

Adaptive Latency Reduction in IoT-Enabled Smart Grid Architectures

Dr. Arisara Chanthaburi¹, Wei-Hao Lin²

^{1,2} King Mongkut's Institute of Technology Ladkrabang, Thailand / National Taiwan University, Taiwan

Abstract

The transition to a decentralized power grid, driven by the high penetration of Distributed Energy Resources (DERs) and massive IoT deployments, necessitates ultra-reliable, low-latency communication for real-time control and stability. Time-critical applications, particularly Fault Detection, Isolation, and Restoration (FDIR) and Wide-Area Monitoring Systems (WAMS), are fundamentally compromised by variable network latency. Existing solutions, including static-edge computing architectures, fail to provide adaptive performance, as they are typically oblivious to the dynamic state of both the physical grid and the communication network. This paper introduces the Adaptive Latency Reduction Smart Grid Architecture (ALR-SGA), a novel distributed intelligence framework. ALR-SGA leverages a two-part mechanism: (1) a Graph Neural Network (GNN) model that performs topology-aware prediction of network congestion by fusing traffic data with physical grid-state information, and (2) a distributed Deep Reinforcement Learning (DRL) agent at each edge node that makes autonomous data routing and processing decisions. The DRL agent's reward function is uniquely tied to grid stability metrics, enabling it to prioritize data streams based on their real-time impact on grid operations. Through co-simulation of the IEEE 39-bus system (MATPOWER) and a detailed communication network (ns-3), we demonstrate that ALR-SGA significantly outperforms benchmark models. Results show a reduction in end-to-end latency for critical FDIR messages by over 35% during high-congestion fault scenarios, providing a viable architecture for next-generation, resilient grid control.

1. Introduction

The traditional electric power grid is undergoing a fundamental paradigm shift, evolving from a centralized, unidirectional system into a decentralized, bi-directional, and data-intensive cyber-physical system [1]. This transformation is characterized by the large-scale integration of Distributed Energy Resources (DERs) such as solar, wind, and battery storage, alongside the massive deployment of Internet of Things (IoT) devices, including Phasor Measurement Units (PMUs), smart meters, and remote actuators [2]. This new architecture unlocks unprecedented potential for efficiency and flexibility, but it also creates an absolute dependency on high-speed, reliable, and intelligent communication networks for monitoring and control.

In this decentralized paradigm, low-latency data is not merely a performance metric but a prerequisite for grid stability. Many critical grid operations are executed within strict, sub-second time windows. For instance, Fault Detection, Isolation, and Restoration (FDIR) protocols require the detection and isolation of a fault in milliseconds to prevent cascading failures [3]. Similarly, Wide-Area Monitoring Systems (WAMS) and associated frequency regulation services rely on the continuous, time-synchronized stream of PMU data to maintain grid-wide stability [4]. Any significant delay or jitter in the delivery of this critical data can lead to inaccurate state estimation, delayed control actions, and, in the worst case, systemic collapse.

While the limitations of centralized cloud computing for real-time applications are well-established, the initial-generation edge and fog computing architectures also present significant-yet-subtle shortcomings [5]. These solutions reduce propagation delay by moving computation closer to the data source; however, they largely operate on static, pre-defined resource allocation and Quality of Service (QoS) policies. Such a static approach is fundamentally misaligned with the stochastic nature of the smart grid. These systems are ill-equipped to handle the *simultaneous* occurrence of a physical grid event (e.g., a transmission line fault, which suddenly generates a high-priority burst of PMU data) and a network-level event (e.g., traffic congestion from non-critical meter data). Existing edge platforms lack the requisite intelligence to autonomously identify and prioritize data based on its *contextual importance* to the physical grid's state.

To address this critical gap, this paper introduces the Adaptive Latency Reduction Smart Grid Architecture (ALR-SGA). ALR-SGA is a novel, fully distributed framework that imbues edge nodes with autonomous decision-making capabilities, making them *grid-state-aware*. Our architecture is built upon a synergistic combination of advanced predictive and decision-making models. First, we employ a Graph Neural Network (GNN) to model the cyber-physical topology, enabling it to *predict* incipient network congestion by correlating network traffic patterns with physical grid telemetry. This predictive state information is then fed into a distributed Deep Reinforcement Learning (DRL) agent located at each edge node. The DRL agent is trained to make optimal, real-time decisions regarding data routing (e.g., which neighbor node to offload to) and processing. Crucially, the DRL's reward function is engineered to optimize for physical grid stability metrics, not just network throughput, thereby ensuring

that critical FDIR and WAMS data is prioritized during contingencies.

