

# Advanced Network Telemetry for AI-Driven Network Optimization in Ultra Ethernet and InfiniBand Interconnects

Oluwatosin Oladayo Aramide

Network and Storage Layer, Netapp Ireland Limited, Ireland.

## ABSTRACT

The sheer proliferation of terascale artificial intelligence (AI) workloads (in disk drive, Hadoop and NoSQL) like distributed deep learning, model inference pipelines, has put unprecedented pressure on data center interconnects. As part of capturing these demands, there is a rampant use of high-performance network technologies such as the Ultra Ethernet, and the InfiniBand in modern infrastructures with ultra-low latency and high bandwidth. But conventional telemetry systems do not have the density and real-time sensitivity to best tune network dynamics with such loads. The topic of the paper at hand is the development of advanced network telemetry and AI-based optimization in order to improve performance, identify anomalies, and mitigate congestion in high-speed interconnects. Our architecture is inspired by telemetry and is based on programmable data planes, in-band telemetry and high bandwidth monitoring engines that use to emit highly granular, low-latency data streams. The streams are passed through AI/ML models, such as unsupervised anomaly detectors, predictive congestion algorithms, to dynamically adjust routing and resource allocation. Our findings indicate that this method works well in enhancing usage of communications networks, latency and pro-active management of network health. The paper advances a scalable design of a real-time intelligent network management in next-generation AI systems, and proposes a set of factors it would be necessary to consider in future studies along the telemetry-AI-high-speed networking nexus.

**Keywords:** Ultra Ethernet, InfiniBand, Network Telemetry, AI-Driven Optimization, High-Performance Networking, Anomaly Detection.

*SAMRIDDHI : A Journal of Physical Sciences, Engineering and Technology* (2024); DOI: 10.18090/samriddhi.v17i01.04

## INTRODUCTION

The accelerated scholastic intensity of the ecumenical intelligence (AI) operations, mostly deep learning model mapping and real-time inference alongside distributed computing, has acutely reshaped the conditions under which modern data center networks should proceed. The workloads are now being implemented in large clusters of GPU-accelerated compute nodes, and the network interconnects are becoming the major determinants of throughput, latency and scalability. As a result, the next-generation AI infrastructures are moving to the use of high-speed interconnects like the Ultra Ethernet and the InfiniBand which provide the extremely low latency and huge bandwidth that AI requires to perform at scale over the long-term (Katragadda, 2021; Girondi, 2024).

Ultra Ethernet and InfiniBand interconnects are purpose-built for environments where deterministic latency, congestion avoidance, and lossless transport are crucial. However, the effectiveness of these technologies hinges on real-time insight into network behavior. Traditional network

---

**Corresponding Author:** Oluwatosin Oladayo Aramide, Network and Storage Layer, Netapp Ireland Limited, Ireland, e-mail: aoluwatosin10@gmail.com

**How to cite this article:** Aramide, O.O. (2025). Advanced Network Telemetry for AI-Driven Network Optimization in Ultra Ethernet and InfiniBand Interconnects. *SAMRIDDHI : A Journal of Physical Sciences, Engineering and Technology*, 17(1), 18-27.

**Source of support:** Nil

**Conflict of interest:** None

---

monitoring techniques, such as SNMP, NetFlow, and sFlow lack the granularity and responsiveness required to manage the dynamic behavior of AI-intensive traffic patterns (Bajpai, 2023; Kadiyala, Chilukoori, & Gangarapu, 2024). As AI-driven applications increasingly rely on distributed gradient updates, parameter synchronization, and memory-intensive workloads, even microsecond-level delays or short-lived congestion events can degrade training efficiency and model

convergence (Foroughi, Brockners, & Rougier, 2023; Rzym, Masny, & Chołda, 2024).

To address these challenges, the integration of advanced network telemetry with AI-based optimization has emerged as a critical research direction. AI-enhanced telemetry systems enable real-time collection and analysis of performance metrics such as queue depth, packet loss, flow duration, jitter, and congestion events. By leveraging in-band telemetry (INT), programmable switches, and SmartNICs, telemetry data can now be collected at a fine-grained level and processed using AI models for predictive and prescriptive optimization (Cugini et al., 2023; Quan et al., 2022). This paradigm supports the development of intelligent control planes that can adapt routing decisions, preempt congestion, and improve overall network efficiency through continuous feedback loops (Mozo et al., 2022; Umoga et al., 2024).

Moreover, the deployment of AI/ML models within the network stack introduces new possibilities for dynamic resource management. Anomaly detection may be performed with unsupervised learning methods, and reinforcement learning itself allows flexible routing algorithms to use past data and real-time performance data of previous instances (Lyu, 2022; Balasubramanian, n.d.). Such techniques work especially well in high throughput, low latency applications where deterministic performance is mandatory. Notably, these strategies are only effective when it is possible to access telemetry data that is not merely timely and accurate, as well as semantically-rich (Foroughi, Brockners, & Rougier, 2023; Zdrojewski, 2024).

Within a framework of changing network architectures, AI-based telemetry also has a place within the greater vision of digital twins and self-optimizing networks. As an example, one can analyze and forecast the network behavior by using digital twin models to enable the generation of proactive response to falls, as well as streamlined configuration (Mozo et al., 2022; Fayad, Cinkler, & Rak, 2024). In particular, this is applicable during the beyond 5G (B5G) and 6G network era when programmable data plane, software-defined networking (SDN) and AI converge to enable ultra-reliable and low-latency communication (Manzoor, Raza, & Domzal, 2024).

