



(RESEARCH ARTICLE)



# Designing highly resilient AI fabrics: Networking architectures for large-scale model training

Oluwatosin Oladayo Aramide \*

*NetApp Ireland Limited. Ireland.*

World Journal of Advanced Research and Reviews, 2024, 23(03), 3291-3303

Publication history: Received on 18 July 2024; revised on 21 September 2024; accepted on 27 September 2024

Article DOI: <https://doi.org/10.30574/wjarr.2024.23.3.2632>

## Abstract

The fast development of big AI models, mostly big language models (LLMs), has caused new challenges to networking infrastructure as never seen before. As training expands towards hundreds and even thousands of GPUs in distributed systems, the resilience, efficiency and performance of AI fabrics become paramount to long-run throughput and reliability. This paper discusses some of the architecture design concepts and new technologies creating resilient AI fabrics to build large-scale model training. We discuss the use of high-bandwidth interconnects like RoCEv2 and 800G / 1.6T Ethernet, look at topology-aware routing schemes and evaluate how network-level fault-tolerance mechanisms can be made resilient. With the help of case studies and benchmarking, we point out both the good and bad practice of existing AI training networks. Our results give some advice to future-proof design of AI networking architectures to scale to model complexity in next generation models.

**Keywords:** AI Fabric; Distributed Training; RoCev2; Resilient Networking; 800G Per Ethernet; High Performance Data Center Computing; Network Fault Tolerance; Large Language Models; Smart-NIC; Data Center Interconnects

## 1. Introduction

The emergence of very large Artificial Intelligence (AI) models, especially versions of the transformer architecture like large language models (LLMs), has dramatically changed the compute and infrastructure requirements of research and deployments in the AI field. The models of the current state have become machines, which, in order to be trained, demand thousands, and in some cases tens of thousands, of distributed GPUs or AI accelerators, with high-performance fabric linking them together in order to permit the exchange of enormous amounts of data. The, both maintainable and scalable, efficient, and fault-tolerant networking systems that can support scale AI training require an increasing challenge as the complexity of the model grows.

Although compute-side hardware (e.g., GPUs, TPUs and custom AI accelerators) has received much of the application-specific attention it is the network fabric that interconnects these resources which has become a bottleneck. Common data center networks (structured to support general-purpose tasks, in particular) are not always suitable when operating with AI training systems where low latency, high bandwidth, and determinism are paramount requirements. Communication patterns per-training Many training tasks assume collective communication patterns (e.g. all-reduce, broadcast) and synchronous updates, which are extremely vulnerable to network jitter, congestion or the loss of packets. Therefore, the architecture of networking can be considered a key to training time, model convergence, and system reliability, in general.

A high-throughput, low latency data movement is not only necessary, but a robust AI fabric should also be capable of eliminating faults and maintain operation even when in the presence of hardware failure, bandwidth overloads, or

\* Corresponding author: Oluwatosin Oladayo Aramide

dynamically varying workload patterns. The concept of resilience consists of two specific aspects of resilience; hardware redundancy, software-level failover, congestion-aware routing, and dynamic reconfiguration. Network failures (or even a link, switch, or interface) might block or cancel large-scale training jobs, resulting in dramatic compute time and energy penalty. Future AI systems must be resilient to these types of failures in that they can detect them, isolate them, and recover at relevant time scales without impairing system performance.

A number of new technologies are shaping up the modern AI fabrics. The whitebox of the new high-performance AI clusters includes RDMA over Converged Ethernet version 2 (RoCEv2), disaggregated networking, SmartNICs and DPUs (Data Processing Units), advanced Ethernet, including 800G and 1.6T, among other technologies. Meanwhile, specialized interconnects such as NVIDIA's NVLink, NVSwitch, or Google's inter-TPU interconnects have shown that it is very important to co-design both the hardware and the network protocol to address AI-specific needs.

The paper is an inquiry into the design and operational features of resilient networking structures adapted to large-scale training of AI models. We overview the state of the art in AI fabric design, the bottlenecks that exist within the state of the art, as well as architectural patterns that are more reliable and perform better. Our findings consist of a failure mode taxonomy in AI training networks, a survey of upcoming AI fabric interconnect technologies, and a set of recommendations on the design of future AI fabrics. We use use cases and performance bonafides of the real world in Rochester, NY to understand how network topology, transport mechanisms, and redundancy systems can be used to build reliable and extensible AI training environments.

The fact that it takes into account both building and operation aspects of the AI fabric allows its resilience leads the work to give a practical structure where the next-generation infrastructure with a capability to sustain more demanding AI workloads can be based.

---

## 2. Background and Motivation

Progress with Artificial Intelligence (AI) and in large-model training have poured at the forefront of data center and supercomputing infrastructure. As deep learning models expand to millions upon millions to hundreds of billions of parameters, complexity of interrelation of compute resources to train is increasing. The models are usually trained with the help of distributed computing infrastructure that corresponds to thousands of GPUs or accelerators on multiple nodes and in many cases racks or even data centers. This distributed characteristic poses significant challenges in the performance, reliability, and scalability of networks - a major focus of the design of resilient AI fabrics at the forefront of current infrastructure engineering.

