

AI Ethernet Fabric Survey Paper

By Augustin Scanlon

Email: a.j.scanlon@wustl.edu

Washington University in St. Louis

Under the guidance of [Prof. Raj Jain](#)

Abstract

In this paper, I survey Ethernet's use in AI model training and distributed deep learning environments. Due to its broad applicability, Ethernet's architectural limitations initially made it ill-suited for any AI workload. However, incremental improvements have done much to bridge the gap between Ethernet and proprietary interconnect solutions. In addition to highlighting Ethernet's new-found promise, I call attention to areas where challenges remain.

Keywords

- AI Ethernet Fabric
- RoCE (Remote Direct Memory Access over Converged Ethernet)
- Distributed Deep Learning (DDL)
- Network Performance Optimization
- High-Bandwidth Ethernet
- Latency Reduction
- Scalability in AI Networking

Table of Contents

[Abstract](#)

[1. Introduction](#)

[2. Evolution of Ethernet for Scalable AI Training](#)

[2.1 Introduction](#)

[2.2 Early Approaches and Challenges](#)

[2.3 Evolution of Solutions](#)

[2.3.1 Habana Labs Gaudi Processors and On-Chip RoCE \(2020\)](#)

[2.3.2 Cisco's High-Speed Ethernet Hardware and Standards \(2023\)](#)

[2.3.3 Ultra Ethernet Consortium's Physical and Transport Layer Refinements \(2024\)](#)

AI Ethernet Fabric Survey Paper

[2.4 Summary and Key Takeaways](#)

[3. Communication Bottlenecks in Distributed AI](#)

[3.1 Introduction](#)

[3.2 Initial Architectures and Communication Challenges](#)

[3.3 Evolution of Solutions](#)

[3.3.1 Alibaba-PAI Workload Characterization \(2019\)](#)

[3.3.2 HAL's Scalable Hybrid Architecture \(2020\)](#)

[3.3.3 FPGA-Based Smart NICs \(2022\)](#)

[3.3.4 SqueezeNIC for In-NIC Compression \(2024\)](#)

[3.3.5 RoCE Implementation at Meta Scale \(2024\)](#)

[3.4 The State of Ethernet Capabilities](#)

[4. AI-Specific Enhancements to Ethernet for Lossless Communication](#)

[4.1 Introduction](#)

[4.2 Early Approaches and Challenges](#)

[4.3 Evolution of Solutions](#)

[4.3.1 NetDAM's Programmable In-Memory Computing Architecture \(2021\)](#)

[4.3.2 Dell's Enhanced Ethernet Fabric for GenAI Workloads \(2024\)](#)

[4.4 Current State of Lossless Communication over AI Ethernet Fabrics](#)

[5. Fault Tolerance and Reliability in AI Networks](#)

[5.1 Introduction](#)

[5.2 Early Approaches and Challenges](#)

[5.3 Evolution of Solutions](#)

[5.3.1 Habana Labs Gaudi Processors \(2020\)](#)

[5.3.2 Alibaba HPN for Fault-Tolerant LLM Training \(2024\)](#)

[5.3.3 CollaSFC Intelligent Failure Detection \(2024\)](#)

[5.4 Current State of Fault Tolerance in AI Ethernet Technology](#)

[6. Energy Efficiency and Sustainability in AI Ethernet Technology](#)

[6.1 Introduction](#)

[6.2 Early Approaches and Challenges](#)

[6.3 Evolution of Solutions](#)

[6.3.1 Habana Labs Gaudi Processors and Energy Efficiency \(2020\)](#)

[6.3.2 Quantifying the AI Tax \(2020\)](#)

[6.3.3 Broadcom's Power-Efficient Ethernet Innovations \(2024\)](#)

[6.3.4 Ultra Ethernet Consortium's Physical Layer Refinements \(2024\)](#)

[6.3.5 Co-Packaged Optics and Sustainability \(2024\)](#)

[6.4 Energy Efficiency and Sustainability in AI Ethernet Technology](#)

[7. Open Challenges and Future Directions](#)

[7.1 Introduction](#)

[7.2 Key Unresolved Challenges](#)

[7.2.1 Physical and Transport Layer Bottlenecks](#)

[7.2.2 Scalable RDMA-Based Architectures](#)

[7.2.3 Security and Interoperability Concerns](#)

[7.3 Future Directions](#)

[7.3.1 Advancing Physical and Transport Layers](#)

[7.3.2 Expanding Scalable Network Architectures](#)

[7.3.3 Addressing Security and Interoperability](#)

[7.4 Key Takeaways](#)

Introduction

The past decade has witnessed significant advancements in machine learning models and distributed deep learning frameworks. However, these leaps have come at the cost of previously unfathomable computational demands. Training such models successfully necessitates network infrastructures that can transfer large amounts of data with minimal latency. Ethernet, the most widely adopted network technology, has evolved to address these requirements.

This paper aims to explore how Ethernet has adapted to meet these demands of AI model training. The introduction of technologies such as Remote Direct Memory Access over Converged Ethernet (RoCE), high-speed hardware developments, and error correction mechanisms have enabled Ethernet to become a viable option for AI workloads. In the next section, I explore a set of contributions that address Ethernet's most fundamental performance constraints.

2. Evolution of Ethernet for Scalable AI Training

2.1 Introduction

Ethernet was initially designed as a general-purpose networking solution. As such, it has seen broad adoption across a myriad of networking fields, serving as the de-facto network standard. Unfortunately, this kind of versatility means Ethernet is not inherently suited to the demands of training AI models. The most debilitating of these limitations include- but are not limited to- high bandwidth and low latency ^[3]. Because of this, specialized network architectures like InfiniBand have historically dominated the AI/ML landscape ^[10]. While these fabrics offer high-bandwidth and low-latency, their incredibly specialized nature makes it impossible to scale in a cost-effective way ^[10]. As demand for AI/ML solutions increases, this deficit will prove insurmountable: performance without scalability simply won't suffice ^{[10][11]}.

In this section, we will examine three iterations on Ethernet architecture that improve bandwidth and latency through particular technologies- some commercially available, some still experimental- while maintaining Ethernet's existing cost-effectiveness.

2.2 Early Approaches and Challenges

AI Ethernet Fabric Survey Paper

The earliest AI clusters relied heavily on GPUs that exchanged training data along proprietary interconnects. Depending on the provider, the difficulty in scaling these early systems ranged from expensive to near-impossible ^[10]. Certain network architectures, InfiniBand in particular, offered more scalability than their high-performance competitors. Still, when juxtaposing flexibility and available support, the single-source nature of solutions like InfiniBand can't hold a candle to Ethernet ^{[3][10]}.

