



AUTHORS

Ron Westfall

Research Director | The Futurum Group

David Nicholson

Chief Technology Advisor | The Futurum Group

Daniel Newman

Chief Analyst | The Futurum Group

AUGUST 2024

IN PARTNERSHIP WITH





Executive Summary

Ethernet outperforms all other networking technologies in terms of efficiency, as measured by the lowest job completion time (JCT) for AI training and inference. It has become the go-to choice for expansive AI clusters, offering the unique capability to link over a million AI accelerators essential for scaling AI networks. The creation of AI systems is an intricate task that requires interoperability, standardization, and technology prowess among various suppliers, spotlighting the necessity of a shared, open platform to unite these components. Ethernet uniquely delivers networking solutions that ensure peak performance and scalability with cost-effectiveness and agility, due primarily to its foundation on open standards and its extensive ecosystem support.

In contrast, InfiniBand technology is commonly deployed in high-performance computing (HPC) scenarios, where attaining low unloaded latency is the primary focus. However, from our view, when it comes to AI training and inference clusters, Ethernet outperforms InfiniBand. Specifically, Ethernet can offer several advantages, including greater bandwidth, higher radix, multivendor standard solutions, lower cost, power-efficient economics, improved throughput, and overall better performance. As such, for large-scale AI workloads that demand high bandwidth, reliable transport, and predictable tail latency, Ethernet is the preferred choice over InfiniBand.

We explore why the Broadcom portfolio provides Ethernet solutions and capabilities that assure lower costs and address network congestion challenges by optimizing network efficiency, especially for scaling and optimizing fast expanding AI workloads. Broadcom solutions, such as the Switch Scheduled Fabric based on Jericho 3-AI switch, Endpoint Scheduled Fabric based on the Tomahawk switch family, and Broadcom's 400G Ethernet NICs, enable the AI infrastructure and fabric vital to making the network the compute platform essential to fulfilling the unique distributed computing demands of AI.





The Evolving Ethernet Landscape

Al is playing an integral role in driving ecosystem-wide digital transformation across virtually every industry vertical. The demand for compute-intensive Al applications, such as generative Al, is driving demand for Al clustering with immediacy. With demands for generative Al growing swiftly, handling massive data volumes for training large language models (LLMs) requires even more robust server clusters.

These demands require ultra-high radix networking devices that enable the AI fabric to support larger clusters, scaling to tens or hundreds of thousands of GPUs. In 2024, AI clusters have already grown from using 4,000 GPUs to swiftly scaling up to 32,000 GPUs, with the expectation of scaling up to 1 million plus GPUs on the horizon.

Enabling AI infrastructure is fundamentally a distributed computing challenge since no one GPU is big enough to run a massive job. Over two decades ago, Google built one of the planet's largest distributed computing systems using Ethernet to interconnect many hundreds of thousands of CPUs. Essentially, in today's computing landscape, the network is the computer. The network is the only way to execute the massively distributed computing required for today's workloads.

However, attaining optimal network efficiency across AI workload environments is challenging. <u>Up to 50% of the time required to train AI models</u> (according to West Gate Networks) is spent transferring data between GPUs. This results in making GPU resources decidedly more costly since they can be idle up to 50% of the time.

As such, attaining optimal network efficiency is paramount to meet the unparalleled challenge of scaling and managing AI workloads. The process of handling AI workloads involves a continuous loop of compute, synchronize, and communication operations between the various compute nodes. As the processing continues, there is a corresponding exponential increase in the data exchanged between nodes, putting a premium on using networking technology that can deliver improvements in overall network efficiency. This puts a premium on using capabilities that enable efficient bandwidth optimization for AI workloads, especially involving natural language processing that rely on moving large datasets swiftly. This also involves delivering latency reduction techniques that meet the unique requirements of AI workloads as AI back-end networks need



a high-speed, scalable, low-latency fabric. This escalates the demand for implementing solutions that can help substantially improve network efficiency across AI workload fabrics:

- Accelerated Servers: Al servers frequently use GPUs as they excel at parallel processing. GPUs can support
 and scale thousands of tasks simultaneously, making them best suited for training complex Al models. Through
 accelerated servers, users can apply and administer massive datasets that iterate their designs with more
 immediacy.
- Networking for AI Capabilities: Network efficiency gains are augmented by factors such as AI-focused adaptive
 routing, hardware-assisted failure recovery, and end-to-end (E2E) congestion control that are integral to supporting
 AI clusters and AI workload optimization.
- CPO: Co-packaged Optics (CPO) is a heterogeneous integration of optics and silicon on a single packaged substrate targeted at addressing fast evolving bandwidth and power challenges. CPO brings together a vast array of expertise in switch application-specific integrated circuits (ASICs), digital signal processors (DSPs), optical transceivers, as well as packaging and testing to make data center and cloud infrastructure ready for AI workload challenges.
- LPO: Linear Pluggable Optics (LPO) has garnered market interest as a step in the path toward CPO. LPO modules
 remove the DSP from pluggable optics, providing a lower cost and power solution since the DSP can account for
 around half the power consumption and a significant portion of pluggable optics costs. The trade-off is that LPO
 modules require a high-performance switch SerDes and careful PCB signal integrity design.

Why Ethernet Meets the AI Scaling Challenge

Ethernet is essential to fulfilling the unique demands of AI networking including the scaling of AI workloads. Already Ethernet has demonstrated its ability to address the most demanding scaling environments across networks, as evidenced by seven of the top eight hyperscalers using Ethernet for scaling their AI workloads. Ethernet has become the top choice for AI clusters due to built-in benefits such as E2E congestion management, standards-backing, cost-effectiveness, and fabric management:

- E2E Congestion Management: Ethernet fabrics are designed with congestion management capabilities to intelligently control latency and loss, which plays an integral role in minimizing JCT across AI training clusters. As such, large-scale Ethernet-enabled training clusters are ready to support 32,000 GPUs, each with 800Gbps, in a single cluster, and on the horizon scale up to around 1.2 million GPUs.
- Robust Capabilities and Reliability: Broadcom's Ethernet switches provide dynamic load balancing, key to
 optimizing network traffic across multiple paths or links; multi-tenancy that assures different tenants (customers)
 share the same infrastructure securely by keeping their data and resources isolated; and failure recovery advances,
 such as zero impact failover (ZIF), to ensure minimal disruption to job completion.
- Compatibility and Versatility: Ethernet's proven compatibility and versatility with a vast array of devices makes it the topmost selection for AI networking for providing reliable data transmission, traffic management, and error detection throughout the most demanding AI training and inferencing scenarios. Ethernet technology is widely used and mass-produced, yielding lower deployment and maintenance expenses, while InfiniBand requires specialized hardware and software that result in higher costs.



Fabric Management: Ethernet eases the creation and management of fabrics using a network controller that
configures devices with policies based on packet header tags, catalyzing the fabric management process and
ensuring overall network efficiency.

We find that Ethernet especially delivers enduring competitive advantages over alternative proprietary technology such as NVIDIA's InfiniBand. Key considerations encompass:

- Lower Costs: Ethernet fabrics with RoCEv2 support are delivering proven scalability at lower cost-per-bit compared to InfiniBand across swiftly evolving AI cluster environments. Ethernet fabrics with RDMA over Converged Ethernet (ROCE) protocol support are optimized for AI clusters, especially RoCEv2, by enabling direct memory access between systems, minimizing XPU involvement and reducing latency. Emerging custom XPU technology, including GPUs, CPUs, and TPUs, all use Ethernet to meet AI workload requirements. Through advances in Ethernet speeds, such as 400 Gbps and 800 Gbps, the cost of networking infrastructure has decreased.
- Predictable Latency: Across environments that require predictable latency and reliable transport, Ethernet has
 demonstrated that it can <u>outperform</u> InfiniBand when configured with similar SerDes speed and switch bandwidth.
 This includes providing more consistent and predictable performance during critical tail latency scenarios whereby
 the latency experienced by a small fraction of requests can take significantly longer than the average.
- Avoid Proprietary Traps: Ethernet's well-established standards backing and imprimatur assure customers, including
 cloud service providers (i.e., hyperscalers), operators, and enterprises, that they can have choice and flexibility in
 scaling their AI/ML workloads. InfiniBand's proprietary nature offers no such confidence.





Why Broadcom Ethernet Solutions Are Essential to Al Workload Optimization

To meet the unique demands of AI workload scaling, we review why Broadcom E2E networking solutions are integral to ensuring the successful implementation of the largest Ethernet installations worldwide. Underpinning its E2E networking solutions, Broadcom offers two scheduled fabric solutions for AI networks. The scheduling of traffic can originate either at the switch or at the endpoint. For switch-scheduled fabrics, the intelligence for the scheduled fabrics resides within the switches. For endpoint-scheduled fabrics, the scheduled fabric is established at the endpoints, such as NICs or GPUs. Either solution is implemented according to customer choice.

