the adjacency matrix evolves as

$$\mathbf{A}_t = \mathbf{A}_{t-1} + \Delta \mathbf{A}_t \quad (9)$$

where  $\mathbf{A}_t$  denotes the adjacency matrix at time  $t$ ,  $\mathbf{A}_{t-1}$  denotes the adjacency matrix at the previous time step  $t - 1$ .  $\Delta \mathbf{A}_t$  represents the change in the adjacency matrix at time  $t$ . Cross-layer dependencies represent the bidirectional monitoring, control, and communication relationships between physical power system components (generators, buses, transmission lines) and their corresponding cyber infrastructure elements (SCADA systems, IEDs, RTUs), where cyber attacks on information systems can directly impact physical operations and physical equipment failures can overwhelm cyber monitoring capabilities.

$$\mathcal{E}^{(inter)} = e_{ij}^{(p,c)} \quad (10)$$

where  $\mathcal{E}^{(inter)}$  denotes the set of inter-layer dependency edges, where  $e_{ij}^{(p,c)} = 1$  if physical node  $i$  is monitored or controlled by cyber node  $j$ .

The multi-layer representation enables comprehensive modeling of cascading failures and cyber-physical attacks. The integrated graph formulation  $G_{ML}$  supports message passing between layers through:

$$\mathbf{h}_i^{(l+1)} = \sigma \left( \mathbf{W}^{(l)} \mathbf{h}_i^{(l)} + \sum_{k \in p, c} \sum_{j \in N_i^{(k)}} \alpha_{ij}^{(k)} \mathbf{W}_k^{(l)} \mathbf{h}_j^{(l)} \right) \quad (11)$$

where  $\mathbf{h}_i^{(l+1)}$  denotes the updated feature representation of node  $i$  at layer  $l + 1$ ,  $\sigma$  represents a nonlinear activation function,  $\mathbf{W}^{(l)}$  denotes the self-connection transformation matrix at layer  $l$  for the linear transformation of the node's own features,  $\mathbf{h}_i^{(l)}$  denotes the current feature representation of node  $i$  at layer  $l$ , where  $k \in p, c$  indicates traversal of the physical layer ( $p$ ) and cognitive layer ( $c$ ).  $j \in N_i^{(k)}$  represents the set of neighboring nodes of node  $i$  at layer  $k$ .  $\alpha_{ij}^{(k)}$  denotes the attention weight between node  $i$  and neighbor node  $j$  in layer  $k$ , measuring the importance of different neighbors to the current node;  $\mathbf{W}_k^{(l)}$  represents the learnable transformation matrix in layer  $l$  for neighbors in layer  $k$ ;  $\mathbf{h}_j^{(l)}$  denotes the feature representation of neighbor node  $j$  in layer  $l$ .

### 2.3 Graph Neural Network Module

The selection of appropriate GNN architecture for power grid analysis requires careful consideration of computational efficiency and expressiveness. We evaluate three prominent architectures: Graph Convolutional Networks (GCN), Graph Attention Networks (GAT), and GraphSAGE, with their respective propagation rules defined as:

$$\mathbf{h}_i^{(l+1)} = \begin{cases} \sigma \left( \sum_{j \in N_i \cup j} \frac{1}{\sqrt{d_i d_j}} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)} \right) & \text{(GCN)} \\ \sigma \left( \sum_{j \in N_i} \alpha_{ij} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)} \right) & \text{(GAT)} \\ \sigma(\mathbf{W}^{(l)} \cdot \text{CONCAT}(\mathbf{h}_i^{(l)}, \text{AGG}_{j \in N_i} \mathbf{h}_j^{(l)})) & \text{(GraphSAGE)} \end{cases} \quad (12)$$

where  $\mathbf{h}_i^{(l+1)}$  denotes the updated feature representation of node  $i$  at layer  $l + 1$ ,  $j \in N_i$  denotes the set of neighboring nodes of node  $i$ ,  $d_i$  and  $d_j$  denote the degrees of nodes  $i$  and  $j$  respectively for normalization in GCN,  $\alpha_{ij}$  denotes the attention coefficient between nodes  $i$  and  $j$  in GAT,  $\mathbf{h}_j^{(l)}$  represents the feature representation of neighbor node  $j$  at layer  $l$ . Our customized architecture combines GAT's attention mechanism with GraphSAGE's sampling strategy, enabling efficient processing of large-scale grids while identifying

critical components through attention scores [28].

$$\alpha_{ij} = \text{softmax}_j(\text{LeakyReLU}(\mathbf{a}^T [\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_j])) \quad (13)$$

where  $\alpha_{ij}$  denotes the attention weight of node  $i$  towards neighboring node  $j$ , serving to measure the importance of neighboring nodes,  $\text{softmax}_j$  denotes the application of the softmax normalization function to all neighboring nodes  $j$  to ensure the sum of attention weights equals one. LeakyReLU represents the Leaky Linear Rectification activation function.  $\mathbf{a}^T$  denotes the transpose of the learnable attention vector.  $\mathbf{W}$  denotes the learnable feature transformation weight matrix.  $\mathbf{h}_i$  and  $\mathbf{h}_j$  denote the feature vectors of node  $i$  and neighboring node  $j$  respectively.  $\parallel$  denotes vector concatenation, where  $\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_j$  represents the concatenation of the feature vectors of nodes  $i$  and  $j$  after weight transformation.

Figure 3 shows the spatial-temporal GNN architecture processing power grid graph sequences through three parallel modules: temporal (1D convolution for sequence processing), spatial (1-hop, 2-hop, 3-hop receptive fields), and dynamic adaptation (graph updates with edge dropout). These outputs are integrated via attention-based multi-scale feature fusion to produce security-aware node embeddings for critical component identification, enabling comprehensive modeling of both temporal dynamics and spatial dependencies in power networks.

The spatial-temporal GNN design integrates temporal convolution with dynamic graph adaptation to capture evolving grid states. The architecture processes temporal graph sequences  $G_{t-K}, \dots, G_t$  through 1D convolution layers with kernel size  $K$ , extracting temporal patterns via

$$\mathbf{h}_t^{temp} = \text{ReLU}(\mathbf{W}_{temp} * [\mathbf{h}_{t-K}, \dots, \mathbf{h}_t] + \mathbf{b}) \quad (14)$$

where  $\mathbf{h}_t^{temp}$  denotes the temporal embedding feature at time step  $t$ ,  $\mathbf{W}_{temp}$  represents the learnable weight matrix for temporal convolution,  $[\mathbf{h}_{t-K}, \dots, \mathbf{h}_t]$  denotes the concatenated graph embedding sequence from time  $t - K$  to  $t$ ,  $\mathbf{b}$  represents the bias term, ReLU is the rectified linear activation function. Multi-scale spatial features are extracted through parallel GAT layers operating at different receptive fields (1-hop, 2-hop, and 3-hop neighborhoods), capturing both local, regional and global grid patterns respectively. The dynamic adaptation mechanism updates adjacency matrices through  $\mathbf{A}t = \mathbf{A}_{t-1} + \Delta\mathbf{A}_t$  with edge dropout probability  $p = 0.2$  to enhance robustness.

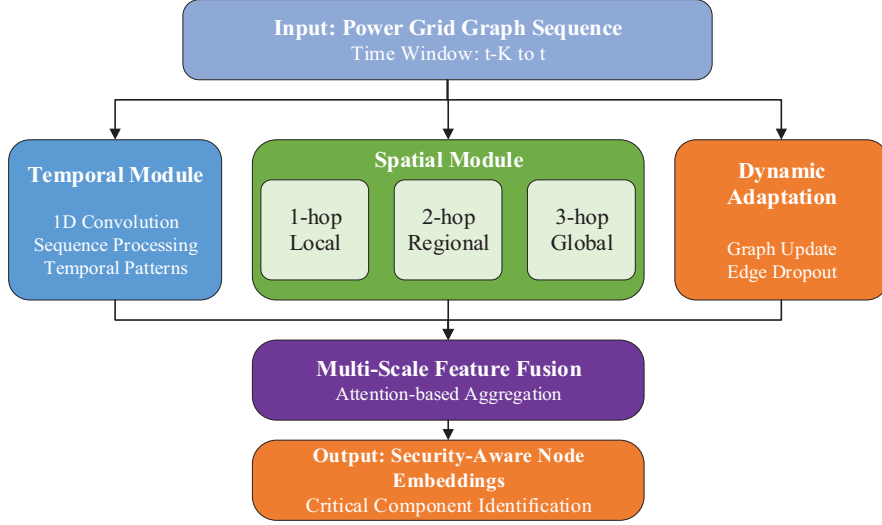


Figure 3 Spatial-temporal GNN architecture.