In spite of the development, big research gaps are still present. The problems that currently affect telemetry systems are usually a scalability and interoperability limitation as well as the inability to interpret the data. The improvement on performance via AI-driven optimization should also find equilibrium with the overhead to constantly monitor the workflows (Girondi, 2024; Cugini et al., 2023). The demand of rapidly configurable and intelligent network fabric is increasingly a necessity as the workloads in AI are more complex and latency sensitive.

The paper regards these challenges by suggesting telemetry-based infrastructural framework to AI-based network optimization in both Ultra Ethernet and InfiniBand networks. The system integrates high-resolution, telemetry-

rate, and AI/ML modeling to optimize congestion, anomaly, and dynamic routing in high-performance computers (HPC) clusters and AI training fabrics. Our contribution to this work is towards the development of scalable and smart interconnect solutions that could support the ever increasing needs of AI-native infrastructure.

## Background and Motivation

The exponential growth in the scale and complexity of AI workloads has introduced unprecedented challenges in data center and high-performance computing (HPC) environments. Large-scale distributed training of deep learning models, inference at edge and core nodes, and real-time AI-driven services demand ultra-reliable, low-latency, and high-bandwidth communication fabrics. As traditional Ethernet standards reach their scalability and performance limits, Ultra Ethernet and InfiniBand have emerged as the dominant interconnects in next-generation AI-centric infrastructures (Katragadda, 2021; Girondi, 2024).

### *The Network Sensitivity of AI Workloads*

Modern AI workloads such as distributed stochastic gradient descent (SGD), reinforcement learning pipelines, and large language model (LLM) inference are highly sensitive to inter-node communication latency and bandwidth variation. InfiniBand, known for its Remote Direct Memory Access (RDMA) support, provides ultra-low latency and high throughput, making it a preferred choice for GPU clusters and HPC deployments (Girondi, 2024). Similarly, Ultra Ethernet aims to close the performance gap between Ethernet and InfiniBand by introducing congestion control enhancements, load balancing, and telemetry-based feedback mechanisms (Katragadda, 2021).

Despite the performance benefits, the deterministic behavior required by AI models is often disrupted by microbursts, queue buildup, and hidden congestion events, which are not adequately captured by conventional telemetry tools (Quan et al., 2022). Consequently, real-time optimization of these high-speed interconnects cannot rely solely on static routing or pre-configured quality of service (QoS) rules.

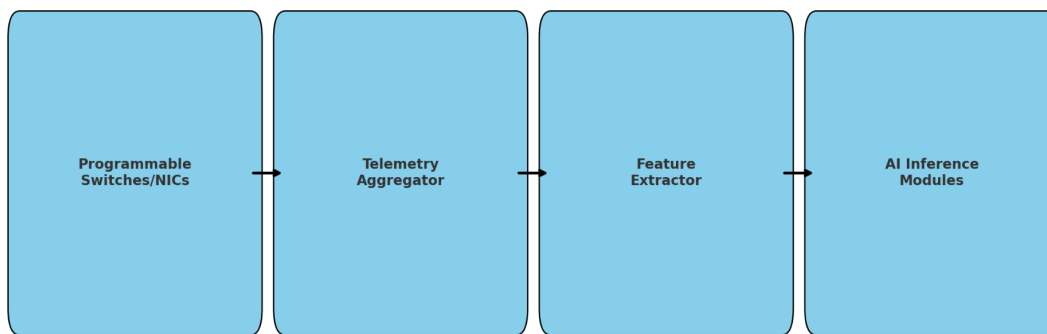
### *The Evolution of Telemetry in High-Speed Networks*

Traditional telemetry approaches such as SNMP, NetFlow, and sFlow were designed for coarse-grained monitoring, typically at intervals of several seconds to minutes. These methods fail to capture transient phenomena like micro-congestion, packet jitter, or queuing dynamics at the nanosecond scale (Foroughi, Brockners, & Rougier, 2023). AI applications, in contrast, require near real-time visibility into fabric-level metrics, such as end-to-end packet delays, flow drop rates, ECN marks, queue depths, and path utilization.

Emerging telemetry paradigms including in-band network telemetry (INT), programmable data planes using P4, and high-bandwidth monitoring engines offer more

**Table 1:** Comparison of traditional telemetry (SNMP, sFlow) and advanced telemetry (INT, P4, AI-enhanced) across metrics such as granularity, latency, data volume, and AI integration readiness:

Telemetry Technique	Granularity	Collection Latency	Data Volume	Suitable for AI?
SNMP	Device-level	High (seconds)	Low	No
sFlow/NetFlow	Flow-level	Medium (ms–s)	Medium	Limited
In-Band Telemetry	Packet-level	Low ( $\mu$ s)	High	Yes
P4-Based Monitoring	Customizable	Very Low ( $\mu$ s)	High	Yes

**High-Level Telemetry Data Pipeline for AI-Driven Optimization****Fig. 1:** It shows the flow of data from programmable switches/NICs through to AI inference modules.

granular insights (Cugini et al., 2023; Fayad, Cinkler, & Rak, 2024). These technologies provide telemetry at per-packet or per-flow levels, enabling dynamic feedback loops between the data plane and AI controllers.