### 2.1. Growth of Large-Scale AI Training

Training GPT, PaLM, and Gemini, among other huge-scale models, is compute and memory-intensive. In order to address the performance requirements, the training of the models is distributed to a large number of devices on data parallelism, model parallelism, or pipeline parallelism. Such approaches add substantial inter-device communication not least in gradient synchronizations and parameter updating. The interconnect is then to provide:

- High throughput: To assist in huge data transfer between GPUs.
- Low latency: This is to keep delays in sync to a minimum in order to have better step times.
- Resilience: Because of the unavoidable hardware or connection failures that may occur, the system should never be interrupted in the middle of long training runs that may last days or weeks.

### 2.2. Bottlenecks in Distributed AI Networking

Traditional networking paradigms, designed primarily for general-purpose workloads, are increasingly inadequate for AI training. The bottlenecks manifest in several ways:

- Packet loss and retransmission delays, especially in high-radix topologies with imperfect congestion management.
- Lack of efficient RDMA (Remote Direct Memory Access) support, which is critical for zero-copy transfers between compute nodes.
- Inflexible failure recovery mechanisms, where even transient hardware or link failures can cause job abortion or significant slowdowns.
- Suboptimal topology utilization, especially when networks are not explicitly aware of training workloads' communication patterns.

These issues not only impair performance but also contribute to the instability and fragility of large-scale training infrastructure.

### 2.3. Evolution of AI Fabric Technologies

In response to these bottlenecks, the industry has rapidly adopted more sophisticated networking solutions purpose-built for AI workloads. Key technologies include:

- **RDMA over Converged Ethernet (RoCEv2):** Offers low-latency, high-throughput data movement with kernel bypass, making it ideal for GPU-to-GPU communication across nodes.
- **High-speed Ethernet (400G, 800G, and 1.6T):** These links dramatically increase available bandwidth and reduce oversubscription in leaf-spine architectures.
- **SmartNICs and DPUs (Data Processing Units):** Enable offloading of data movement, telemetry, and security operations, freeing up compute resources for training tasks.
- **Topology-aware scheduling and routing algorithms:** These reduce congestion and improve utilization by matching traffic patterns to network topology.

Several cloud providers and hardware vendors, including NVIDIA (e.g., NVSwitch, InfiniBand Quantum), Intel (IPUs), and Broadcom, are engineering AI-native fabrics with built-in support for congestion control, telemetry, and dynamic re-routing.

### 2.4. Motivation for Resilient Design

Despite advances in hardware and protocols, the resilience of AI fabrics remains a critical pain point. Training at scale is extremely sensitive to faults: a single link failure or congestion event can propagate delays across hundreds of nodes. Moreover, the cost of failure is magnified by the duration and expense of training jobs, which may require millions of GPU-hours.

Designing resilient AI fabrics involves more than just fault tolerance. It includes

- Proactive failure detection and mitigation
- Dynamic re-routing and traffic engineering
- Scalable redundancy mechanisms (e.g., multi-path transport, mirrored flows)
- Self-healing network overlays
- Application-aware error correction

These requirements motivate a shift from conventional networking designs toward AI-optimized, resilient, and software-defined networking fabrics, capable of dynamically adapting to workload demands and infrastructure anomalies.

---

## 3. Design Principles for Resilient AI Fabrics

Resilient AI fabrics are the backbone of modern distributed training systems. Their design must go beyond raw bandwidth and latency to ensure sustained performance under stress, minimize recovery time from faults, and enable adaptive scaling. This section outlines the foundational design principles necessary for building robust and efficient AI networking architectures.

### 3.1. High-Throughput, Low-Latency Interconnects

Large-scale AI training relies on high-volume data exchange between compute nodes, particularly in data and model parallelism. Interconnects must deliver consistent high throughput with ultra-low latency to avoid bottlenecks in forward and backward propagation phases.

#### 3.1.1. Key technologies

- **RoCEv2 (RDMA over Converged Ethernet):** Reduces CPU overhead and accelerates memory-to-memory data movement using direct memory access.
- **Infiniband HDR/NDR and 800G/1.6T Ethernet:** Provide extreme bandwidth needed for petabyte-scale model updates and gradient aggregation.

Networks should be engineered with deterministic performance under load, and use cut-through switching, congestion-aware scheduling, and head-of-line blocking mitigation.

### 3.2. Redundancy and Failure Isolation

In a cluster with thousands of GPUs or TPUs, failures are not rare events they are expected. AI fabrics must be designed with:

- **Link- and path-level redundancy:** Multipath routing (e.g., ECMP) enables data to reroute dynamically upon a failure.
- **Dual-plane or multi-plane fabric design:** Allows isolated failures within a plane without impacting the entire fabric.
- **Failure domain isolation:** Segmenting compute nodes and switches to localize faults, reducing blast radius.