The primary contributions of this work are threefold:

1. **A Novel Grid-State-Aware Architecture:** The design of the ALR-SGA, a fully distributed edge intelligence framework that uniquely combines GNN-based predictive models with DRL-based autonomous agents for dynamic data routing.
2. **Contextual Prioritization Mechanism:** A system where the *criticality* of IoT data is not static but is dynamically assessed and prioritized by the DRL agent based on real-time grid events, ensuring stability-critical data bypasses congestion.
3. **Rigorous Co-Simulation Validation:** A comprehensive performance evaluation using a high-fidelity co-simulation (MATPOWER and ns-3). We provide quantitative analysis demonstrating ALR-SGA's superior latency reduction and reliability for critical FDIR messages compared to static-edge and cloud-based benchmarks.

The remainder of this paper is organized as follows: Section 2 reviews related work in edge computing for smart grids and AI-based network control. Section 3 details the system model and formal problem formulation. Section 4 presents the ALR-SGA architecture in detail. Section 5 describes the co-simulation setup and discusses the performance results. Finally, Section 6 concludes the paper and outlines directions for future research.

2. Related Work

The evolution of smart grid communication architectures has been the subject of extensive research; however, the intersection of ultra-low latency networking and physics-aware control remains a contested area. This section reviews existing literature across three thematic pillars: IoT-Edge orchestration, low-latency protocols, and AI-driven network control, highlighting specific limitations that motivate the proposed ALR-SGA framework.

2.1 Limitations of Static IoT-Edge Orchestration

The migration from cloud-centric to edge-centric processing is well-documented as a necessity for reducing propagation delay in smart grids. Recent surveys indicate that while fog and edge computing successfully offload computation from the core network, the orchestration strategies employed remain largely static or heuristic-based [1]. Standard edge load-balancing algorithms, such as Round-Robin or Min-Min, typically optimize for generic metrics like CPU utilization or throughput without "grid-state awareness." For instance, existing architectures often treat a data packet from a smart meter (billing data) with the same priority as a packet from a Phasor Measurement Unit (PMU) indicating a voltage sag, provided they arrive at the same time. Research by *Cao et al. (2023)* demonstrates that while dynamic task offloading can reduce general latency, it frequently fails under "bursty" traffic conditions typical of grid faults because the offloading policies do not account for the underlying physical urgency of the data [2]. The lack of cross-layer visibility—where the network layer is blind to the power system's physical state—remains a critical deficiency in current edge orchestration models.

2.2 Protocol Constraints: IEC 61850 and SDN

IEC 61850, particularly the GOOSE (Generic Object Oriented Substation Events) protocol, is the industry standard for substation automation, requiring end-to-end latency under 4ms for trip signals. However, mapping GOOSE over wide-area IP networks (Routed-GOOSE) introduces significant jitter and non-deterministic delays that standard TCP/IP stacks cannot mitigate [3]. To address this, Software-Defined Networking (SDN) has been proposed to provide deterministic QoS by separating the control and data planes. *Mozayani and Vali (2024)* demonstrated that SDN could improve link recovery times in smart grids using centralized machine learning controllers [4]. However, the centralized nature of SDN controllers creates a single point of failure and a latency bottleneck during rapid, multi-point topology changes. In highly volatile grid scenarios, where communication topology must change milliseconds after a physical line fault, the round-trip time (RTT) required for an SDN controller to install new flow rules is often prohibitive.

2.3 AI in Network Control: The "Energy vs. Latency" Bias

The application of Reinforcement Learning (RL) in smart grids has surged, yet a review of the literature reveals a distinct bias toward *energy* optimization over *latency* control. The vast majority of Deep Reinforcement Learning (DRL) applications focus on Demand Response (DR), economic dispatch, or battery arbitrage, where the decision timescales are in minutes or hours [5]. While some recent works, such as *Munikoti et al. (2023)*, have explored the intersection of Graph Neural Networks (GNN) and RL for grid topology analysis, these are primarily utilized for identifying power flow vulnerabilities rather than optimizing real-time data packet routing [6]. There is a scarcity of research that applies DRL specifically to the *millisecond-level* routing of FDIR data, and even fewer that utilize GNNs to predict communication congestion based on physical grid anomalies.

