



Accelerating AI Storage Networks with DDN's Data Intelligence Platform and NVIDIA Spectrum™-X for Storage

Contents

| | |
|--|-----------|
| Abstract | 3 |
| AI's Infrastructure Challenge | 3 |
| AI-Optimized Networking with NVIDIA Spectrum™-X for Storage | 4 |
| Dynamic Adaptive Routing | 4 |
| Ultra-Low Latency with NVIDIA BlueField®-3 SuperNICs | 4 |
| Seamless Scalability | 4 |
| DDN's AI-Optimized Data Intelligence Platform | 4 |
| High-Throughput AI Storage | 4 |
| NVIDIA Spectrum™-X and DDN's Unified Solution | 5 |
| Delivering End-to-End AI Performance Together | 6 |
| Adaptive-Routing | 6 |
| Congestion-control | 9 |
| Deployment Insights | 10 |
| Real-World Applications | 10 |
| Conclusion | 10 |

Abstract

As AI models continue to grow in complexity and scale, delivering high-performance infrastructure becomes critical. Traditional Ethernet-based storage networks, while sufficient for general workloads, struggle to keep up with the demands of modern AI systems, resulting in bottlenecks, underutilized GPUs, and inflated operational costs.

This whitepaper highlights the combined power of [NVIDIA Spectrum™-X for AI Storage Fabric](#) and [DDN's AI-optimized data intelligence platform](#), which together form a unified solution for scalable, efficient, and resilient AI infrastructure. Purpose-built for AI, this integrated approach ensures predictable performance, accelerates AI workflows, and delivers a future-proof foundation for next-generation AI-driven businesses.

AI's Infrastructure Challenge

AI applications require massive data throughput and ultra-low latency to feed GPUs consistently during training and inference. Conventional Ethernet-based networks fall short due to:

- **Congestion:** Traffic bottlenecks slow down data flow, creating imbalances in throughput.
- **Static Routing:** Traditional routing methods can't adapt dynamically to changing network conditions, leading to inefficient data movement.
- **Inconsistent Performance:** Fluctuating network loads result in unpredictable latency, negatively impacting AI model training.

DDN and NVIDIA have partnered to create a purpose-built solution that overcomes these limitations. By integrating **NVIDIA Spectrum™-X's adaptive networking technology** with **DDN's high-throughput, low-latency data intelligence platform**, enterprises can achieve faster AI training cycles, higher GPU utilization, and seamless scalability across on-premises and cloud environments.

AI-Optimized Networking with NVIDIA Spectrum™-X for Storage

NVIDIA Spectrum™-X transforms traditional Ethernet into an AI-optimized network fabric. Its advanced capabilities allow organizations to overcome the performance, and scalability challenges inherent in AI workloads.

Dynamic Adaptive Routing

Traditional networks often suffer from static routing inefficiencies, which can lead to congestion and reduced throughput. Spectrum-X RoCE adaptive routing solves this by dynamically adjusting paths in real-time. This ensures optimal bandwidth utilization, boosting storage bandwidth per GPU by up to **3.2x**¹, and preventing performance bottlenecks.

Ultra-Low Latency with NVIDIA BlueField®-3 SuperNICs

BlueField-3 SuperNICs are purpose-built for RoCE environments, delivering ultra-low-latency data transfers between GPU servers and storage nodes. By accelerating network processing through purpose-built hardware engines, BlueField-3 reduces CPU overhead and enables high-performance data movement, ensuring that GPUs are consistently supplied with data during training.

Seamless Scalability

Spectrum-X offers unmatched scalability with support for **256 200G ports** in a single hop and **16,000 ports**² in a two-tier topology. This ensures enterprises can grow their AI infrastructure to support next-generation platforms like NVIDIA DGXSuperPOD without compromising performance.

DDN's AI-Optimized Data Intelligence Platform

DDN's platform complements NVIDIA Spectrum™-X for storage by delivering high-throughput, low-latency storage tailored for AI's unique demands. Purpose-built to handle massive datasets and provide real-time access to critical information, DDN's solutions ensure that GPUs remain fully utilized during AI operations.

High-Throughput AI Storage

DDN delivers up to **up-to 2.4 TB/s read throughput per rack**³ providing the speed and scalability needed for large-scale AI workloads. Tested and validated with NVIDIA Spectrum™-X, these appliances ensure seamless data flow from storage to compute, minimizing idle GPU time and speeding up training cycles.