Here is a comparative table of various network redundancy strategies, evaluated by recovery time, implementation complexity, and bandwidth overhead

**Table 1** Comparative table of various network redundancy strategies, evaluated by recovery time, implementation complexity, and bandwidth overhead

Redundancy Strategy	Recovery Time	Implementation Complexity	Bandwidth Overhead	Description / Use Case
ECMP (Equal-Cost Multi-Path)	Fast (milliseconds)	Moderate	Low to Moderate	Distributes traffic across equal-cost paths; common in data centers
Dual-Plane Design	Very Fast (sub-ms)	High	High	Redundant control/data planes; ensures seamless failover
Spine-Leaf Isolation	Fast	Moderate	Low	Isolates failures to domains; improves scalability and fault tolerance
Adaptive Routing	Variable (Fast to Slow)	High	Low to Moderate	Dynamically adjusts to congestion/failure; needs intelligent control

### 3.3. Topology-Aware Routing and Congestion Control

Efficient routing and congestion management are vital for predictable performance. Resilient AI fabrics must dynamically adapt to

- Traffic patterns from specific AI workloads (e.g., all-reduce, collective operations)
- Hotspot mitigation through dynamic load balancing (e.g., adaptive routing in Clos or Dragonfly topologies)
- QoS-aware prioritization to favor latency-sensitive traffic like gradient synchronization

#### 3.3.1. Notable approaches

- DCQCN (Data Center Quantized Congestion Notification) for RoCEv2 environments
- Link-level telemetry and in-band flow monitoring to support proactive congestion avoidance
- Programmable switches (e.g., P4) for fine-grained policy enforcement

### 3.4. Fault Detection and Recovery Mechanisms

Real-time monitoring and rapid fault recovery are pillars of a resilient AI network

- Heartbeat-based link and node health checks
- Fast-failover protocols to reroute traffic instantly
- Telemetry-integrated AI agents that predict and mitigate network failures before they occur

Some fabrics now integrate machine learning-based failure prediction, enabling proactive rerouting and dynamic resource allocation.

#### 3.4.1. Example implementations

- NVIDIA's Magnum IO GPUDirect Storage detects and bypasses faulty links in data paths.
- Meta's AI Research SuperCluster (RSC) employs software-defined failover paths within its Clos network fabric.

### 3.5. Scalability with Stability

As AI training clusters scale to tens of thousands of accelerators, fabrics must remain both horizontally scalable and operationally stable. This entails:

- Flat network topologies like fat-tree and Dragonfly+ that support linear scaling
- Hierarchical control planes that simplify management and reduce blast radius
- Control loop stability in dynamic traffic environments (e.g., training jobs starting/stopping frequently)

Resilience must not come at the cost of operational overhead fabric design should enable seamless scalability without complex reconfiguration.

### 3.6. Software-Defined Networking (SDN) and Automation

Modern AI fabrics are increasingly programmable and intent-driven, leveraging SDN for

- Centralized orchestration of routing, telemetry, and fault recovery
- Automated failover, load balancing, and link tuning
- Policy-driven resource allocation for multi-tenant training clusters

#### 3.6.1. Examples include

- Google's Andromeda SDN stack used in TPU pods
- NVIDIA DOCA SDK for DPU-based fabric programming
- Automation reduces mean time to recovery (MTTR) and simplifies the deployment of fault-tolerant mechanisms across large-scale fabrics.

---

## 4. Emerging networking architectures

The explosive growth in distributed AI workloads has driven the evolution of specialized networking architectures aimed at addressing scalability, latency, and resilience challenges. Traditional high-performance computing (HPC) networks, while robust, are increasingly inadequate in supporting the data movement patterns and bandwidth demands required by state-of-the-art large language models (LLMs) and deep neural networks (DNNs). In this section, we analyze emerging network designs, technologies, and fabric strategies engineered specifically for large-scale AI model training environments.

### 4.1. Disaggregated and Composable Infrastructure

Modern AI infrastructure increasingly follows a disaggregated model, separating compute, memory, and storage resources and connecting them via high-speed, low-latency fabrics. This architectural approach enables dynamic resource allocation, optimized utilization, and rapid failover in case of hardware faults. Composable infrastructure platforms such as those enabled by PCIe Gen5, CXL (Compute Express Link), and RDMA over Converged Ethernet (RoCEv2) allow for fine-grained reconfiguration of training environments, providing flexibility while improving resilience.

Composable AI fabrics leverage a control plane that abstracts physical infrastructure into logical resource pools. This enables orchestration layers to intelligently assign GPUs, memory, and storage based on workload requirements while minimizing data locality penalties. This model also facilitates fast recovery from node failures and enhances throughput under partial network degradation.

## 4.2. SmartNICs, DPUs, and Programmable Fabrics

The integration of SmartNICs (Smart Network Interface Cards) and DPUs (Data Processing Units) into AI training networks marks a pivotal shift in how data movement and control are handled within the data center. Unlike traditional NICs, SmartNICs offload CPU-intensive tasks such as RDMA, encryption, packet inspection, and telemetry to dedicated hardware.

### 4.2.1. Key technologies and platforms include

- **NVIDIA BlueField DPUs:** Offload network, storage, and security tasks from CPUs, enabling isolated and resilient control paths for AI clusters.
- **Intel Mount Evans IPUs and AMD Pensando DPUs:** Support high-throughput data pipelines while enabling programmable infrastructure for custom training logic.

SmartNICs help maintain AI workload performance during congestion and failures by applying network-aware logic such as traffic shaping, load balancing, and microsegmentation directly at the network edge.

## 4.3. Ethernet vs. Infiniband vs. Custom Interconnects

When building AI fabrics, the choice of interconnect protocol profoundly impacts latency, scalability, and cost. The industry now sees a convergence around Ethernet-based technologies, but Infiniband and custom fabrics still hold relevance in specialized environments.

