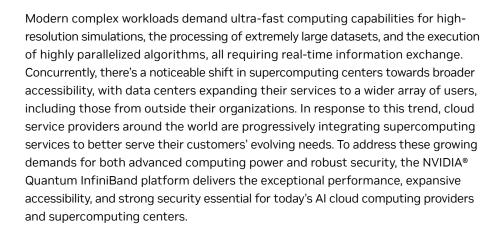


# **NVIDIA Quantum InfiniBand Product Guide**

High-performance networking platform for HPC, AI, and cloud data centers.



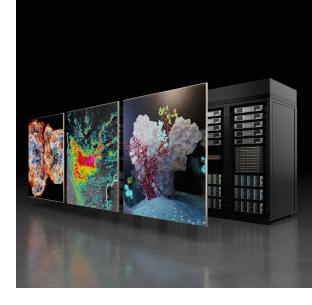


The NVIDIA Quantum InfiniBand switch portfolio enables compute clusters to operate at any scale while reducing capital expenses, operational costs, and infrastructure complexity.

NVIDIA Quantum InfiniBand switch systems deliver the highest performance and port density available. Innovative In-Network Computing capabilities, including NVIDIA Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)™ and advanced management features, such as self-healing network capabilities, quality of service, and enhanced congestion control, provide a performance boost for industrial, Al, and scientific applications. This includes hardware-based acceleration of collective communication operations that are used extensively with Al systems. Select switches with NVIDIA InfiniBand router capabilities enable the scaling-out of InfiniBand clusters to a vast number of nodes. This ensures the preservation of peak performance and reliability required for research, simulations, Al, and cloud applications data routing and processing.

#### **BlueField Data Processing Units**

NVIDIA BlueField® data processing units (DPUs) integrate powerful computing, high-speed networking, and extensive programmability, enabling organizations to build software-defined, hardware-accelerated infrastructures from cloud to core data center to edge. BlueField excels in offloading and accelerating tasks related to networking, storage and security. It isolates functions in software-defined networking, storage, safeguarding, and management, profoundly improving data



center performance, efficiency, and security. In addition to providing optimal baremetal performance and native support for multi-node tenant isolation, BlueField's innovative architecture and capabilities set a new benchmark in optimizing data center operations and paving the way for advanced supercomputing infrastructures.

#### **NVIDIA ConnectX Network Adapters**

NVIDIA ConnectX® network adapters provide a broad set of software-defined, hardware-accelerated networking, storage, and security capabilities, which enable organizations to modernize and secure their IT infrastructures. ConnectX provides ultra-low latency and extreme throughput, and innovative In-Network Computing engines such as MPI Tag Matching and All-to-All hardware engines, addressing traditional enterprise needs up to the world's most-demanding Al, scientific computing, and hyperscale cloud data center workloads.

## InfiniBand Long-Haul Systems

NVIDIA MetroX® systems extend the reach of InfiniBand to up to 40 kilometers, enabling native RDMA connectivity between remote data centers - for high availability and disaster recovery - from edge infrastructures to data centers, or between data centers and remote storage. Delivering up to 100Gb/s data throughout on long-haul ports and 200Gb/s on standard ports, MetroX lets users easily migrate application jobs from one InfiniBand center to another, or combine the compute power of multiple remote data centers together for higher overall performance and scalability.

#### Skyway Gateway to Ethernet

NVIDIA Skyway™ enables scalable and efficient connectivity from high-performance, low-latency InfiniBand data centers to external Ethernet networks and infrastructures. Supporting high availability and load balancing, with standard IP routing protocols, Skyway is a simple and cost-effective option to empower data centers to achieve the lowest interconnect latency.

#### **LinkX Cables and Transceivers**

**NVIDIA LinkX®** cables and transceivers maximize the performance of InfiniBand networks to deliver high-bandwidth, low-latency, highly reliable, and robust connectivity. Providing superior performance, LinkX products undergo rigorous testing to ensure the highest quality.

#### **NVIDIA UFM Network Fabric Management**

**NVIDIA UFM®** (Unified Fabric Management) revolutionizes data center networking management. Supporting scale-out InfiniBand data centers, UFM combines enhanced, real-time network telemetry with Al-powered cyber intelligence and analytics to realize higher utilization of fabric resources and a competitive advantage. UFM includes fabric diagnostics, monitoring, alerting, provisioning, and advanced features such as congestion monitoring, and fabric segmentation and isolation.

## **Software for Optimal Performance**

**NVIDIA HPC-X®** is a full-featured, tested, and packaged software toolkit that enables MPI and SHMEM/PGAS programming languages to achieve high performance, scalability, and efficiency. HPC-X leverages In-Network Computing to increase application and network performance, reducing latencies and increasing throughput for improved data processing and faster access to data.

