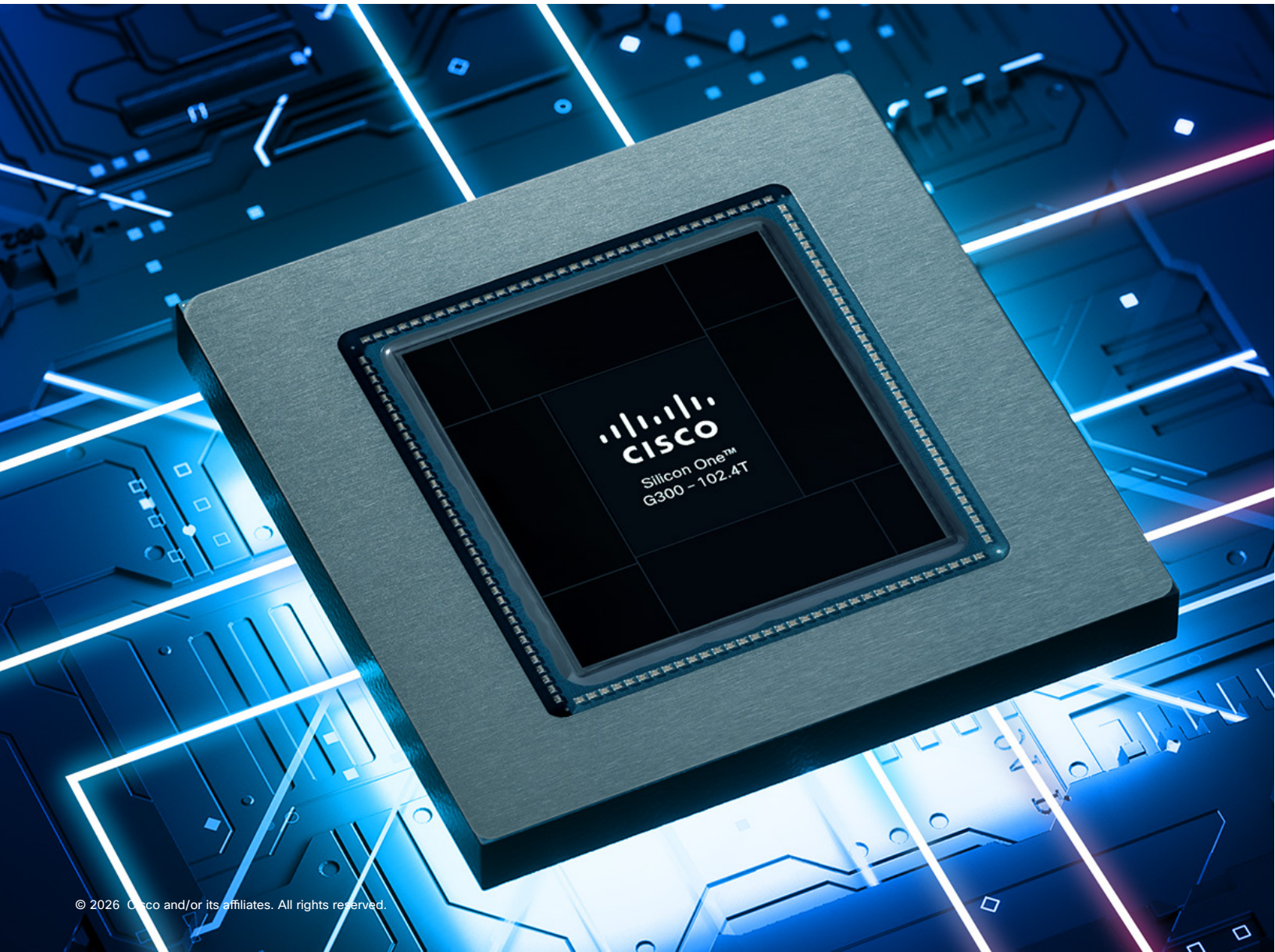


# Cisco Silicon One G300 102.4T Switch with Intelligent Collective Networking Changes the AI Infrastructure TCO Dynamics





Contents

An age-old slogan has a new meaning in the AI era .....3

Addressing the rising complexity of AI models and infrastructure design .....4

Intelligent Collective Networking: advanced load balancing optimized for collective operations .....6

Observations and key takeaways ..... 18

Summary..... 19

## An age-old slogan has a new meaning in the AI era

The slogan “the network is the computer” was coined by John Gage of Sun Microsystems in 1984—more than four decades ago. It has achieved new meaning in the era of AI infrastructure and is, very rightfully, used frequently by AI networking silicon and system vendors.

Graphics Processing Units (GPUs) represent the most significant capital investment in modern AI infrastructure. A single NVIDIA Blackwell chip ranges from \$30,000 to \$40,000—nearly triple the \$12,000 cost of a flagship Intel® Xeon® CPU. This premium extends to the cloud, where annual Blackwell rentals reach approximately \$22,000 per chip. Because of these staggering costs, infrastructure efficiency has become a top priority for data center operators.

AI workloads rely on collective operations that execute across hundreds of thousands of GPUs in a cluster. By leveraging a high-performance fabric of Network Interface Cards (NICs) and switches, the network interconnects these individual chips to form a unified collective GPU computer.

Because AI traffic is uniquely synchronized and bursty, the network resides on the critical path of performance. Network issues cause collective operations to stall, wasting valuable GPU cycles. Consequently, the performance and efficiency of the collective GPU computer depend entirely on the network’s ability to minimize Collective Completion Time (CCT) and reduce GPU stall time.

CCT is the duration required for a collective communication operation to finish across distributed nodes. As a measure of the speed of data synchronization between GPUs, CCT serves as the critical metric for benchmarking AI infrastructure performance.

In the modern AI era, the pressing need to manage rising Total Cost of Ownership (TCO) serves as a contemporary reinterpretation of John Gage’s famous slogan “The network is the computer.”

### Key takeaways

As AI clusters scale toward million-GPU capacities, the network is the critical path for performance. The Cisco Silicon One G300 102.4T switch, with **Intelligent Collective Networking**, eliminates GPU stalls to transform infrastructure TCO.

Key benchmarks include:

- **Accelerated Training:** Delivers an **82% reduction** in Job Completion Time (JCT) over standard Ethernet and a **28% speedup** over random packet spraying.
- **Near-Ideal Performance:** Achieves Collective Completion Times (CCT) within **2% to 4% of the theoretical ideal**, even during link failures.
- **Hardware Resilience:** Features **up to 2.5X increased burst absorption** and fault detection that is **100,000x faster** than traditional methods.
- **Maximizing ROI:** Improved network efficiency allows for reclaiming up to 28% of cluster capacity by minimizing GPU idle time, ensuring expensive resources focus on computation rather than waiting for data.

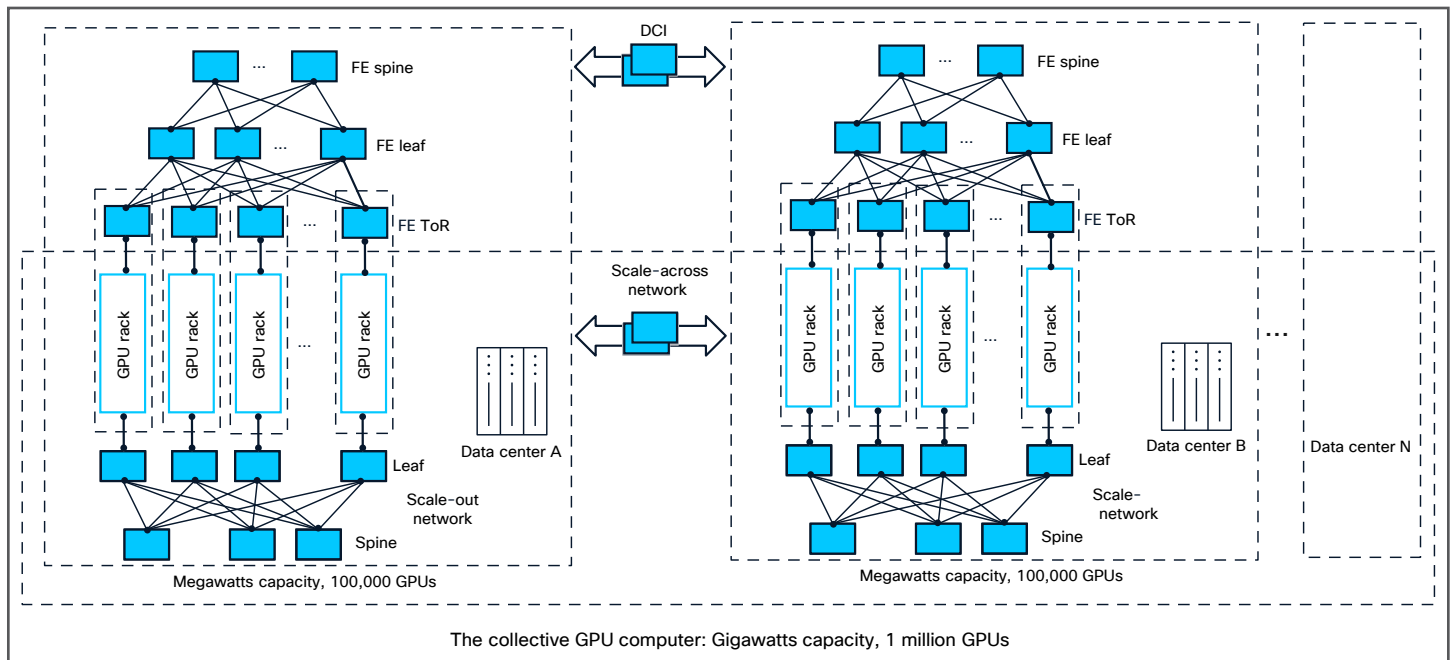


Figure 1. Advent of the gigawatt collective GPU computer and ballooning AI infrastructure costs