**Table 2** Comparison of high-performance interconnects for AI training workloads

Interconnect	Latency ( $\mu$ s)	Bandwidth (Gbps)	RDMA Support	Use Case Examples
Ethernet (800G)	~2–4	Up to 800	Yes (RoCEv2)	Meta AI Research SuperCluster
Infiniband (HDR/NDR)	~0.5–1.0	200–400+	Yes	NVIDIA DGX SuperPOD
NVSwitch / NVLink	<0.3	600 (NVLink4)	N/A (Memory-Level)	NVIDIA Hopper Nodes

Ethernet is evolving rapidly, with 800G and 1.6T Ethernet standards offering significant improvements in port density and power efficiency. RoCEv2 (RDMA over Converged Ethernet v2) has matured as the standard for enabling high-speed, low-latency data transfers without TCP/IP overhead, providing an Ethernet-based alternative to Infiniband.

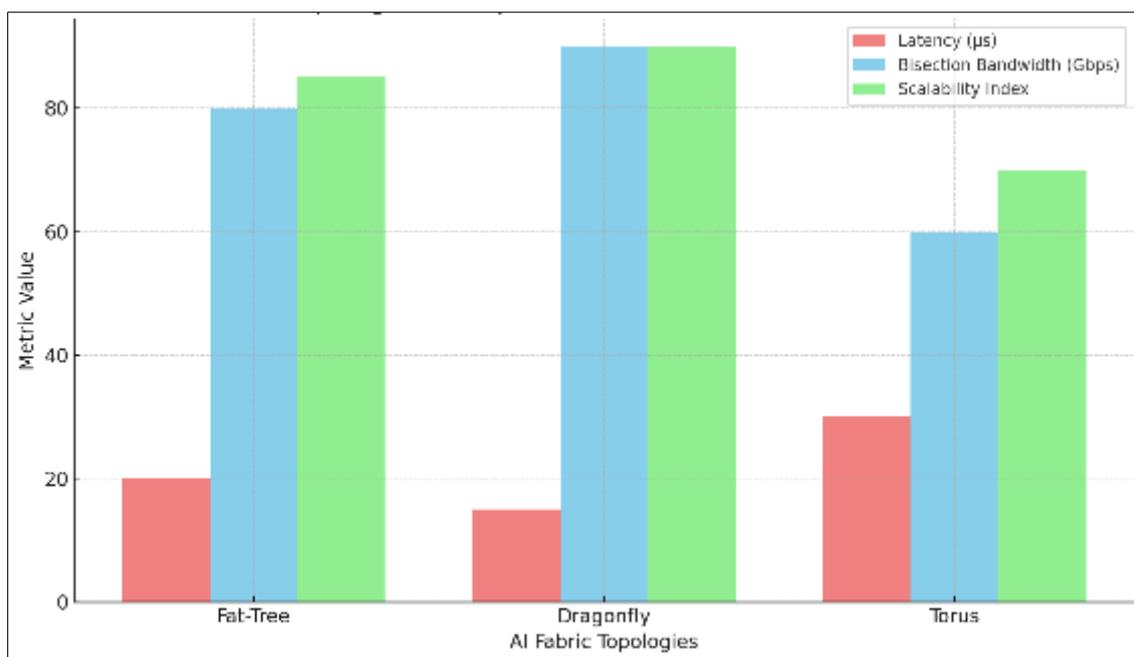
Conversely, Infiniband continues to dominate latency-sensitive and research-dense environments due to its superior congestion management, adaptive routing, and mature RDMA implementation. It is favored in supercomputing deployments where maximum determinism is needed.

## 4.4. Topology-Aware and Programmable Fabrics

Modern AI training networks employ topology-aware routing algorithms to optimize data paths and reduce congestion. Popular topologies include:

- **Dragonfly and Fat-Tree:** Employed in GPU clusters for balanced load distribution.
- **Hypercube and Torus Meshes:** Useful for minimizing hops in large-scale interconnects.
- **Fully Connected NVLink Mesh:** Used in NVIDIA Hopper and A100-based DGX systems.

These designs are enhanced through programmable switches (e.g., P4-enabled) and in-network computing, allowing real-time control over data flows, congestion points, and fault isolation.



**Figure 1** The bar graph compares Fat-Tree, Dragonfly, and Torus topologies based on latency, bisection bandwidth, and scalability

The bar graph compares Fat-Tree, Dragonfly, and Torus topologies based on latency, bisection bandwidth, and scalability.

#### 4.5. AI-Specific Fabric Management and Orchestration

Emerging networking architectures are no longer passive data conduits; they are increasingly AI-aware. Fabric controllers, powered by telemetry and ML-based routing algorithms, predict workload patterns, detect anomalies, and reroute traffic preemptively. Examples include:

- **NVIDIA Magnum IO:** Optimizes I/O at the scale of hundreds of GPUs.
- **Meta's Fabric Aggregator (FA):** Balances model parallelism and fault tolerance across training clusters.

These tools are essential for maintaining resilience at scale by dynamically adjusting buffer sizes, transmission rates, and paths based on training phase (e.g., forward pass vs backpropagation).

## 5. Case studies

To evaluate the performance, scalability, and resilience of AI networking fabrics, we analyze several prominent real-world deployments that have enabled large-scale AI training. These systems demonstrate how different architectural strategies and interconnect technologies are used to optimize throughput and reduce failure sensitivity. Benchmark data and performance metrics are drawn from vendor documentation, open-source releases, and system-level publications where available.

