

*ServerMO Engineering - Proprietary & Confidential*

# Architecting High-Availability AI Clusters: Overcoming Network Bottlenecks with 100Gbps Backbone and Nvidia H100 GPUs

A Technical Whitepaper by ServerMO Engineering Version: 1.2 (Objective Architecture Revision) | Date: April 2026

## 1. Abstract

The infrastructure requirements for Large Language Models (LLMs) and distributed deep learning frequently test the limits of standard virtualized data centers. While public cloud environments are often utilized for variable, stateless web applications, scaling sustained, I/O-heavy AI workloads introduces complex challenges related to the "interconnect wall," storage throughput, and unpredictable data movement costs.

This paper provides an architectural overview of the ServerMO bare-metal framework. We examine the engineering rationale behind integrating up to 100Gbps unmetered networking, RDMA over Converged Ethernet (RoCE v2), and AMD EPYC Genoa platforms to mitigate specific data movement bottlenecks. We also detail the thermal engineering methodologies required to maintain sustained GPU performance across our global footprint.

## 2. The Infrastructure Dilemma: Virtualized Clouds vs. Dedicated Bare Metal

As enterprises transition workloads to AI-centric models, the architectural trade-offs between managed virtual environments and dedicated bare-metal infrastructure must be evaluated objectively based on workload profiles.

### 2.1 Data Gravity and Egress Economics

In the AI lifecycle, data possesses "gravity." Training a model requires shifting massive datasets constantly from storage to compute nodes, and often across regions.

- **The Cloud Trade-off:** The convenience of virtualized clouds often comes with metered data movement. Outbound data transfer (egress) typically incurs fees between \$0.05 and \$0.09 per GB. For sustained workloads moving 500TB monthly, egress bills can introduce significant variable costs that disrupt budget forecasting.
- **The Bare Metal Alternative:** ServerMO targets this specific bottleneck by offering 1Gbps, 10Gbps, 20Gbps, 40Gbps, and 100Gbps unmetered uplink ports. While this shifts infrastructure management to a more direct model, it converts variable network costs into a fixed operational expense, significantly reducing the financial friction of data mobility regardless of scale.

## 2.2 Latency, Jitter, and Virtualization Overhead

Distributed AI training relies heavily on synchronous processes. In collective communications (like NCCL's Ring AllReduce), GPUs must synchronize gradients across nodes before proceeding.

- **The Hypervisor Impact:** Virtualized clouds utilize hypervisors to pool resources. However, this multi-tenancy introduces "noisy neighbor" effects. Depending on the workload and framework tuning, even minor variations in network packet arrival (jitter) can lead to measurable drops in overall training efficiency.
- **The Solution:** Bare-metal infrastructure removes the virtualization layer entirely, granting direct access to the NIC and PCIe lanes. This hardware-level access minimizes latency variation, providing the predictable network environment critical for tightly coupled HPC workloads.

## 3. Network Architecture: High-Bandwidth RoCE v2 Fabric

High-throughput AI clusters require a fabric explicitly engineered to reduce CPU overhead during data transfers.

### 3.1 Intra-Cluster RDMA (Bypassing the CPU)

Standard TCP/IP networking requires OS kernel interrupts to process packets, which introduces latency. ServerMO implements **RoCE v2**, enabling GPUs to read/write directly to the memory of other GPUs across the network. While global latency remains constrained by physical distance, this architecture drops *intra-cluster* (node-to-node within the same switch fabric) latency to sub-microsecond levels.

## 3.2 BGP Path Selection and Global Routing

ServerMO operates a carrier-neutral network with direct peering at major Internet Exchange Points (IXPs) using a blend of NTT, Lumen, Arelion, and Cogent.

- **Traffic Engineering:** While public internet routing is inherently unpredictable, our network engineers utilize custom BGP community strings to influence shortest-path routing. For example, transatlantic links (e.g., Dallas to London) are engineered to target a highly consistent RTT of 88ms to 91ms under normal network conditions.

## 3.3 Edge Security and DDoS Mitigation

AI clusters are high-value targets for both corporate espionage and volumetric cyberattacks. No network provider can guarantee absolute immunity, but proactive risk mitigation is critical.

- **The Architecture:** ServerMO includes 250Gbps DDoS protection as a standard feature across all bare-metal deployments, embedded directly at our edge scrubbing centers.
- **The Mechanism:** By utilizing deep packet inspection (DPI) at the edge, we mitigate the impact of volumetric Layer 3/4 and Layer 7 attacks before they saturate core uplinks. This capability strongly supports our 99.99% uptime SLA without introducing perceptible latency to legitimate traffic.

# 4. Hardware Benchmarks: AMD EPYC Genoa and PCIe Gen 5

Storage throughput is a primary factor in preventing GPU starvation. ServerMO standardizes AI nodes on the AMD EPYC 9004 (Genoa) platform to maximize PCIe Gen 5.0 lane utilization.

## 4.1 Benchmark Methodology and Synthetic Results

To provide baseline transparency, ServerMO conducts internal benchmarks using industry-standard tools like [fio](#).