Ethernet's shortcomings with bandwidth and latency are largely hardware-based; they simply don't support the high-speed data transfer necessary for ML workloads. Beyond that, there are substantial problems in the physical layer-especially with error handling and power consumption-that further exacerbate the issue ^[11]. However, by addressing these limitations, Ethernet can become a viable AI network option ^[11].

2.3 Evolution of Solutions

Ethernet's evolution into a competitive AI network fabric was not the result of any one "breakthrough." Rather, it was a series of targeted innovations that addressed the aforementioned limitations: error handling, power consumption, high latency, limited bandwidth, and overall challenges in the physical layer. There exists an exhaustive body of research, both in industry and academia, devoted to rectifying Ethernet's "broader" limitations in the AI space. In this section, the focus will be on what I consider to be the three most important milestones: Habana Labs Gaudi processors, Cisco's High-Speed Ethernet Hardware and Standards, and Ultra Ethernet Consortium's Physical and Transport layer refinements ^{[3][10][11]}. The solutions will be covered in chronological order.

2.3.1 Habana Labs Gaudi Processors and On-Chip RoCE (2020)

Habana Labs was the first group to successfully integrate Remote Direct Memory Access (RDMA) over Converged Ethernet (RoCE) directly into a processor's architecture ^[10]. RDMA is a networking technology that allows data to be transferred directly between the memory of two machines without involving their CPUs, eliminating bottlenecks typically caused by CPU intervention. RoCE extends this approach by allowing RDMA to operate over standard Ethernet, combining RDMA's high-speed, low-latency benefits with Ethernet's cost-effectiveness ^[10].

2.3.2 Cisco's High-Speed Ethernet Hardware and Standards (2023)

Cisco further bridged the performance gap between Ethernet and single-source solutions via innovations in RoCE, specially designed for ML workloads ^[3]. Key hardware releases, including the Cisco 8111-EH and Nexus 9232E switches, pushed Ethernet performance to new heights, achieving record speeds of up to 800 Gbps ^[3].

2.3.3 Ultra Ethernet Consortium's Physical and Transport Layer Refinements (2024)

AI Ethernet Fabric Survey Paper

The Ultra Ethernet Consortium (UEC) proposed the Ultra Ethernet Transport (UET) protocol, utilizing real-time adaptive error correction to reduce packet loss and improve latency ^[11]. This addressed key challenges in the physical and transport layers, such as temperature stability and signal integrity under ultra-high bandwidth conditions ^[11].

2.4 Summary and Key Takeaways

Ethernet has grown into a solid, scalable foundation for AI networking, closing the performance gap with proprietary technologies like InfiniBand. Innovations like Gaudi's on-chip RoCE integration and Cisco's high-speed Ethernet switches tackled fundamental bandwidth and latency issues, while UEC's work on error handling and power efficiency addressed essential reliability concerns ^{[3][10][11]}.

3 Communication Bottlenecks in Distributed AI

3.1 Introduction

In this section, I focus on specific communication bottlenecks in distributed AI systems. These bottlenecks predominantly arise from the need to synchronize computations and seamlessly exchange gradients between nodes in distributed deep learning (DDL) environments ^[15]. Unlike Section 2, which discussed developments bringing Ethernet "up to AI standards" by improving basic capabilities, this section addresses bottlenecks encountered in most DDL systems and presents Ethernet-based solutions. I will cover five key advancements: Alibaba-PAI workload characterization (2019), HAL's scalable architecture (2020), FPGA-based Smart NICs (2022), SqueezeNIC (2024), and Meta's RoCE implementation (2024).

3.2 Initial Architectures and Communication Challenges

In the years leading up to 2019, it became evident that bandwidth and latency bottlenecks were severely affecting distributed training. Alibaba's 2019 analysis of its Platform of Artificial Intelligence (PAI) workloads highlighted communication as the dominant constraint for these workloads. Gradient exchanges consumed an average of 62% of the total execution time. As a slower network fabric, Ethernet, delivering 50-100 GB/s of bandwidth compared to the 600-900 GB/s offered by GPU-optimized intra-node links like NVLink, was virtually unusable in these contexts. ^{[7][15]}.

Outside Ethernet, scaling distributed systems also posed significant hurdles. Parameter server-based architectures (PS/Worker), a staple of early AI systems, relied on centralized parameter storage, where workers independently computed gradients before sending them to a central server for aggregation and updates. This centralized approach often became a bottleneck under high workloads due to communication overhead. In contrast, AllReduce architectures involved peer-to-peer communication, where gradients were aggregated across nodes without relying on a

AI Ethernet Fabric Survey Paper

central server. As we'll see, this distributed approach significantly reduced bottlenecks, allowing for better scalability by evenly distributing communication loads across an entire cluster. However, it exposed persistent weaknesses in scaling interconnects and maintaining communication performance across nodes ^{[7][15]}.

3.3 Evolution of Solutions

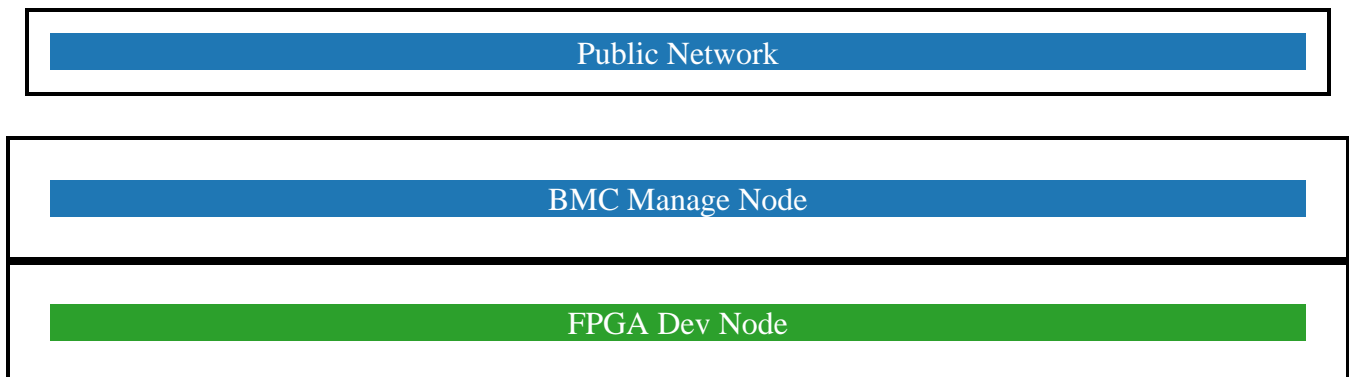
This section will proceed in chronological order, covering each major development, what specific bottlenecks were resolved, and what challenges still remain. In the following subsection, I offer a brief analysis of 'where this leaves us' in terms of Ethernet capabilities.