2.4 Research Gap

Synthesizing the above, a clear gap exists for a **physics-aware, fully distributed network control architecture**. Current state-of-the-art solutions are either: (1) *Static* (unable to adapt to sudden fault-induced traffic bursts), (2) *Centralized* (SDN approaches that introduce bottleneck latency), or (3) *Energy-Focused* (AI models optimizing for cost rather than speed). ALR-SGA addresses this gap by embedding a GNN-based predictor and a DRL-based router directly at the edge, creating a system that is both decentralized and explicitly driven by physical grid stability metrics.

3. System Model and Problem Formulation

In this section, we model the IoT-enabled smart grid as a dynamic cyber-physical system. We define the network topology, the heterogeneous traffic characteristics of grid devices, and formulate the latency minimization problem as a constrained optimization task.

3.1 Network and Communication Model

We consider a hierarchical, three-tier smart grid communication architecture consisting of an IoT device layer, an Edge/Fog layer, and a Cloud layer. The network is modeled as a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} represents the set of nodes and \mathcal{E} represents the set of communication links.

The node set \mathcal{V} is the union of three disjoint sets: $\mathcal{V} = \mathcal{V}_{IoT} \cup \mathcal{V}_{Edge} \cup \{v_{Cloud}\}$.

- $\mathcal{V}_{IoT} = \{1, \dots, N\}$ represents sensors and actuators (e.g., PMUs, Smart Meters).
- $\mathcal{V}_{Edge} = \{1, \dots, M\}$ represents edge computing nodes co-located with substations or aggregators.
- v_{Cloud} represents the central control center.

Each link $(i, j) \in \mathcal{E}$ is characterized by a bandwidth W_{ij} and a propagation latency d_{ij}^{prop} . The achievable transmission rate $R_{ij}(t)$ at time slot t is modeled according to the Shannon-Hartley theorem, considering time-varying channel gain $h_{ij}(t)$ and noise power σ^2 :

$$R_{ij}(t) = W_{ij} \log_2 \left(1 + \frac{P_i^{tx} |h_{ij}(t)|^2}{\sigma^2} \right)$$

where P_i^{tx} is the transmission power of node i (Mao et al., 2017).

3.2 Heterogeneous Traffic Model

The grid traffic is heterogeneous, characterized by distinct requirements for latency and reliability. We classify the traffic generated by node $i \in \mathcal{V}_{IoT}$ into two distinct classes:

Class A: Periodic Monitoring Data (Delay-Tolerant)

Generated by Smart Meters for billing and load profiling. These tasks arrive with a predictable generation rate λ_i^P . A task $\tau_{i,k}^P$ is defined as a tuple $\{L_k, C_k, T_k^{max}\}$, where L_k is the data size (bits), C_k is the computational cycles required, and T_k^{max} is the loose deadline (typically minutes).

Class B: Event-Driven Critical Data (Delay-Sensitive)

Generated by PMUs or protection relays during grid anomalies (e.g., voltage sags, line faults). This traffic follows a bursty arrival process (often modeled as a Poisson process during

normal operation, but spiking during faults) (Zhou et al., 2019). A critical task $\tau_{i,k}^C$ has a strict deadline $T_k^{crit} \ll T_k^{max}$ (typically ≤ 20 ms for GOOSE messages).

We define a priority indicator variable $\phi_k \in \{0, 1\}$, where $\phi_k = 1$ denotes Class B (Critical) traffic.

3.3 Computation and Latency Model

For a given task τ_k generated at node i , the total latency T_k^{total} depends on the offloading decision. The system can process the task locally, offload to a neighboring edge node, or offload to the cloud.

If processed at Edge Node j with CPU frequency f_j , the processing delay is:

$$T_{k,j}^{proc} = \frac{C_k}{f_j}$$

The transmission delay to move data from node i to node j is:

$$T_{i,j}^{trans} = \frac{L_k}{R_{ij}(t)}$$

The total latency for task k routed through path \mathcal{P} is the sum of transmission, propagation, queueing, and processing delays:

$$T_k^{total} = \sum_{(u,v) \in \mathcal{P}} \left(\frac{L_k}{R_{uv}(t)} + d_{uv}^{prop} + d_u^{queue}(t) \right) + T_{v,dest}^{proc}$$

Where $d_u^{queue}(t)$ represents the waiting time in the buffer of node u , which is highly dynamic and dependent on current network congestion (Liu et al., 2019).

3.4 Problem Formulation

The objective of the ALR-SGA architecture is to minimize the long-term system cost, defined as a weighted sum of latency and packet loss, with a heavy penalty for violating critical deadlines.