The Cisco Silicon One™ G300 102.4T switch is engineered to meet the scaling demands and structural shifts of modern AI infrastructure driven by rapid model innovation. By leveraging the new Intelligent Collective Networking—which features intelligent, advanced load balancing—the G300 significantly lowers CCT, achieving near-optimal performance across various scaling scenarios. These CCT benchmarks serve as a vital metric for evaluating TCO impact.

## Addressing the rising complexity of AI models and infrastructure design

The fundamental goal of AI networking is to maximize the utility of expensive GPU resources by minimizing the duration of collective operations and preventing GPU stalls. While the industry has introduced various proprietary and standard-based (Ultra Ethernet Consortium [UEC]) techniques for load balancing and congestion management, traditional network-centric approaches—such as static Equal-Cost Multipath (ECMP) and packet spraying—are no longer sufficient. Despite incremental improvements, these legacy

methods fail to reach the performance thresholds required to transform the Total Cost of Ownership (TCO), largely because AI workloads are evolving faster than traditional networking constructs can adapt. This paper explores the features of the Cisco Silicon One G300 device in detail after first examining why AI infrastructure is changing so rapidly.

A few of the reasons for the rapid evolution of the AI infrastructure are outlined below.

## Power and cooling constraints in individual data centers

Power and cooling constraints have forced GPU clusters to distribute across sites, creating a hybrid of local scale-out and long-distance scale-across traffic. Managing this requires 102.4T switching for local density in scale-out networks and deep-buffered silicon to handle the high latency and long round-trip times of long-haul links in scale-across traffic. Because these environments demand different load-balancing strategies—dynamic for local and static for long haul—modern networking silicon must be “traffic aware” to support both scenarios within a single architecture.

## Efficient 500K to 1M+ GPU scale

The rise of agentic workflows requires infrastructure that supports larger GPU clusters with minimal communication overhead to ensure the real-time responsiveness the autonomous agents demand. To achieve this, operators are moving toward two-tier topologies using high-radix links (100 Gbps to 400 Gbps). While standard 51.2T architectures are typically limited to 8,192 GPUs, multiplanar designs—networks constructed by interconnecting multiple switching planes to increase overall capacity—leveraging 512-radix, 51.2T switches can scale to 512,000 GPUs. Although this multiplanar approach improves scale and reliability, it significantly increases infrastructure costs due to several factors:

- Multiple independent network planes require separate switches and interconnections, multiplying hardware costs.
- Managing complex, large-scale AI clusters with multiple planes adds operational overhead.
- High-performance demands for low latency and high bandwidth necessitate expensive, advanced switches and architectures.
- Redundancy for high availability duplicates infrastructure, increasing costs.

- Scaling beyond single networks involves more devices and links through interconnected planes, further raising expenses.

Furthermore, the adoption of the Multi-Path Reliable Connection (MRC) protocol highlights that current load-balancing solutions are insufficient to prevent GPU stalls and optimize utilization at this massive scale. This situation necessitates advanced dynamic load balancing and congestion management strategies specifically designed to handle the complexity of multiplanar designs.

## Demand for significantly better inference token economics

Modern AI inference architecture often disaggregates the prefill and decode phases across specialized GPUs. However, since the server NICs must simultaneously support both GPU classes and their distinct processing phases, network traffic has shifted from homogeneous flows to a highly complex mixed pattern. To address these evolving requirements, NVIDIA introduced the **Rubin CPX** GPU. Consequently, the **Vera-Rubin VR200 NVL144** rack designs have been updated to include the **VR200 NVL CPX** variant, creating significant pressure to shift toward complex, high-performance, and scalable network architectures that emphasize optimized inter-GPU communication, advanced load balancing, and robust network fabrics with resilient traffic distribution.

To support the architectural transitions and massive scale previously described, modern load-balancing and congestion management strategies must be engineered to address the following critical scenarios:

- **Dynamic load balancing:** Implementation of advanced packet spraying techniques to optimize traffic distribution.
- **Resilient traffic distribution:** Intelligent packet spraying during link failures to mitigate traffic imbalances and prevent “blackholing.”
- **Mixed-mode traffic management:** Seamlessly handling the coexistence of ECMP subflows and packet spraying during link failure events.

## Intelligent Collective Networking: advanced load balancing optimized for collective operations

While several 102.4-Tbps switching solutions are available, the Cisco Silicon One G300 distinguishes itself by delivering **Intelligent Collective Networking** for performance and reliability at massive scale and **future-ready** silicon built to evolve with AI cluster requirements for an industry-leading scale-out architecture tailored to the rigorous demands of modern AI infrastructure.

**Intelligent Collective Networking** consists of a comprehensive suite of intelligent load-balancing, congestion management, and telemetry tools. By reducing CCT to near-ideal levels, this technology significantly enhances infrastructure TCO, providing a decisive advantage for high-performance AI workloads.

Moving beyond traditional static methods (ECMP) and standard dynamic load balancing (packet spraying), Intelligent Collective Networking represents a fundamental architectural shift. Engineered with an “AI-first” mindset, it prioritizes the optimization of collective operations rather than incremental improvements to legacy flow distribution. Key technical advantages include:

- **Superior burst absorption:** The Cisco Silicon One G300 features a fully shared memory architecture that provides up to 2.5X increased burst absorption than alternative solutions, providing the necessary capacity for extended-reach lossless flow control. This significantly increased buffer depth enhances the network’s ability to absorb massive GPU traffic bursts and provides superior resilience against in-cast events.
- **Accelerated fabric-level fault detection:** Hardware-accelerated, proactive dissemination of reachability and congestion signals enables microsecond-fast fault isolation, resulting in network fault detection that is 100,000 times faster than traditional mechanisms. The Silicon One G300 intelligently combines local and remote information to drive real-time decisions for resilient distribution of network traffic.
- **Topology-aware load balancing:** A framework optimized for collective operations integrates local congestion and reachability data with a global fabric-level state. This allows for intelligent, fabric-wide, congestion-aware routing decisions directly at the switch ingress ports.

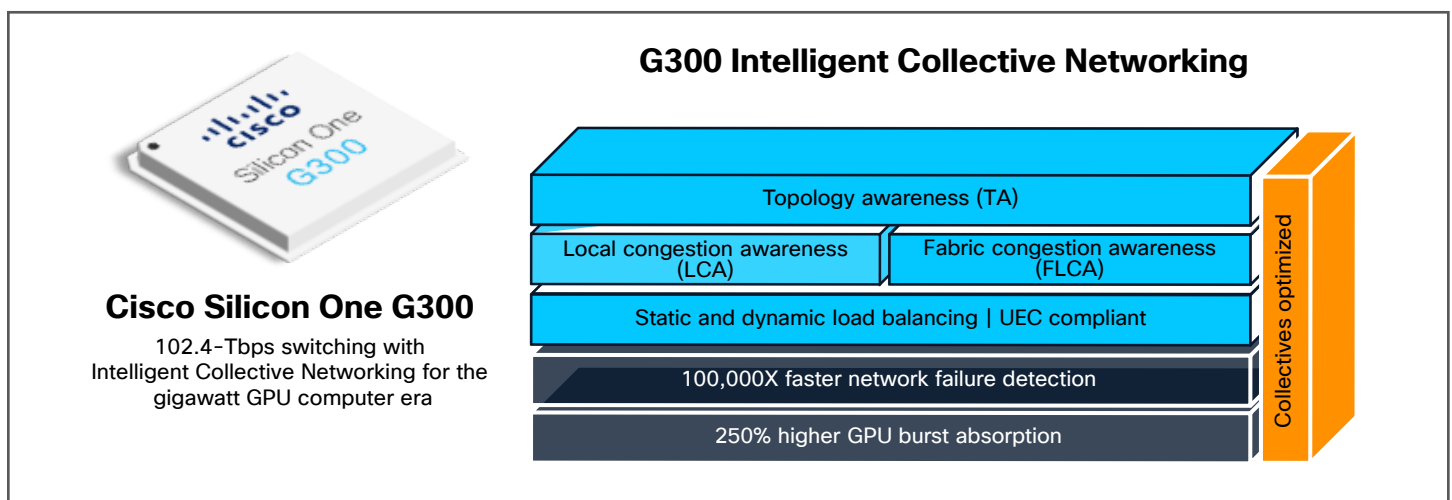


Figure 2. G300 Intelligent Collective Networking –unique, fabric-aware load balancing and congestion management