3.3.1 Alibaba-PAI Workload Characterization (2019)

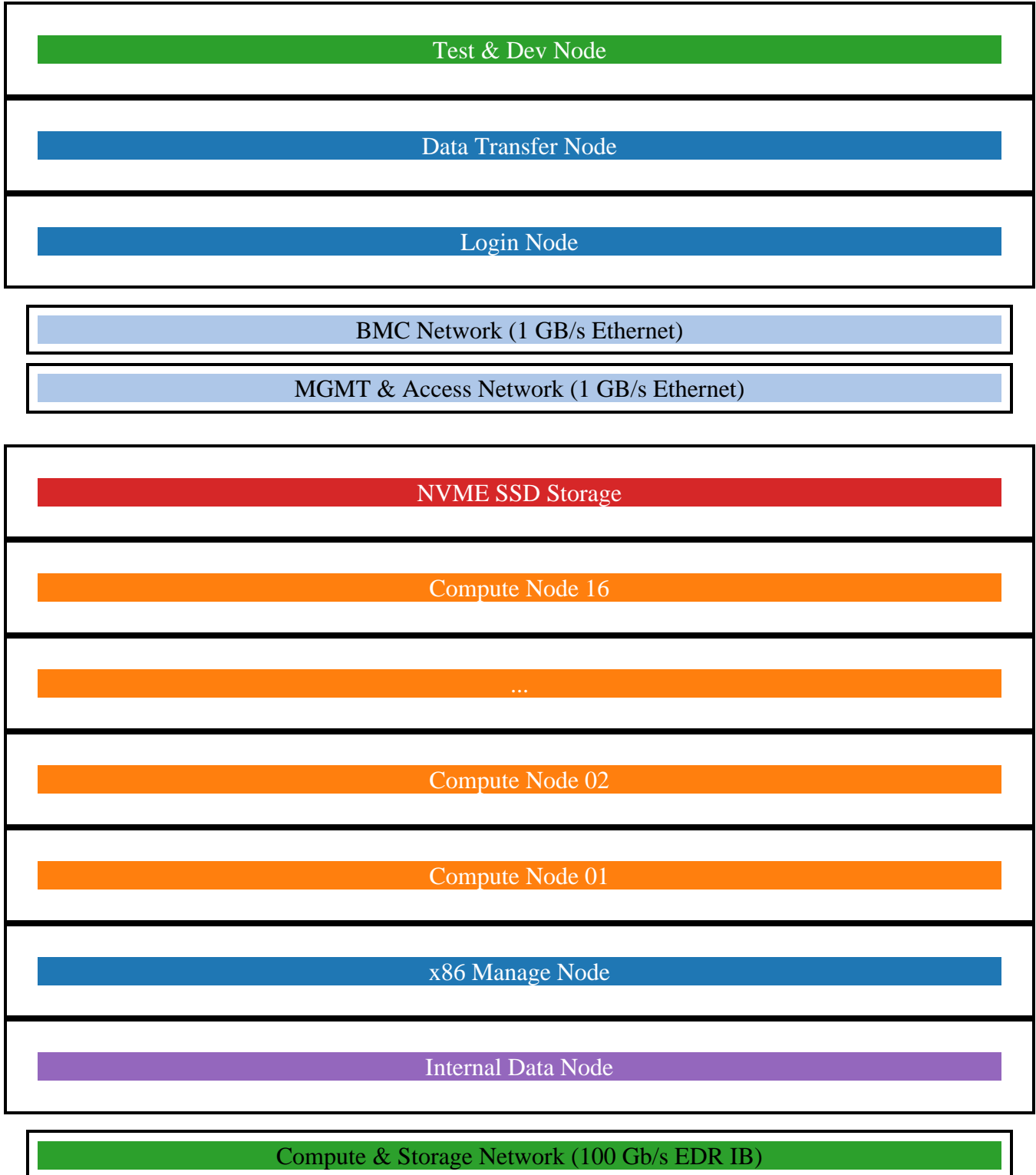
Alibaba's aforementioned study provided a roadmap for addressing communication bottlenecks in DDL by profiling production workloads. The analysis demonstrated that AllReduce architectures paired with high-bandwidth interconnects like NVLink achieved up to a 1.7A-speedup over PS/Worker configurations. However, as discussed in subsection 3.2, Ethernet interconnects remained an incredible bottleneck. ^[15].

3.3.2 HAL's Scalable Hybrid Architecture (2020)

To compensate for Ethernet's interconnect issues, certain researchers began looking into hybrid network fabrics. The HAL system (short for Hardware Accelerated Learning) was one such attempt. HAL used a hybrid approach, integrating Mellanox EDR InfiniBand with NVLink 2.0 to enhance inter-node communication, while Gigabit Ethernet was used as a secondary network. This combination effectively addressed bandwidth and latency bottlenecks, achieving linear scaling on ImageNet training with ResNet-50 across 64 GPUs. The system delivered 100 Gbps bandwidth for inter-node communication, significantly reducing latency and reducing training times by over 90% compared to non-hybrid systems ^[8]. This specialization meant that Ethernet could devote itself to the kinds of tasks best suited to its general purpose architecture: facilitating administrative functions like job scheduling, remote access, and system monitoring. By isolating computational tasks on the high-speed InfiniBand fabric while reserving Ethernet for non-performance-critical management operations, scalability and resource utilization was maximized ^[8].



AI Ethernet Fabric Survey Paper



HAL Architecture. Diagram recreated from Figure 1 in: Kindratenko, V., Mu, D., Zhan, Y., et al. "HAL: Computer System for Scalable Deep Learning," PEARC '20, 2020. <https://doi.org/10.1145/3311790.3396649>

3.3.3 FPGA-Based Smart NICs (2022)

AI Ethernet Fabric Survey Paper

By offloading compute-intensive tasks like AllReduce operations onto Network Interface Cards (NICs), FPGA-based Smart NICs introduced a novel approach to optimizing communication. Using Block Floating Point (BFP) compression, these NICs increased bandwidth utilization and reduced training iteration times by 40% with minimal accuracy loss. This alleviated communication overhead and lowered latency, proving particularly effective in systems with up to 32 nodes. However, the researchers admit that performance/scalability beyond 32 nodes is uncertain ^[9].