Let $\mathbf{x}_k \in \mathcal{V}_{Edge} \cup \{v_{Cloud}\}$ be the offloading decision variable for task k . We define a cost function J_k for task k :

$$J_k = \alpha \cdot T_k^{total} + \beta \cdot \mathbb{I}(T_k^{total} > T_k^{deadline}) \cdot \phi_k$$

Here, α is the weight for latency, and β is a large penalty factor for deadline violations of critical tasks $\phi_k=1$. $\mathbb{I}(\cdot)$ is the indicator function.

The global optimization problem is to find the optimal policy π that maps the system state $S(t)$ to action $A(t)$ (routing and offloading decisions) to minimize the expected average cost:

$$\min_{\pi} \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T \sum_{k \in \mathcal{K}(t)} J_k(t)$$

Subject to:

1. **Stability Constraint:** $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T Q_i(t) < \infty, \forall i \in \mathcal{V}$ (Queue stability).
2. **Energy Constraint:** $E_i(t) \leq E_i^{max}$ (Battery/Power limits for IoT nodes).
3. **Bandwidth Constraint:** $\sum_k R_{i,j}^k(t) \leq W_{ij}$.

This problem is non-convex and NP-hard due to the combinatorial nature of discrete routing decisions and the stochasticity of channel conditions and traffic arrivals (Ren et al., 2020). Conventional optimization techniques are too slow for real-time grid control; thus, we propose a DRL-based approach in the subsequent section.

4. Proposed ALR-SGA Methodology

The core of the ALR-SGA framework is a hierarchical control loop that operates at the network edge. This framework is composed of two coupled subsystems: (1) A **Topology-Aware Congestion Prediction Module** based on Graph Neural Networks (GNN), and (2) A **Distributed Decision Module** based on Deep Reinforcement Learning (DRL). This section details the prediction module, which provides the necessary state-space anticipation for the DRL agent.

4.1 Module I: Topology-Aware Congestion Prediction via GNN

Traditional edge routing protocols (e.g., OSPF, RIP) rely on historical averages or instantaneous snapshots of link metrics, often failing to capture the complex, non-Euclidean spatial dependencies of a smart grid network (Jiang et al., 2021). To address this, we employ a Message Passing Neural Network (MPNN) architecture. This GNN variant is designed to learn a low-dimensional embedding of the network state that encodes both the communication traffic

patterns and the underlying physical grid stability.

4.1.1 Feature Engineering (Cyber-Physical Fusion)

We define the input feature vector $\mathbf{x}_i(t)$ for each node $i \in \mathcal{V}$ at time t as a concatenation of cyber-layer and physical-layer metrics. This fusion is critical; a voltage sag in the physical grid is a precursor to a burst of PMU traffic in the cyber network.

$$\mathbf{x}_i(t) = [q_i(t), \rho_i(t), V_i^{mag}(t), \Delta f_i(t)]$$

Where:

- $q_i(t)$: Current packet queue length at node i .
- $\rho_i(t)$: Average link utilization of connected edges.
- $V_i^{mag}(t)$: Local bus voltage magnitude (p.u.).
- $\Delta f_i(t)$: Local frequency deviation (Hz).

4.1.2 Graph Convolution and Aggregation

The GNN operates on the graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ defined in Section 3. At each GNN layer l (where $l = 1, \dots, L$), node i aggregates information from its immediate neighbors $\mathcal{N}(i)$. We adopt the GraphSAGE inductive learning framework (Hamilton et al., 2017) to ensure scalability to large dynamic graphs.

The aggregation step for layer l is defined as:

$$\mathbf{h}_{\mathcal{N}(i)}^{(l)} = \text{AGGREGATE}^{(l)} \left(\{\mathbf{h}_j^{(l-1)}, \forall j \in \mathcal{N}(i)\} \right)$$

Here, $\mathbf{h}_j^{(l-1)}$ is the hidden state (embedding) of neighbor j from the previous layer. $\mathbf{h}_i^{(0)} = \mathbf{x}_i(t)$ is the initial input feature vector. We utilize a *mean-pooling* aggregator to capture the average congestion level of the local neighborhood.

Following aggregation, the node updates its own embedding:

$$\mathbf{h}_i^{(l)} = \sigma \left(\mathbf{W}^{(l)} \cdot [\mathbf{h}_i^{(l-1)} || \mathbf{h}_{\mathcal{N}(i)}^{(l)}] \right)$$

Where:

- \parallel denotes the concatenation operation.
- $\mathbf{W}^{(l)}$ is the trainable weight matrix for layer l .
- $\sigma(\cdot)$ is a non-linear activation function (ReLU).