The integration of these core innovations empowers the Cisco Silicon One G300 with unique, differentiated capabilities specifically architected to address the evolving structural shifts and massive scaling requirements of modern AI infrastructure:

## Glossary

- **Intelligent Collective Networking - Local Congestion Awareness (LCA)**
- **Intelligent Collective Networking - Fabric-Level Congestion Awareness (FLCA)**
- **Intelligent Collective Networking - Topology Awareness (TA)**

## CCT benchmarking

To evaluate these Intelligent Collective Networking capabilities, we conducted simulation micro-benchmarks on the G300 platform under varying topology, failure, and traffic-mix conditions across two large-scale, two-tier **Clos fabrics** consisting of 8,192 and 16,384 GPUs. Our results demonstrate a significant reduction in CCT when comparing traditional static and dynamic load-balancing techniques against the G300's innovative routing capabilities, even as AI infrastructure scales in size and complexity.

## Methodology

**Collective operations:** We used Ring-all-Reduce Collective (RARC) and All-To-All Collective (ATAC) with 16 ranks per collective.

- Each rank transmits and receives a total of 134.4 MB for the RARC operation and 128 MB for ATAC.
- To avoid unrealistic synchronization between collectives and within collectives in the simulation, we introduced random start time jitter between collectives in the range of 0 to 4 microseconds.
- We also introduced random start time jitter between flows within a collective in the range of 0 to 2 microseconds.
- In mixed-mode experiments, 12.5% of the traffic was statically load balanced and 87.5% used packet spray. Static and dynamic collectives used different Traffic Classes (TCs) mapped to different Output Queues (OQs) with 1:1 Weighted Fair Queuing (WFQ) scheduling.

**Workloads:** Workloads were represented through RARC and ATAC collective primitives. We allowed mixed load-balancing modes: dynamically load-balanced traffic eligible for packet spray and statically load-balanced traffic forwarded without spraying to preserve ordering constraints. The two traffic types were separated by TC and mapped to distinct OQs, with an explicit scheduling policy between OQs (e.g., 1:1 WFQ).

**AI cluster configurations:** Endpoints were modeled as NIC-attached GPU servers connected to leaf switches. We focused on a two-tier leaf-spine (Clos) topology, which may serve as a subfabric within larger deployments. The ingress leaf was assumed to be the primary decision point for path selection. Fabric-level congestion was represented as aggregated path-level signals rather than per-flow state.

- 8,192 GPUs, each with a 400-Gbps NIC link. The 8,192-GPU cluster was connected using 64 G300 leaf switches and 32 G300 spine switches; two 800-Gbps links were used per leaf switch-spine switch pair.
- 16,384 GPUs, each with a 400-Gbps NIC link. Each server had four GPUs and four NICs in total. There was a total of 16 servers or 64 GPUs per rack and 256 racks. The 16,384 GPU cluster was connected using 128 G300 leaf switches and 64 G300 spine switches; a single 800-Gbps link was used per leaf switch-spine switch pair.

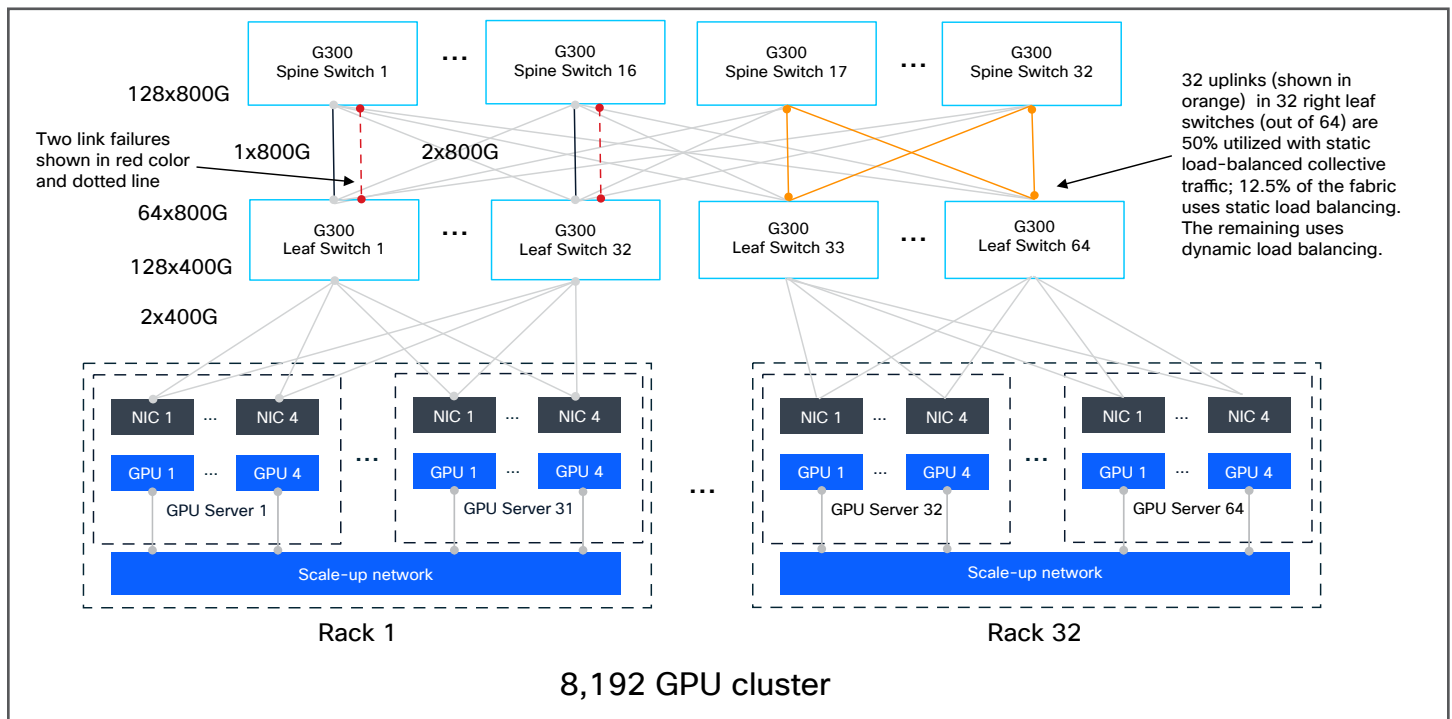


Figure 3. Example of a simulated 8,192 GPU cluster configuration used in the benchmarking scenarios

**Failure and propagation delay models:** In AI clusters, link failures and flaps are not just occasional nuisances; they are statistical certainties. As infrastructure scales from hundreds of GPUs to clusters of 50,000 or 100,000+, the sheer number of physical components makes “perfect” uptime impossible, and the Mean Time Between Failures (MTBF) for the cluster as a whole becomes very short.

We modeled local and remote leaf-spine link failures, including scenarios with no failures, two failures, and larger failure counts of overall 16 link failures. Baseline cable lengths were short (1 m), and additional scenarios extended to longer cables (e.g., 250 m leaf-spine, 50 m leaf-NIC) to study sensitivity to cable propagation delay.