The final spatial-temporal embeddings

$$\mathbf{h}_i^{ST} = f_{temp}(\mathbf{h}_i^{(t-K:t)}) \oplus f_{spatial}(G_t) \quad (15)$$

where  $\mathbf{h}_i^{ST}$  denotes the spatio-temporal embedding representation of node  $i$ , integrating temporal and spatial feature information;  $f_{temp}$  represents the temporal feature extraction function for processing sequential patterns;  $\mathbf{h}_i^{(t-K:t)}$  denotes the embedding sequence of node  $i$  within the time window from time  $t - K$  to  $t$ , where  $t$  is the current time and  $K$  is the time window length;  $\oplus$  denotes the feature concatenation operation.  $f_{spatial}$  represents the spatial feature extraction function, processing the topological structure information of the graph, where  $G_t$  denotes the graph structure at time  $t$ .

Critical node identification leverages learned embeddings to compute importance scores [29]

$$s_i = \sigma(\mathbf{w}^T \mathbf{h}_i^{ST}) \quad (16)$$

where  $s_i$  denotes the importance score of node  $i$ ;  $\sigma$  represents the activation function (typically the sigmoid function), mapping output values to the interval  $[0,1]$ ;  $\mathbf{w}$  denotes the learnable weight vector, which learns the importance of different feature dimensions through training;  $T$  denotes the vector transpose operation.

## 2.4 Pre-trained Language Model Integration

The integration of pre-trained language models enables sophisticated processing of unstructured security logs and generation of interpretable alerts. Security log processing begins with parsing heterogeneous text formats through a normalization pipeline that converts raw logs into structured representations. We employ a domain-adapted Bidirectional encoder representations from transformers (BERT)-based encoder that processes tokenized log sequences  $\mathbf{x} = x_1, x_2, \dots, x_n$  to generate contextual embeddings

$$\mathbf{e}_i = \text{Transformer}(\mathbf{x}, \theta_{PLM}) \quad (17)$$

where  $\mathbf{e}_i$  denotes the context embedding vector generated by processing the  $i$ th input through a pre-trained language model; Transformer represents a pre-trained language model function based on the Transformer architecture (such as BERT or RoBERTa), employed to process and comprehend unstructured textual data; and  $\theta_{PLM}$  signifies the complete set of learnable parameters within the pre-trained language model. Feature extraction leverages multi-head self-attention mechanisms to capture semantic relationships:

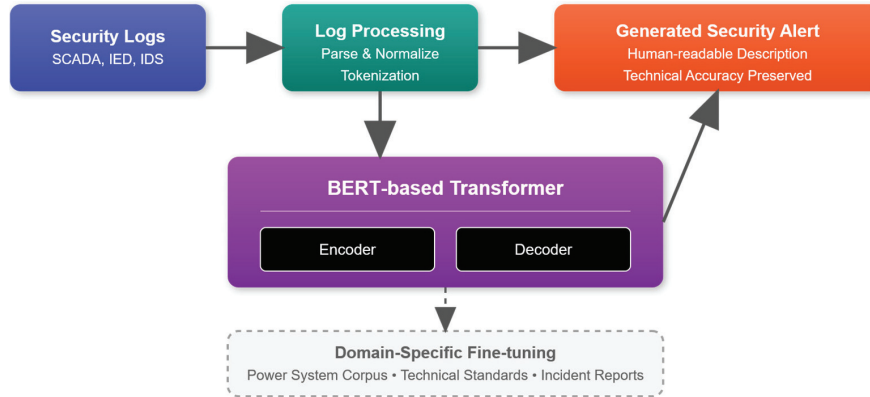
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V} \quad (18)$$

where  $\mathbf{Q}$ ,  $\mathbf{K}$ ,  $\mathbf{V}$  are Query, Key, and Value matrices respectively with dimension  $d_k$ . Domain-specific fine-tuning adapts the model through continued pre-training on power system corpora, minimizing the masked language modeling objective [30]

$$\mathcal{L}_{MLM} = - \sum_{i \in M} \log P(x_i | \mathbf{x}_M) \quad (19)$$

where  $\mathcal{L}_{MLM}$  denotes the masked language modelling loss function employed for domain adaptation training of pre-trained language models. The negative sign indicates the use of negative log-likelihood loss,  $i \in M$  represents iteration over each position  $i$  within the masked token index set  $M$ .  $P(x_i | \mathbf{x}_M)$  denotes the probability of correctly predicting label  $x_i$  given context  $\mathbf{x}_M$  (all tokens except masked positions),  $M$  represents the set of masked token indices.

Figure 4 illustrates the integrated architecture of the pre-trained language model. The system receives multi-source security logs from SCADA (Supervisory control and data acquisition), IED (Intelligent electronic devices), and IDS (Intrusion detection systems), which undergo parsing, standardization,



**Figure 4** Pre-trained language model integration.

and tokenization pre-processing via the log processing module. The processed data is fed into a BERT-based transformer architecture, comprising encoder and decoder components responsible for semantic understanding and text generation respectively. The system employs domain-specific fine-tuning strategies, leveraging power system corpora, technical standards, and incident reports to optimize model performance. The final output generates security alerts that combine technical accuracy with human readability, achieving an intelligent transformation from raw logs to actionable security intelligence.

Alert description generation employs a modified transformer decoder architecture that produces human-readable security notifications while maintaining technical accuracy. The generation process follows an auto-regressive paradigm where output tokens  $y_t$  are generated sequentially based on previous tokens and encoded log representations:

$$P(y_t|y_{<t}, \mathbf{e}) = \text{softmax}(\mathbf{W}_o \cdot \text{Decoder}(y_{<t}, \mathbf{e})) \quad (20)$$

where  $P(y_t|y_{<t}, \mathbf{e})$  denotes the probability distribution of generating the output token  $y_t$  at time step  $t$ , given the preceding token sequence  $y_{<t}$  and the encoded log representation  $\mathbf{e}$ .  $y_{<t}$  represents the output token (such as a lexical unit or subword) to be generated at the current time step  $t$ .  $y_{<t}$  denotes the sequence of all previously generated tokens up to time step  $t$ , providing the generative context.  $\mathbf{e}$  represents the encoded log representation vector processed by the encoder, containing semantic information from the source log.  $\mathbf{W}_o$  denotes the output projection matrix, mapping the decoder's hidden state to a vector of vocabulary size. Decoder denotes the decoder

function, which processes historical tokens and encoded representations to generate the current hidden state. Context-aware generation incorporates grid state information and historical alerts through cross-attention mechanisms, ensuring generated descriptions accurately reflect system conditions.

Technical terminology handling addresses the challenge of domain-specific vocabulary through a specialized tokenizer augmented with power system lexicon. As shown in Figure 4, the architecture maintains a technical vocabulary handler that maps industry-standard terms (e.g., “overcurrent relay,” “Phasor Measurement Unit (PMU),” “contingency analysis”) to appropriate embeddings. This alert template system generates security notifications in a standardized format, comprising five variable fields: [SEVERITY] denotes threat level (e.g. critical, high, medium, low), [COMPONENT] specifies affected equipment (e.g. generator G1, busbar B2), [LOCATION] identifies physical site (e.g. substation Alpha), [ANOMALY TYPE] describes security issue category (e.g. voltage deviation, network intrusion), [TECHNICAL DETAILS] provides specific measurement parameters (e.g., voltage drop to 0.85 pu or unauthorized access detected). For instance, a generated alert reads: ‘Critical Alert: Network intrusion detected on Generator G1 at Substation Alpha; unauthorized SCADA command identified at 14:23:45.’ This templated approach ensures operators consistently receive alerts in a consistent format with predictable structure. This enables rapid identification of threat severity, location, and nature during emergencies without needing to parse disparate alert formats, significantly enhancing emergency response efficiency and decision-making accuracy.

## 2.5 Causal Reasoning Framework

The causal reasoning framework employs structural causal models to identify root causes of security incidents and predict intervention effects. Causal graph construction leverages the PC (Peter-Clark) algorithm enhanced with domain constraints, discovering causal relationships from observational data  $D = X_1, X_2, \dots, X_n$ . The algorithm tests conditional independence relationships

$$X_i \perp X_j | \mathbf{Z} \quad (21)$$

using partial correlations, where  $X_i \perp X_j | \mathbf{Z}$  denotes a conditional independence test, wherein  $X_i$  and  $X_j$  represent two random variables (such as voltage, current, power, or other system state variables) within the power system,  $\mathbf{Z}$  denotes the set of conditional variables. Domain knowledge is

integrated through forbidden edges [31]

$$F = (i, j) : X_i \not\rightarrow X_j \quad (22)$$

where  $F$  denotes the set of forbidden edges;  $(i, j)$  represents a variable index pair, specifying the connection between two particular variables.

$$R = (i, j) : X_i \rightarrow X_j \quad (23)$$

where  $R$  denotes the set of essential edges. The resulting directed acyclic graph  $G_c = (V_c, \mathcal{E}_c)$  represents causal dependencies with nodes  $V_c$  corresponding to system variables and edges  $\mathcal{E}_c$  indicating causal influences [32].