4.1.3 Spatiotemporal Prediction Output

After L layers of message passing, the final embedding $\mathbf{h}_i^{(L)}$ captures the structural information of the L-hop neighborhood. This embedding is passed through a Multi-Layer Perceptron (MLP) readout function to predict the latency state for the next time step $t + 1$

$$\hat{y}_i(t + 1) = \text{MLP}(\mathbf{h}_i^{(L)})$$

Here, $\hat{y}_i(t + 1)$ represents the predicted nodal processing latency and local congestion level. This predictive state allows the system to anticipate bottlenecks before packets are dropped. The loss function for training the GNN is the Mean Squared Error (MSE) between the predicted latency \hat{y} and the ground-truth latency y observed during offline training phases:

$$\mathcal{L}_{GNN} = \frac{1}{N} \sum_{i=1}^N (y_i(t + 1) - \hat{y}_i(t + 1))^2$$

The output $\hat{y}_i(t + 1)$, combined with the local queue state, forms the augmented state space \mathcal{S} for the DRL agent described in the next section. This specific coupling enables the "Adaptive" nature of the ALR-SGA, distinguishing it from static optimization methods (Rusek et al., 2019).

4.2 Module II: Distributed DRL-based Resource Allocation

While the GNN module provides *situational awareness* (predicting future congestion), the decision-making capability is governed by the **Distributed Resource Controller (DRC)**. We model the resource allocation problem as a Decentralized Partially Observable Markov Decision Process (Dec-POMDP), where each edge node acts as an independent agent learning a policy π_θ to maximize a global reward signal centered on grid stability.

4.2.1 MDP Formulation

For an agent at edge node i , the decision process at time slot t is defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \gamma)$:

1. State Space (\mathcal{S}):

The state s_t^i captures both the local observations and the predictive insights from the GNN module defined in Section 4.1.

$$s_t^i = \{\mathcal{K}_i(t), E_i(t), \mathbf{q}_{\mathcal{N}(i)}(t), \hat{\mathbf{y}}_{\mathcal{N}(i)}(t+1)\}$$

Where:

- $\mathcal{K}_i(t)$: The set of active tasks in the local buffer, characterized by their priority ϕ_k (Critical/Non-Critical).
- $E_i(t)$: The current residual energy of the edge node.
- $\mathbf{q}_{\mathcal{N}(i)}(t)$: The current queue lengths of neighboring nodes.
- $\hat{\mathbf{y}}_{\mathcal{N}(i)}(t+1)$ The **predicted latency** of neighbors for the *next* time step, derived from the GNN module. This predictive term allows the agent to avoid routing data to nodes that are about to become congested (Chen et al., 2021).

2. Action Space (\mathcal{A}):

The action a_t^i determines the offloading strategy for the head-of-line task k . The action space is discrete:

$$a_t^i \in \{0, 1, \dots, |\mathcal{N}(i)|, |\mathcal{N}(i)| + 1\}$$

- $a_t^i = 0$: **Local Execution.** Process task k on node i .
- $a_t^i \in \{1, \dots, |\mathcal{N}(i)|\}$ **Edge Offloading.** Transmit task k to neighbor j .
- $a_t^i = |\mathcal{N}(i)| + 1$: **Cloud Offloading.** Transmit task k to the central cloud (fallback for high-compute tasks).

3. Grid-Stability-Aware Reward Function (\mathcal{R}):

The design of the reward function is the most critical component for achieving "grid-state awareness." Unlike standard offloading schemes that minimize average latency, our reward function r_t^i is non-linear and priority-weighted:

$$r_t^i = - [\omega_1 \cdot (T_k^{total}) \cdot (1 + \eta\phi_k) + \omega_2 \cdot P_{drop} + \omega_3 \cdot E_{cost}]$$

Here:

- T_k^{total} : The realized end-to-end latency of the task.
- $\phi_k \in \{0, 1\}$: The priority indicator (1 for critical FDIR/PMU data).
- η : A large penalty factor (e.g., $\eta=10$) applied strictly to critical data. This forces the agent to learn that **delaying a critical packet is significantly worse than delaying a routine metering packet**.
- P_{drop} : Penalty for packet loss or deadline violation.
- E_{cost} : Energy consumption cost.
- $\omega_1, \omega_2, \omega_3$: Weighting coefficients.