# **Software for Optimal Performance**

Switch	Configuration	Advanced Features	Size
QM9700 NVIDIA Quantum-2 Switch Series	<ul> <li>64 400Gb/s ports that can be split into 128 200Gb/s ports</li> <li>51.2Tb/s total throughput</li> <li>Ultra-low switch latency</li> </ul>	<ul> <li>Accelerated In-Network Computing</li> <li>Limitless scalability</li> <li>3rd generation of NVIDIA SHARP (SHARPv3)</li> <li>Enhanced management with out-of-the box bring-up for up to 2,000 nodes</li> <li>Internally managed and externally managed flavors</li> </ul>	1U
QM8700 NVIDIA Quantum Fixed Configuration Switch Series  NVIDIA Skyway InfiniBand to Ethernet Gateway	<ul> <li>40 200Gb/s ports or 80 100Gb/s ports</li> <li>16Tb/s aggregate switch throughput</li> <li>Ultra-low switch latency</li> <li>8 200Gb/s/100Gb/s InfiniBand ports</li> <li>8 200/100Gb/s Ethernet ports</li> <li>1.6Tb/s aggregate switch throughput</li> </ul>	<ul> <li>Internally managed and externally managed flavors</li> <li>Self-healing networking</li> <li>NVIDIA SHARPv2: In-network collective offloads support of low-latency and streaming aggregation for Al applications</li> <li>Adaptive routing, congestion control, and QoS</li> <li>Industry-leading InfiniBand to Ethernet gateway</li> <li>Future-ready architecture</li> </ul>	1U
NVIDIA MetroX-3 XC Switch Systems	<ul> <li>2 100Gb/s InfiniBand QSFP112 long-haul ports</li> <li>2 100Gb/s InfiniBand QSFP112 local ports</li> </ul>	<ul> <li>Connectivity over long distances and DWDM</li> <li>Adaptive routing and congestion control</li> <li>Self-healing networking</li> </ul>	1U
NVIDIA MetroX-2 Switch Systems	<ul> <li>2 100Gb/s InfiniBand QSFP56 long-haul ports</li> <li>8 200Gb/s InfiniBand QSFP56 local ports</li> </ul>	<ul> <li>Adaptive routing and congestion control</li> <li>Self-healing networking</li> </ul>	1U

# **NVIDIA DPUs**

DPU	Speed	Connectors	Bus	Capabilities	Form Factor
BlueField-3	Up to 400Gb/s	QSFP112	PCle Gen5 x16 2x PCle Gen4 x16	BlueField-3 SoC > 16 Arm A78 cores	PCle stand-up
	·			> 400Gb/s InfiniBand	
				ConnectX-7 hardware offload	
BlueField-2	Up to	QSFP56	PCle Gen3/4 x16	BlueField-2 SoC	PCIe stand-up
	200Gb/s		2x PCle Gen3 x16	> 168 Arm A78 cores	
				> 200Gb/s InfiniBand	
				ConnectX-6 Dx hardware offloads	

### **NVIDIA Network Adapters**

Adapters	Speed	Connectors	Bus	RDMA Message Rate (Mmps)	Advanced Features	Form Factor
ConnectX-7	Up to 400Gb/s	OSFP or QSFP112	PCIe Gen5 x16 2x PCIe Gen4 x16	370	<ul> <li>MPI All-to-All offloads</li> <li>Enhanced congestion control</li> <li>Secure boot with hardware root of trust</li> <li>NVIDIA Multi-Host up to 4x hosts (specific OPNs)</li> <li>NVMe-oF target offload</li> <li>NVIDIA SHARPv3 support</li> </ul>	<ul><li>PCle stand-up</li><li>PCle Socket Direct</li><li>OCP 3.0</li></ul>
ConnectX-6	Up to 200Gb/s	QSFP56	PCIe Gen3/4 x16 2x PCIe Gen3 x16	215	<ul> <li>MPI tag-matching offload</li> <li>Block-level XTS-AES hardware encryption</li> <li>Secure firmware update</li> <li>NVIDIA Multi-Host up to 4x hosts</li> <li>NVMe-oF target offload</li> <li>NVIDIA SHARPv2 support</li> </ul>	<ul><li>PCle stand-up</li><li>PCle Socket Direct</li><li>OCP 3.0</li></ul>

#### InfiniBand Interconnect

Direct Attach Copper (DAC)	Active Copper Cables (ACC)	Active Optical Cables (AOC)
> HDR max reach: 2m	> NDR max reach: 5m	> HDR max reach: 150m
> NDR switch-to-switch max reach: 2m		> NDR-to-HDR max reach: 30m
NDR switch-to-network adapter/DPU max reach: 3m		

#### **Optical Transceivers**

- > Maximum reaches:
  - 2x 400Gb/s: Single-mode twin port 100m, 500m, 2km
- 2x 400Gb/s: Multimode twin-port 50m
- 400Gb/s: Single-mode 100m, multimode 50m
- 200Gb/s: Single-mode 2km, multimode 100m
- 100Gb/s: Single-mode multimode 100m

# Ready to Get Started?

To learn more about NVIDIA Quantum InfiniBand Networking Solutions, visit: <a href="mailto:nvidia.com/networking/products/infiniband">nvidia.com/networking/products/infiniband</a>