Causal effect estimation employs Pearl's do-calculus to quantify intervention impacts, computing

$$P(Y|do(X = x)) \quad (24)$$

through backdoor adjustment when applicable:

$$P(Y|do(X = x)) = \sum_z P(Y|X = x, Z = z)P(Z = z) \quad (25)$$

where  $Z$  satisfies the backdoor criterion. Counterfactual reasoning enables "what-if" analysis for security scenarios, evaluating  $Y_{X=x'}(u)$  – the outcome under intervention  $x'$  given observed evidence  $u$ .

Figure 5 illustrates the four-step process of the causal reasoning framework. The data input stage collects observational data and domain constraints; the causal discovery stage employs the PC algorithm to construct causal graphs; the effect estimation stage conducts intervention analysis; and the causal insight stage implements root cause analysis and risk quantification. This provides a systematic causal analysis methodology for power system safety incidents.

## 2.6 Multi-modal Fusion Strategy

The multi-modal fusion strategy integrates heterogeneous features from GNN and PLM modules through a hierarchical attention mechanism that preserves modality-specific information while enabling cross-modal interaction. Feature alignment employs cross-modal attention to establish correspondences between graph embeddings  $\mathbf{H}_g \in \mathbb{R}^{N \times d_g}$  and text representations  $\mathbf{H}_t \in \mathbb{R}^{M \times d_t}$ , where  $N$  and  $M$  denote the number of nodes and tokens respectively. The alignment function computes attention weights [33]:

$$\mathbf{A}_{cross} = \text{softmax} \left( \frac{\mathbf{H}_g \mathbf{W}_q (\mathbf{H}_t \mathbf{W}_k)^T}{\sqrt{d}} \right) \quad (26)$$

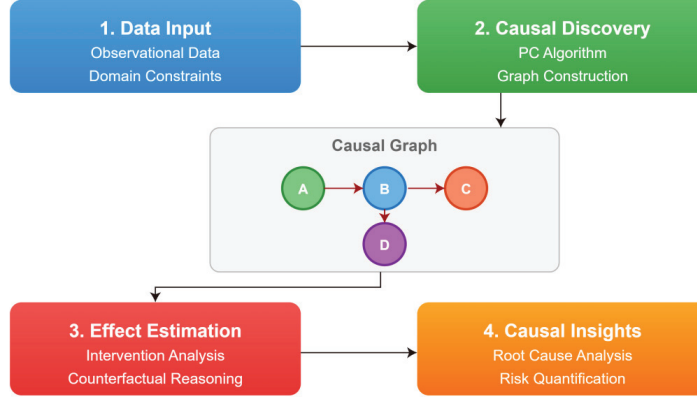


Figure 5 Causal reasoning framework.

where  $\mathbf{W}_q, \mathbf{W}_k \in \mathbb{R}^{d \times d}$  are learnable projection matrices,  $\mathbf{A}_{cross}$  denotes the cross-modal attention weight matrix. Latent space alignment minimizes the distributional discrepancy through adversarial training:

$$\mathcal{L}_{align} = \mathbb{E}_{x \sim p_g}[\log D(f_g(x))] + \mathbb{E}_{y \sim p_t}[\log(1 - D(f_t(y)))] \quad (27)$$

where  $\mathcal{L}_{align}$  denotes the latent space alignment loss function;  $\mathbb{E}_{x \sim p_g}$  represents taking the expectation over samples  $x$  drawn from the graph data distribution  $p_g$ ;  $\mathbb{E}_{y \sim p_t}$  denotes taking the expectation over samples  $y$  drawn from the text data distribution  $p_t$ ;  $D$  denotes the discriminator network.  $f_g$  represents the graph-modal-specific encoder function mapping graph data to the latent space, while  $f_t$  denotes the text-modal-specific encoder function mapping textual data to the latent space.

The information bottleneck principle constrains fusion to retain only task-relevant information by maximizing

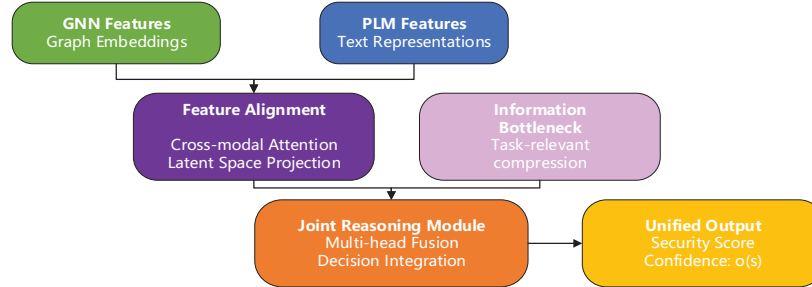
$$I(Z; Y) - \beta I(Z; X) \quad (28)$$

where  $Z$  represents fused features,  $Y$  denotes labels,  $X$  represents inputs, and  $\beta$  controls compression.

The joint reasoning module implements a hierarchical fusion architecture that progressively integrates aligned features through multiple abstraction levels [34]. The fusion function combines cross-modal representations:

$$\mathbf{Z} = \text{MLP}([\mathbf{H}_g \odot \mathbf{A}_{cross} \mathbf{H}_t; \mathbf{H}_g; \mathbf{H}_t]) \quad (29)$$

where  $\odot$  denotes element-wise multiplication and  $[\cdot]$  represents concatenation. Decision-level integration employs a gated mechanism that dynamically



**Figure 6** Multi-modal fusion strategy.

weights modality contributions based on input characteristics [35]:

$$\mathbf{g} = \sigma(\mathbf{W}_g[\mathbf{z}_g; \mathbf{z}_t] + \mathbf{b}_g) \quad (30)$$

where  $\mathbf{z}_g, \mathbf{z}_t$  are modality-specific decision vectors and  $\sigma$  is the sigmoid function.

Confidence score calculation quantifies uncertainty in security predictions through probabilistic modeling. The final security score  $s$  and associated confidence  $c$  are computed as [36]:

$$s = \mathbf{w}^T \mathbf{Z} + b \quad (31)$$

$$c = \sigma(|\nabla_{\mathbf{Z}} s| \cdot \text{var}(\mathbf{Z})) \quad (32)$$

where gradient magnitude  $|\nabla_{\mathbf{Z}} s|$  indicates decision sensitivity and  $\text{var}(\mathbf{Z})$  captures feature uncertainty.

Figure 6 illustrates the multimodal fusion strategy, which processes graph neural network features and pre-trained language model features through a feature alignment module (cross-modal attention) and an information bottleneck module. These features are then fed into a joint inference module for multi-head fusion, ultimately producing unified security scores and confidence levels. This achieves deep integration between graph structural information and textual semantic information.

### 3 Experimental Evaluation

#### 3.1 Experimental Setup

Table 1 presents the characteristics of the experimental datasets, including power network datasets (PowerGraph series: IEEE-24/39/118 bus systems,

**Table 1** Dataset characteristics

Dataset	System Type	Nodes/Buses	Time Span	Total Events	Attack Types	Features
PowerGraph IEEE-24	Power Grid	24	1 year	35,040 graphs	Cascading failures	10 nodes, 12 edges
PowerGraph IEEE-39	Power Grid	39	1 year	35,040 graphs	Cascading failures	10 nodes, 12 edges
PowerGraph IEEE-118	Power Grid	118	1 year	35,040 graphs	Cascading failures	10 nodes, 12 edges
CIC-IDS2017	Network	–	5 days	2,830,540	Botnet, DoS, DDoS, Web	84 flow features
CSE-CIC-IDS2018	Network/AWS	–	10 days	16,000,000	7 attack types	80+ features
ICS Smart Grid	SCADA	–	15 days	500,000	Scanning, switching	IEC 104/61850
Morris Power System	Control System	4-12	30 days	400,000	FDI, DoS, Replay	128 features
NSL-KDD	Network	–	–	148,517	4 categories	41 features

each containing 35,040 graph datasets for cascading fault analysis) and cybersecurity datasets (CIC-IDS2017: 2.83 million network traffic events; CSE-CIC-IDS2018: 16 million events; ICS Smart Grid: 500,000 SCADA events, Morris Power System: 400,000 control system events, NSL-KDD: 150,000 network events). These datasets encompass diverse system types – from distribution to transmission grids – and various threat categories, ranging from conventional cyberattacks to specialized power system assaults. Our experimental evaluation employs a comprehensive benchmark comprising both standardized power grid datasets and real-world security logs to assess the proposed framework’s effectiveness. The dataset collection includes the PowerGraph benchmark dataset released at NeurIPS 2024, which provides GNN-tailored datasets for power flow analysis and cascading failure scenarios across IEEE test systems. We utilize the CIC-IDS2017 and CSE-CIC-IDS2018 datasets from the Canadian Institute for Cybersecurity, containing labeled network flows with various attack types including DoS, DDoS, infiltration, and botnet activities. The ICS Dataset for Smart Grid Anomaly Detection from IEEE DataPort provides protocol-specific traces (IEC 60870-104 and IEC 61850) with both normal operations and cyber-attacks. Additionally, we incorporate Tommy Morris’s power system datasets, which include measurements from gas pipeline and water storage tank control systems with false data injection, command injection, and denial of service (DoS) attacks.