4.2.2 Optimization Algorithm: Proximal Policy Optimization (PPO)

To solve for the optimal policy π_{θ}^* , we utilize Proximal Policy Optimization (PPO), a policy gradient method selected for its stability and sample efficiency compared to DQN (Schulman et al., 2017). PPO avoids large, destabilizing policy updates by clipping the objective function.

The agent maximizes the clipped surrogate objective:

$$L^{CLIP}(\theta) = \mathbb{E}_t \left[\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right]$$

Where \hat{A}_t is the generalized advantage estimation (GAE), which measures how much better an action is compared to the average action in that state. The "clip" function ensures the new policy does not deviate excessively from the old policy, ensuring stable convergence even in the highly volatile environment of a faulted power grid (Wang et al., 2020).

By integrating the GNN's future predictions into the PPO state space, the ALR-

SGA agent transitions from *reactive* routing (responding to current congestion) to *proactive* routing (avoiding future bottlenecks), thereby securing the sub-cycle latency required for FDIR.

5. Performance Evaluation

To rigorously validate the ALR-SGA framework, we developed a high-fidelity Cyber-Physical System (CPS) co-simulation platform. This platform synchronizes the continuous dynamics of the power grid with the discrete-event behavior of the communication network, allowing us to capture the cascading effects of latency on grid stability.

5.1 Experimental Setup

The co-simulation architecture integrates three distinct simulation environments via a ZeroMQ-based middleware for real-time data exchange:

1. **Physical Layer (Power Grid):** We utilize the **IEEE 39-Bus System** (New England Power System), simulated in **MATPOWER** and **OpenDSS**. The system comprises 10 generators and 29 load buses. We introduce stochastic load profiles based on real-world data traces from the NYISO dataset to simulate realistic demand fluctuations (Zimmerman et al., 2011).
2. **Cyber Layer (Communication Network):** The communication topology is modeled in **NS-3** (Network Simulator 3). We overlay a mesh communication network where each IEEE 39-bus node is equipped with an edge computing server. The links are modeled as fiber-optic connections with bandwidths varying between 100 Mbps and 1 Gbps to simulate heterogeneity. We utilize the UDP protocol for critical PMU data streams (Class B) to minimize overhead, consistent with IEC 61850-90-5 standards.
3. **Decision Layer (AI Agent):** The ALR-SGA agents (GNN predictor and PPO router) are implemented in **PyTorch**. The agents interact with the NS-3 environment every transmission time interval (TTI = 10ms).

5.2 Baselines and Scenarios

We compare the performance of ALR-SGA against three state-of-the-art benchmark algorithms:

1. **Cloud-Only Processing:** A centralized approach where all data is offloaded to a remote cloud server. This represents the legacy SCADA architecture.
2. **Static Edge (Greedy):** A heuristic approach where tasks are offloaded to the edge node with the highest currently available CPU capacity, without considering network congestion or future grid states.
3. **Round-Robin (RR):** A load-balancing scheme that distributes tasks sequentially among neighboring nodes.

Simulation Scenario: We simulate a "Cascading Fault" scenario. At $t=50s$, a three-phase short circuit is introduced at Bus 16. This triggers a burst of high-priority protection messages (Class B) from adjacent PMUs, simultaneously with heavy background traffic (Class A) from smart meters.

5.3 Results and Discussion

5.3.1 End-to-End Latency Analysis

Figure 3 illustrates the Cumulative Distribution Function (CDF) of the end-to-end latency for critical Class B (FDIR) packets during the fault interval ($t=50s$ to $t=60s$).

The results demonstrate that ALR-SGA achieves the lowest tail latency. Specifically, the 99th percentile latency for ALR-SGA is 18.4 ms, which is well within the 20ms requirement for critical GOOSE messages (Standard IEC 61850). In contrast, the Static Edge and Cloud-Only approaches exhibit 99th percentile latencies of 42.1 ms and 105.3 ms, respectively.

The Cloud-Only model suffers from inherent propagation delays, rendering it unsuitable for real-time protection. The Static Edge model fails because it greedily offloads tasks to computationally free nodes, often routing traffic through links that are already congested. ALR-SGA outperforms these baselines by utilizing the GNN module to anticipate link congestion induced by the fault event and proactively rerouting packets through less congested paths before the buffer overflow occurs.

5.3.2 Packet Delivery Ratio (PDR) under Congestion

Figure 4 depicts the Packet Delivery Ratio (PDR) as the network traffic load increases from 10 Mbps to 100 Mbps.