- **Test Environment:** Dual AMD EPYC 9654 processors, 1.5TB DDR5 ECC RAM, 4x Enterprise NVMe SSDs (Software RAID 10).
- **Observed Synthetic Performance:** Under ideal lab conditions (e.g., 4K block size, high queue depth), the array delivers peak synthetic throughput of **3.2 Million Random Read IOPS**.
- *Engineering Note:* Real-world application I/O (such as PyTorch dataloaders or checkpointing) will yield lower, yet highly competitive, throughput depending on file sizes, batching, and framework efficiency.

## 4.2 Architectural Use-Case Summary

| Technical Metric           | Typical Virtualized Cloud               | ServerMO Bare Metal Architecture                           |
|----------------------------|---|--|
| <b>Optimal Use Case</b>    | Variable traffic, stateless web apps    | <b>Sustained AI training, Heavy I/O, Predictable loads</b> |
| <b>Network Protocol</b>    | Virtualized Networking (TCP/IP focus)   | <b>RoCE v2 / RDMA (Hardware-level bypass)</b>              |
| <b>Storage Performance</b> | Governed by instance size/volume limits | <b>Raw PCIe Gen 5 speeds (Hardware limited)</b>            |
| <b>Egress Cost</b>         | Metered variable cost                   | <b>Unmetered (Fixed port cost)</b>                         |

## 5. Thermal Engineering: Managing the 10kW Rack Challenge

Nvidia H100 SXM5 nodes present severe thermal challenges, pulling up to 10kW per 8-GPU chassis. Standard Computer Room Air Conditioning (CRAC) units generally struggle to efficiently cool densities beyond 15kW per rack.

### 5.1 Hybrid Liquid Cooling Topology

To safely support 50kW+ rack densities, ServerMO utilizes a hybrid cooling approach:

- **Direct-to-Chip (D2C):** Liquid cold plates mounted directly to the GPUs and CPUs capture the majority of the thermal load at the silicon source.
- **Rear Door Heat Exchangers (RDHx):** Active liquid-to-air radiators replace standard rack doors to neutralize the remaining exhaust heat before it enters the ambient room environment.

## 5.2 Efficiency (PUE) and Clock Stability

This topology achieves an internal Power Usage Effectiveness (PUE) of 1.15. More importantly, it mitigates the risk of catastrophic thermal throttling. The engineering reality is that this system ensures GPUs can reliably sustain their base and boost compute clocks throughout multi-week, high-intensity training runs without heat-induced degradation.

### Cooling Performance Comparison

| Cooling Topology            | Max Rack Power Density | PUE Score   | Thermal Throttling Risk                        |
|-----------------------------|------------------------|-------------|--|
| Standard Air Cooling (CRAC) | ~15 kW                 | 1.50 - 1.80 | High (under sustained heavy load)              |
| <b>ServerMO D2C + RDHx</b>  | <b>50 kW+</b>          | <b>1.15</b> | <b>Negligible (Sustains base/boost clocks)</b> |

## 6. Enterprise Implementation: Generative AI Case Study

*Note: The following metrics represent a highly optimized, specific client deployment scenario after extensive framework tuning. Actual results will vary depending on model architecture, dataset constraints, and network configurations.*

**The Scenario:** A generative AI studio required infrastructure for steady-state, high-resolution video model training. **The Challenge:** The client experienced significant variable egress fees and efficiency drops due to inter-node network jitter on a public cloud. **The Deployment:** ServerMO provisioned a 32-node H100 bare-metal cluster in Dallas, interconnected via 100Gbps RoCE v2. **Observed Outcomes:**

- **Performance:** Training iteration speed saw measurable improvements due to the reduction of hypervisor network jitter in their highly synchronous workload.
- **Cost Efficiency:** By converting variable egress fees to a fixed unmetered uplink, overall infrastructure costs were heavily optimized for their specific continuous-training use case.

- **Hardware Utilization:** Sustained GPU utilization remained high, validating the hybrid cooling topology's ability to maintain hardware stability.

## 7. The Horizon: High-Density Trunking and PCIe Gen 6

As LLM parameter counts scale, network bandwidth requirements grow linearly. ServerMO is provisioning high-density 100Gbps trunking fabrics in primary hubs (Dallas, Frankfurt, Singapore). Pairing this core capacity with upcoming PCIe Gen 6 architectures ensures our infrastructure roadmap aligns with the extreme I/O demands of future AI innovations.

## 8. Conclusion

While virtualized cloud environments offer tools for variable scaling, the architectural requirements shift significantly for the specific domain of sustained AI model training. By focusing on raw bare-metal performance, predictable fixed-cost data movement, and advanced thermal engineering, ServerMO provides an infrastructure foundation designed specifically for the heavy computational realities of enterprise AI.

*© 2026 ServerMO. All rights reserved. | Confidential - Proprietary Information Disclaimer: Technical specifications and synthetic benchmarks are based on internal ServerMO testing environments. For architecture inquiries, contact ServerMO Network Operations.*