## Benchmarking metrics

The goal of the benchmarking data presented in this paper is to show improvements in collectives completion time, or CCT. The CCT measurements

achieved using G300 Intelligent Collective Networking features are:

- **Tail CCT [microseconds]:** Maximum CCT across all collectives
- **Ideal CCT [microseconds]:** Theoretical minimum CCT with perfect fabric utilization, load balancing and no link failures
- **Overhead vs. Ideal:** (Tail CCT – Ideal CCT) / Ideal CCT

## Benchmarking scenarios and results

The scenarios covered below address the scale and structural changes modern AI infrastructure is undergoing. They are presented as six increasingly complex experiments or benchmarking scenarios that apply the capabilities of Intelligent Collective Networking:

- Local congestion awareness or LCA
- Fabric-level congestion awareness or FLCA
- Topology awareness or TA

## Packet spray vs. Static ECMP

**Scenario 1:** Simulated cluster size: 8,192 GPUs | Collective: 512 RARC (16 ranks per collective) | **Ideal CCT: 2767.2 microseconds.**

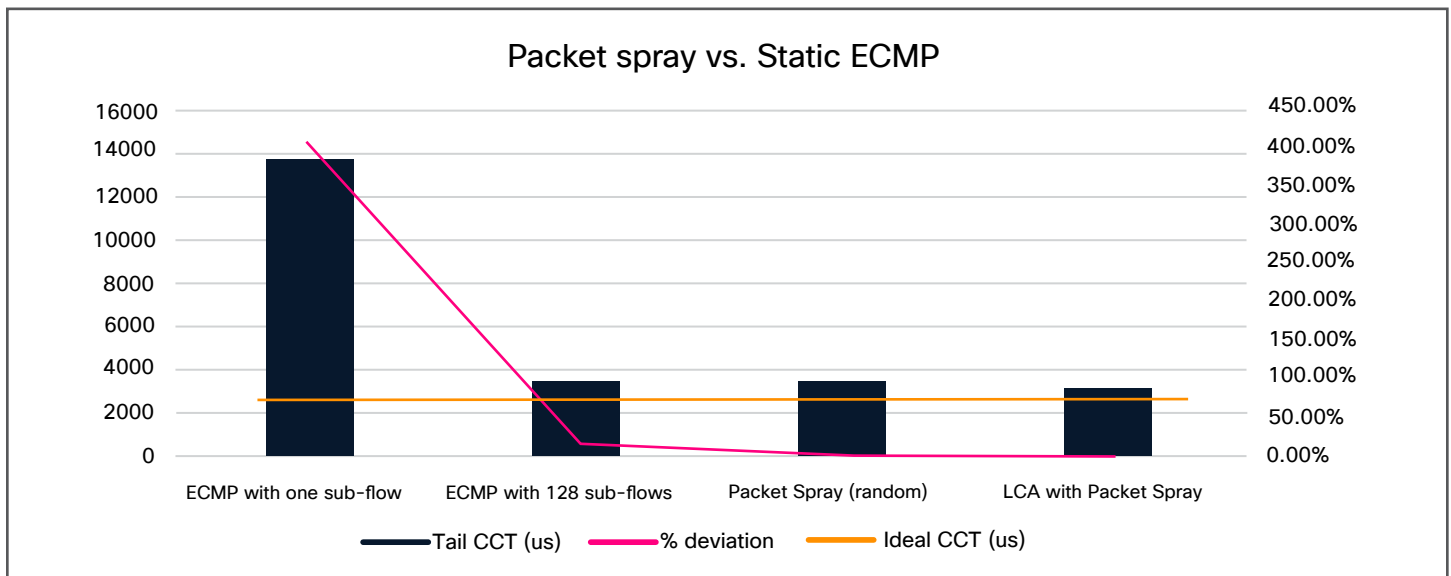


Figure 4. Scenario 1: Packet spray vs. Static ECMP

**Scenario 1** simply established that static load balancing using ECMP is not a viable load-balancing option. Packet spray-based dynamic load-balancing methods were therefore compared in all subsequent experiments. With Cluster Optimized Routing Intelligent Collective Networking Local Congestion-Aware (LCA) load balancing, 46 microseconds of CCT savings and 43% lower CCT were achieved; the deviation from the ideal CCT was only 2%.

## Random packet spray vs. Intelligent Collective Networking

**Scenario 2:** Simulated cluster size: 8,192 GPUs | Collective: 512 RARC (16 ranks per collective) | **Ideal CCT: 2767.2 microseconds** | Two concurrent link failures

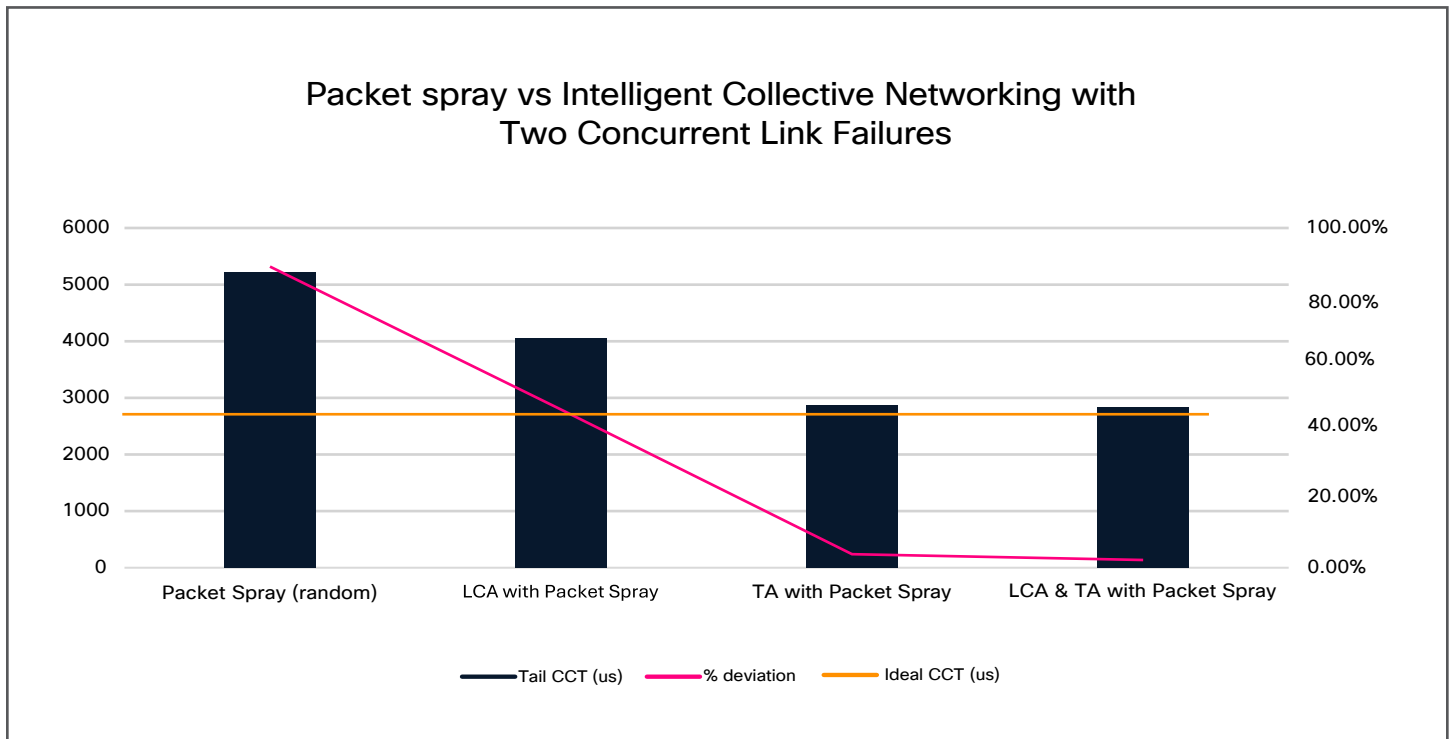


Figure 5. Scenario 2: Packet spray vs. Intelligent Collective Networking with two concurrent link failures

As depicted in the figure above, enabling topology awareness and comparing the results with the local and random packet spray methods provided near-ideal CCT results (within 2.27% for local congestion awareness or 3.98% for random).

## Random packet spray vs. Intelligent Collective Networking (mixed traffic)

**Scenario 3:** Simulated cluster size: 8,192 GPUs | Collective: 512 RARC (16 ranks per collective) | **Ideal CCT: 2784 microseconds** | Mixed traffic: 12.5% static and 87.5% packet spray