As the network load surpasses 60 Mbps, the PDR of the Round-Robin and Static Edge schemes degrades rapidly, dropping below 85%. This degradation occurs because these algorithms treat all data packets equally. ALR-SGA, however, maintains a **PDR of >98.5% for Critical Class B data** even under heavy congestion (90 Mbps). This resilience is a direct result of the DRL reward function (R), which heavily penalizes the dropping of critical packets ($\phi_k=1$). The agent learns to intentionally drop or delay non-critical Class A data (smart meter readings) to preserve bandwidth for stability-critical messages, a behavior we define as "Grid-State Awareness."

5.3.3 Convergence and Stability

The training convergence of the proposed PPO-based agent is compared against a standard Deep Q-Network (DQN) in Figure 5. The PPO algorithm stabilizes after approximately 1,200 episodes, whereas the DQN agent exhibits significant oscillation and requires over 2,500 episodes to reach a comparable reward level. The faster convergence of ALR-SGA is attributed to the inclusion of the GNN-predicted state $\hat{y}(t+1)$. By providing the agent with a "look-ahead" capability regarding network latency, the complexity

of the policy search space is effectively reduced, allowing the agent to learn optimal routing policies with fewer interactions with the environment.

6. Conclusion

The transition toward a decentralized, renewable-intensive smart grid is fundamentally contingent upon the reliability and responsiveness of its underlying communication infrastructure. As this paper has demonstrated, traditional static edge computing models and centralized cloud architectures are insufficient for meeting the sub-cycle latency requirements of critical grid applications like FDIR and WAMS.

In this work, we introduced the **Adaptive Latency Reduction Smart Grid Architecture (ALR-SGA)**, a novel distributed intelligence framework that bridges the gap between physical grid dynamics and cyber-network orchestration. By fusing Graph Neural Networks (GNN) for topology-aware congestion prediction with Deep Reinforcement Learning (DRL) for autonomous routing, ALR-SGA enables edge nodes to anticipate network bottlenecks and prioritize traffic based on physical grid stability.

Our comprehensive co-simulation, integrating MATPOWER and ns-3, provided quantitative validation of the proposed architecture. The results indicate that ALR-SGA achieves a **35% reduction in end-to-end latency** for critical protection messages during fault scenarios compared to static edge baselines. Furthermore, the system maintains a **Packet Delivery Ratio (PDR) exceeding 98%** for high-priority data, even under severe network congestion. These findings suggest that physics-aware, distributed AI is not merely an enhancement but a necessity for the resilient operation of next-generation cyber-physical power systems.

7. Future Work

While this study provides a robust theoretical and simulation-based framework, several avenues remain for future research to bridge the gap to practical deployment:

1. **Hardware-in-the-Loop (HIL) Implementation:** The current evaluation relies on software co-simulation. Future work will focus on porting the ALR-SGA inference models (GNN and DRL policy networks) onto real-world edge hardware, specifically **Field-Programmable Gate Arrays (FPGAs)** or edge TPUs. This will allow for the evaluation of inference latency and energy consumption constraints in a realistic hardware environment.
2. **Adversarial Robustness:** As the grid becomes increasingly dependent on AI-driven control, it becomes susceptible to adversarial attacks. A malicious actor could inject subtle noise into the traffic data (adversarial examples) to fool the GNN into misclassifying congestion states, leading to deliberate routing failures. Future research must investigate **adversarial training techniques** to enhance the robustness of the GNN and DRL models against such cyber-physical attacks.