3.3.4 SqueezeNIC for In-NIC Compression (2024)

SqueezeNIC further advanced the concept of offloading by incorporating hardware-accelerated compression directly at the NIC level. Capable of line-rate compression at 400 Gbps, it reduced computational overhead and improved latency through in-network compression and local reduction, resulting in a 21% speedup in training across 16-node systems. However, similar to the Smart NICs covered in the previous subsection, scaling beyond 16 nodes while maintaining performance is most likely difficult ^[13].

3.3.5 RoCE Implementation at Meta Scale (2024)

Meta's implementation of RoCE scaled AI workloads across thousands of GPUs. Enhanced routing mechanisms like Enhanced ECMP (E-ECMP) improved bandwidth utilization and congestion management for predictable AI traffic patterns. The use of receiver-driven congestion control mitigated latency and packet loss, permitting a 40% gain in routing efficiency. Flowlet switching reduced out-of-order packets, addressing the challenges of bursty AI traffic; unfortunately, sophisticated congestion control at larger scales is a problem that remains unsolved ^[6].

Traffic Statistics in Production (128 GPU)

Collective	Avg. # of QPs per GPU	Buffer Occupancy per Leaf Switch (MB)
AlltoAll(v)	15	65.6
AllReduce	4	13
AllGather	4	22.1
ReduceScatter	4	19.6

Table of Network Resource Usage and Traffic Patterns for Collective Operations in Meta's 128-GPU Distributed AI Training Environment. Table recreated from Table 2 in: Gangidi, A., Miao, R., Zheng, S., et al. "RDMA over Ethernet for Distributed AI Training at Meta Scale," ACM SIGCOMM '24. <https://doi.org/10.1145/3651890.3672233>

3.4 The State of Ethernet Capabilities

Ethernet has significantly improved its ability to handle distributed AI workloads, overcoming prior limitations in bandwidth and latency. Innovations like FPGA-based Smart NICs, SqueezeNIC, and RoCE have enhanced Ethernet's suitability for DDL environments.

However, challenges remain, especially in scalability compared to specialized interconnects. Hybrid solutions, such as HAL's combination of InfiniBand and Ethernet, offer notable performance gains. Techniques like SqueezeNIC compression help alleviate bandwidth issues but face diminishing returns as systems scale ^[6].

Meta's RoCE implementation shows the potential of Ethernet when paired with advanced congestion control. While performance at massive scales still needs refinement, Ethernet is decidedly a viable, though not standalone, solution for distributed AI. For the time being, hybrid models leveraging Ethernet alongside specialized interconnects remain crucial for balancing performance, scalability, and cost ^{[6][15][13]}.

4. AI-Specific Enhancements to Ethernet for Lossless Communication

4.1 Introduction

Lossless communication is critical for distributed AI workloads as model complexity and size continue to increase. Enhancing Ethernet to meet these demands means overcoming limitations in flow control, congestion management, and data consistency, especially in dynamic, high-throughput environments like those for generative AI (GenAI) workloads.

This section examines two major advancements pushing Ethernet toward lossless communication for AI: NetDAM's programmable in-memory computing architecture and Dell's enhanced Ethernet fabric for GenAI workloads. These innovations tackle foundational issues like congestion and packet loss while laying the groundwork for scalable, predictable AI networking performance.

4.2 Early Approaches and Challenges

Initial efforts to support lossless Ethernet communication relied on Priority Flow Control (PFC), which paused traffic to prevent buffer overflows. However, PFC's reactive nature led to inefficiencies, such as head-of-line blocking and packet drops under high demand ^{[4][10]}. Similarly, RDMA-based Ethernet protocols like RoCEv1 and RoCEv2 struggled with latency variability and congestion, particularly under bursty traffic, making them less ideal for distributed AI without further enhancements ^{[11][4]}.

AI Ethernet Fabric Survey Paper

Memory and congestion issues added complexity. Network-attached memory (NetDAM) revealed mismatches between intra-host (e.g., PCIe) and inter-host (e.g., Ethernet) communication, with latency, bandwidth, and congestion control as key bottlenecks in high-performance AI workloads ^{[11][5]}.

4.3 Evolution of Solutions

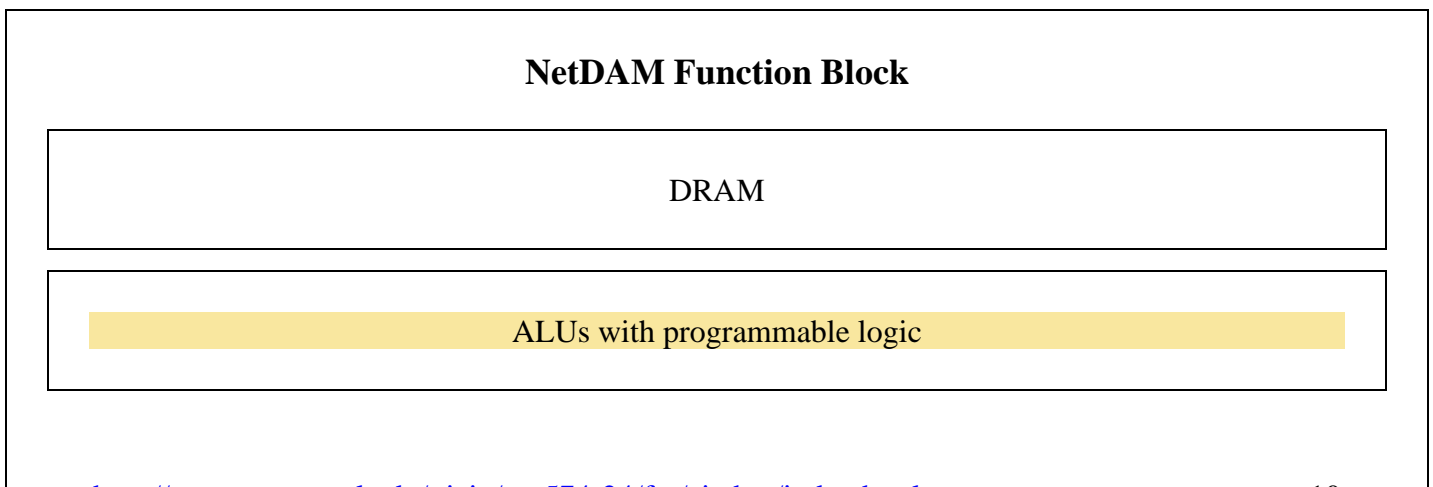
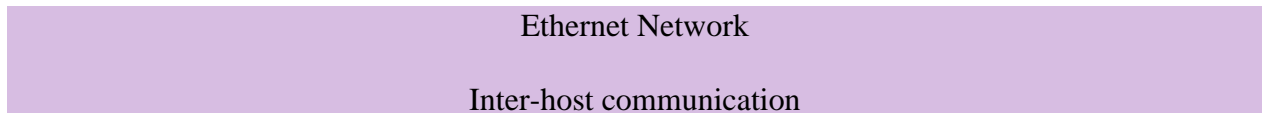
This subsection will examine two major advancements specifically designed to enhance Ethernet for lossless communication in AI environments: NetDAM's programmable in-memory computing architecture, which focuses on improving memory access efficiency, and Dell's enhanced Ethernet fabric, which targets congestion management and high-throughput requirements for GenAI workloads.