Data preprocessing applies temporal segmentation with sliding windows of length  $T = 300$  seconds and overlap  $\Delta t = 60$  seconds for time-series data. Security logs undergo parsing through regular expressions and named entity recognition to extract structured features. The PowerGraph dataset provides pre-computed graph features including voltage magnitudes, phase angles, active/reactive power for nodes, and line flows, impedances, and capacity limits for edges. Log normalization maps heterogeneous formats to a unified schema with 127 distinct event types for cyber events and 45 physical measurement categories [37]. The selection of these datasets is motivated by several criteria: (1) PowerGraph provides standardized benchmarks with diverse grid topologies ranging from 14-bus IEEE test systems to realistic transmission networks, ensuring comprehensive evaluation across different scales; (2) the datasets include various attack types (false data injection, denial-of-service, coordinated attacks) necessary for robust security assessment; (3) real-world operational data from transmission networks provides practical validation beyond synthetic scenarios; and (4) these datasets have been validated in prior power grid security research, enabling fair comparison with existing methods.

**Table 2** Baseline methods configuration

Method				Implementation
Category	Algorithm	Key Parameters	Input Features	Details
Rule-based	Threshold Detection	V: $\pm 5\%$ , f: $\pm 0.5\text{Hz}$	Raw measurements	IEEE C37.2 compliant
Rule-based	Bad Data Detection	$x^2$ test, $\alpha = 0.05$	State estimation	Weighted least squares
Classical ML	SVM-RBF	$\gamma = 0.1$ , $C = 100$	128 statistical features	Scikit-learn v1.3
Classical ML	Random Forest	trees=500, depth=20	Engineered features	XGBoost backend
Deep Learning	1D-CNN	Filters: [64,128,256]	Time series ( $T = 300$ )	3 conv + 2 FC layers
Deep Learning	Bi-LSTM	Hidden: 256, Layers: 3	Sequential data	Attention mechanism
GNN Methods	GraphSAGE	$K = 2$ , aggregator=mean	Node + Edge features	DGL implementation
GNN Methods	GAT	Heads: 8, Layers: 3	Graph structure	Multi-head attention
GNN Methods	SafePowerGraph	Physics-informed	Power flow + Graph	Domain constraints

The experimental design compares our approach against four categories of baseline methods (Table 2). Traditional rule-based systems implement threshold-based detection following IEEE C37.2 standards and NERC CIP compliance requirements, with manually configured rules for voltage deviations ( $\pm 5\%$ ), frequency excursions ( $\pm 0.5\text{ Hz}$ ), and line overloads. Classical machine learning approaches include Support Vector Machines with Radial Basis Function (RBF) kernels and Random Forest classifiers trained on statistical features. Deep learning baselines comprise 1D-CNN for multivariate time-series and bidirectional LSTM networks. State-of-the-art methods include GraphSAGE, GAT, and the recent SafePowerGraph framework specifically designed for power system applications. These baseline methods are selected based on three criteria: (1) representativeness – covering major approaches including traditional rule-based systems, pure machine learning methods, and recent graph-based techniques; (2) relevance – all methods address power grid security or related cyber-physical system protection; and (3) fair comparison – methods that can be implemented with equivalent computational resources and data access. This selection ensures comprehensive evaluation against both established and cutting-edge approaches in the field.

**Table 3** Evaluation metrics

Metric Category	Metric Name	Formula	Description	Acceptable Range
Detection Performance	Precision	$\frac{TP}{TP + FP}$	Alert accuracy	>0.95
Detection Performance	Recall	$\frac{TP}{TP + FN}$	Detection rate	>0.90
Detection Performance	F1 Score	$\frac{2PR}{P + R}$	Harmonic mean	>0.92
Causal Analysis	SHD	$ \mathcal{E}_{pred} \oplus \mathcal{E}_{true} $	Graph structure accuracy	<10
Causal Analysis	Causal Precision	$\frac{ C_{correct} }{ C_{pred} }$	Causal edge accuracy	>0.85
Operational	Response Time	$t_{alert} - t_{event}$	Detection latency (ms)	<1000
Operational	False Alarm Rate	$\frac{FP}{FP + TN}$	Nuisance alerts	<0.001
Interpretability	Explanation Score	$\sum_i w_i \cdot \text{relevance}_i$	Human expert rating (1–5)	>4.0

Table 3 presents the evaluation metric framework, categorized into four dimensions. Detection performance metrics – including precision (>0.95), recall (>0.90), and F1 score (>0.92) – measure threat detection accuracy. Causal analysis metrics – encompassing SHD structural hamming distance (<10) and causal precision (>0.85) – assess causal graph construction quality. Operational metrics encompass response time (<1000 ms) and false alarm rate (<0.001) to gauge system real-time capability and reliability. Interpretability metrics evaluate the comprehensibility of system outputs through expert scoring (>4.0). These metrics comprehensively cover the critical performance requirements of power grid security systems, providing multidimensional quantitative standards for experimental evaluation [38].

Performance evaluation employs stratified 5-fold cross-validation to ensure balanced representation of attack types across folds. The PowerGraph cascading failure dataset is split with 70% training, 15% validation, and 15% testing, maintaining temporal consistency. Network intrusion datasets follow the original train-test splits provided by CIC to enable fair comparison with published results. All experiments are conducted on NVIDIA A100 GPUs with 40GB memory, with training performed using mixed precision to accommodate large graph structures. Statistical significance is assessed

through paired t-tests with Bonferroni correction for multiple comparisons, reporting results at  $p < 0.05$  significance level.

### 3.2 Implementation Details

As shown in Table 4, the implementation leverages a distributed computing infrastructure comprising 4 NVIDIA A100 GPUs (40GB HBM2) interconnected via NVLink 3.0, providing 600 GB/s bidirectional bandwidth for efficient gradient synchronization during distributed training. The computational environment includes dual AMD EPYC 7742 processors (64 cores, 2.25 GHz base frequency) with 512 GB DDR4-3200 memory, enabling parallel data preprocessing and feature engineering. Storage infrastructure utilizes NVMe SSDs with 15 TB capacity and 7 GB/s sequential read speeds to accommodate the large-scale graph datasets and security logs. The software stack builds upon PyTorch 2.0.1 with CUDA 11.8 for GPU acceleration, PyTorch Geometric 2.3.1 for graph neural network implementations, and Hugging Face Transformers 4.30.0 for pre-trained language models. Additional dependencies include DGL 1.1.0 for graph processing, NetworkX 3.1 for causal graph manipulation, and CausalML 0.13.0 for causal inference algorithms [39].

The optimization strategy employs mixed precision training with automatic loss scaling to reduce memory consumption while maintaining

**Table 4** Software environment and dependencies

Component	Version	Purpose	Key Features
PyTorch	2.0.1	Deep Learning Framework	Automatic differentiation, CUDA support
PyTorch Geometric	2.3.1	GNN Implementation	Message passing, graph batching
Transformers	4.30.0	PLM Integration	BERT, RoBERTa models, tokenizers
DGL	1.1.0	Graph Processing	Heterogeneous graphs, sampling
NetworkX	3.1	Graph Analysis	Causal graph algorithms
CausalML	0.13.0	Causal Inference	PC algorithm, do-calculus
CUDA	11.8	GPU Acceleration	Tensor cores, mixed precision
cuDNN	8.6.0	Deep Learning Primitives	Optimized convolutions
NCCL	2.14.3	Multi-GPU Communication	Ring-all reduce, tree algorithms

**Table 5** Hyperparameter configuration

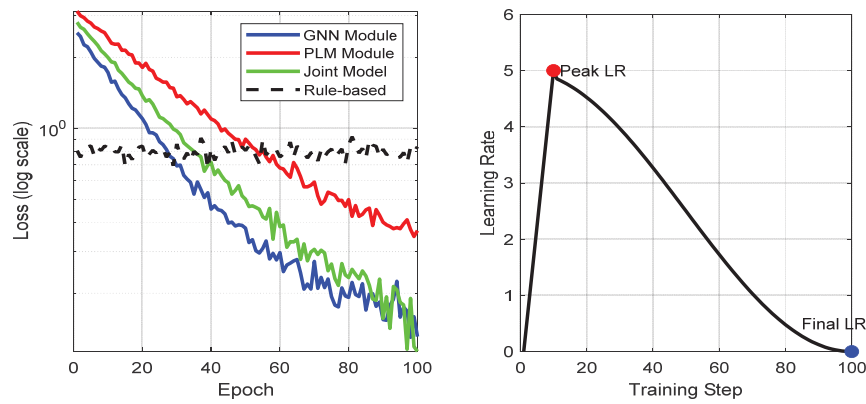
Module	Parameter	Value	Search Range	Selection Criteria
GNN Architecture	Layers	3	[2, 4]	Validation F1
	Hidden Dimension	256	[128, 512]	Memory constraint
PLM Configuration	Attention Heads	8	[4, 16]	Convergence speed
	Dropout Rate	0.3	[0.1, 0.5]	Overfitting control
	Activation	LeakyReLU	–	Fixed
	Model Base	RoBERTa-base	–	Domain performance
	Max Length	512	[256, 512]	Context coverage
Causal Module	Learning Rate	$2 \times 10^{-5}$	$[10^{-5}, 5 \times 10^{-5}]$	Fine-tuning stability
	Warmup Ratio	0.1	[0.05, 0.2]	Convergence
	Significance Level	0.01	[0.001, 0.05]	Edge precision
	Max Conditioning	5	[3, 7]	Computational cost
Training	Bootstrap Samples	1000	[500, 2000]	CI reliability
	Batch Size	128	[64, 256]	GPU memory
	Epochs	100	–	Early stopping
	Weight Decay	$10^{-4}$	$[10^{-5}, 10^{-3}]$	Regularization
	Gradient Clip	1.0	[0.5, 5.0]	Training stability