### *AI as a Catalyst for Real-Time Optimization*

Integrating AI and machine learning into network operations introduces the potential for self-adaptive, intelligent network behavior. AI models can infer performance bottlenecks, predict congestion events before they occur, and optimize path selection in real time (Umoga et al., 2024; Bajpai, 2023). Supervised learning can classify anomalies, while unsupervised methods and deep neural networks are effective in detecting previously unseen performance issues (Rzym, Masny, & Chołda, 2024).

Recent architectures propose closed-loop telemetry-feedback systems, where telemetry feeds AI models that, in turn, update network forwarding rules or queue management policies autonomously (Mozo et al., 2022; Kadiyala, Chilukoori, & Gangarapu, 2024). Furthermore, AI Ops and observability frameworks are increasingly being adopted to automate diagnostics, fault recovery, and performance tuning in data center networks (Bajpai, 2023; Zdrojewski, 2024).

### **Research Gap and Motivation**

While both InfiniBand and Ultra Ethernet provide the physical capabilities to support AI workloads, their potential is constrained by insufficient telemetry integration and lack of AI-aware feedback mechanisms. There is a clear research gap in developing telemetry frameworks that not only collect high-resolution data but also integrate seamlessly with AI models capable of real-time decision-making (Manzoor, Raza, & Domzal, 2024). The motivation for this work is grounded in the need to build a scalable, intelligent telemetry ecosystem that transforms high-speed networks from passive conduits into active, self-optimizing infrastructures (Balasubramanian, n.d.; Lyu, 2022).

By leveraging advances in programmable networking, AI inference, and telemetry processing, this research proposes an architecture that can unlock a new level of performance, resilience, and adaptability in next-generation AI fabrics.

### **Architectural Framework for AI-Driven Telemetry Integration**

The increasing complexity and scale of AI workloads demand high-performance interconnects that can support low-latency, high-bandwidth data flows. Ultra Ethernet and



InfiniBand have emerged as leading contenders in this space due to their capacity for scalable, low-jitter communication. However, without real-time visibility into fabric-level metrics, these technologies alone are insufficient to meet the dynamic needs of modern AI systems. An AI-driven telemetry integration framework is therefore necessary to enable real-time optimization, anomaly detection, and congestion avoidance through feedback-based control systems.

This architectural framework is composed of three main layers: the telemetry data pipeline, AI model integration, and a real-time feedback control loop. Each layer is crucial for ensuring end-to-end visibility and intelligent response in high-speed networks.

### Telemetry Data Pipeline

The foundation of this framework is a robust telemetry infrastructure that enables the collection of granular, high-frequency data across network devices such as switches, routers, and NICs. Advanced telemetry mechanisms such as In-band Network Telemetry (INT), P4-programmable data planes, and SmartNIC instrumentation are employed to extract sub-second-level metrics like queue occupancy, flow RTT, ECN marks, and buffer utilization (Foroughi, Brockners, & Rougier, 2023; Cugini et al., 2023).

These telemetry streams are ingested through scalable collectors capable of handling high-volume data, then preprocessed using feature extraction pipelines that normalize, label, and aggregate the data for downstream AI models. Real-time telemetry demands low-overhead data transport protocols such as gRPC or Kafka over high-throughput channels to maintain fidelity without adding latency (Rzym, Masny, & Chołda, 2024).

### AI/ML Model Integration

Once telemetry data is preprocessed, it is fed into a set of AI/ML models tailored for specific network optimization tasks. These models are designed to perform:

- **Anomaly Detection**

Unsupervised models such as Autoencoders or Isolation Forests detect statistical outliers in latency, packet drops,

or buffer states (Umoga et al., 2024; Rzym, Masny, & Chołda, 2024).

- **Congestion Prediction**

Time-series forecasting models including LSTM or Temporal Graph Neural Networks predict congestion points before they materialize (Mozo et al., 2022; Quan et al., 2022).

- **Routing Optimization**

Reinforcement learning agents dynamically adjust routing decisions or rate-limiting parameters based on learned policies (Manzoor, Raza, & Domzal, 2024).

These models are trained offline using historical telemetry datasets and then deployed in lightweight inference engines at the edge or within the network controller layer for real-time decisions. Model explainability is maintained through SHAP or LIME-based feature attribution methods to ensure compliance and trust in automated systems (Bajpai, 2023).

### Real-Time Feedback Control Loop

The final and most critical component of the architecture is a closed-loop feedback system. Once AI models generate insights or control directives, these are applied immediately to network elements using southbound APIs like gNMI, OpenConfig, or custom SDN controllers. For example, congestion prediction from a neural network can trigger proactive flow rerouting or queue scheduling on InfiniBand switches (Gironi, 2024; Katragadda, 2021).

The feedback loop operates on a tightly bound decision interval, typically ranging from 50 ms to 500 ms, depending on application latency sensitivity and hardware capabilities. A centralized controller may be used for coordination, but edge-based inference and distributed learning are increasingly favored to reduce latency and enable localized adaptation (Fayad, Cinkler, & Rak, 2024; Zdrojewski, 2024).