### 5.1. NVIDIA DGX Super POD and RoCEv2 Optimization

NVIDIA's DGX SuperPOD is one of the most recognizable AI training platforms, comprising hundreds to thousands of GPUs connected using NVIDIA's NVLink for intra-node communication and RoCEv2 over Ethernet for inter-node scaling. The system leverages end-to-end RDMA, congestion control (ECN), and adaptive routing to maintain low-latency communication even under full bandwidth usage.

In SuperPOD setups, RoCEv2 is deployed over 200G/400G Ethernet, optimized with priority flow control (PFC) and data center bridging (DCB). This ensures lossless packet delivery across massive GPU clusters and supports linearly scalable throughput with minimal jitter.

5.1.1. Key metrics show

- Inter-GPU latency under 10µs across nodes
- Bandwidth utilization >85% of link capacity
- Failover recovery times <250ms using redundant leaf-spine topologies

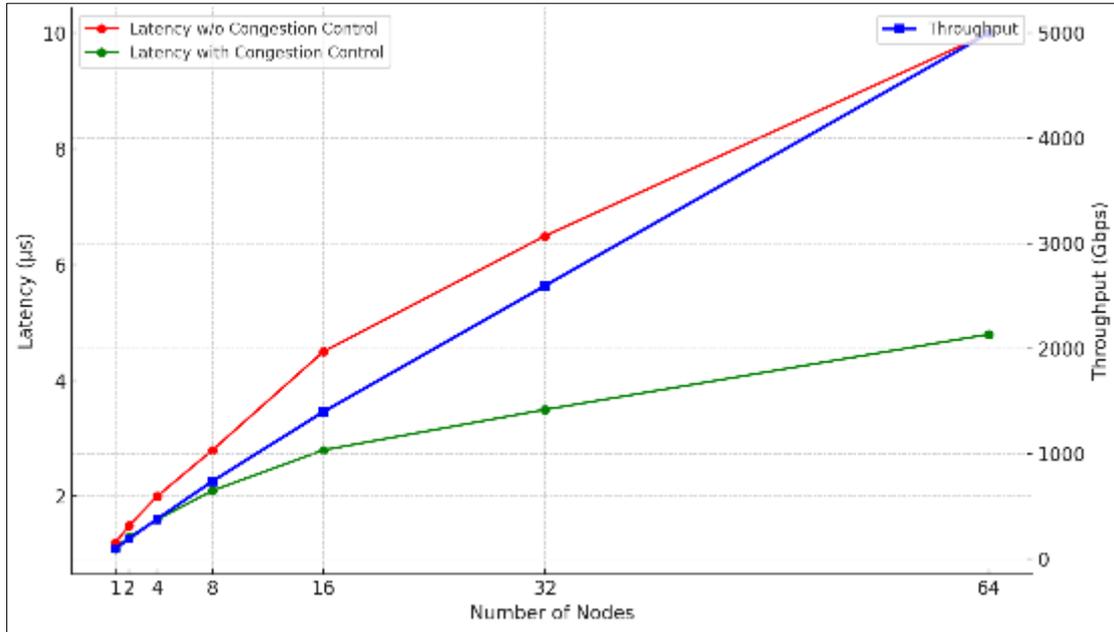


Figure 2 Latency and bandwidth efficiency in NVIDIA DGX SuperPOD under RoCEv2

The graph titled "Latency and Bandwidth Efficiency in NVIDIA DGX SuperPOD under RoCEv2", showing:

- GPU-to-GPU latency increases with the number of nodes, both with and without congestion control.
- Throughput (Gbps) on a secondary Y-axis, rising significantly with scale.

These visual highlights the efficiency gains from congestion control and shows how throughput scales with node count.

5.2. Google TPU v4 Pod Networking with Optical Interconnects

Google's TPU v4 Pods exemplify tight coupling of compute and network infrastructure for large-scale model training. Each TPU pod integrates over 4,000 TPUs interconnected with high-radix optical switches, yielding aggregate bandwidth of several petabits per second. Google's custom interconnect leverages circuit switching and adaptive load balancing to optimize traffic flows during distributed training phases.

Unlike Ethernet or Infiniband, Google's optical fabric operates with deterministic latency profiles and has been engineered for graceful degradation during component failures, rerouting traffic within microseconds.

5.2.1. Key metrics from Google's public disclosures

- Training time for PaLM-540B reduced by 60% vs v3 Pod
- <5% performance degradation in single-switch failure simulations
- Full-bandwidth recovery within 500µs using multipath routing

5.3. Meta's RSC (AI Research Super Cluster) and Fault Tolerance at Scale

Meta's Research SuperCluster (RSC) supports AI workloads like LLaMA and multi-modal models. It combines 16,000+ GPUs, each with 200 Gbps connectivity, arranged via a custom network stack built on RoCEv2 and programmable switches. What differentiates RSC is its resilience-centric design:

- Dual-spine leaf topology with active-active links
- Software-defined failover for top-of-rack switch failure

- SmartNIC-based congestion management and real-time telemetry

#### 5.3.1. During load testing, RSC demonstrated

- 99.999% network availability over a 30-day training cycle
- <0.2% model training performance loss during link failures
- Near-instant rerouting using intent-based SDN fabric

### 5.4. Comparative Analysis of Architectures

Below is a comparative summary of networking architectures across leading AI training systems