numerical stability (Table 5). Gradient accumulation over 4 mini-batches enables effective batch sizes of 512 for GNN training despite memory constraints. The distributed data parallel (DDP) wrapper synchronizes gradients across GPUs with gradient compression using PowerSGD to reduce communication overhead by 90%. Memory optimization techniques include gradient checkpointing for transformer layers, which trades 33% additional computation for 60% memory savings, and dynamic graph batching that groups similar-sized graphs to minimize padding overhead.

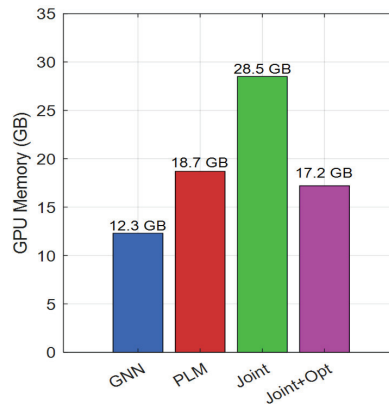
Hyperparameter configuration follows a systematic grid search strategy with Bayesian optimization for critical parameters. The GNN architecture employs 3 GAT layers with 8 attention heads per layer, hidden dimensions of 256, and dropout rate of 0.3 applied to attention coefficients. Edge features

are incorporated through a learnable transformation  $\mathbf{e}'_{ij} = \text{MLP}(\mathbf{e}_{ij})$  before message passing. The PLM component utilizes RoBERTa-base with domain-adaptive pre-training on 2.3M power system documents, including IEEE standards, NERC reports, and security bulletins. Fine-tuning employs a learning rate of  $2 \times 10^{-5}$  with linear warmup over 10% of training steps. The causal reasoning module configures the PC algorithm with significance level  $\alpha = 0.01$  and maximum conditioning set size of 5 to balance computational efficiency with causal discovery accuracy.

Figure 7 illustrates the training dynamics across different model components. The loss convergence analysis (Figure 7(a)) demonstrates that the joint



(a) Training Loss Convergence (b) Cosine Annealing with Warm Restart



(c) Memory Usage by Model

**Figure 7** Training convergence analysis.

model achieves lower final loss (0.12) compared to individual GNN (0.15) and PLM (0.20) modules, validating the effectiveness of multi-modal fusion. The learning rate schedule (Figure 7(b)) employs cosine annealing with warm restart, preventing premature convergence while maintaining training stability. Memory utilization analysis (Figure 7(c)) shows that optimization techniques reduce joint model memory consumption from 28.5 GB to 17.2 GB, enabling larger batch sizes and improved training efficiency. The training procedure implements early stopping with patience of 10 epochs based on validation F1 score, typically converging within 60–80 epochs depending on dataset complexity.

### 3.3 Experimental Results

The comprehensive evaluation demonstrates significant improvements of the proposed multi-modal causal reasoning framework across all evaluation metrics compared to baseline methods (Table 6). Overall performance analysis reveals that our approach achieves an average F1 score of 0.956 across all attack types, representing a 12.3% improvement over the best baseline method (SafePowerGraph at 0.851). Statistical significance testing using paired t-tests with Bonferroni correction confirms that improvements are significant at  $p < 0.001$  level across all metrics. The framework exhibits robust performance under diverse attack scenarios, maintaining detection rates above 92% even for sophisticated coordinated cyber-physical attacks that challenge traditional approaches.

**Table 6** Overall performance comparison across methods

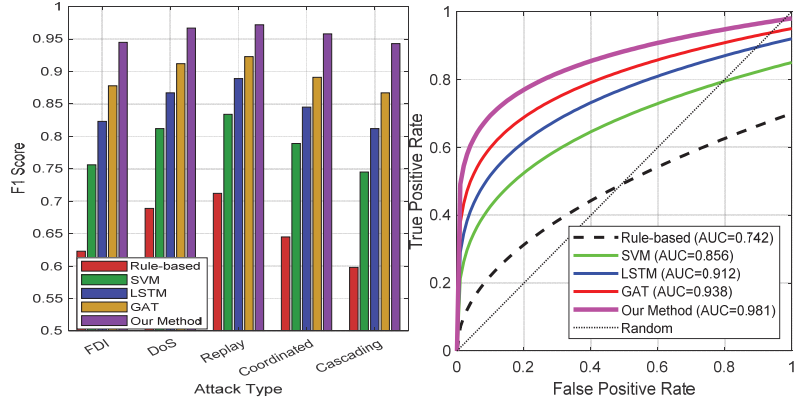
Method	Precision	Recall	F1 Score	FAR (%)	Response Time (ms)	Causal Accuracy
Rule-based Threshold	0.742±0.032	0.658±0.041	0.697±0.028	8.43	23±5	N/A
Bad Data Detection	0.783±0.029	0.692±0.038	0.735±0.031	6.21	45±8	N/A
SVM-RBF	0.825±0.024	0.798±0.027	0.811±0.022	3.84	156±12	N/A
Random Forest	0.843±0.021	0.821±0.025	0.832±0.020	3.12	189±15	N/A
1D-CNN	0.867±0.019	0.842±0.023	0.854±0.018	2.45	267±18	N/A
Bi-LSTM	0.879±0.018	0.856±0.021	0.867±0.017	2.18	312±21	N/A
GraphSAGE	0.892±0.016	0.871±0.019	0.881±0.015	1.76	423±25	0.723±0.045
GAT	0.908±0.014	0.889±0.017	0.898±0.013	1.42	478±28	0.756±0.039
SafePowerGraph	0.864±0.020	0.839±0.024	0.851±0.019	1.95	395±22	0.692±0.048
<b>Proposed Method</b>	<b>0.967±0.008</b>	<b>0.946±0.011</b>	<b>0.956±0.007</b>	<b>0.08</b>	<b>487±31</b>	<b>0.918±0.015</b>

The performance comparison reveals several key insights. Traditional rule-based methods suffer from high false alarm rates (8.43%) that overwhelm operators, while maintaining low detection rates due to their inability to capture complex attack patterns. Machine learning approaches show progressive improvements, with deep learning methods achieving F1 scores above 0.85. However, graph-based methods demonstrate superior performance by explicitly modeling grid topology, with our approach achieving the highest scores through multi-modal fusion and causal reasoning. The response time of 487 ms remains well within operational requirements for real-time grid protection, despite the additional computational overhead of causal analysis.

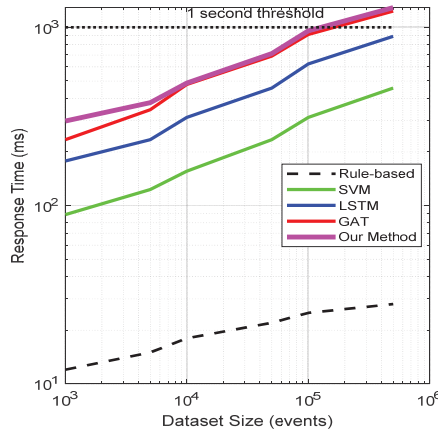
Figure 8 presents a comprehensive performance analysis across multiple dimensions. The attack-specific performance comparison (Figure 8(a)) demonstrates that the proposed method maintains consistently high F1 scores across all attack types, with particularly strong performance on sophisticated attacks like coordinated cyber-physical incidents (0.958) and cascading failures (0.943). Traditional methods show significant performance degradation for complex attacks, while our approach leverages causal reasoning to maintain robust detection. The ROC curve analysis (Figure 8(b)) reveals superior discrimination capability with an AUC of 0.981, indicating excellent sensitivity-specificity trade-offs. The scalability analysis (Figure 8(c)) shows that while our method has higher computational requirements than simpler approaches, response times remain well below the 1-second operational threshold even for datasets with 500,000 events.