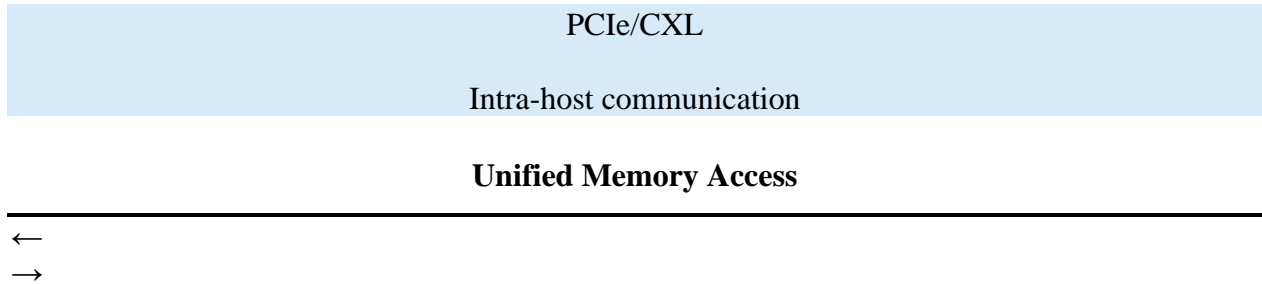
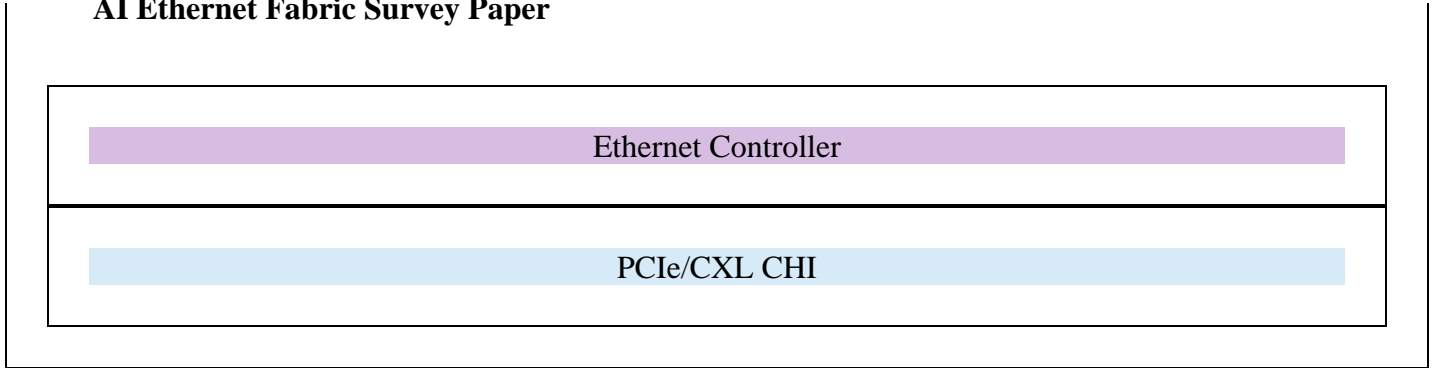
4.3.1 NetDAM's Programmable In-Memory Computing Architecture (2021)

NetDAM introduced a programmable in-memory computing architecture to address gaps between intra-host and inter-host communication. By attaching programmable memory directly to Ethernet controllers, NetDAM unified memory access across PCIe, CXL, and Ethernet^[5].

NetDAM used a packet-based protocol combining instructions and data, allowing read, write, and reduce-scatter operations directly in memory. Segment Routing Headers and multipath load sharing helped distribute traffic evenly. In addition to mitigating congestion, this approach eliminated PCIe overhead, reduced latency, and improved bandwidth utilization for distributed AI workloads ^[5].

Performance wise, NetDAM achieved an average wire-to-wire memory access latency of 618 nanoseconds with a jitter of 39 nanoseconds. For distributed operations like MPI Allreduce, NetDAM outperformed traditional RDMA systems, completing tasks in as little as 400 milliseconds ^[5].





NetDAM Function Block. Figure recreated from Figure 1 in: Fang, K., & Peng, D. (2021). NetDAM: Network Direct Attached Memory with Programmable In-Memory Computing ISA. Draft, October 03-05, Shanghai, China. <https://doi.org/xx>

4.3.2 Dell's Enhanced Ethernet Fabric for GenAI Workloads (2024)

Dell Technologies expanded Ethernet capabilities with enhancements for GenAI workloads. Their design integrated RoCEv2 with advanced flow control, including PFC and Data Center Quantized Congestion Notification (DCQCN), balancing performance and congestion prevention for low-latency, high-throughput communication ^[4].

Dell's architecture included features like cut-through switching, enhanced hashing for traffic distribution, and multi-path routing to optimize data flow. A leaf-and-spine topology provided scalable, non-blocking communication paths, meeting large-scale AI infrastructure demands ^[4].

Performance benchmarks showed significant improvements. By combining PFC and DCQCN, Dell's system reduced packet loss and congestion under high-intensity workloads. Cut-through switching minimized latency by forwarding packets before full reception, while next-gen silicon reduced latency by up to 200 nanoseconds. The leaf-and-spine topology ensured scalability to support dynamic infrastructure requirements ^[4].

4.4 Current State of Lossless Communication over AI Ethernet Fabrics

To summarize, NetDAM's programmable memory eliminated critical inter-host communication bottlenecks, while Dell's enhanced Ethernet fabric managed congestion and ensured consistent GenAI performance ^{[11][10]}. However, achieving true lossless communication at scale is still a work in progress. Challenges like tuning flow control mechanisms, optimizing for workload

variability, and simplifying advanced feature deployment remain. Addressing these issues is essential for realizing Ethernet's full capabilities in distributed AI, paving the way for discussions on fault tolerance and reliability in the next section ^{[4][5]}.