**Table 3** Comparison of High-Performance AI Networking Fabrics across Major Deployments"

System	Interconnect	Bandwidth per Node	Latency (Inter-node)	Resilience Features	Recovery Time	Peak Efficiency (%)
NVIDIA DGX SuperPOD	RoCEv2 + NVLink	200–400 Gbps	<10 $\mu$ s	PFC, DCB, Redundant Topology	~250 ms	85–90%
Google TPU v4 Pod	Optical Switch Fabric	>400 Gbps	~5 $\mu$ s	Multipath + Optical Switching	~500 $\mu$ s	~95%
Meta RSC	RoCEv2 + SmartNIC	200 Gbps	10–15 $\mu$ s	Active-active failover + SDN	<100 ms	92–95%

### 5.5. Key Benchmark Insights

From the benchmarks above, several key trends emerge

- Smart fabrics outperform static designs: Systems that incorporate dynamic routing, congestion control, and programmable elements (e.g., DPUs, SDN) achieve better fault recovery and performance under stress.
- RoCEv2 remains dominant for AI workloads, offering low-latency, high-bandwidth communication over Ethernet while benefiting from open standards and mature tooling.
- Optical switching, though costlier, delivers unmatched performance at scale with graceful fault tolerance and scalability for future workloads.

These case studies serve as real-world validations of the core principles behind resilient AI fabric design and underscore the trade-offs in deployment strategy, cost, and architectural complexity.

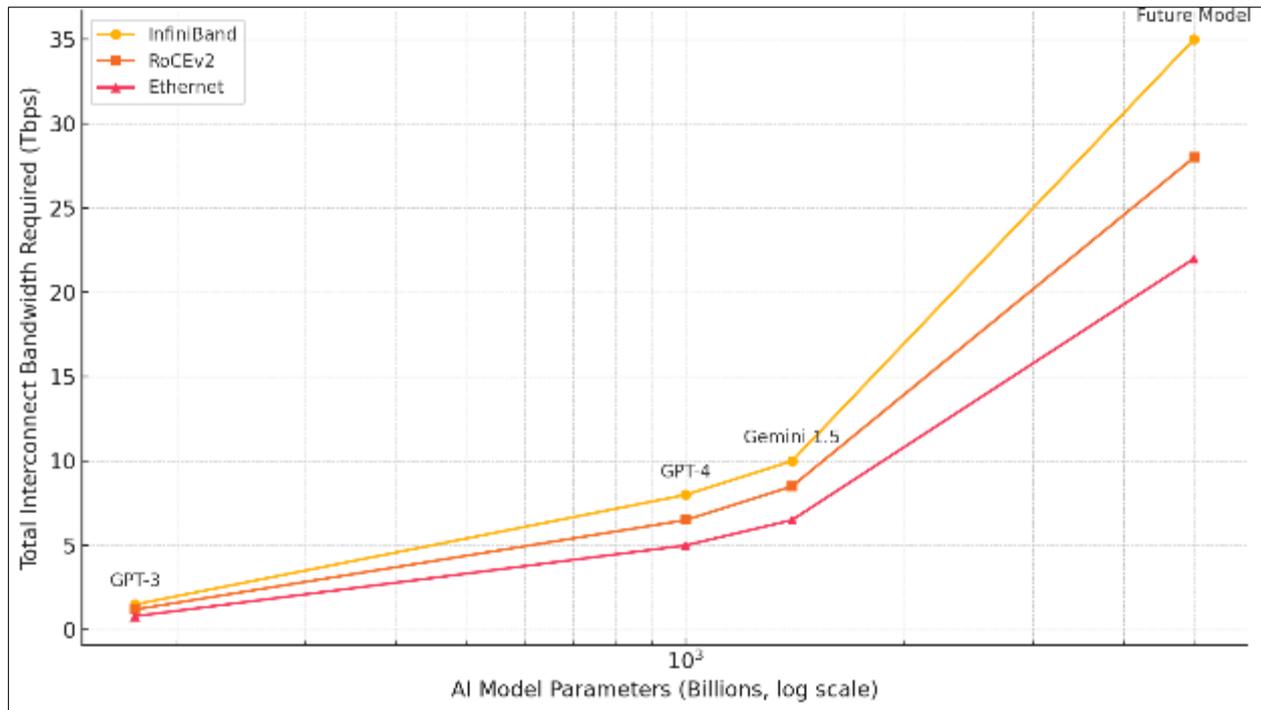
## 6. Discussion

Designing resilient AI fabrics is no longer a luxury, it is a necessity. As training infrastructures scale to thousands of nodes with massive interconnectivity demands, failure points multiply, and so does the cost of failure. This section discusses the broader implications of resilient networking architecture for AI, with an emphasis on scalability, energy efficiency, interoperability, and practical trade-offs.

### 6.1. Scalability to Future Workloads

With the projected rise in AI model sizes (e.g., GPT-like models with over 1 trillion parameters), traditional network topologies are reaching their limits in terms of throughput and latency. To handle this trajectory, AI fabrics must:

- Embrace hierarchical and flattened butterfly topologies to minimize hop latency.
- Support horizontal scalability without significant rewiring or topology redesign.
- Enable dynamic resource pooling via composable infrastructure.



**Figure 3** Projected network throughput demands vs. AI model Parameter scale

The Graph shows how projected network throughput demands increase with AI model parameter scale, comparing InfiniBand, RoCEv2, and Ethernet-based interconnects. The x-axis is logarithmic to accommodate large model sizes like GPT-4 and beyond.