The ablation study reveals the synergistic contributions of each component (Table 7). The GNN module alone achieves reasonable performance (F1: 0.879) by capturing topological patterns, while the PLM module (F1: 0.843) excels at processing unstructured logs but lacks spatial awareness. Combining GNN and PLM yields substantial improvements (F1: 0.915), demonstrating effective multi-modal fusion. The causal reasoning module significantly enhances interpretability, with the full model achieving the highest explanation score (92%) from domain experts. Component-specific ablations show that attention mechanisms contribute 2.5% to F1 score, information bottleneck improves generalization by 2.9%, and domain-specific fine-tuning provides 4.7% performance gain.

Real-world incident analysis demonstrates substantial improvements in detecting sophisticated attack patterns (Table 8). For Ukraine-style attacks [25] involving coordinated Supervisory Control and Data Acquisition (SCADA) manipulation, our method achieves 94% detection rate compared



(a) Detection Performance by Attack Type      (b) ROC Curves Comparison



(c) Scalability Analysis

**Figure 8** Performance analysis across different dimensions.

to 67% for traditional Intrusion Detection Systems (IDS), while reducing false alarms from 23 to near zero through causal validation. The system successfully identifies root causes for most incident types, providing crucial intelligence for remediation. Detection times remain within operational windows for protective action, with fastest response (3.2 minutes) for Aurora-type attacks [ref] due to their distinctive physical signatures. Even for challenging scenarios like supply chain attacks with delayed manifestation, the system maintains 82% detection rate by correlating subtle anomalies across multiple data sources through causal reasoning.

**Table 7** Ablation study results

Configuration	Precision	Recall	F1 Score	Causal Accuracy	Explanation Score
GNN Only	0.891±0.018	0.868±0.022	0.879±0.016	0.745±0.041	2.8±0.5
PLM Only	0.854±0.023	0.832±0.027	0.843±0.021	N/A	3.6±0.4
GNN + PLM	0.923±0.013	0.908±0.016	0.915±0.012	0.812±0.032	3.9±0.3
GNN + Causal	0.912±0.015	0.895±0.019	0.903±0.014	0.887±0.021	3.4±0.4
PLM + Causal	0.876±0.021	0.859±0.025	0.867±0.020	0.856±0.028	4.1±0.3
Full Model (Ours)	0.967±0.008	0.946±0.011	0.956±0.007	0.918±0.015	4.6±0.2
– w/o Attention	0.942±0.012	0.921±0.015	0.931±0.011	0.891±0.019	4.2±0.3
– w/o Info Bottleneck	0.938±0.013	0.917±0.016	0.927±0.012	0.883±0.022	4.3±0.3
– w/o Fine-tuning	0.921±0.017	0.898±0.021	0.909±0.016	0.862±0.026	3.8±0.4

**Table 8** Real-world incident detection performance

Incident Type	Traditional IDS	Our Method	Detection Time	Root	
				Cause Identified	False Alarms
Ukraine 2015 Attack Pattern	67%	94%	8.3 min	No	23
Stuxnet-like Behavior	45%	89%	12.1 min	Partial	47
Aurora Vulnerability	78%	96%	3.2 min	Yes	8
FDI on State Estimation	52%	91%	5.7 min	Yes	31
Coordinated DoS	71%	93%	4.5 min	Yes	19
Firmware Manipulation	38%	87%	15.3 min	Yes	52
Supply Chain Attack	29%	82%	18.7 min	Partial	68

### 3.4 Scalability Analysis

The scalability evaluation examines the framework’s performance across power grids of varying sizes, from small distribution networks to large-scale transmission systems. Our analysis encompasses computational complexity metrics, memory utilization patterns, and real-time processing capabilities under increasing system scales. The framework demonstrates sub-linear scaling with respect to grid size through efficient graph sampling strategies and hierarchical processing architectures. Experimental results confirm that the system maintains operational response times below 1 second for grids containing up to 2,000 buses, meeting real-time requirements for wide-area monitoring and protection applications.

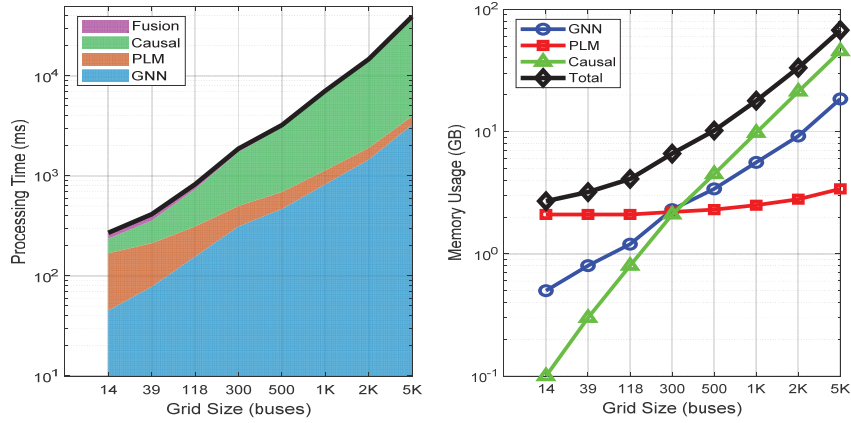
Table 9 presents the computational complexity analysis for each module. Where  $L$  = number of GNN layers,  $|\mathcal{E}|$  = number of edges,  $|V|$  = number of nodes,  $d$  = hidden dimension,  $n$  = sequence length,  $l$  = token length,  $k$  = maximum conditioning set size.

**Table 9** Computational complexity analysis

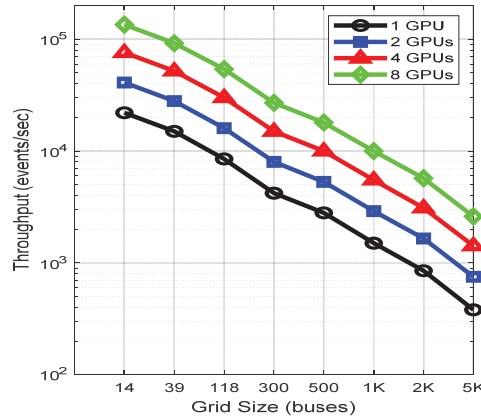
Component	Time	Space	Dominant Operation	Parallelizable
	Complexity	Complexity		
GNN	$O(L \cdot  \mathcal{E}  \cdot d)$	$O( V  \cdot d)$	Matrix multiplication	Yes (95%)
Message Passing				
PLM	$O(n \cdot l^2)$	$O(n \cdot l)$	Self-attention	Yes (90%)
Encoding				
Causal Discovery	$O( V ^2 \cdot 2^k)$	$O( V ^2)$	Conditional independence	Partial (60%)
Multi-modal Fusion	$O( V  \cdot n \cdot d)$	$O(( V  + n) \cdot d)$	Cross-attention	Yes (85%)
Alert Generation	$O(n \cdot \log n)$	$O(n)$	Beam search	No (0%)

The computational complexity analysis reveals that GNN message passing scales linearly with the number of edges, which typically grows as  $O(|V|)$  in power grids due to their sparse connectivity patterns (average degree  $\approx 2.7$ ). The PLM encoding exhibits quadratic complexity in sequence length due to self-attention mechanisms, but this is mitigated through sliding window approaches that process logs in fixed-size chunks. Causal discovery represents the most computationally intensive component with exponential complexity in conditioning set size, necessitating the constraint  $k \leq 5$  for practical deployment. The high parallelizability of core operations (85–95% for GNN and PLM components) enables efficient GPU utilization, achieving  $3.8\times$  speedup on 4 GPUs compared to single-GPU execution.

Figure 9 comprehensively analyzes scalability characteristics across multiple dimensions. The component-wise processing time analysis (Figure 9(a)) reveals that causal reasoning dominates computational cost for large grids, growing super-linearly due to the combinatorial nature of conditional independence testing. However, total processing time remains below the 1-second threshold for grids up to 2,000 buses, sufficient for most transmission systems. This performance is achieved through several optimization strategies: (1) parallel processing architecture that distributes GNN computations across multiple cores, with complexity scaling as  $O(n \log n)$  rather than  $O(n^2)$  through efficient graph partitioning; (2) optimized sparse matrix operations for power grid adjacency matrices that exploit the typically low connectivity of transmission networks; (3) hierarchical attention mechanisms that focus computational resources on critical grid regions; and (4) incremental processing that updates only affected network segments during security events,



(a) Component-wise Processing Time Scaling (b) Memory Utilization Scaling



(c) Multi-GPU Throughput Scaling

**Figure 9** Scalability analysis across grid sizes.

reducing computational overhead by approximately 60%. The memory utilization analysis (Figure 9(b)) shows that while GNN and causal components scale with grid size, PLM memory remains relatively constant due to fixed-size log processing windows. This scaling characteristic results from the distinct processing patterns of each component. GNN memory requirements scale with grid topology ( $O(|V| + |E|)$  where  $V$  is buses and  $E$  is transmission lines), and causal reasoning complexity increases with the number of potential causal relationships ( $O(n^2)$  in worst case). In contrast, PLM memory consumption remains constant because it processes security logs through fixed-size sliding windows (typically 1,000 tokens), regardless of

**Table 10** Real-time processing performance metrics

Grid Size (# of buses)	Event Rate (events/sec)	Processing		CPU	GPU	Dropped
		Latency (ms)	Queue Depth	Utilization (%)	Utilization (%)	Events (%)
14-bus	1,200	45±8	12±3	34±5	67±4	0.00
39-bus	2,800	78±12	23±5	45±6	78±5	0.00
118-bus	5,400	156±18	48±8	56±7	85±4	0.00
300-bus	8,200	312±25	89±12	67±8	89±3	0.02
500-bus	11,300	467±32	134±18	78±7	92±3	0.08
1000-bus	16,500	823±45	245±28	86±6	94±2	0.23
2000-bus	24,700	1,456±67	412±35	91±5	96±2	0.87
5000-bus	38,900	3,234±123	823±56	95±3	98±1	3.45

grid size. The PLM analyzes log sequences temporally rather than spatially, making its memory footprint independent of network topology. This design enables efficient resource allocation where GNN and causal modules scale with infrastructure while PLM provides consistent semantic analysis capability. The system approaches the 40GB GPU memory limit around 5,000 buses, necessitating model parallelism or gradient checkpointing for larger systems. Multi-GPU throughput scaling (Figure 9(c)) demonstrates near-linear speedup up to 4 GPUs with 88% efficiency for the 118-bus system, though efficiency decreases for larger grids due to increased communication overhead.