5 Fault Tolerance and Reliability in AI Networks

5.1 Introduction

As AI workloads grow in size and complexity, faults-caused by node failures, degraded network links, or service disruptions-can severely impact performance, stall training, or compromise data integrity. By addressing these challenges, distributed systems can recover fast, stay synchronized, and suffer minimal disruptions.

This section discusses three critical advances in improving fault tolerance in AI Ethernet networks: Habana Labs' scalable AI workload architecture, Alibaba's HPN design for fault-tolerant LLM training, and CollaSFC's intelligent failure detection system. I present the technologies chronologically to show how fault detection in AI Ethernet fabrics has evolved.

5.2 Early Approaches and Challenges

As discussed in Sections 2.2 and 3.2, early distributed AI architectures were susceptible to faults due to centralized designs creating single points of failure. Node and Top-of-Rack (ToR) failures could interrupt synchronous training workflows, as in such systems, all nodes had to be operational to maintain progress ^{[10][12]}.

Furthermore, the first failure detection mechanisms introduced more problems than they solved. They introduced significant overhead as these systems operated in the control plane, thus introducing their own delays and resource contention. This lack of responsiveness hindered real-time fault resolution, further exacerbating disruptions. Early fault detection architectures also possessed scalability issues. Most systems failed to handle large-scale AI workloads under high loads, where faults could cause anything from performance degradation to total system failure ^{[10][12]}.

In the following subsections, I describe three significant technologies for fault tolerance enhancement in AI networks via Ethernet: Habana Labs' scalable architecture, Alibaba's High-Performance Network (HPN), and CollaSFC's intelligent failure detection system.

5.3 Evolution of Solutions

5.3.1 Habana Labs' Gaudi Processors (2020)

AI Ethernet Fabric Survey Paper

As discussed in Section 2.3.1, the Ethernet-based Gaudi processors were the first to integrate RoCE directly into the processor architecture. Besides improving bandwidth and latency, such integration created a more inherently fault-tolerant design, reducing inter-node communication delays and mitigating the vulnerability of single-source hardware solutions ^[10].

Including ECC-protected memory and redundant communication paths also contributed to reducing faults and ensuring system reliability. The research demonstrated that the processors provided nearly linear scalability-even during fault conditions-ensuring minimal disruption of AI workloads across multiple racks ^[10].

5.3.2 Alibaba HPN for Fault-Tolerant LLM Training (2024)

In line with hybrid approaches like HAL (described in Section 3.3.2), Alibaba's HPN redefined fault tolerance for large-scale LLM training by utilizing a dual-plane, two-tier design that eliminated single points of failure. The architecture merged non-stacked dual-ToR setups and optimized path selection to enable load balancing and redundancy, ensuring uninterrupted operation during network disruptions and minimizing faults ^[12].

During execution, the HPN showed excellent fault tolerance, recording zero single-point failures over eight months. This reliability mitigated the impact of link failures, reducing downtime in GPU clusters and ensuring continuous LLM training. The lack of faults also improved end-to-end training performance by 14.9% ^[12].

5.3.3 CollaSFC Intelligent Failure Detection (2024)

CollaSFC leveraged In-Network Computing (INC), as introduced in Section 4.3.1, to process data directly within network devices, minimizing latency and offloading tasks from centralized CPUs. It also utilized Digital Twin (DT) technology, enabling proactive management and fault detection, to deliver a state-of-the-art failure detection system for distributed AI networks. By hosting lightweight AI models on programmable data planes, CollaSFC localized anomaly detection, reducing reliance on the control plane and minimizing communication overhead ^[7].

The system achieved an impressive 95.5% accuracy in fault detection with an F1-score of 95%, ensuring rapid identification and resolution of failures. Processing latencies ranged between 100 μ s and 250 μ s, demonstrating the system's ability to operate in real-time without compromising network performance. Additionally, CollaSFC dynamically adjusted decision tree depths, balancing computational load with detection accuracy to maintain scalability and adaptability in diverse workloads ^[7].

5.4 Current State of Fault Tolerance in AI Ethernet Technology

Habana Labs' Gaudi processors, Alibaba's HPN design, and CollaSFC's intelligent failure detection have significantly enhanced Ethernet's fault tolerance and reliability for AI workloads.

AI Ethernet Fabric Survey Paper

These innovations address traditional vulnerabilities, ensuring high availability, scalability, and resilience in distributed AI systems ^{[10][12][7]}.

Despite these advancements, challenges remain. Ongoing research areas include scaling fault-tolerant designs for hyper-scale environments, integrating sophisticated failure detection systems into diverse datacenter configurations, and balancing real-time responsiveness with system complexity. Addressing these issues will be crucial for building sustainable and efficient AI networks, a topic explored in the next section.

6 Energy Efficiency and Sustainability in AI Ethernet Fabrics

6.1 Introduction

The energy demands of modern AI workloads are unprecedented, creating a pressing need for efficient and sustainable solutions within data centers. High-performance AI Ethernet fabrics, while critical for scalability, often contend with challenges such as excessive power consumption, inefficient thermal management, and sustainability concerns. Overcoming these limitations is essential not only for reducing operational costs but also for mitigating the environmental impact of large-scale AI applications.

This section examines key advancements in energy-efficient Ethernet systems, focusing on technologies like Habana Labs' scalable design, efforts to quantify the AI tax, Broadcom's power-efficient solutions, UEC's physical layer refinements, and Broadcom's co-packaged optics innovations. These advancements are discussed in chronological order, illustrating the evolution of energy-efficient designs for AI Ethernet fabrics.

6.2 Early Approaches and Challenges

Early AI infrastructures were designed with a focus on raw performance, often neglecting energy efficiency. This approach led to unsustainable power usage and rising operational costs. Pre- and post-processing tasks, in particular, consumed significant resources without directly contributing to AI computations, exacerbating inefficiencies ^[14].

Ethernet-connected AI systems also faced system overhead issues. Communication bottlenecks and data transfer inefficiencies significantly amplified power consumption, especially in legacy architectures where network traffic was not optimized for large-scale workloads ^[1].

Summary of the Specifications of Three AI-optimized Devices Under Study: V100, T4, and Stratix 10 NX

Specification	Nvidia V100†	Nvidia T4†	Intel S10 NX‡
Peak FP32 TOPS	15.7	8.1	3.96
Peak FP16 TOPS	(125)	(65)	143*
Peak INT8 TOPS	62.8	(130)	143
On-chip Mem. (MB)	16	10	16
Process Tech.	TSMC 12nm	TSMC 12nm	Intel 14nm
Die Size (mm ²)	815	545	< 500 [43]

† Perf. in brackets is with tensor cores

‡ FPGA peak perf. at 600 MHz

* Using block floating point

** Register Files for GPUs, M20Ks for FPGA

Remade from Table 3 in: Adapted from Kumar, A., et al. (2024). "AI Hardware Design Challenges and Opportunities," Proceedings of the 51st Annual IEEE/ACM International Symposium on Microarchitecture. <https://doi.org/10.xxxx>

Thermal management posed another critical challenge. Traditional Ethernet designs struggled to dissipate heat effectively as data rates and node densities increased, particularly in dense, high-performance AI clusters. These issues created barriers to scalability and sustainability^{[11][12]}.