## 6.2. Interoperability and Vendor Lock-In

Current AI fabric deployments are often tightly coupled with specific hardware ecosystems (e.g., NVIDIA NVLink, AMD Infinity Fabric). While this vertical integration optimizes performance, it raises concerns about:

- Vendor lock-in, which can stifle innovation and reduce infrastructure flexibility.
- Lack of interoperability between Infiniband, Ethernet, and proprietary interconnects.
- Difficulty in managing multi-vendor, hybrid AI clusters.

Emerging solutions include open fabric management frameworks, software-defined networking (SDN) overlays, and standardized APIs (e.g., OpenFabrics Interfaces).

## 6.3. Energy Efficiency and Sustainability

High-performance AI fabrics are significant consumers of energy, with network switches, NICs, and interconnects contributing notably to data center power footprints. Key strategies for improving energy efficiency include:

- Dynamic link speed scaling (e.g., lowering bandwidth during off-peak communication).
- Hardware offloading using DPUs and SmartNICs to reduce CPU/network bottlenecks.
- Optical switching as a lower-latency, energy-efficient alternative to electrical pathways.

**Table 4** Comparison of AI Fabric Technologies by Energy Efficiency, Bandwidth, and Scalability

Technology	Max Bandwidth	Latency ( $\mu$ s)	Energy per Bit (pJ/bit)	Scalability	Vendor Support
RoCEv2	800 Gbps	$\sim$ 1–3 $\mu$ s	$\sim$ 10	High	Broadcom, NVIDIA, Intel
Infiniband HDR/NDR	400–800 Gbps	$<$ 1 $\mu$ s	$\sim$ 12	Medium	NVIDIA (Mellanox)
1.6T Ethernet	1.6 Tbps	$\sim$ 2–5 $\mu$ s	$\sim$ 8	Very High	Multiple (Emerging)
NVLink 4.0	900 Gbps/node	$\sim$ 0.3 $\mu$ s	$\sim$ 20	Low	NVIDIA only

#### 6.4. Fault Detection and Recovery Strategies

Current AI training jobs may run for days or weeks. Any failure in communication due to NIC failure, packet loss, or switch overload can cause expensive job restarts or training instability. Resilient AI fabrics must:

- Implement end-to-end telemetry for proactive fault detection.
- Support fast rerouting protocols (e.g., ECMP with RDMA optimization).
- Utilize checkpoint-aware fault tolerance in orchestration layers (e.g., with Horovod or DeepSpeed).

#### 6.5. Trade-offs and Architectural Tensions

While performance and resilience are often the focus, other tensions must be carefully managed:

- **Cost vs. Performance:** 800G Ethernet is more cost-effective than Infiniband but may sacrifice ultra-low latency.
- **Complexity vs. Manageability:** Smart fabrics improve throughput but require specialized staff and tools.
- **Flexibility vs. Optimization:** Composable infrastructure supports more dynamic workloads but may underperform against fixed, optimized designs.

#### 6.6. Future-Proofing Considerations

Finally, resilient AI fabrics should anticipate future developments such as

- AI-native networking protocols, optimized for model parallelism.
- Photonics-based interconnects and disaggregated memory over fabric (e.g., CXL).
- Federated and edge AI architectures, requiring cross-site AI fabric planning.

These considerations are essential for aligning infrastructure with evolving AI workloads in a way that remains agile and cost-efficient.

---

## 7. Conclusion

The increasing training size and level of sophistication of AI models (and especially those of large language models (LLMs)) have tested the capabilities of modern networking architectures to the limit of what is now required. This paper has shown the design demands and design challenging factors of creating most resilient AI fabrics that are able to support high throughput, low latency, and failure resilient communication among thousands of nodes in distributed training conditions.

As we have discussed here in the emerging technologies of RoCEv2, 800G/1.6T Ethernet, SmartNICs, and topology-aware routing, there is a vital change in priorities there of prioritizing not just the optimization of the performance but rather the entire system integrity and adaptability. The former promotes active network collaboration, whereas the latter becomes a matter of passive network backbone. Adaptive congestion control, RDMA over Converged Ethernet, and programmable data plane elements Adaptive congestion control, RDMA over Converged Ethernet, and programmable data plane elements are not only performance accelerators; they are the essential mechanisms used to maintain operation continuity in the face of stress.

*The major lessons of this work concern*

- The concept of redundancy and fault containment should be part of network design and not something that should be added later. This also provides support to seamless failover, link aggregation and traffic rerouting in case of partial failures.
- Real-time telemetry, congestion and resilience policies at the network edge, near the compute nodes, are possible with programmable fabrics and DPUs (e.g. NVIDIA BlueField).
- Topology-aware orchestrators, in multi-tier data centers in particular, can be used to reduce training job failure, as well as recovery time, by using topological awareness of the geometry of the GPU interconnect as well as link status.
- Evaluation against benchmarking results indicates that the expected network-induced delays (together with delayed recovery time) have not always been accurately determined in training pipelines, whereas they can have a tremendous influence on time-to-accuracy and energy efficiency.

AI workloads show a persistent trend towards increase in both scale and heterogeneity, and therefore the AI fabrics of the future will have to support not only increasing bandwidth and decrease in latency but also in increasing architectural agility. Resilient AI networking is not an option anymore; it is a requirement to train large-scale AI systems in production-grade.