Key to the Broadcom E2E portfolio, consisting of both the switch-scheduled fabric and endpoint-scheduled fabric solutions, are the following product offerings:

- Switch-Scheduled Fabric based on Jericho3-AI: The Broadcom switched-scheduled fabric solution uses the Jericho3-AI family, optimized for AI clusters. The Jericho3-AI switch is purpose-built for meeting the specific demands of AI workloads. AI workloads have unique attributes such as a low number of massive, long-lived flows, all beginning concurrently upon completion of an AI computation cycle. To assure top performance for the specific requirements of AI workloads, Jericho3-AI delivers:
 - Ultra-high radix: Enables Jericho3-Al fabric to scale connectivity to 32,000 GPUs, each with 800 Gbps, in a single cluster.
 - Zero-Impact failover (ZIF): Provides sub-10ns automatic path convergence, assuring no impact to JCTs.
 - Perfect load balancing: Distributes traffic over all the links of the fabric, ensuring maximum network utilization under the most intense network loads.
 - Minimal congestion: Uninterrupted operation through E2E traffic scheduling, which eliminates flow collisions and jitters.
- Endpoint-Scheduled Fabric based on Tomahawk 5: The Broadcom endpoint-scheduled fabric solution uses the Tomahawk 5 family, providing 51.2 Tbps of switching capacity on a single chip, which we identify as a top runner among merchant switch chips available today. Additional key benefits include:
 - High radix: Enables single-hop connectivity between, for example, 256 high-performance AI/ML accelerators, each with 200 Gbps of network bandwidth.
 - AI/ML workload virtualization: Delivers efficient use of shared infrastructure in large data centers by offering capabilities such as single-pass VxLAN (Virtual Extensible LAN) routing and bridging for AI/ML workloads.
 - Configuration flexibility: Supports multitude of link configurations, including up to 64x 800GbE, 128x 400GbE, or 256x 200GbE links.



• 400G RoCE/RDMA Ethernet NICs: Broadcom's recently introduced 400G PCle Gen 5.0 Ethernet adapters offer breakthrough RoCEv2 performance, assisted by smart congestion control engine and peer-direct capabilities for the largest AI clusters. These new adapters can enable the most power and thermally efficient design in the market. Combined with the device's ability to drive passive copper cables up to five meters or ultra-low power linear pluggable optics transceivers, the adapter delivers higher rack density using mainstream air-cooled GPU data center technology.

Broadcom E2E open congestion control algorithm streamlines RoCEv2 implementations in data centers to help achieve the low latencies critical to scaling AI workloads.

For both the switch-scheduled fabric solution and the endpoint-scheduled fabric solutions, the endpoint can be one of the following four options: Broadcom NIC, customer NIC, merchant NIC, or native Ethernet interface from the GPU/accelerator.

The Broadcom Tomahawk family improves AI performance through endpoint-scheduled fabric solutions. Through this approach, the NIC manages traffic scheduling while Tomahawk switches efficiently forward traffic. As noted, the available endpoint options include Broadcom NICs alongside third-party NICs as well as GPU/accelerator-based native Ethernet interfaces. Broadcom's Tomahawk 5 switch series, which can provide accelerated job completion and greater bandwidth, plays an integral role at optimizing AI cluster performance.

The Tomahawk 5 family endpoint-scheduled fabric solution delivers optimized AI load balancing features such as Cognitive Routing to enable efficient handling of large, low-entropy flows typical of East–West AI/ML workloads. Broadcom Ethernet reliability features deliver time-to-failure reactions that the company reports as providing a 30X improvement over comparable InfiniBand implementations.

Multi-tenancy plays a vital role in the implementation of AI training networks. This includes Broadcom's support for VxLAN as an extension to Layer 2 VLANs, designed to enhance VLAN functionality with broader flexibility and greater scalability. Benefits consist of using VxLAN to extend Layer 2 segments across the data, enabling tenant workloads to span physical pods, key to multitenant environments.

VxLAN capabilities support key use cases such as multi-tenancy in cloud environments, enabling cloud providers to support multiple tenants with isolated virtual networks that scale beyond VLANs and efficiently manage tenant-specific L2 domains as well as assuring software-defined datacenters which enable the creation of virtual networks for many VMs while preserving the datacenter network architecture.