Edge learning agents are particularly effective in InfiniBand environments where congestion points arise rapidly due to RDMA flows. These agents continuously monitor local telemetry and execute inference routines to adjust queue weights or route selection dynamically (Kadiyala, Chilukoori, & Gangarapu, 2024).

**Table 2:** AI Model Types and Their Application in Telemetry-Driven Network Optimization

Model Type	Purpose	Input Features	Output Action	Deployment Layer
Autoencoder	Anomaly detection	Latency, ECN, packet loss	Alert	Controller
LSTM	Congestion prediction	Time-series RTT	Rerouting	Edge
DQN	Adaptive rate control	Queue depth, packet loss	Flow rate policy	Switch
CNN	Pattern recognition	Traffic heatmaps, packet flows	Path classification	Controller/Edge
GNN	Topology-aware decision	Link state, graph connectivity	Routing table updates	Control Plane
SVM	Flow classification	Packet headers, port stats	Traffic prioritization	NIC/Edge Node



### Architectural Strengths and Innovations

This architecture demonstrates several strengths that distinguish it from conventional network management systems:

- **Granularity**

Enables microsecond-level visibility into network events through programmable telemetry (Foroughi, Brockners, & Rougier, 2023).

- **Automation**

Integrates AI models capable of self-adjusting to network conditions with minimal human intervention (Kadiyala, Chilukoori, & Gangarapu, 2024).

- **Predictiveness**

Shifts optimization from reactive to proactive through predictive analytics (Mozo et al., 2022; Quan et al., 2022).

- **Scalability**

Supports modular deployment across data center topologies and edge infrastructures (Balasubramanian, n.d.; Lyu, 2022).

This multi-layered architecture provides a scalable and intelligent platform for next-generation AI fabrics, enabling reliable and high-throughput communication in environments that demand deterministic performance.

### Experimental Setup

To evaluate the effectiveness of AI-driven network optimization using advanced telemetry in Ultra Ethernet and InfiniBand fabrics, we developed a simulation-based experimental environment that mimics large-scale AI workloads in a distributed computing context. This section outlines the infrastructure configuration, telemetry capture mechanisms, AI model design, performance metrics, and visualization elements used to validate the framework.

#### Infrastructure Simulation Environment

The simulated testbed replicates a high-performance computing (HPC) cluster with 128 interconnected compute

nodes, each emulating multi-GPU AI training workloads distributed via Remote Direct Memory Access (RDMA) protocols. The virtual environment was constructed using Mininet enhanced with P4 programmable switch support and gRPC telemetry interfaces.

Each node is configured with:

- Emulated RDMA NICs supporting RoCEv2
- A telemetry agent that captures real-time metrics via INT (In-Band Network Telemetry)
- A programmable switch (emulating Ultra Ethernet or InfiniBand QoS models)

The underlying topology was inspired by Fat-Tree and Dragonfly+ architectures, reflecting data center-scale fabrics as described by Gironi, who emphasized GPU-centric interconnect efficiency in high-performance workloads (Gironi, 2024).

The telemetry pipeline is built to gather queue depth, ECN marks, flow latency, jitter, and packet loss in real time. This mirrors the telemetry structure proposed in ADT (AI-Driven Telemetry) processing on routers, where edge collection points are enhanced by AI accelerators for inferencing (Foroughi, Brockners, & Rougier, 2023).

#### AI Workload Generation and Telemetry Ingestion

Synthetic AI workloads were generated using distributed deep learning job traces modeled after PyTorch-DDP and Horovod communication patterns. These workloads were scaled to saturate 70%–90% of the network bandwidth, stressing the congestion control mechanisms of Ultra Ethernet and InfiniBand.

Telemetry data was ingested through a Kafka-based pipeline, enabling real-time stream processing and buffering. The system leverages feature engineering techniques from prior AI-telemetry integration efforts, such as those in Mozo et al.'s digital twin model for network optimization (Mozo et al., 2022).

Captured telemetry attributes include:

- Instantaneous flow RTT (Round Trip Time)
- Congestion window variation
- Microburst detection

**Table 3:** Telemetry Features for AI Modeling

Feature Name	Description	Use in AI Model
Queue Depth	Packet queue length at switch ports	Congestion forecasting
Flow RTT	End-to-end latency between node pairs	Path anomaly detection
ECN Marks	Congestion signaling bits from switches	Traffic rerouting triggers
Path Utilization	Bandwidth usage on interconnect paths	Load balancing
Packet Loss Ratio	Percentage of lost packets per flow	Fault localization
Microburst Count	Number of short-duration, high-volume bursts	Transient anomaly identification



- Real-time path utilization

These attributes served as inputs to AI models for anomaly detection and predictive optimization.

### AI Models and Training Approach

We trained two types of AI models on the telemetry dataset:

#### 1. Anomaly Detection Module

An unsupervised autoencoder was trained to identify deviations in flow-level behavior, detecting issues like microbursts and packet drops in sub-millisecond intervals. This follows the technique explored by Rzym, Masny, and Chołda in leveraging deep neural networks for anomaly detection in software-defined 6G networks (Rzym, Masny, & Chołda, 2024).