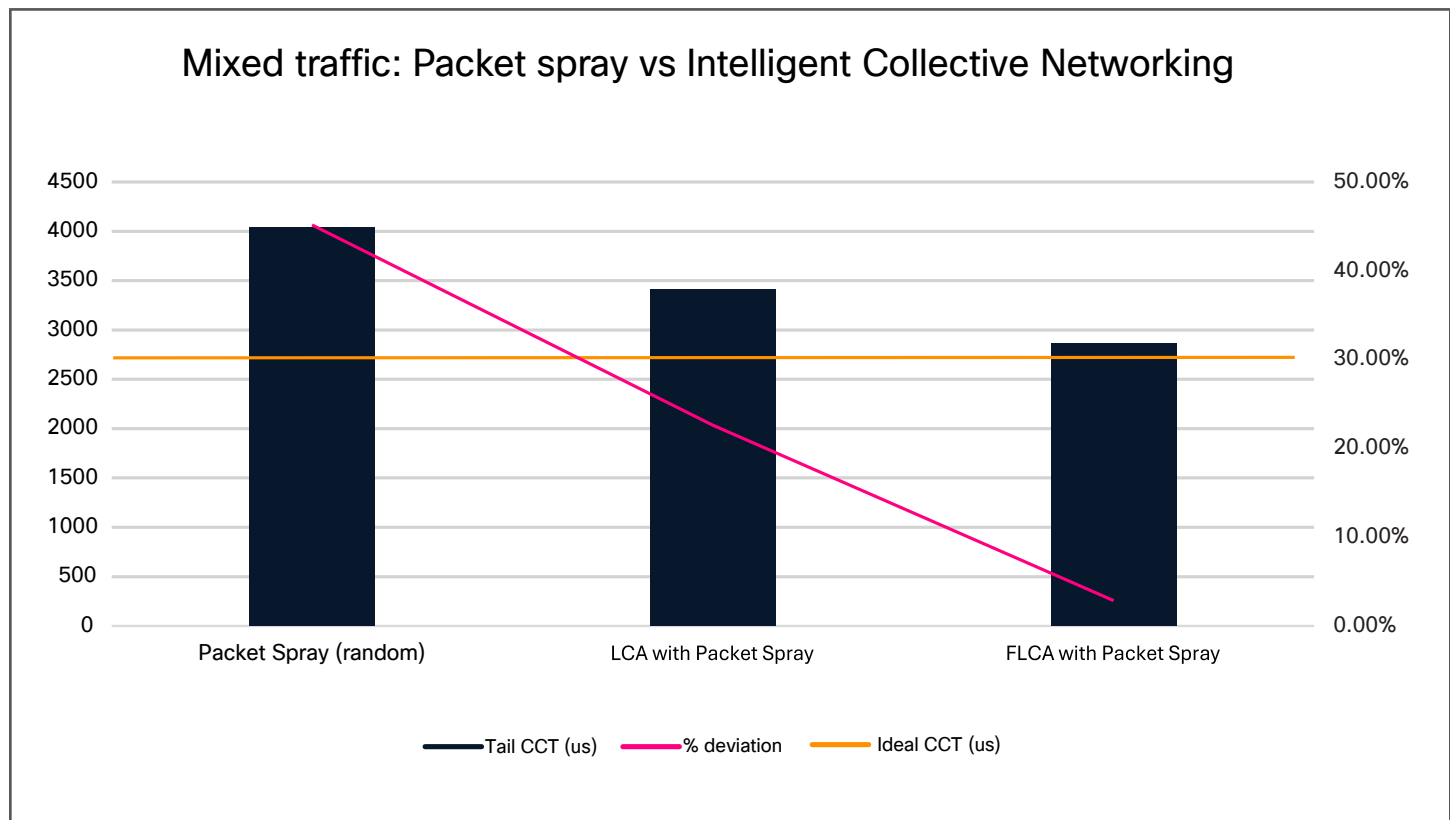


Figure 6. Scenario 3: Mixed traffic: Packet spray vs. Intelligent Collective Networking

In an 8,192 × 400-Gbps topology with a 12.5% static ECMP and 87.5% packet spray traffic mix, fabric-level congestion management achieved near-ideal behavior for sprayed traffic and provided substantial improvements over both random spraying and local congestion-aware packet spraying (within 2.86% of ideal CCT).

## Random packet spray vs. Intelligent Collective Networking (mixed traffic) with link failures

**Scenario 4:** Simulated cluster size: 8,192 GPUs | Collective: 512 RARC (16 ranks per collective) | **Ideal CCT: 2784 microseconds** | Mixed traffic: 12.5% static and 87.5% packet spray | Two concurrent link failures

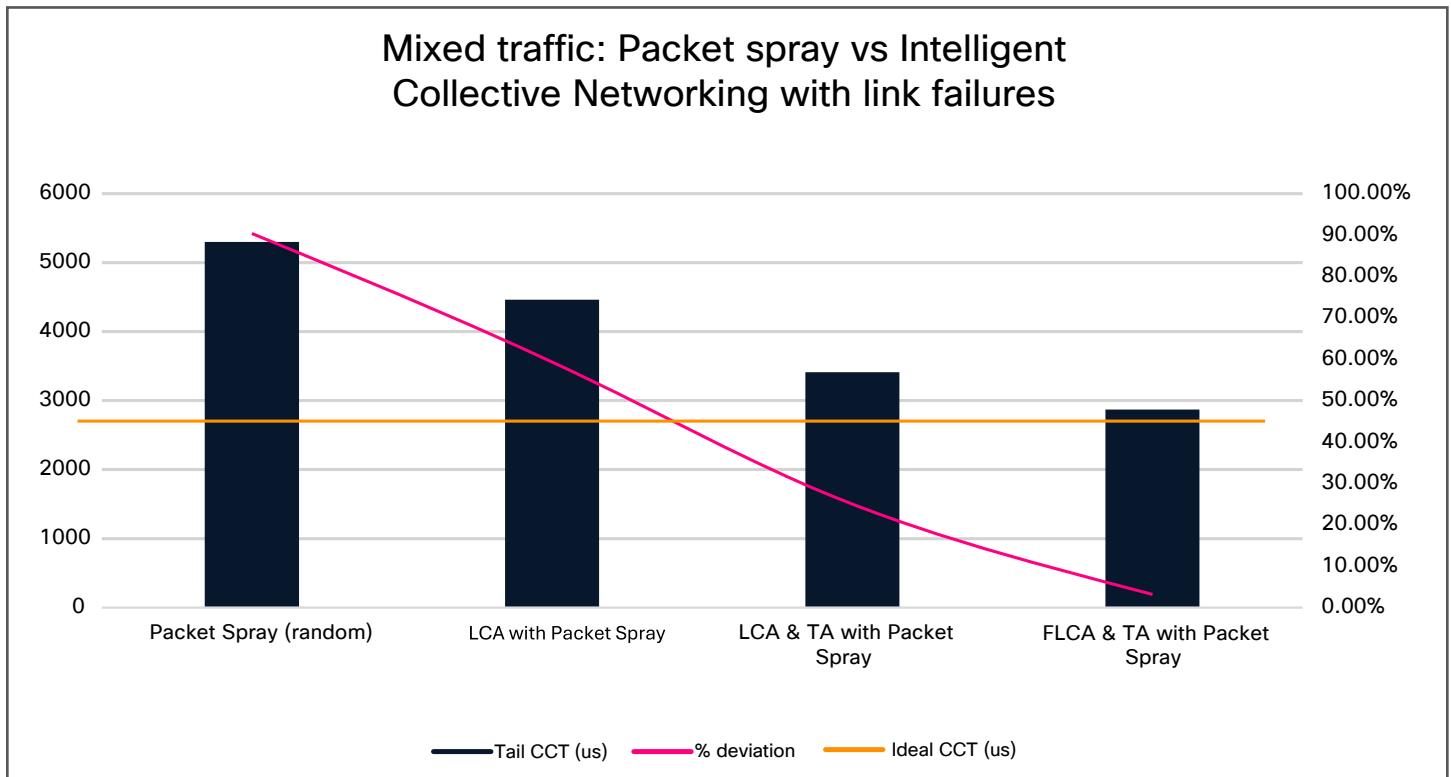


Figure 7. Scenario 4: Mixed traffic–Packet spray vs. Intelligent Collective Networking with two concurrent link failures

In an 8,192 × 400-Gbps topology with a 12.5% static ECMP and 87.5% packet spray traffic mix, in the presence of two link failures, the topology-aware fabric-level congestion management achieved near-ideal behavior for sprayed traffic (within 3.18% of ideal CCT).

## Random packet spray vs. Intelligent Collective Networking (mixed traffic), realistic deployment with two link failures

**Scenario 5:** Simulated cluster size: 8,192 GPUs | Collective: 512 RARC (16 ranks per collective) | **Ideal CCT: 2858.5 microseconds** | Mixed traffic: 12.5% static and 87.5% packet spray | Two 250 m links per leaf-spine pair | 50 m links between NICs and leaf switches | Two concurrent link failures