With reliability and consistency, Broadcom has supported multi-tenancy with VxLAN and IPinIP tunnels, enabling multiple customers to co-exist on a single network with exacting isolation. Built-in ZIF ensures sub-10ns automatic path convergence, avoiding disruption to JCTs. Plus, a deep set of summarized telemetry metrics enhance endpoint congestion management protocols. Underpinning these key Broadcom portfolio capabilities is the company's proven track record in working with the extensive ecosystem of hardware, software, and orchestration vendors in providing Ethernet-based solutions that fulfill the distinct requirements of AI workload optimization.

From our view, the ability of the Broadcom E2E portfolio providing support for either the switch-scheduled fabric solution, or the endpoint-scheduled fabric solution, delivers all the essential features required to attain top-performing AI training clusters, including especially minimizing JCTs. This includes high bandwidth, efficient E2E congestion management, load balancing, fabric management, telemetry capabilities, and cost advantages over InfiniBand.



Furthermore, Ethernet has a diverse ecosystem of numerous silicon vendors, OEMs, ODMs, cables, optics, software, and continuous innovations. The recently launched Ultra Ethernet Consortium (UEC) aims to standardize features for high-performant networks for large-scale AI/ML and HPC networks, furthering Ethernet technology's deployment agility and democratizing an already vibrant ecosystem. The Broadcom scheduled fabric solutions are aligned with the vision of UEC and are optimized to provide superior AI/ML networking performance.



Broadcom E2E Ethernet Portfolio: Major Cost Improvements

With Broadcom E2E portfolio capabilities, the Jericho3-AI fabric provides at least 10% shorter JCTs in comparison to alternative networking solutions for key AI benchmarks such as AII-to-AII. From our view, this performance improvement has a multiplicative effect on reducing the cost of running AI workloads as it indicates that expensive AI accelerators are used 10% more efficiently. As such, the total cost of ownership (TCO) is lowered through the network-wide performance gains delivered by Jericho3-AI. In sum, the network cost justifies itself.

Broadcom E2E Ethernet Portfolio: Sharp Power Advantages

Integral to Broadcom's E2E Ethernet portfolio advantages are the immense power savings that are essential to supporting the economic and societal justification for AI networking buildouts and the overall scaling of AI workloads. Notably, the Jericho 3-AI provides 40% power reduction savings in comparison to the previous generation. CPO and LPO deliver up to 25% to 33% system power reduction versus conventional pluggable optics. Broadcom can spotlight how its high-performance 400G suite consumes 25% less power than competition. Moreover, Broadcom 400G PCIe Gen 5.0 Ethernet adapters take advantage of 5nm process technology to deliver major power and thermal efficiency advances. The Jericho3-AI fabric offers 26 petabits per second of Ethernet bandwidth, almost four times the bandwidth of the previous generation, while simultaneously delivering 40% lower power per gigabit.

Moreover, Broadcom actively supports Software for Open Networking in the Cloud (SONiC), fulfilling customer demand for the disaggregated OS capabilities integral to cloud data center optimization of AI/ML workloads. SONiC works on over 100 different platforms, enabling cloud service providers to share the same workload stack across hardware from multiple vendors. This directly aligns with Broadcom's collaboration with a large ecosystem of ODM, OEM, Network Operating System (NOS), and orchestration/monitoring partners.





Path Forward – UEC's Vital Role in Advancing Al Training and Inference

We find that Al's destiny is firmly embedded in Ethernet. To that end, Broadcom's active role in Ultra Ethernet Consortium (UEC) shows there is a path forward for the Al ecosystem to optimize Ethernet to ensure high-performance networking in Al/ML workload environments. UEC, hosted by the Linux Foundation, prioritizes delivering an Ethernet-based, open, and interoperable high-performance full communications stack architecture that surpasses current specialized solutions, including particularly proprietary InfiniBand.

From our view, UEC has a solid foundation with a growing membership now numbering over 90 industry leaders including Arista, Eviden, HPE, Meta, and Microsoft. UEC emphasizes delivering the higher scale, bandwidth density, multi-pathing, fast congestion reaction, and low tail latency that AI workloads require. Underscoring meeting this goal is UEC's focus on achieving TCO, developer and user friendliness, performance, and functionality gains.

Key to the UEC mission is to advance beyond existing protocols that can only address some of the unique demands of scaling AI workloads. Moreover, UEC intends to coordinate with established standards development organizations to help ease adoption of UEC technology with the widely deployed family of Ethernet-based products and solutions.