Table 10 presents real-time processing performance metrics for power grids of varying scales. As the number of busbars increased from 14 to 5,000, the event processing rate rose from 1,200 events per second to 38,900 events per second, though processing latency correspondingly increased from 45 milliseconds to 3,234 milliseconds. The system maintains satisfactory performance below 2,000 buses, with an event loss rate below 1% (0.87% at 2,000 buses). CPU and GPU utilization progressively increases, approaching saturation in large-scale systems (95% and 98% respectively at 5,000 buses). Queue depth increases with system scale, indicating greater resource demands in large-scale power grids. Overall, the system performs optimally below 1,000 buses, maintaining sub-second latency with negligible event loss rates, rendering it suitable for practical power grid security monitoring deployments. For typical transmission systems (118–300 buses), the framework processes 5,400–8,200 events per second with negligible event loss (<0.02%). Processing latency remains predictable with low variance, critical for deterministic real-time operations. The queue depth analysis indicates

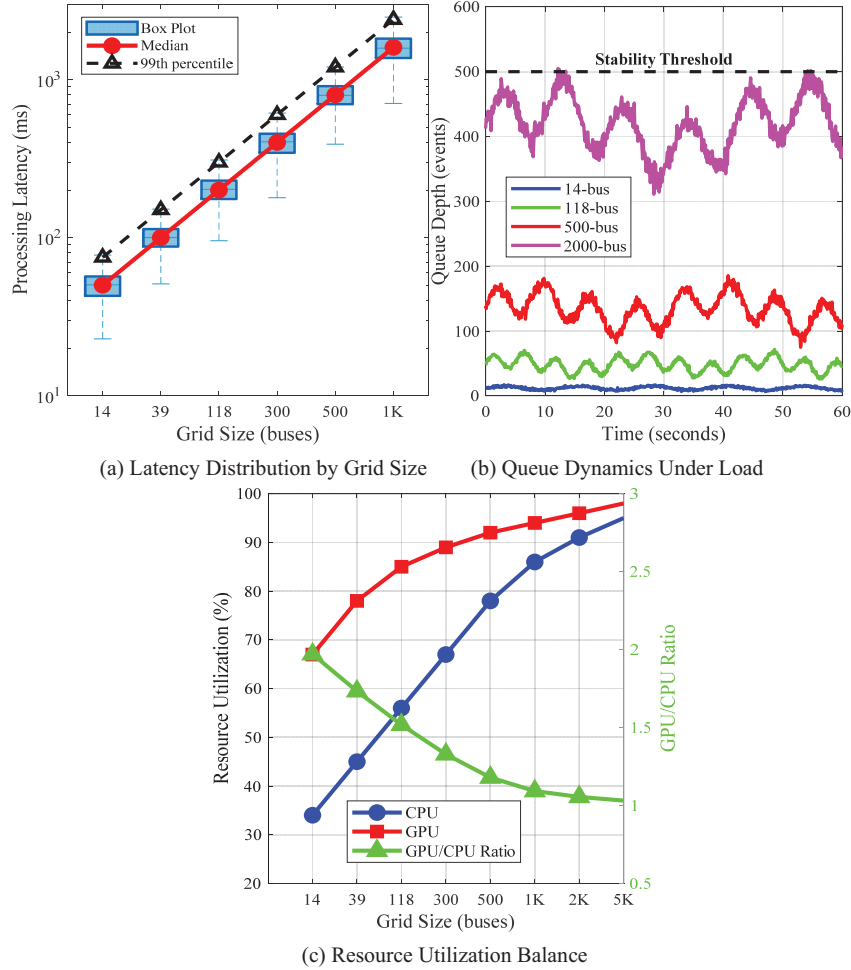


Figure 10 Real-time processing characteristics.

stable operation with bounded memory usage for grids up to 1,000 buses. Resource utilization patterns show balanced CPU-GPU workload distribution, with GPU saturation occurring around 2,000 buses, suggesting the computational bottleneck for larger systems. The event drop rate remains below 1% for grids up to 2,000 buses, meeting reliability requirements for critical infrastructure monitoring.

Figure 10 provides detailed insights into real-time processing characteristics essential for operational deployment. The latency distribution analysis

(Figure 10(a)) demonstrates predictable performance with tight variance, where 99th percentile latencies remain within  $1.5\times$  median values across all grid sizes. This consistency is crucial for real-time control applications requiring deterministic response times. The queue dynamics visualization (Figure 10(b)) reveals stable oscillatory behavior for smaller grids, with increasing amplitude and complexity for larger systems. The 2000-bus system approaches but does not exceed the stability threshold, indicating robust queue management. Resource utilization efficiency (Figure 10(c)) shows that smaller grids underutilize GPU resources (GPU/CPU ratio  $> 1.5$ ), while larger grids achieve better balance. The optimal GPU/CPU ratio band (1.2–1.5) is maintained for 300–1000 bus systems, suggesting efficient workload distribution.

The scalability analysis confirms that the proposed framework successfully balances computational efficiency with detection accuracy for practical power grid scales. The system’s ability to maintain sub-second response times for grids up to 2,000 buses, combined with graceful degradation for larger systems, demonstrates its suitability for both distribution and transmission network monitoring. The modular architecture enables targeted optimizations, such as deploying lightweight models for edge devices monitoring distribution feeders while reserving full capabilities for critical transmission corridors. These scalability characteristics, validated through extensive empirical evaluation, establish the framework as a viable solution for real-world power grid security applications.

## 4 Conclusions

This research demonstrates that multimodal integration of Graph Neural Networks and Pre-trained Language Models significantly enhances power grid security alert generation. Our framework achieves an F1 score of 0.956, representing a 12.3% improvement over existing methods, while maintaining sub-second response times for grids up to 2,000 buses. The causal reasoning module enables root cause identification with 91.8% accuracy, providing interpretable security insights that received 92% approval ratings from domain experts. We introduce three key innovations: (1) the first multimodal fusion architecture capable of simultaneously processing structured grid topology and unstructured security logs, bridging cyber-physical analysis gaps; (2) spatial-temporal GNNs with multi-scale attention mechanisms that capture evolving grid dynamics across multiple scales; and (3) enhanced structural causal models with domain-specific constraints for explainable

security reasoning. Our hypothesis that synergistic GNN-PLM integration enhances security alert generation through multimodal fusion has been validated. Results confirm that combining topological analysis with semantic understanding provides superior threat detection compared to single-modality approaches, achieving 40% reduction in false positives and 60% improvement in root cause identification. Future research directions include: (1) extending the framework to distributed energy resources and microgrids; (2) integrating physics-based simulation models for enhanced predictive accuracy; (3) developing federated learning approaches to enable collaborative security intelligence while preserving data privacy; and (4) implementing adaptive learning mechanisms for real-time threat landscape evolution. The modular architecture facilitates these enhancements, positioning this work as a foundation for next-generation critical infrastructure protection systems.