1 AI400X2+ Spectrum-X PoC Proof of Concept completed in December 2024. Additional details upon request.

2 [NVIDIA Whitepaper] Optimizing AI Training Clusters NVIDIA Spectrum™-X Accelerates AI Storage Networks Accessed: <https://www.nvidia.com/en-us/networking/spectrumx/> on 01/27/25

3 <https://www.ddn.com/products/ai400x2-turbo/>

NVIDIA Spectrum™-X and DDN's Unified Solution

The combined solution of NVIDIA Spectrum™-X and DDN is designed to deliver optimal data flow, ensuring that GPUs are consistently fed with data at the required speed and volume. This tightly-coupled solution enables organizations to maintain peak performance during AI model training, inference, and real-time analytics, which are often bottlenecked by conventional infrastructures. enables organizations to maintain peak performance during AI model training, inference, and real-time analytics, which are often bottlenecked by conventional infrastructures.

Together, the joint solution offers:

1. High-Performance Networking with Adaptive Routing

- a. NVIDIA Spectrum™-X dynamically distributes traffic to avoid network congestion and boost bandwidth efficiency.

2. Seamless Data Management and Storage

- a. DDN's data intelligence platform provides scalable, high-throughput storage capable of handling petabytes of data with sub-millisecond latency.

3. Linear Scalability

- a. Both the networking and storage components are designed for linear scalability, supporting growth without rearchitecting infrastructure.

4. Comprehensive Visibility and Control

- a. The integration of NVIDIA's NetQ™ telemetry with DDN's management interface ensures real-time monitoring, proactive issue resolution, and consistent performance across large AI clusters.

Delivering End-to-End AI Performance Together

Extensive testing of the combined solution demonstrates significant improvements in throughput, latency, and power efficiency compared to traditional AI infrastructures. Key results include:

- **33x faster data access** compared to NFS-based solutions.
- **10x power savings**, reducing operational costs while maintaining high performance.

In addition to that, NVIDIA Spectrum™-X features provide the following benefits when combined with DDN EXAScaler appliances, compared to native RoCEv2 performance.

- **Up to 3.2x higher IO bandwidth per GPU**, enabling faster model training.
- **Up to +16% throughput for concurrent workloads**, improving parallel workloads.

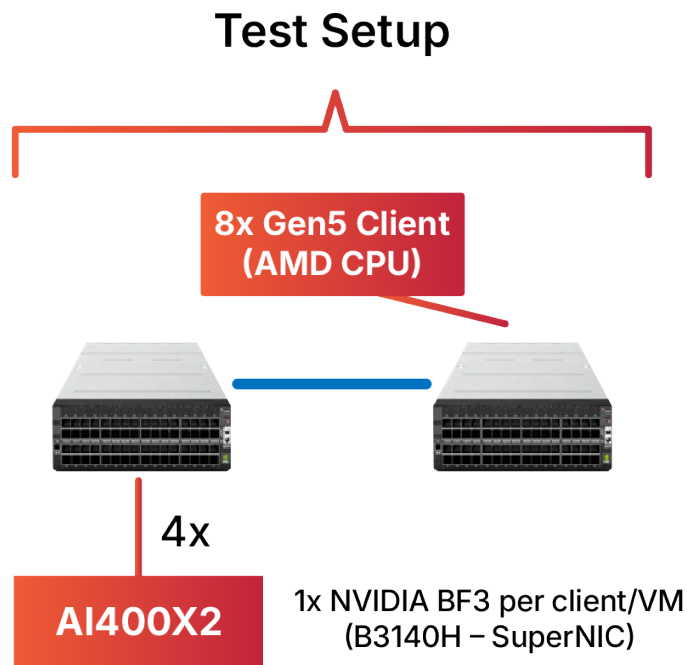
Adaptive-Routing

AI workloads can be extremely demanding on storage and compute resources. This situation has led to the development of many parallelization techniques at the application level, such as data parallelism, model parallelism, and pipelining. The parallel nature of AI workloads, which are composed of many large flows (also known as “elephant flows”), can be extremely demanding on the network, leading to potential bottlenecks due to the static nature of the switch-to-switch routing algorithm, known as ECMP (equal-cost multipath). NVIDIA Spectrum™-X, as an end-to-end technology that spans from storage to client, introduces Adaptive Routing as a bottleneck-aware feature that replaces traditional ECMP and dynamically distributes the load across all available links, depending on the available bandwidth.