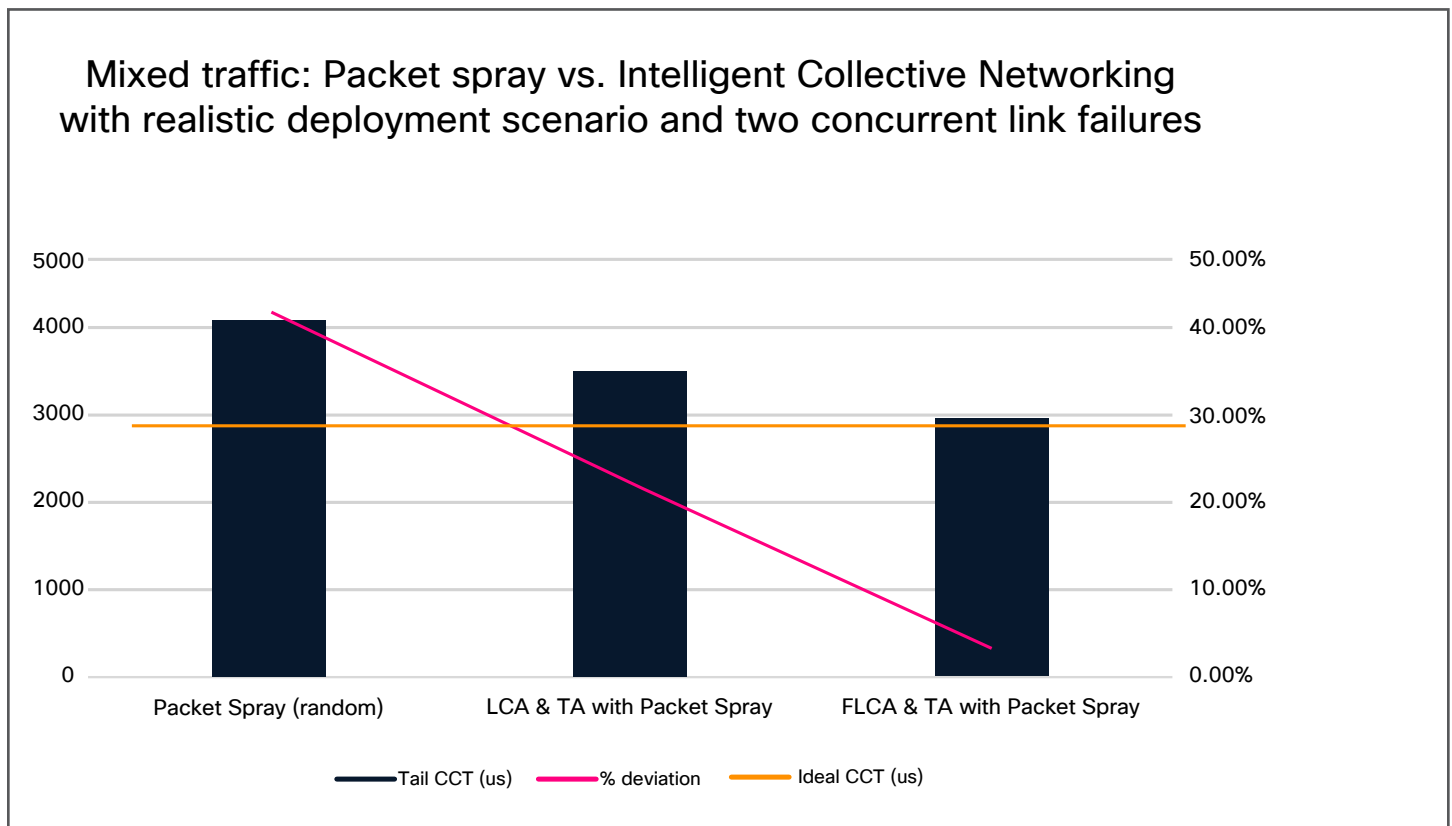


Figure 8. Scenario 5: Mixed traffic–Packet spray vs. Intelligent Collective Networking with realistic deployment scenario and two concurrent link failures

As depicted above, in an 8,192 × 400-Gbps realistic deployment scenario (two 250 m links per leaf-spine pair and 50 m links between NICs and leaf switches) topology with a 12.5% static ECMP and 87.5% packet spray traffic mix, in the presence of two link failures, the topology-aware fabric-level congestion management achieved near-ideal behavior for sprayed traffic (within 3.66% of ideal CCT).

## Random packet spray vs. Intelligent Collective Networking (mixed traffic), realistic deployment with 16 link failures

**Scenario 6:** Simulated cluster size: 16,384 GPUs | Collective: 512 RARC (16 ranks per collective) | **Ideal CCT: 2858.5 microseconds** | Mixed traffic: 12.5% static and 87.5% packet spray | Two 250 m links per leaf-spine pair | 50 m links between NICs and leaf switches | 16 concurrent link failures

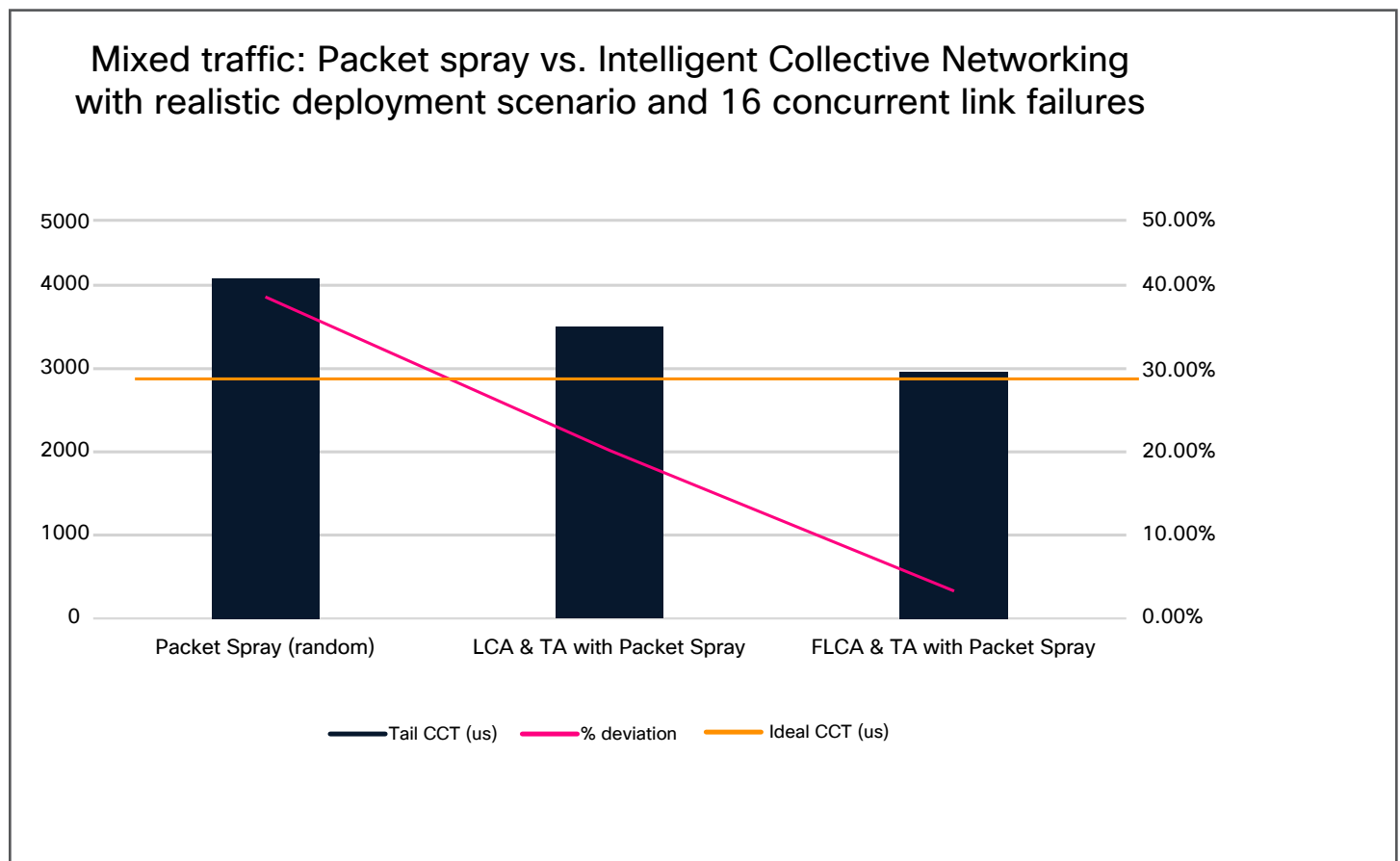


Figure 9. Scenario 6: Mixed traffic–Packet spray vs. Intelligent Collective Networking with realistic deployment scenario and 16 concurrent link failures

With 16 link failures in an 8,192 × 400-Gbps realistic deployment scenario (two 250 m links per leaf-spine pair and 50 m links between NICs and leaf switches) topology with a 12.5% static ECMP and 87.5% packet spray traffic mix, the topology-aware fabric-level congestion management achieved near-ideal behavior for sprayed traffic (within 4.54% of ideal CCT).

## Predictive modeling of JCT improvements in distributed training

While the packet-level simulations demonstrate significant reductions in CCT using the advanced features of the Cisco Silicon One G300, quantifying the value for AI practitioners requires translating these network gains into end-to-end Job Completion Time (JCT). This section presents an analytical framework linking fabric performance to application speedup, instantiated using empirical profiles from a Llama 3.1-8B training workload.

### Methodology and analytical model

To bridge the gap between network simulation and application performance, we modeled the total JCT as the aggregate of **N** training iterations. We focused on a ZeRO-2 based Fully Sharded Data Parallel (FSDP) configuration.