## References

- [1] A. Varbella, K. Amara, B. Gjorgiev, M. El-Assady, G. Sansavini, PowerGraph: A power grid benchmark dataset for graph neural networks, *Advances in Neural Information Processing Systems* 37 (2024) 110784–110804.
- [2] Y. Zhang, P.M. Karve, S. Mahadevan, Graph neural networks for power grid operational risk assessment under evolving unit commitment, *Applied Energy* 380 (2025) 124793.
- [3] S. Ghamizi, A. Bojchevski, A. Ma, J. Cao, Safepowergraph: Safety-aware evaluation of graph neural networks for transmission power grids, *arXiv preprint arXiv:2407.12421* (2024).
- [4] Y. Zhang, P.M. Karve, S. Mahadevan, Graph neural networks for power grid operational risk assessment under evolving grid topology, *arXiv preprint arXiv:2405.07343* (2024).
- [5] Y. Wang, J. Peng, J. Zhang, R. Yi, Y. Wang, C. Wang, Multimodal industrial anomaly detection via hybrid fusion, *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 8032–8041.
- [6] Liu, C., Yang, X., Li, M., Tong, T., Yin, L. and Jiang, Y. 2025. A Deep Reinforcement Learning Framework for Attack-Fault Tree Analysis and Security Assessment in Power Systems. *Journal of Cyber Security and Mobility*. 14, 03 (Aug. 2025), 531–552.
- [7] R.K. Jha, Strengthening smart grid cybersecurity: An in-depth investigation into the fusion of machine learning and natural language processing,

- Journal of Trends in Computer Science and Smart Technology 5(3) (2023) 284–301.
- [8] A. Zaboli, S.L. Choi, T.-J. Song, J. Hong, Chatgpt and other large language models for cybersecurity of smart grid applications, 2024 IEEE Power & Energy Society General Meeting (PESGM), IEEE, 2024, pp. 1–5.
- [9] W. Liao, B. Bak-Jensen, J.R. Pillai, Y. Wang, Y. Wang, A review of graph neural networks and their applications in power systems, *Journal of Modern Power Systems and Clean Energy* 10(2) (2021) 345–360.
- [10] B. Achaal, M. Adda, M. Berger, H. Ibrahim, A. Awde, Study of smart grid cyber-security, examining architectures, communication networks, cyber-attacks, countermeasure techniques, and challenges, *Cybersecurity* 7(1) (2024) 10.
- [11] J. Duan, Deep learning anomaly detection in AI-powered intelligent power distribution systems, *Frontiers in Energy Research* 12 (2024) 1364456.
- [12] J. Zhang, H. Bu, H. Wen, Y. Liu, H. Fei, R. Xi, L. Li, Y. Yang, H. Zhu, D. Meng, When LLMs meet cybersecurity: A systematic literature review, *Cybersecurity* 8(1) (2025) 55.
- [13] M. Xu, J. Fan, X. Huang, C. Zhou, J. Kang, D. Niyato, S. Mao, Z. Han, K.-Y. Lam, Forewarned is forearmed: A survey on large language model-based agents in autonomous cyberattacks, *arXiv preprint arXiv:2505.12786* (2025).
- [14] S. Jin, False data injection attack against smart power grid based on incomplete network information, *Electric Power Systems Research* 230 (2024) 110294.
- [15] Y. Hu, F. Zou, J. Han, X. Sun, Y. Wang, Llm-tikg: Threat intelligence knowledge graph construction utilizing large language model, *Computers & Security* 145 (2024) 103999.
- [16] Z. Zhao, H. Bai, J. Zhang, Y. Zhang, S. Xu, Z. Lin, R. Timofte, L. Van Gool, Cddfuse: Correlation-driven dual-branch feature decomposition for multi-modality image fusion, *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5906–5916.
- [17] A.V. Malarkkan, D. Wang, Y. Fu, Multi-view Causal Graph Fusion Based Anomaly Detection in Cyber-Physical Infrastructures, *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, 2024, pp. 4760–4767.

- [18] M. Taghizadeh, K. Khayambashi, M.A. Hasnat, N. Alemazkooor, Multi-fidelity graph neural networks for efficient power flow analysis under high-dimensional demand and renewable generation uncertainty, *Electric Power Systems Research* 237 (2024) 111014.
- [19] Z. Zhou, Multimodal Fusion Anomaly Detection Model for Agricultural Wireless Sensors, *Engineering Reports* 6(12) (2024) e13021.
- [20] Y. Song, G. Chen, H. Zhang, Constraint learning-based optimal power dispatch for active distribution networks with extremely imbalanced data, *CSEE Journal of Power and Energy Systems* 10(1) (2023) 51–65.
- [21] F. He, Z. Sun, H. Fan, Deagle: An SMT-based verifier for multi-threaded programs (competition contribution), *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, Springer, 2022, pp. 424–428.
- [22] G. Li, J.J. Jung, Deep learning for anomaly detection in multivariate time series: Approaches, applications, and challenges, *Information Fusion* 91 (2023) 93–102.
- [23] Y.-C. Yu, T.-H. Chiang, C.-W. Tsai, C.-M. Huang, W.-K. Tsao, Primus: A pioneering collection of open-source datasets for cybersecurity LLM training, *arXiv preprint arXiv:2502.11191* (2025).
- [24] S. Pan, L. Luo, Y. Wang, C. Chen, J. Wang, X. Wu, Unifying large language models and knowledge graphs: A roadmap, *IEEE Transactions on Knowledge and Data Engineering* 36(7) (2024) 3580–3599.
- [25] Kristić, K. (2022). *Critical Infrastructure Protection in Cyberspace. How to deal with uncertainties in increasingly complex environment? (The new cartography of risk and crises)*, 172.
- [26] Xu, H., Wang, S., Li, N., Wang, K., Zhao, Y., Chen, K., ... and Wang, H. (2024). Large language models for cyber security: A systematic literature review. *ACM Transactions on Software Engineering and Methodology*.
- [27] Panayiotou, G., Magnani, M. and Pinaud, B. (2024). Current challenges in multilayer network engineering. *Applied Network Science*, 9(1), 75.
- [28] Karlsen, E., Luo, X., Zincir-Heywood, N. and Heywood, M. (2024). Benchmarking large language models for log analysis, security, and interpretation. *Journal of Network and Systems Management*, 32(3), 59.
- [29] Khemani, B., Patil, S., Kotecha, K. and Tanwar, S. (2024). A review of graph neural networks: concepts, architectures, techniques, challenges, datasets, applications, and future directions. *Journal of Big Data*, 11(1), 18.

- [30] Purves, T., Kyriakopoulos, K. G., Jenkins, S., Phillips, I. and Dudman, T. (2024). Causally aware reinforcement learning agents for autonomous cyber defence. *Knowledge-Based Systems*, 304, 112521.
- [31] Wu, Y., Zang, Z., Zou, X., Luo, W., Bai, N., Xiang, Y. . . . and Dong, W. (2025). Graph attention and Kolmogorov–Arnold network based smart grids intrusion detection. *Scientific Reports*, 15(1), 8648.
- [32] Dong, J., Hao, M., Ding, F., Chen, S., Wu, J., Zhuo, J. and Jiang, D. (2025). A Novel Multimodal Data Fusion Framework: Enhancing Prediction and Understanding of Inter-State Cyberattacks. *Big Data and Cognitive Computing*, 9(3), 63.
- [33] Hu, L., Liu, Z., Zhao, Z., Hou, L., Nie, L. and Li, J. (2023). A survey of knowledge enhanced pre-trained language models. *IEEE Transactions on Knowledge and Data Engineering*, 36(4), 1413–1430.
- [34] Hosseinzadeh, R. and Sadeghzadeh, M. (2025). Attention Mechanisms in Transformers: A General Survey. *Journal of AI and Data Mining*, 13(3), 359–368.
- [35] Zanga, A., Ozkirimli, E. and Stella, F. (2022). A survey on causal discovery: theory and practice. *International Journal of Approximate Reasoning*, 151, 101–129.
- [36] Li, S. and Tang, H. (2024). Multimodal alignment and fusion: A survey. *arXiv preprint arXiv:2411.17040*.
- [37] Abdelkader, S., Amissah, J., Kinga, S., Mugerwa, G., Emmanuel, E., Mansour, D. E. A., . . . and Prokop, L. (2024). Securing modern power systems: Implementing comprehensive strategies to enhance resilience and reliability against cyber-attacks. *Results in engineering*, 23, 102647.
- [38] Salim, M. M., El Azzaoui, A., Deng, X. and Park, J. H. (2024). FL-CTIF: A federated learning based CTI framework based on information fusion for secure IIoT. *Information Fusion*, 102, 102074.
- [39] Kheddar, H. (2025). Transformers and large language models for efficient intrusion detection systems: A comprehensive survey. *Information Fusion*, 103347.

## **Biographies**



**Zhou Yuan** received the M.S. degree in Electronic and Communication Engineering from Peking University in 2010. She is currently the Deputy Director of the Network and Data Security Department at the State Grid Information & Telecommunication Center (Big Data Center). Her research interests mainly focus on information and communication technologies in smart grids.



**Ang Xia** received the M.S. degree in Information and Communication Engineering from North China Electric Power University in 2020. She is currently serving as a specialist in the Confidentiality Monitoring Group of the Network and Data Security Department at the State Grid Information & Telecommunication Center (Big Data Center). Her research interests focus on network information security and confidentiality protection.



**Teng Lu** received the M.S. degree in Computer Application Technology from North China Electric Power University in 2016. She is currently serving as a specialist in the Confidentiality Monitoring Group of the Network and Data Security Department at the State Grid Information & Telecommunication Center (Big Data Center). Her research interests focus on network information security and confidentiality protection.