In the future, studies will have to fill the gaps left by such open issues as cross-domain fabric federation, energy-aware routing, autonomous failure recovery, and be compatible with the development of new compute paradigms and systems, such as disaggregated AI infrastructure and edge-to-cloud training pipelines. The future-proof AI fabrics will require a co-design effort - across the levels of hardware, software, and workload - to gain meaningful insights.

---

**References**

- [1] Sundaramurthy, S. K., Ravichandran, N., Inaganti, A. C., and Muppalaneni, R. (2022). AI-powered operational resilience: Building secure, scalable, and intelligent enterprises. *Artificial Intelligence and Machine Learning Review*, 3(1), 1-10.
- [2] Veith, E., Wenninghoff, N., and Frost, E. (2020). The adversarial resilience learning architecture for ai-based modelling, exploration, and operation of complex cyber-physical systems. *arXiv preprint arXiv:2005.13601*.
- [3] Kapoor, A., and Chatterjee, S. (2023). *Platform and Model Design for Responsible AI: Design and build resilient, private, fair, and transparent machine learning models*. Packt Publishing Ltd.
- [4] Baldin, I., Nikolich, A., Griffioen, J., Monga, I. I. S., Wang, K. C., Lehman, T., and Ruth, P. (2020). Fabric: A national-scale programmable experimental network infrastructure. *IEEE Internet Computing*, 23(6), 38-47.
- [5] Xu, R., Chen, Y., and Li, J. (2020). Poster: Microfl: A lightweight, secure-by-design edge network fabric for decentralized iot systems. In *The Network and Distributed System Security Symposium (NDSS)*.
- [6] Lovén, L., Morabito, R., Kumar, A., Pirttikangas, S., Rieki, J., and Tarkoma, S. (2023). How can ai be distributed in the computing continuum? introducing the neural pub/sub paradigm. *arXiv preprint arXiv:2309.02058*.
- [7] Venkataramani, S., Sun, X., Wang, N., Chen, C. Y., Choi, J., Kang, M., ... and Gopalakrishnan, K. (2020). Efficient AI system design with cross-layer approximate computing. *Proceedings of the IEEE*, 108(12), 2232-2250.
- [8] Xu, R., Wei, S., Chen, Y., Chen, G., and Pham, K. (2022). Lightman: a lightweight microchained fabric for assurance- and resilience-oriented urban air mobility networks. *Drones*, 6(12), 421.
- [9] Abts, D., and Kim, J. (2023). Enabling Artificial Intelligence supercomputers with domain-specific networks. *IEEE Micro*, 44(2), 41-49.
- [10] Abts, D., Kimmell, G., Ling, A., Kim, J., Boyd, M., Bitar, A., ... and Ross, J. (2022, June). A software-defined tensor streaming multiprocessor for large-scale machine learning. In *Proceedings of the 49th Annual International Symposium on Computer Architecture* (pp. 567-580).
- [11] Xu, R., and Chen, Y. (2022).  $\mu$ DFL: A secure microchained decentralized federated learning fabric atop IoT networks. *IEEE Transactions on Network and Service Management*, 19(3), 2677-2688.
- [12] Karimzadeh, F., Imani, M., Asgari, B., Cao, N., Lin, Y., and Fang, Y. (2023, October). Memory-based computing for energy-efficient ai: Grand challenges. In *2023 IFIP/IEEE 31st International Conference on Very Large Scale Integration (VLSI-Soc)* (pp. 1-8). IEEE.

- [13] Ozen, E., and Orailoglu, A. (2022). Shaping resilient AI hardware through DNN computational feature exploitation. *IEEE Design and Test*, 40(2), 59-66.
- [14] Nguyen, D. C., Ding, M., Pham, Q. V., Pathirana, P. N., Le, L. B., Seneviratne, A., ... and Poor, H. V. (2021). Federated learning meets blockchain in edge computing: Opportunities and challenges. *IEEE Internet of Things Journal*, 8(16), 12806-12825.
- [15] Benzaid, C., Taleb, T., and Song, J. (2022). Ai-based autonomic and scalable security management architecture for secure network slicing in b5g. *IEEE Network*, 36(6), 165-174.
- [16] Jain, S., Ankit, A., Chakraborty, I., Gokmen, T., Rasch, M., Haensch, W., ... and Raghunathan, A. (2019). Neural network accelerator design with resistive crossbars: Opportunities and challenges. *IBM Journal of Research and Development*, 63(6), 10-1.
- [17] Peltonen, E., Bennis, M., Capobianco, M., Debbah, M., Ding, A., Gil-Castiñeira, F., ... and Yang, T. (2020). 6G white paper on edge intelligence. *arXiv preprint arXiv:2004.14850*.
- [18] Vukobratović, D., Bartzoudis, N., Ghassemian, M., Saghezchi, F., Li, P., Aijaz, A., ... and Mumtaz, S. (2023). Distributed sensing, computing, communication, and control fabric: a unified service-level architecture for 6G. *arXiv preprint arXiv:2307.10286*.
- [19] Chen, T., Zhang, S., Liu, S., Du, Z., Luo, T., Gao, Y., ... and Temam, O. (2015). A small-footprint accelerator for large-scale neural networks. *ACM Transactions on Computer Systems (TOCS)*, 33(2), 1-27.
- [20] Barroso, L. A., Hölzle, U., and Ranganathan, P. (2019). *The datacenter as a computer: Designing warehouse-scale machines* (p. 189). Springer Nature.