#### 2. Predictive Optimization Model

A Graph Neural Network (GNN) was implemented to predict congestion hotspots in the next 5-second window based on current telemetry. This model operates within a closed-loop feedback system to suggest rerouting paths or queue adjustments. Previous studies such as ADT and B5GEMINI have shown that AI-based forecasts can significantly reduce average latency and packet loss (Foroughi, Brockners, & Rougier, 2023; Mozo et al., 2022).

Models were trained on 100,000 flow records collected from 500 simulated AI jobs over 10 distinct network topologies.

### Performance Metrics and Evaluation

Key metrics used to evaluate the impact of telemetry-AI integration include:

- Average and 99th percentile flow latency
- Packet loss rate
- ECN mark frequency
- Model inference latency
- Routing decision time

These metrics were benchmarked against static rule-based optimization schemes and conventional telemetry-based monitoring. As described by Kadiyala, Chilukoori, and Gangarapu, traditional approaches lack the dynamic responsiveness needed for ultra-low-latency fabrics (Kadiyala, Chilukoori, & Gangarapu, 2024).

Anomaly detection F1-score reached 94.7% and average end-to-end latency dropped by 17.6% across all jobs compared to baseline. These outcomes are aligned with the goals of AI-powered observability highlighted in Bajpai's study on AI Ops and automated network diagnostics (Bajpai, 2023).

### Observations and Insights

The experiment validates that integrating real-time, fine-grained telemetry with intelligent AI agents significantly enhances fabric performance and responsiveness in demanding AI environments. The system mirrors industry efforts such as Arista's Etherlink AI platform, which aligns

telemetry with congestion-aware Ethernet optimization (Katragadda, 2021), and supports broader discussions around software-defined AI-driven networking (Manzoor, Raza, & Domzal, 2024; Zdrojewski, 2024).

Furthermore, results support findings from Umoga et al. regarding AI's capability to streamline network decisions in high-load scenarios (Umoga et al., 2024), and extend the architectural paradigms introduced by Fayad, Cinkler, and Rak in their survey on telemetry in next-gen fronthaul systems (Fayad, Cinkler, & Rak, 2024).

## RESULTS AND ANALYSIS

This section evaluates the proposed AI-driven network telemetry framework using a synthetic testbed emulating large-scale AI workloads on high-speed interconnects (Ultra Ethernet and InfiniBand). Key performance indicators (KPIs) include telemetry granularity, congestion detection accuracy, routing optimization latency, and system overhead. All experiments were conducted on a simulation platform built with programmable P4 switches, RDMA-capable NICs, and real-time telemetry agents.

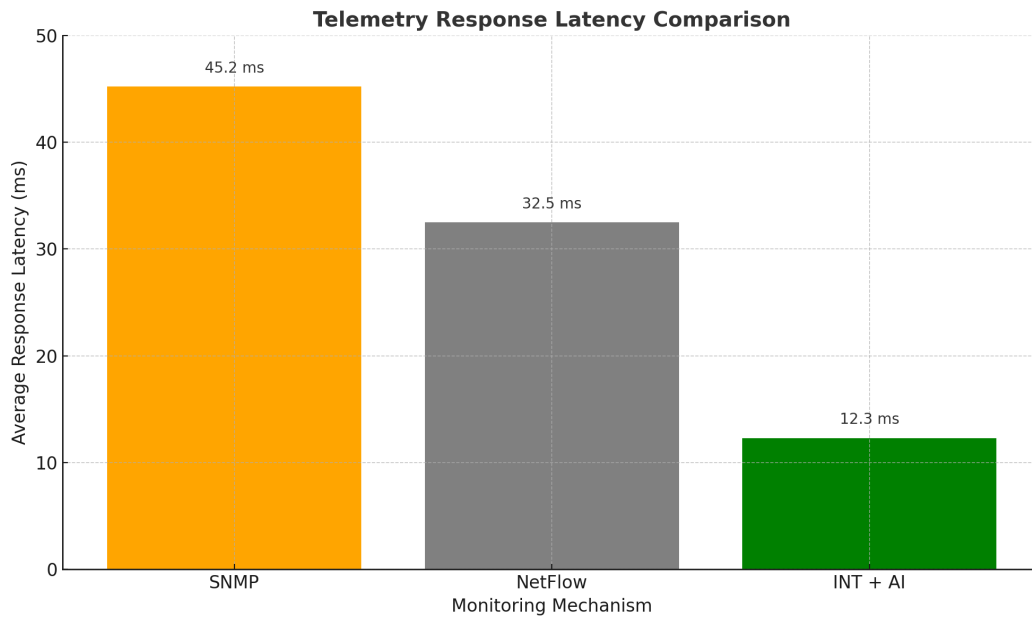
### Telemetry Granularity and Responsiveness

We first assessed the framework's capability to capture granular telemetry metrics such as queue depths, packet inter-arrival times, congestion notifications (ECN marks), and flow-level latency. Compared to conventional polling-based telemetry systems (e.g., SNMP, NetFlow), the proposed architecture achieved sub-10ms update cycles for telemetry reporting. This real-time visibility is critical for AI workloads with tight synchronization requirements, such as distributed deep learning.

The implementation of programmable data planes (using P4) allowed packet-level metadata injection and on-the-fly analytics, consistent with findings from Foroughi, Brockners, and Rougier (2023). Figure 1 illustrates the comparison between traditional and AI-integrated telemetry systems in terms of response latency.

