

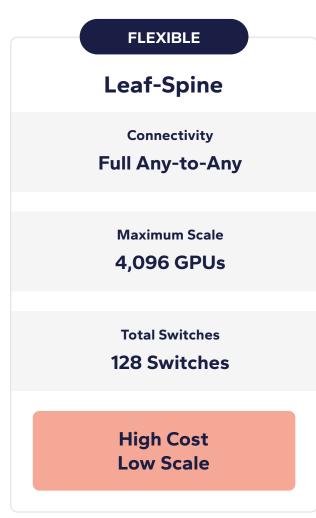
Al Network Topology Analysis

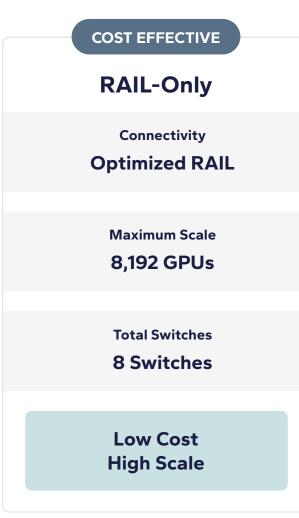
Scaling Considerations for Training & Inference

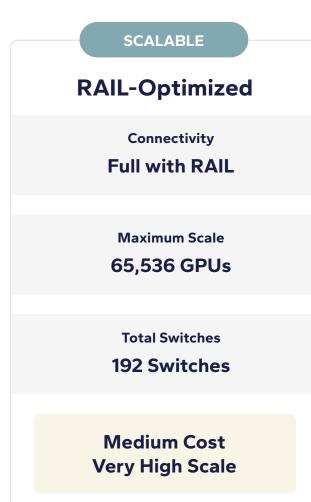
The Infrastructure Challenge

Modern Al workloads demand unprecedented network bandwidth. Traditional CLOS architectures are being pushed beyond their breaking point, requiring purpose-built topologies for Al training and inference.

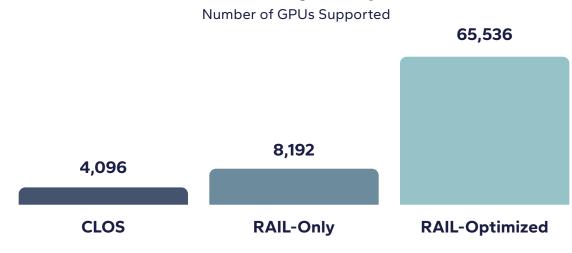
Three Network Topologies Compared







GPU Scalability Comparison



RAIL-Only delivers 2x more servers than CLOS with only 6% of the switches (8 vs 128).

Mixture-of-Experts (MoE) Impact

Massive Performance Gains

Llama 4 MoE models deliver 2-3x better tokens/sec than dense models of similar size, enabling faster inference and training.

Cost Optimization

With MoE models, the 5-10% performance trade-off of RAIL-Only vs RAIL-Optimized becomes economically justified by massive cost savings.

Better Resource Utilization

Sparse activation patterns mean fewer GPUs active per token, reducing power consumption and thermal load.

Dell Open Networking: Cost Advantages

Potential reduction in network infrastructure costs moving from CLOS to RAIL-Only architecture. moving from CLOS to RAIL-Only architecture.

Fewer Switches Required

RAIL-Only uses 8 switches vs 128 for CLOS at comparable scale, dramatically reducing CapEx on switching hardware.

Lower Optics Costs

Optical transceivers account for most network costs. RAIL-Only requires far fewer interconnects between switches.

Reduced Power & Cooling

Fewer active network components translate to lower OpEx for power consumption and data center cooling requirements.

Simplified Managemen

Less complex topology means easier deployment, monitoring, and maintenance, reducing operational overhead.

Decision Framework

Choose RAIL-Only

when building dedicated Al infrastructure primarily for inference, MoE workloads, or when TCO optimization is critical for competitive advantage.

Consider RAIL-Optimized

when requiring maximum flexibility for experimental architectures, diverse legacy models, or extreme scalability for large-scale Al training.

Consider CLOS Networks

when building heterogeneous infrastructure with limited Al workloads or when simplified routing and management are preferred.

Ready to Optimize Your AI Infrastructure?

Purpose-built architectures that precisely align with modern ML workloads define the next generation of Al infrastructure.

For more, see the full report on the Signal65 website.