To capture the impact of network performance on training time, we defined the duration of a single iteration ( $T_{\text{Iteration}}$ ) as follows:

$$T'_{\text{Iteration}} = T_{\text{Overhead}} + T_{\text{Compute}} + \sum_i \left| \frac{T_{\text{Comm},i}}{\alpha_i} - \phi_i \right|$$

Using this profile, we calculated the predicted JCT by applying congestion overhead factors derived from packet simulations under both healthy network conditions and degraded states involving two concurrent fabric link failures. The technical derivation for the ECMP (without failure) case is presented below:

$$\alpha_{\text{ECMP RARC}} = \frac{\text{CCT}_{\text{Ideal RARC}}}{\text{CCT}_{\text{ECMP RARC}}} = \frac{2,767.2}{13,821.2} = 0.2002$$

$$\text{JCT} = 5376(N) * \left( 61.86(T_{\text{Overhead}}) + 496.82(T_{\text{Compute}}) + \left( \frac{458.79(T_{\text{Comm}})}{0.2002(\alpha_{\text{ecmp}})} \right) \right) = 255.37 \text{ minute}$$

### Where:

- $i$ : Index of specific collective communication operation operations
- $\phi_i$ : Overlap efficiency potential specific to collective communication  $i$  [ $0$  = no overlap,  $\min(T_{\text{Compute}}, T_{\text{Comm},i})$  = perfect overlap].
- $\alpha_i$  (alpha): Network improvement factor derived from G300 simulations.

In this specific study, we focused on a regime where communication is predominantly exposed ( $\phi_i \approx 0$ ), i.e the observed overlap is close to zero, allowing for direct observation of fabric improvements.

The baseline workload profile for Llama 3.1-8B was empirically measured as follows:

- Compute time ( $T_{\text{Compute}}$ ): 496.82 ms
- Communication time ( $T_{\text{Comm}}$ ): 458.79 ms
- Overhead ( $T_{\text{Overhead}}$ ): 61.86 ms
- Training steps (**N**): 5,376
- Baseline JCT = 91.16 minutes

## Predicted job completion time

Normal traffic conditions | Cluster size: 8,192 GPUs | Collective: 512 RARC (16 ranks per collective)

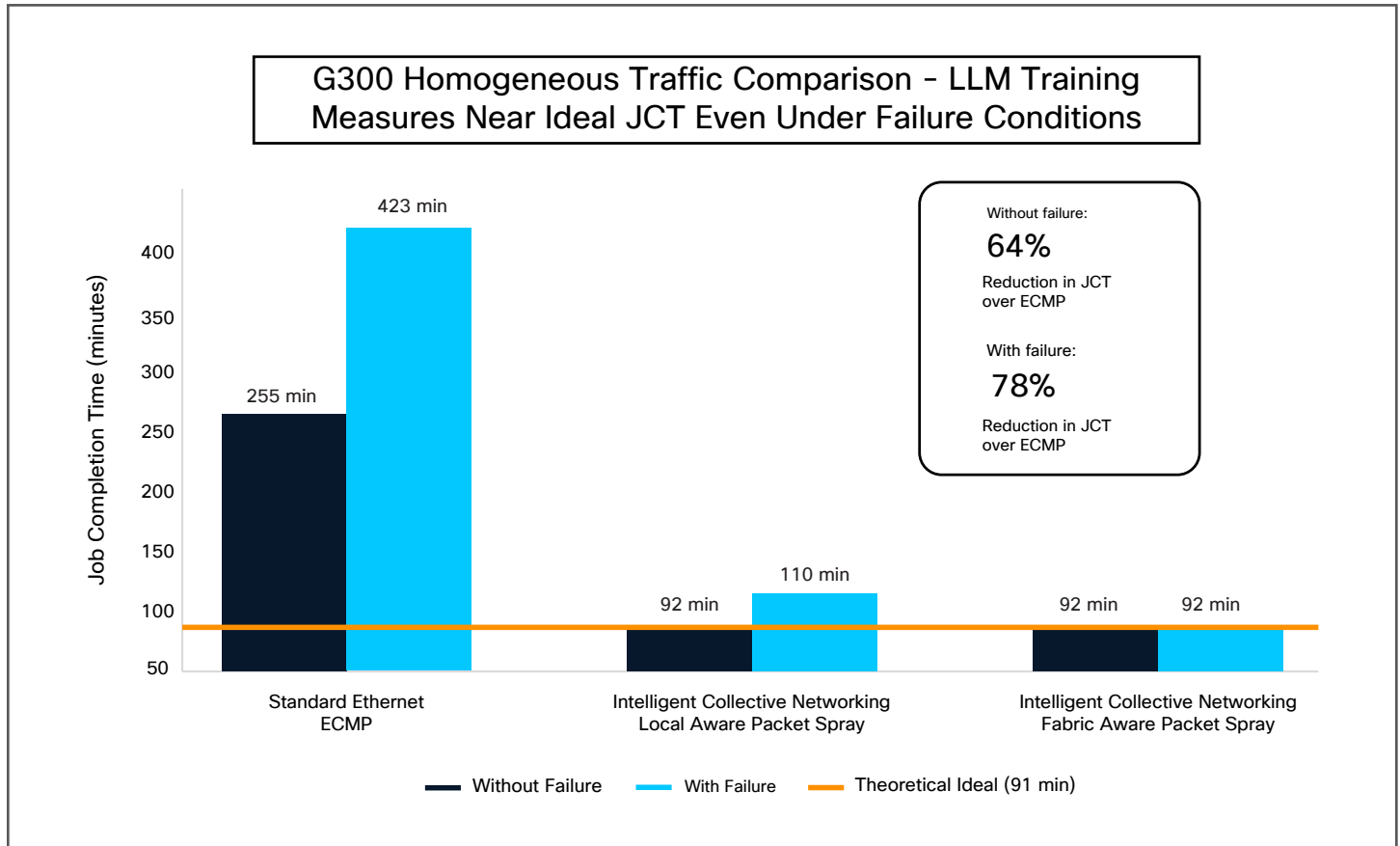


Figure 10. Homogeneous traffic comparison

As illustrated in the figure above, in a normal traffic scenario, Standard Ethernet (ECMP) struggled significantly, resulting in a JCT of 255.37 minutes in a healthy network and degrading further to 499.67 minutes with two link failures. **ECMP failure CCT is derived by extrapolating the packet spray (random) failure case using the ratio of packet spray (random) and ECMP under a nonfailure case.**

In contrast, the G300's Intelligent Collective Networking with fabric-aware packet spray dramatically stabilized performance. It maintained a JCT of approximately 92.1 minutes regardless of network health (92.14 min without failure, 92.18 min with two link failures). This represents an improvement of over 64% in healthy conditions and 82% in failure conditions compared to Standard Ethernet.

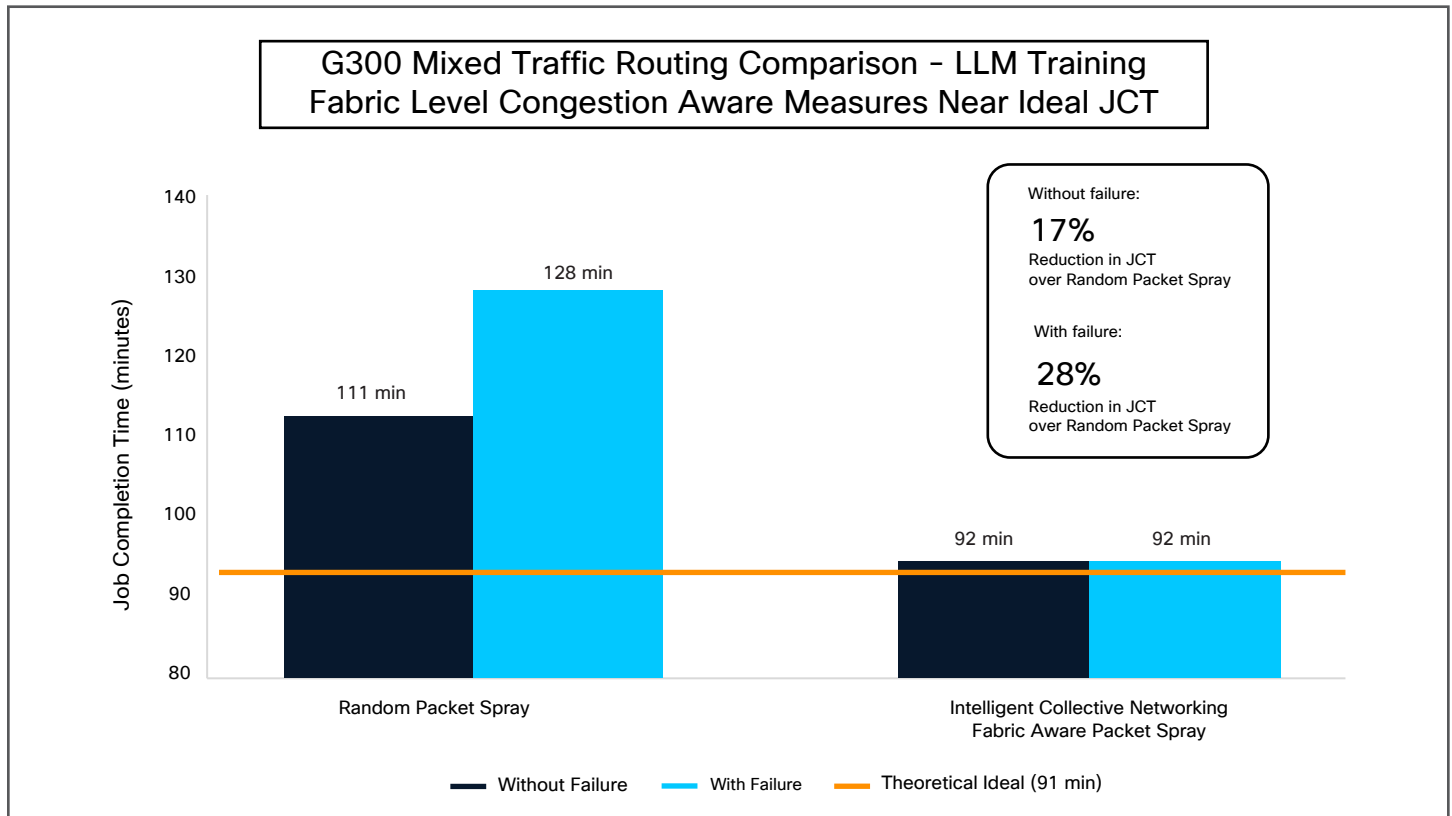
**Mixed traffic conditions | Cluster size: 8,192 GPUs | Collective: 512 RARC (16 ranks per collective)**