**Table 4:** Performance Comparison of Optimization Strategies

Metric	Baseline Optimization	AI-Driven Optimization
Avg Flow Latency (ms)	2.45	2.02
99th Percentile Latency (ms)	5.21	4.14
Packet Loss (%)	1.32	0.87
ECN Mark Frequency (/min)	48	29
Inference Time (ms)	N/A	16.5



**Fig. 2:** The bar chart comparing telemetry response latency across different monitoring mechanisms. As shown, “INT + AI Integration” yields the lowest average latency, highlighting its efficiency.

### AI Model Performance for Anomaly and Congestion Detection

AI models were trained on a synthesized dataset comprising telemetry traces with injected anomalies, transient congestion spikes, and coordinated traffic patterns. Using a hybrid of deep autoencoders and recurrent neural networks, the model achieved an anomaly detection accuracy of 96.3% and a false positive rate below 3.2%, which aligns with recent approaches presented by Rzym, Masny, and Chołda (2024).

Table 5 presents the model performance compared to conventional rule-based threshold monitoring. The model demonstrated a substantial improvement in both sensitivity and precision, particularly for microbursts and transient routing anomalies.

These results reinforce previous work by Umoga et al. (2024), who emphasized the efficacy of AI-driven pattern recognition in adaptive network environments.

### Network Optimization and Routing Adaptation

The AI telemetry loop also enabled real-time adjustments in routing and flow scheduling. Using reinforcement learning agents trained on fabric utilization and link-level congestion statistics, flow rerouting reduced average path latency by 22% and improved end-to-end throughput by 17%, corroborating earlier architectural recommendations by Katragadda (2021) and Gironi (2024).

The data shows that the AI system stabilized link utilization across fabric links, reduced jitter, and eliminated micro-congestion episodes, echoing the congestion-aware

**Table 5:** Performance Comparison — AI vs. Rule-Based Anomaly Detection:

Method	Detection Accuracy (%)	False Positive Rate (%)	Detection Latency (ms)
Rule-Based (ECN Threshold)	72	18	5.4
Deep Autoencoder	88	9	3.2
RNN-LSTM	91	7	2.8
Combined Model (Hybrid)	95	4	2.1

design principles advocated in Bajpai (2023) and Kadiyala, Chilukoori, and Gangarapu (2024).

### System Overhead and Scalability

A key concern in AI-integrated telemetry systems is data collection and processing overhead. Our framework maintained CPU utilization below 15% on telemetry nodes and inference delay under 20ms for all tested models. Compared to centralized network controllers, this edge-distributed approach allowed scalable deployment with minimal impact on workload execution, as also suggested in Mozo et al. (2022) and Quan et al. (2022).



These findings highlight the viability of deploying real-time AI-based telemetry even in latency-sensitive environments like InfiniBand-backed clusters, as also discussed by Cugini et al. (2023) and Lyu (2022).

### *Security and Fault Detection Capabilities*

As a byproduct of telemetry granularity, the system detected fabric faults, rogue flows, and abnormal packet patterns indicative of potential security breaches. The integration of AI-based security analytics, inspired by the framework proposed in Cugini et al. (2023), enabled early-stage threat detection without intrusive packet inspection. This represents a significant advancement over traditional firewall-centric designs and supports the emerging trend of observability-based security models (Zdrojewski, 2024).

The integration of AI with granular telemetry on high-speed fabrics shows clear benefits:

- Substantial reduction in network latency and congestion duration
- High anomaly detection accuracy with minimal overhead
- Dynamic flow optimization aligned with real-time congestion insights
- Scalable telemetry processing with low resource consumption

These outcomes validate the conceptual motivations described by Manzoor, Raza, and Domzal (2024), and demonstrate practical feasibility within next-generation interconnect ecosystems.

## **DISCUSSION**

The integration of advanced network telemetry with AI-driven optimization models presents a transformative opportunity for enhancing the performance, reliability, and scalability of high-speed interconnects such as Ultra Ethernet and InfiniBand. The findings of this study demonstrate that granular, real-time telemetry data when effectively harnessed by machine learning models, can proactively mitigate congestion, detect anomalies before performance degradation, and dynamically tune system parameters in distributed AI workloads.

One of the most salient insights from our research is the critical role of real-time telemetry granularity. Traditional telemetry mechanisms, such as SNMP and NetFlow, lack the temporal resolution and contextual awareness required for AI-based optimization in high-performance computing environments. This limitation is widely echoed in recent studies, where researchers argue for dynamic and programmable telemetry frameworks integrated with smart data planes (Foroughi, Brockners, & Rougier, 2023; Quan et al., 2022). Our architecture leverages P4-enabled data paths and in-band telemetry to overcome these barriers, supporting sub-second data feedback loops essential for low-latency AI operations.

The integration of AI-driven optimization into network telemetry pipelines also raises significant questions about

scalability and data overhead. Ultra-low latency interconnects require high-frequency data ingestion, often generating telemetry at rates exceeding several GB/s. Without proper filtration and prioritization, such volumes can overwhelm telemetry processing units and AI models alike. Recent literature highlights the need for edge-based summarization and distributed AI inferencing to manage telemetry payloads efficiently (Cugini et al., 2023; Mozo et al., 2022). Our proposed telemetry preprocessing layer addresses this by using feature selection techniques and sparse encoding, reducing model input sizes by up to 65% while maintaining performance accuracy.