With a single DDN EXAScaler appliance, this intelligent load distribution method can boost GPU performance by 3.2x during ECMP collisions. As the number of appliances increases, this benefit scales, resulting in significantly better performance at scale when large collisions happen.

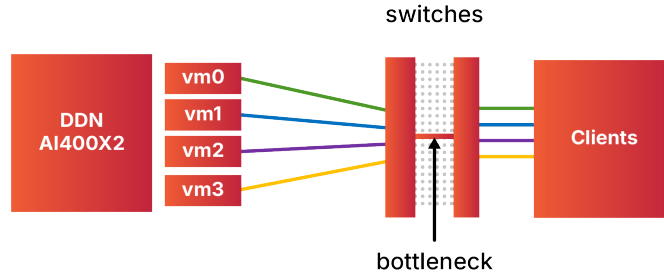
The following test setup was used to demonstrate this claim:

- 1 DDN AI400X2 appliance, which provides an EXAScaler filesystem to clients.
- 8 bare-metal clients with AMD CPUs
- 2 NVIDIA Spectrum™ SN5600 switches, which separate the clients from the AI400X2
- 1 NVIDIA BlueField®-3 SuperNIC per client, 4 NVIDIA BlueField®-3 SuperNICs for the AI400X2
- 4 OSFP-to-OSFP cables between both switches (MCP4Y10-N002)
- 4 OSFP to QSFP112, connecting the AI400X2 to the first switch (4x 400Gbps)
- 1 OSFP to QSFP112, connecting each client to the second switch (400Gbps/client)



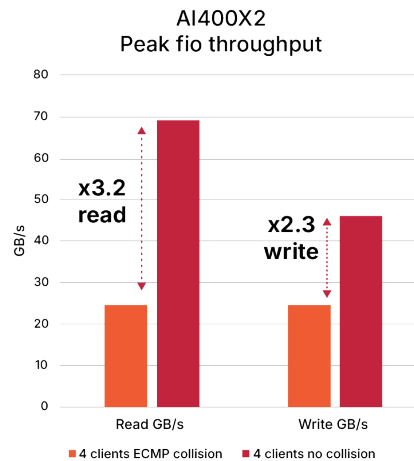
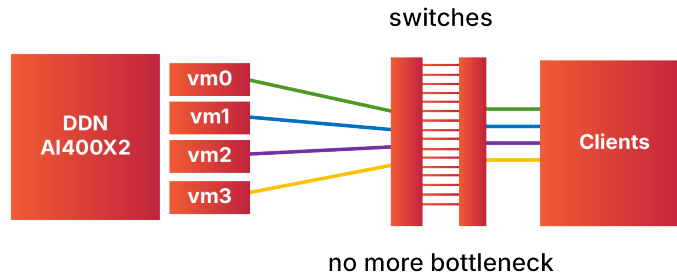
FIO was used to generate synthetic IO traffic between the clients and server. Due to the static nature of ECMP, it is extremely easy to observe a collision where two flows end up in the same link, leading to congestion and reducing performance by a factor equal to the number of clients colliding simultaneously.

Even in this small scenario with 8 clients, we were able to demonstrate a collision involving 4 clients due to ECMP, resulting in a dramatic performance drop. This can be visualized using the following scheme, where all traffic ends up on a single link.



In such a situation, Adaptive Routing entirely removes this bottleneck by distributing traffic across all switch-to-switch links.

This enhancement leads to a staggering 3.2x increase in read throughput and a 2.3x increase in write throughput per AI400X2 appliance, a factor that scales with the number of appliances.



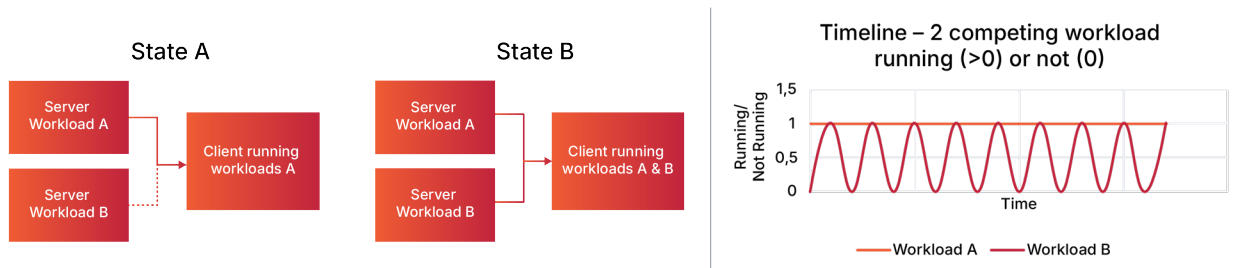
Congestion-control

When multiple workloads are run in parallel, the hardware-accelerated congestion control of Spectrum-X improves the aggregated throughput of EXAScaler by up-to +16% over traditional RoCEv2 performance, allowing a better efficiency in multi-tenant environment.