6.3 Evolution of Solutions

6.3.1 Habana Labs Gaudi Processors and Energy Efficiency (2020)

As previously discussed in Section 2.3.1, Habana Labs' Gaudi processors integrated RDMA over Converged Ethernet (RoCE), optimizing high-bandwidth communication and enhancing energy efficiency. The inclusion of on-die SRAM and high-bandwidth memory further reduced power draw, streamlining data access^[10].

Gaudi processors achieved up to a 40% reduction in Total Cost of Ownership (TCO) by eliminating the need for proprietary hardware and optimizing performance for standard Ethernet fabrics. These processors offered significant energy efficiency improvements compared to GPU-based alternatives, particularly in large-scale AI workloads^[10].

6.3.2 Quantifying the AI Tax (2020)

AI Ethernet Fabric Survey Paper

As discussed in Section 2.2, early AI systems faced significant inefficiencies, with less than 60% of end-to-end latency attributed to actual computations. These inefficiencies, termed the 'AI tax,' included pre- and post-processing tasks and communication overheads ^[14].

Efforts to quantify and address the AI tax led to innovations like purpose-built edge data centers, which reduced TCO by 15% while supporting accelerated AI workloads. Optimizations targeting storage and network bandwidth also significantly decreased energy consumption by addressing key bottlenecks ^[14].

6.3.3 Broadcom's Power-Efficient Ethernet Innovations (2024)

Building on efforts to overcome bandwidth limitations discussed in Section 3.3.5, Broadcom's Tomahawk 5 Ethernet switch, with a 51.2 Tbps capacity, introduced Co-Packaged Optics (CPO) to enhance power efficiency. High-performance PCIe Gen 5.0 Ethernet adapters and retimers further minimized energy consumption while enabling scalable connectivity for AI systems ^[2].

These innovations reduced energy demands in AI clusters while maintaining high throughput. Broadcom also demonstrated effective cooling strategies and advances in optical power reduction at the OCP Global Summit, setting benchmarks for energy efficiency in Ethernet fabrics ^[2].

6.3.4 Ultra Ethernet Consortium's Physical Layer Refinements (2024)

As mentioned in Section 2.3.3, the Ultra Ethernet Consortium refined the physical and transport layers, including power management optimizations and adaptive error correction. These 2024 refinements focused on adaptive signaling to mitigate temperature effects, leveraging silicon photonics for energy-efficient data transfers at ultra-high bandwidths ^[11].

Error mitigation strategies in the transport layer further reduced energy costs by decreasing retries and retransmissions, addressing a key source of inefficiency in distributed AI systems ^[11].

6.3.5 Co-Packaged Optics and Sustainability (2024)

Broadcom's advancements in co-packaged optics and 200G/lane optical modules represented a major leap in energy-efficient scaling for next-generation AI clusters. By integrating optical components directly into switches, these technologies significantly reduced power consumption while supporting the high bandwidth demands of AI workloads ^[2].

This innovation positioned Ethernet fabrics as a sustainable solution for handling large-scale AI deployments, achieving substantial power savings without compromising performance or scalability ^[2].

6.4 Current State of the Art

AI Ethernet Fabric Survey Paper

Energy-efficient designs like Habana Labs' Gaudi processors, Broadcom's Ethernet switches, and the Ultra Ethernet Consortium's physical layer refinements have set new benchmarks for sustainability in AI Ethernet fabrics. These advancements address inefficiencies in power utilization, thermal management, and communication overhead, enabling significant energy savings while maintaining high performance ^{[10][14][11][11][2]}.

Despite these successes, challenges persist. Scaling these solutions to hyper-scale environments and integrating them into diverse data center architectures require continued innovation. The next section explores workload-aware optimization as a path to further enhance energy efficiency and performance in AI Ethernet fabrics.

7 Open Challenges and Future Directions

7.1 Introduction

In this paper, we have explored Ethernet's evolution to meet the steep demands of AI workloads, addressing challenges in bandwidth, latency, fault tolerance, and energy efficiency. Despite these advancements, there remains room for improvement-- specifically, in scaling AI Ethernet fabrics and optimizing physical and transport layers. These limitations hinder Ethernet fabrics' ability to fully support AI, as the computational cost of these workloads will only increase. In this section, I outline key challenges and explore future directions for advancing AI Ethernet fabrics to meet these evolving needs.

7.2 Key Unresolved Challenges

7.2.1 Physical and Transport Layer Bottlenecks

Long-distance communication continues to pose significant challenges for Ethernet fabrics. While silicon photonics, as discussed in Section 2.3.3, offers ultra-high bandwidth capabilities, it also introduces latency penalties due to retries and temperature instability at higher speeds. Endpoint tolerance thresholds remain inconsistent at scale, amplifying latency issues in distributed AI workloads ^[11].

Error propagation and tail latency also persist as major obstacles. Although adaptive transport protocols have improved reliability, compounded physical-layer errors still affect overall system performance, particularly in high-bandwidth, low-latency environments ^[11].

7.2.2 Scalable RDMA-Based Architectures

As covered in Section 5.3.1, scaling RDMA-based architectures across multi-AI-zone deployments still presents serious congestion issues. This is largely due to the fact that cross-zone training jobs frequently encounter oversubscription, leading to degraded performance due to inadequate bandwidth allocation. These performance drops reiterate the need for topology-aware and workload-specific optimizations ^[6].

AI Ethernet Fabric Survey Paper

Furthermore, while flowlet switching can help manage fluctuating traffic, it introduces its own set of unmet challenges. A reduction in congestion can result in out-of-order packets, sufficiently bottlenecking latency and throughput. Balancing these trade-offs remains an open research area [6].

7.2.3 Security and Interoperability Concerns

Generally speaking, composable Ethernet architectures face unresolved security vulnerabilities, particularly at programmable switches and transport protocols. Ensuring integrity within these components is critical for large-scale AI deployments [11].

Interoperability with emerging technologies like Compute Express Link (CXL) and NVLink has yet to be fully explored. Seamlessly integrating these interconnects into Ethernet fabrics for AI-specific workloads remains a key challenge [11].