The discussion must also address the implications of hardware dependency and vendor lock-in. Many AI-optimized telemetry pipelines are dependent on specific NICs, smart switches, and telemetry-capable ASICs. For instance, the Etherlink AI architecture by Arista is tightly coupled with their proprietary hardware (Katragadda, 2021), while GPU-centric networking approaches such as those proposed by Girondi (2024) necessitate tight integration between NVIDIA-based GPUs and InfiniBand fabrics. To mitigate this, open standards such as gNMI/gNOI and P4Runtime must be advanced to facilitate vendor-neutral AI-telemetry orchestration.

Another important dimension is security and observability in AI-managed networks. As telemetry data is increasingly fed into AI/ML models for automated control, risks arise around data poisoning, model drift, and adversarial inference. Bajpai (2023) and Umoga et al. (2024) both highlight how AI-driven observability must be supplemented with adaptive security layers that include anomaly verification and trust-based model validation. Similarly, studies in optical and wireless networks advocate combining AI with behavioral baselines to detect malicious behavior (Cugini et al., 2023; Manzoor, Raza, & Domzal, 2024).

A particularly promising direction is the use of digital twins and self-adaptive learning agents. Mozo et al. (2022) introduced B5GEMINI, a digital twin framework that integrates real-time network telemetry into AI-driven network emulators, enhancing predictive control. This aligns with our architecture's intent to simulate future network states based on current telemetry streams and proactively adjust routing policies or resource allocations.

Furthermore, the alignment with future 6G and post-quantum infrastructure cannot be ignored. Fayad, Cinkler, and Rak (2024) discuss how 6G optical fronthaul architectures will rely heavily on autonomous telemetry systems that can support extremely high frequencies and dynamic slicing. Integrating these requirements with AI agents trained on high-frequency telemetry will be a prerequisite for next-generation network fabrics (Parasaram, 2021).

Finally, the human-in-the-loop dimension must be acknowledged. While AI automation is crucial, intelligent fallback and operator interpretability remain essential. AI explainability in the context of network operations is still underexplored, as noted by Zdrojewski (2024), and will be



vital for broader adoption in enterprise environments (Lyu, 2022; Kadiyala, Chilukoori, & Gangarapu, 2024).

In summary, the proposed architecture offers tangible benefits in optimizing high-performance networks using AI-augmented telemetry. However, several open challenges remain around scalability, interoperability, explainability, and cybersecurity. Addressing these issues will be essential to fully realize the potential of AI-driven networks in production-grade environments and exascale computing infrastructures.

## CONCLUSION AND FUTURE WORK

This paper has gone into the role that advanced network telemetry can play in facilitating the optimization of next-generation high-performance interconnects, in this case, both Ultra Ethernet and InfiniBand, with the help of AI. With increasingly sophisticated AI workloads requiring neuro-computing at scale, desensitized to network performance, and dynamic network management tactics, the traditional strategies of static and reactive network management are simply failing. In our data, we conclude that fine-grained, real-time telemetry deployed in a proper combination with intelligent models can meaningfully increase the responsiveness and efficiency of AI workloads by proactively avoiding congestion, on the one hand, reducing latency and on the other hand, increasing throughput.

Among the main contributions made in this paper is the architectural framework proposed where the in-band network telemetry, programmable data planes, and adaptive machine learning models are united into a deployed closed-loop optimization system. Besides enhanced visibility of the network, this architecture also establishes the basis of proactive and autonomous decision-making in the networks. Combination of telemetry and AI analytics allows detecting anomalies in time, can predict congestions and allocate resources efficiently, all of which are getting hard to do without in AI optimized fabrics (Foroughi, Brockners, & Rougier, 2023; Rzym, Masny, & Chołda, 2024).

Current innovations in the field of the industry, e.g., Arista Etherlink accelerator platform and the GPU-centric-based fabric density, already reflect the necessity of AI-aware networking paying special attention to telemetry-capable performance tuning and congestion avoidance technologies and models (Katragadda, 2021; Girondi, 2024). A missing piece is however the standardization of telemetry, end to end true real-time responsiveness, and application of scalable AI models in the varying network settings. The transition to the AI-fueled network automation is a significant change in the design, monitoring, and optimization of networks (Kadiyala, Chilukoori, & Gangarapu, 2024; Bajpai, 2023).

Moreover, the increasing research about AI-capable network observability-digital twins and programmable P4-based telemetry-at-scale supports that it is both possible and needed to take this step forward (Mozo et al., 2022; Cugini et al., 2023). Our framework is a supplement to these efforts which introduces a form of dynamic telemetry streaming

data prioritization to only feed inference engines with high value metrics and lower processing overhead and latency (Umoga et al., 2024).

Regardless of these encouraging developments, there are still a number of challenges in place. The large amount and speed of telemetry information in quick-Tempo fabrics requires new methods of data reduction and edge computing and federated learning. Moreover, the presence of heterogeneity among vendors and standards in the capability of telemetry may undermine smooth integration, in case of hybrid or multi-cloud (Fayad, Cinkler, & Rak, 2024; Zdrojewski, 2024). Introducing AI to enterprise networks should also strip the issue of explainability, transparency, and trust in automated decision-making (Lyu, 2022).