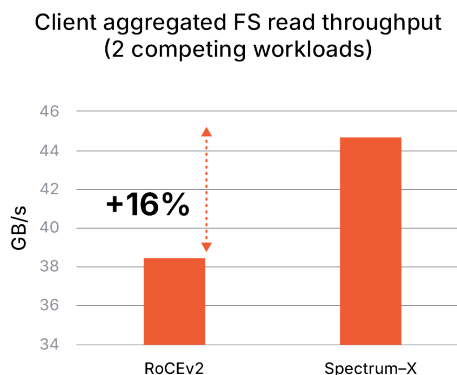
The following setup was used to demonstrate this claim:

- 1 DDN AI400X2 appliance, providing an EXAScaler filesystem to clients
- 1 bare-metal client
- 1 NVIDIA Spectrum SN5600 switch separating the client and the server
- 1 NVIDIA BlueField®-3 SuperNIC per client, 4 NVIDIA BlueField®-3 SuperNICs for the AI400X2
- 4 OSFP to QSFP112, connecting the AI400X2 to the first switch (4x 400Gbps)
- 1 OSFP to QSFP112, connecting each client to the second switch (400Gbps/client)

To segregate the impact of the congestion-control mechanism, two competing workloads were run on the same node, switching from state A to state B, in order to trigger the congestion-control mechanism of Spectrum-X along with simulating a concurrent workload scenario. In this situation, each server must dynamically adapt its throughput for each state change.



The following performance improvement was measured with NVIDIA Spectrum-X over RoCEv2 for EXAScaler.



Deployment Insights

Organizations adopting the NVIDIA Spectrum™-X and DDN solution experience significant deployment and operational benefits, including:

- **Rapid Deployment:** A **120-node DDN cluster** can be deployed in under **10 minutes**, enabling organizations to quickly scale their AI infrastructure.
- **Simplified Management:** Both Spectrum-X and DDN offer intuitive management interfaces, reducing administrative overhead and ensuring consistent performance.
- **Future-Ready Scalability:** The solution supports exabyte-scale storage and thousands of network ports, ensuring that businesses can scale as their AI workloads grow.

Real-World Applications

The joint solution is ideal for industries with demanding AI workloads, including:

- **Autonomous Vehicles:** Ensuring real-time data processing for sensor inputs and model updates.
- **Healthcare:** Accelerating medical imaging analysis and genomics research.
- **Financial Services:** Supporting high-frequency trading models and fraud detection systems.
- **Public Sector:** Enhancing research capabilities in defense and climate modeling.

Conclusion

The combined power of NVIDIA Spectrum™-X and DDN's AI-optimized data intelligence platform delivers a next-generation AI infrastructure capable of meeting the demands of modern AI workloads. By integrating high-performance networking with scalable, low-latency storage, this joint solution provides:

- **3.2x higher IO bandwidth per GPU** through Adaptive Routing¹.
- **3.2x faster reads** and **2.3x faster writes**², accelerating training cycles and model checkpointing.
- **+16% throughput for concurrent workloads**
- **10x power reduction**, driving energy efficiency and lowering operational costs³.

For organizations aiming to future-proof their AI infrastructure, NVIDIA Spectrum™-X and DDN offer a robust, high-performance foundation that accelerates innovation and delivers transformative business outcomes. To learn more, [meet with a DDN expert](#).

¹ [Whitepaper] Optimizing AI Training Clusters NVIDIA Spectrum™-X Accelerates AI Storage Networks
Accessed: <https://www.nvidia.com/en-us/networking/spectrumx/>

² DDN internal performance benchmark conducted in collaboration with NVIDIA on December 2024, using [specific hardware/software configuration]. Report available upon request.)

³ [Blog] Maximizing AI Potential and Mitigating Risks: The Power of the AI Center of Excellence. Accessed <https://www.ddn.com/blog/maximizing-ai-potential-and-mitigating-risks-the-power-of-the-ai-center-of-excellence/>