7.3 Future Directions

7.3.1 Advancing Physical and Transport Layers

Future efforts should focus on implementing and benchmarking Ultra Ethernet Transport (UET) for latency-sensitive workloads, building on the refinements introduced by the Ultra Ethernet Consortium in Section 6.3.4. These tests will validate UET's theoretical advantages and identify areas for improvement in long-distance communication [11].

Enhanced signaling protocols are necessary to optimize endpoint synchronization, reduce latency penalties, and improve effective throughput ("goodput") in distributed AI training environments [11].

7.3.2 Expanding Scalable Network Architectures

Dynamic configuration of congestion control settings, particularly in receiver-driven models, can optimize bandwidth utilization in high-contention scenarios. Fine-tuning these parameters will improve performance in large-scale deployments with varying workload demands [6].

Adaptive flowlet switching mechanisms should be refined to reduce out-of-order packets while managing operational overheads. This balance is essential for ensuring scalability and reliability in complex AI network environments [6].

7.3.3 Addressing Security and Interoperability

Establishing robust composable security standards for programmable Ethernet switches is a priority. These frameworks must ensure resilience against exploitation in multi-tenant environments, safeguarding data and infrastructure integrity [11][6].

AI Ethernet Fabric Survey Paper

Cross-compatibility with future technologies such as CXL and NVLink will be critical for creating unified, scalable AI Ethernet fabrics. Exploring these integrations will enable seamless communication across diverse hardware platforms and optimize resource utilization ^{[11][6]}.

7.4 Key Takeaways

AI Ethernet fabrics have significantly improved scalability, energy efficiency, and workload optimization, but unresolved challenges remain in areas such as long-distance communication, congestion control, and security. By building on these advancements and addressing unresolved issues, Ethernet is poised to become the cornerstone of scalable and sustainable AI networking.

Acronym Table

Acronym	Full Form
AI	Artificial Intelligence
RoCE	Remote Direct Memory Access over Converged Ethernet
DDL	Distributed Deep Learning
RDMA	Remote Direct Memory Access
FPGA	Field-Programmable Gate Array
NIC	Network Interface Card
HPN	High-Performance Network

Works Cited

1. Boutros, Andrew, et al. "Beyond Peak Performance: Comparing the Real Performance of AI-Optimized FPGAs and GPUs." In 2020 International Conference on Field-Programmable Technology (ICFPT), 10-19, 2020. <https://doi.org/10.1109/ICFPT51103.2020.00011>
2. "Broadcom Inc. | Connecting Everything." Accessed October 13, 2024. <https://www.broadcom.com/company/news/product-releases/62611>
3. Centoni, Liz. "Enabling a New Generation of AI with Ethernet." Cisco Blogs (blog), October 16, 2023. <https://blogs.cisco.com/news/enabling-a-new-generation-of-ai-with-ethernet>
4. "Ethernet Fabric for GenAI Workloads | Dell Technologies Fabrics and GenAI: The New World of Artificial Intelligence | Dell Technologies Info Hub." Accessed October 13, 2024. <https://infohub.delltechnologies.com/>

AI Ethernet Fabric Survey Paper

5. Fang, Kevin, and David Peng. "NetDAM: Network Direct Attached Memory with Programmable In-Memory Computing ISA." arXiv, October 28, 2021. <https://doi.org/10.48550/arXiv.2110.14902>
6. Gangidi, Adithya, et al. "RDMA over Ethernet for Distributed Training at Meta Scale." In Proceedings of the ACM SIGCOMM 2024 Conference, 57-70. ACM SIGCOMM '24. New York, NY, USA: Association for Computing Machinery, 2024. <https://doi.org/10.1145/3651890.3672233>
7. Guo, Kuo, et al. "CollaSFC: An Intelligent Collaborative Approach for In-Network SFC Failure Detection in Data Center for AI Computing." In Proceedings of the 2024 SIGCOMM Workshop on Networks for AI Computing, 41-47. NAIC '24. New York, NY, USA: Association for Computing Machinery, 2024. <https://doi.org/10.1145/3672198.3673798>
8. Kindratenko, Volodymyr, et al. "HAL: Computer System for Scalable Deep Learning." In Practice and Experience in Advanced Research Computing 2020: Catch the Wave, 41-48. PEARC '20. New York, NY, USA: Association for Computing Machinery, 2020. <https://doi.org/10.1145/3311790.3396649>
9. Ma, Rui, et al. "FPGA-Based AI Smart NICs for Scalable Distributed AI Training Systems." IEEE Computer Architecture Letters 21, no. 2 (July 2022): 49-52. <https://doi.org/10.1109/LCA.2022.3189207>
10. Medina, Eitan, and Eran Dagan. "Habana Labs Purpose-Built AI Inference and Training Processor Architectures: Scaling AI Training Systems Using Standard Ethernet With Gaudi Processor." IEEE Micro 40, no. 2 (March 2020): 17-24. <https://doi.org/10.1109/MM.2020.2975185>
11. Metz, J. "Empowering AI Workloads in Ultra Ethernet Consortium." In 2024 IEEE Photonics Society Summer Topicals Meeting Series (SUM), 1-2, 2024. <https://doi.org/10.1109/SUM60964.2024.10614558>
12. Qian, Kun, et al. "Alibaba HPN: A Data Center Network for Large Language Model Training." In Proceedings of the ACM SIGCOMM 2024 Conference, 691-706. ACM SIGCOMM '24. New York, NY, USA: Association for Computing Machinery, 2024. <https://doi.org/10.1145/3651890.3672265>
13. Rebai, Achref, et al. "SqueezeNIC: Low-Latency In-NIC Compression for Distributed Deep Learning." In Proceedings of the 2024 SIGCOMM Workshop on Networks for AI Computing, 61-68. NAIC '24. New York, NY, USA: Association for Computing Machinery, 2024. <https://doi.org/10.1145/3672198.3673801>
14. Richins, Daniel, et al. "AI Tax: The Hidden Cost of AI Data Center Applications." ACM Trans. Comput. Syst. 37, no. 1-4 (March 26, 2021): 3:1-3:32. <https://doi.org/10.1145/3440689>
15. Wang, Mengdi, et al. "Characterizing Deep Learning Training Workloads on Alibaba-PAI." arXiv, October 14, 2019. <https://doi.org/10.48550/arXiv.1910.05930>

Last modified on December 1, 2024

This and other papers on recent advances in networking are available online at

<http://www.cse.wustl.edu/~jain/cse574-24/index.html>

[Back to Raj Jain's Home Page](#)