We expect that UEC can prove pivotal in advancing Ethernet innovations related to the Ultra Ethernet Transport (UET) protocol. Specifically, UET will be a new RDMA protocol that is designated to replace ROCE alongside new APIs that replace the Verbs API from the InfiniBand legacy.

Important UET features include multipath spraying, out-of-order packet delivery, and novel rate-control algorithms. Broadcom is contributing intellectual property related to edge-queued datagram service (EQDS) that can provide precise E2E congestion control to UET, a testament to Broadcom's influence across the entire Ethernet ecosystem.





Conclusions and Recommendations

Ethernet is best suited for meeting the demands of rapidly evolving networking and scaling demands of AI clusters. Broadcom E2E networking solutions are integral to ensuring the successful implementation of the largest Ethernet installations worldwide and AI clusters are clearly no exception. The foundation for the competitive advantages of Broadcom's portfolio are the delivery of major cost and power efficiency improvements over alternative solutions. The portfolio-wide merits of Broadcom's E2E networking solution is underpinned by the company's strategic commitment to keeping Ethernet open. As such, AI networking decision makers should prioritize the following considerations:

- Build AI Clusters on Open Ethernet Principles: Using Ethernet on an open basis, reinforced by UEC standards, organizations can control their own journey in the networking and scaling of AI clusters across their network and data centers, avoiding the pitfalls of proprietary alternatives.
- Prioritize Power and Cost Efficiencies: Organizations should prioritize evaluating Broadcom E2E Ethernet solutions, such as switch-scheduled fabric based on Jericho3-AI, endpoint-scheduled fabric based on Tomahawk 5, and 400G RoCE/RDMA Ethernet NICs, to achieve greater levels of power and cost efficiencies throughout their AI cluster implementations by using Broadcom's open hardware and software capabilities. This would enable them to customize and enhance their networks to meet the specific demands of AI networking.
- Ultra Ethernet's Integral Role: Organizations need to take advantage of the extensive community of contributors
 that drive Ethernet Innovation and readiness as exemplified by the ecosystem-wide focus of Ultra Ethernet
 Consortium and Broadcom's integral role. This is important to meet the swiftly evolving demands of AI training and
 inference using next-generation AI interconnect capabilities.



Important Information About this Report

CONTRIBUTORS

Ron Westfall

Research Director | The Futurum Group

David Nicholson

Chief Technology Advisor | The Futurum Group

Daniel Newman

Chief Analyst | The Futurum Group

PUBLISHER

Daniel Newman

CEO | The Futurum Group

INQUIRIES

Contact us if you would like to discuss this report and The Futurum Group will respond promptly.

CITATIONS

This paper can be cited by accredited press and analysts, but must be cited in-context, displaying author's name, author's title, and "The Futurum Group." Non-press and non-analysts must receive prior written permission by The Futurum Group for any citations

LICENSING

This document, including any supporting materials, is owned by The Futurum Group. This publication may not be reproduced, distributed, or shared in any form without the prior written permission of The Futurum Group.

DISCLOSURES

The Futurum Group provides research, analysis, advising, and consulting to many high-tech companies, including those mentioned in this paper. No employees at the firm hold any equity positions with any companies cited in this document.



ABOUT BOARDCOM

Broadcom Inc. (NASDAQ: AVGO) is a global technology leader that designs, develops, and supplies a broad range of semiconductor, enterprise software and security solutions. Broadcom's category-leading product portfolio serves critical markets including cloud, data center, networking, broadband, wireless, storage, industrial, and enterprise software. Our solutions include service provider and enterprise networking and storage, mobile device and broadband connectivity, mainframe, cybersecurity, and private and hybrid cloud infrastructure. Broadcom is a Delaware corporation headquartered in Palo Alto, CA. For more information, go to www.broadcom.com.

For more information on Broadcom AI technologies and solutions please click https://www.broadcom.com/company/events/ai-day



ABOUT THE FUTURUM GROUP

The Futurum Group is an independent research, analysis, and advisory firm, focused on digital innovation and market-disrupting technologies and trends. Every day our analysts, researchers, and advisors help business leaders from around the world anticipate tectonic shifts in their industries and leverage disruptive innovation to either gain or maintain a competitive advantage in their markets.



CONTACT INFORMATION

The Futurum Group LLC | futurumgroup.com | (833) 722-5337 |