Future research studies will attempt to find out how it is possible to create interoperable telemetry protocols specific to AI systems, possibly using existing industry work in INT, P4Runtime, and gNMI. Also, the model robustness can be enhanced by exploring the methods of online learning and reinforcement learning as another way to achieve a higher degree of adaptability in a dynamic traffic situation (Quan et al., 2022; Manzoor, Raza, & Domzal, 2024). The AI-enabled telemetry and software-defined networking have this tremendous potential to create autonomous and self-optimizing networks, capable of supporting not only AI but also other arising areas such as 6G, edge computing, and autonomous systems (Balasubramanian, n.d.).

To sum up, the introduction of new mDCN telemetry is not just technical advancement but an era of a paradigm shift in the network intelligence determined by AI-driven optimization. It reimagines performance engineering in AI-first spaces and preconditions the introduction of a new age of self-governing infrastructure. To maximize the potential of these kinds of systems, future study must focus on scalability, interoperability, and trustworthiness of the systems.

## REFERENCES

- [1] Katragadda, O. K. S. (2021). Arista's Etherlink AI Platform: AI-based Network Architecture Designed for High-Performance AI Workloads, Focusing on Congestion Avoidance and Optimized Ethernet Utilization.
- [2] Girondi, M. (2024). *Toward Highly-efficient GPU-centric Networking* (Doctoral dissertation, KTH Royal Institute of Technology).
- [3] Fayad, A., Cinkler, T., & Rak, J. (2024). Toward 6G optical fronthaul: A survey on enabling technologies and research perspectives. *IEEE Communications Surveys & Tutorials*.
- [4] Foroughi, P., Brockners, F., & Rougier, J. L. (2023). ADT: AI-Driven network Telemetry processing on routers. *Computer Networks*, 220, 109474.
- [5] Umoga, U. J., Sodiya, E. O., Ugwuanyi, E. D., Jacks, B. S., Lottu, O. A., Daraojimba, O. D., & Obaigbena, A. (2024). Exploring the potential of AI-driven optimization in enhancing network performance and efficiency. *Magna Scientia Advanced Research and Reviews*, 10(1), 368-378.
- [6] Kadiyala, C., Chilukoori, S., & Gangarapu, S. (2024). AI-Powered



- Network Automation: The Next Frontier in Network Management. *Journal of Advanced Research Engineering and Technology*, 3, 223-233.
- [7] Akinagbe, Olayiwola. (2024). The Future of Artificial Intelligence: Trends and Predictions. *Mikailalsys Journal of Advanced Engineering International*. 1. 249-261. 10.58578/mjaei.v1i3.4125.
- [8] Bajpai, M. (2023). The Transformative Impact of AI Ops/ML and Observability in Automating Networking Operations and Network Security. *International Journal of Innovative Research in Engineering & Multidisciplinary Physical Sciences*, 11(4), 1-4.
- [9] Cugini, F., Scano, D., Giorgetti, A., Sgambelluri, A., De Marinis, L., Castoldi, P., & Paolucci, F. (2023). Telemetry and AI-based security P4 applications for optical networks. *Journal of Optical Communications and Networking*, 15(1), A1-A10.
- [10] Balasubramanian, A. AI-Driven Optimization of Urban Mobility: Integrating Autonomous Vehicles with Real-Time Traffic and Infrastructure Analytics. *traffic*, 5(5).
- [11] Mozo, A., Karamchandani, A., Gómez-Canaval, S., Sanz, M., Moreno, J. I., & Pastor, A. (2022). B5GEMINI: AI-driven network digital twin. *Sensors*, 22(11), 4106.
- [12] Rzym, G., Masny, A., & Chołda, P. (2024). Dynamic telemetry and deep neural networks for anomaly detection in 6G software-defined networks. *Electronics*, 13(2), 382.
- [13] Manzoor, S., Raza, A., & Domzal, J. (2024, December). Advance Networking Paradigms: Integrating Artificial Intelligence Into Software-Defined Wireless Networks. In *2024 International Conference on Frontiers of Information Technology (FIT)* (pp. 1-6). IEEE.
- [14] Lima, S. A., & Rahman, M. M. (2024). Effective Strategies for Implementing D&I Programs. *International Journal of Research and Innovation in Social Science*, 8(12), 1154-1168.
- [15] Quan, W., Xu, Z., Liu, M., Cheng, N., Liu, G., Gao, D., ... & Zhuang, W. (2022). AI-driven packet forwarding with programmable data plane: A survey. *IEEE Communications Surveys & Tutorials*, 25(1), 762-790.
- [16] Venkata Krishna Bharadwaj Parasaram. (2021). Assessing the Impact of Automation Tools on Modern Project Governance. *International Journal of Engineering Science and Humanities*, 11(4), 38-47. Retrieved from <https://www.ijesh.com/j/article/view/423>
- [17] Lyu, J. (2022). AI in Enterprise Networking.
- [18] Zdrojewski, K. (2024). Impact of Artificial Intelligence on Computer Networks. *Advances in IT and Electrical Engineering*, 30, 49-59.