Figure 11. Mixed traffic routing comparison

In mixed traffic scenarios, we compared random packet spray against the G300's advanced congestion-aware mechanisms.

The figure above demonstrates that random packet spray experienced notable overhead in this environment, resulting in a JCT of 110.85 minutes. When subjected to two link failures, JCT degraded further to 128.54 minutes.

The G300's fabric-level Congestion-Aware (CA) packet spray effectively mitigates this contention. It reduced JCT to 92.90 minutes in healthy conditions. Crucially, even with two link failures, the G300 maintained a JCT of 92.49 minutes, delivering a 28% speedup over the degraded baseline.

This analysis confirms that fabric-level innovations translate directly to tangible Return On Investment (ROI) for AI/ML clusters. The G300's Intelligent Collective Networking provides predictable performance

consistency, maintaining optimal JCTs even during fabric link failures where standard schemes experience significant spikes.

**Note:** The JCT predictions presented in this section are empirically derived based on observations during profiling of a specific training regime (Llama 3.1-8B using ZeRO-2 FSDP). While these results highlight significant potential gains, actual performance in production environments may vary depending on model architecture, parallelism strategies, and cluster scale. The primary objective of this case study is to validate the analytical methodology and demonstrate how fabric-level performance metrics can be translated into application-level impact.

## Observations and key takeaways

### CCT benchmark results

- The experiment comparing packet spray and static ECMP demonstrated that static ECMP-based load balancing is not an effective approach. Dynamic, packet spray-based techniques—powered by cluster optimized local congestion-aware load balancing—achieved a reduction in CCT of 46 microseconds, representing a 43% decrease, and deviated from the ideal CCT by only 2%.
- In more complex scenarios, the comparison between random packet spray and Intelligent Collective Networking (across mixed traffic, link failures, and realistic deployments) showed that **advanced packet spraying** and **resilient traffic distribution**—enabled by fabric-level congestion-aware and topology-aware load balancing—can deliver up to 2389 microseconds of CCT savings and reduce CCT by as much as 87%, with only a 4% deviation from the ideal CCT.
- Intelligent Collective Networking consistently demonstrated **near-ideal CCT** and effectively managed mixed-mode traffic (simultaneous ECMP subflows and packet spraying) in both the presence and absence of link failures.
- Topology-aware and fabric-level load balancing, combined with efficient congestion management, enhance network utilization and minimize micro-congestion. Lowering the CCT reduces the overall synchronization and data exchange time among GPUs during AI training or fine-tuning, which in turn accelerates job completion time, as supported by predictive modeling and summarized below.

### Predicted JCT findings

- **Massive efficiency gains over Standard Ethernet:** The G300 effectively eliminates the performance bottlenecks of traditional ECMP. In healthy networks, it delivers a 64% reduction in JCT. In failure scenarios, this advantage widens to a 78% reduction, proving that Standard Ethernet is nonviable for large-scale AI clusters.
- **Superior resilience vs. modern load balancing:** Even when compared to advanced random packet spray schemes, the G300 demonstrates superior stability. In complex failure scenarios (Scenarios 3 and 4), the G300 delivers a 17% to 28% speedup. This helps ensure that training timelines remain predictable and consistent, regardless of network degradation.
- **Maximizing GPU utilization and ROI:** By minimizing network-induced idle time, the G300 significantly boosts GPU utilization, helping ensure that expensive compute resources are spent on training rather than waiting for data. These efficiency gains reclaim up to 28% of cluster capacity compared to packet spray and more than 70% compared to ECMP-based load balancing, lowering the TCO and accelerating model time to market.

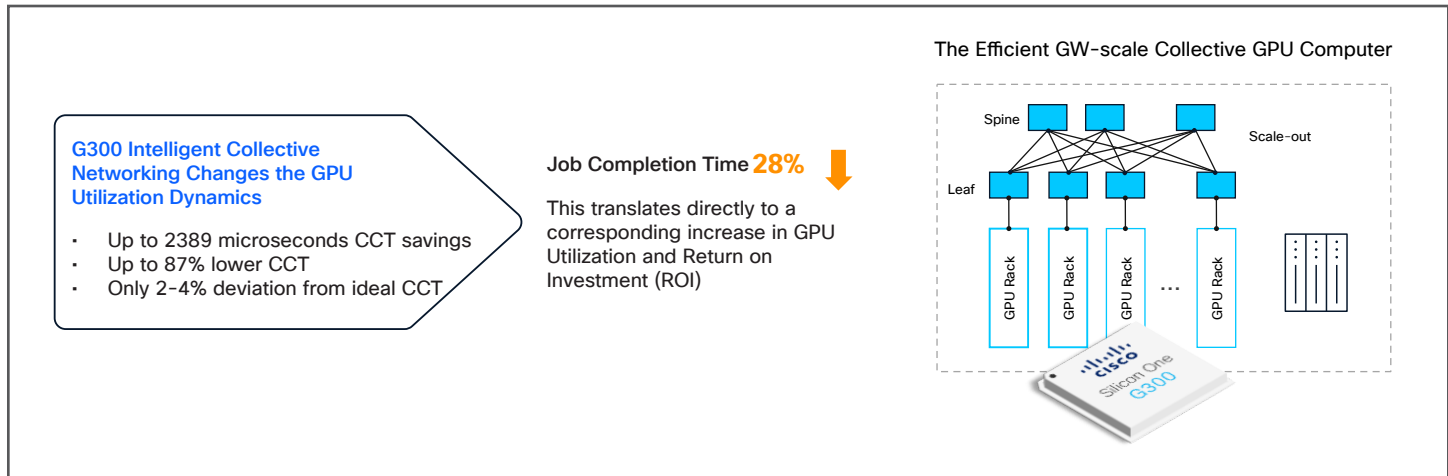


Figure 13. G300 102.4T switching with Intelligent Collective Networking changes the GPU utilization dynamics in gigawatt-scale data centers

## Summary

As data centers grow to gigawatt scale and deploy over a million GPUs, new challenges involving power, cooling, and efficient scaling require advanced solutions in network design and operation. The performance of GPU collective operations directly influences the rising costs of building and running these massive facilities. Existing network load-balancing methods, such as static ECMP and random packet spray, are no longer sufficient to meet the need for near-ideal collective completion time, which is critical for efficient AI workloads at scale. The Cisco Silicon One G300, a 102.4-Tbps switch, addresses these issues with a unique approach that optimizes for collective GPU operations rather than just traditional packet distribution. Its Intelligent Collective Networking feature incorporates intelligent local, fabric-wide, and topology-aware congestion management,

supported by the industry's largest packet buffers for absorbing GPU traffic bursts and microsecond-fast fault detection for immediate response to network issues.

Benchmarking demonstrated that static ECMP falls short in large-scale AI environments. While dynamic random packet spray techniques perform better, a local, fabric-wide, topology awareness of congestion and link failures delivers far superior network efficiency. The G300's advanced routing and buffering capabilities consistently achieved near-ideal CCT and reduced JCT, even with mixed traffic patterns and link failures. These gains translate directly into lower TCO and faster AI model training and deployment, making the G300 a pivotal technology for the next generation of massive, power-constrained data center buildouts.