References

- Bian, D., Kuzlu, M., Pipattanasomporn, M., & Rahman, S. (2015). Analysis of Communication Schemes for Advanced Metering Infrastructure in Smart Grid. *IEEE Transactions on Smart Grid*, 6(5), 2592-2601.
- Cao, H., Hu, Y., & Qu, Y. (2023). Dynamic Task Offloading and Resource Allocation in Mobile Edge Computing for Smart Grids. *IEEE Internet of Things Journal*, 10(14), 12345-12356.
- Chen, X., Zhang, H., Wu, C., Mao, S., & Ji, Y. (2021). Optimized Computation Offloading Performance in Virtual Edge Computing Systems Via Deep Reinforcement Learning. *IEEE Internet of Things Journal*, 8(12), 10053-10067.
- Deng, R., Yang, Z., Chow, M. Y., & Chen, J. (2015). A Survey on Demand Response in Smart Grids: Mathematical Models and Approaches. *IEEE Transactions on Industrial Informatics*, 11(3), 570-582.
- Gadia, S., & Parekh, H. (2024). A Comprehensive Survey on Fog Computing: Architectures, Applications, and Open Challenges. *IEEE Access*, 12, 10450-10478.
- Hamilton, W. L., Ying, Z., & Leskovec, J. (2017). Inductive Representation Learning on Large Graphs. *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS)*, 1024–1034.
- Hussain, S. M. S., & Ustun, T. S. (2019). Integrated software platform for cyber-physical analysis of smart grids. *IEEE Access*, 7, 96321-96333.
- Jiang, J., & Zhu, Y. (2021). Graph Neural Networks for Traffic Prediction in Smart Grids: A Survey. *IEEE Internet of Things Journal*, 8(15), 12034-12050.
- Li, H., Ota, K., & Dong, M. (2018). Learning IoT in Edge: Deep Learning for the Internet of Things with Edge Computing. *IEEE Network*, 32(1), 96-101.
- Liao, W., Bak-Jensen, B., Pillai, J. R., & Yang, Z. (2021). Graph Neural Networks for Power System Security Assessment. *IEEE Transactions on Power Systems*, 36(3), 2250-2260.
- Liu, L., Zhang, X., & Ma, H. (2019). Delay-Optimal Task Offloading in Fog-Enabled IoT Networks: A Deep Reinforcement Learning Approach. *IEEE Internet of Things Journal*, 6(5), 8278-8290.
- Mach, P., & Becvar, Z. (2017). Mobile Edge Computing: A Survey on Architecture and Computation Offloading. *IEEE Communications Surveys & Tutorials*, 19(3), 1628-1656.
- Mao, Y., You, C., Zhang, J., Huang, K., & Letaief, K. B. (2017). A Survey on Mobile Edge

Computing: The Communication Perspective. *IEEE Communications Surveys & Tutorials*, 19(4), 2322-2358.

Min, M., Wan, X., Xiao, L., Chen, Y., & Xia, M. (2019). Learning-Based Edge Computing for Delay Minimization in Smart Grids. *IEEE Access*, 7, 15822-15833.

Mozayani, N., & Vali, H. (2024). A Machine Learning-Based Approach for Link Recovery in Smart Grids Using Software-Defined Networking. *Journal of Resource Management and Decision Engineering*, 3(4), 116-135.

Munikoti, S., Agarwal, K., & Das, L. (2023). Graph Reinforcement Learning in Power Grids: A Survey. *IEEE Transactions on Smart Grid*, 14(5), 3456-3470.

Ren, J., Yu, G., He, Y., & Li, G. Y. (2020). Collaborative Cloud and Edge Computing for Latency Minimization. *IEEE Transactions on Vehicular Technology*, 69(8), 9006-9018.

Rusek, K., Suárez-Varela, J., Almasan, P., Barlet-Ros, P., & Cabellos-Aparicio, A. (2019). RouteNet: Leveraging Graph Neural Networks for Network Modeling and Optimization in SDN. *IEEE Journal on Selected Areas in Communications*, 38(10), 2260-2270.

Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., & Monfardini, G. (2008). The Graph Neural Network Model. *IEEE Transactions on Neural Networks*, 20(1), 61-80.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal Policy Optimization Algorithms. *arXiv preprint arXiv:1707.06347*.

Sutton, R. S., & Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. MIT press.

Ustun, T. S., & Hussain, S. M. S. (2020). Vulnerability and Impact Analysis of the IEC 61850 GOOSE Protocol in the Smart Grid. *Sensors*, 21(4), 1554.

Wang, J., Shao, Y., Ge, Y., & Yu, R. (2020). A Survey of Deep Reinforcement Learning in Smart Grids: A Comprehensive Review. *IEEE Transactions on Smart Grid*, 11(5), 4256-4269.

Zhang, D., Han, X., & Deng, C. (2022). Review on the application of deep reinforcement learning in smart grid energy management. *Energy Reports*, 8, 432-446.

Zhou, X., Wang, R., & Li, T. (2019). Traffic Characterization and Modeling for Smart Grid Communication Networks. *IEEE Transactions on Industrial Informatics*, 15(11), 6152-6161.

Zimmerman, R. D., Murillo-Sanchez, C. E., & Thomas, R. J. (2011). MATPOWER: Steady-State Operations, Planning, and Analysis Tools for Power Systems Research and Education. *IEEE Transactions on Power Systems*, 26(1), 12-19.