



NVIDIA QUANTUM-2 INFINIBAND PLATFORM

Extreme Performance for Exascale AI

The NVIDIA Quantum-2 InfiniBand platform provides AI developers and scientific researchers with the highest networking performance available to take on the world's most challenging problems. New NVIDIA In-Network Computing acceleration engines provide ultra-low latency and double the data throughput, while delivering the scalability and feature-rich capabilities required for supercomputers, artificial intelligence, and hyperscale cloud data centers.

NVIDIA® Quantum-2 enhances and extends In-Network Computing acceleration technology with preconfigured and programmable engines such as NVIDIA Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)™, Message Passing Interface (MPI) tag matching, MPI_Alltoall, and programmable cores, as well as full transport offload with RDMA, GPUDirect RDMA and GPUDirect Storage, which deliver the best cost per node and return on investment (ROI).

World-Leading Performance

The NVIDIA Quantum-2 InfiniBand platform continues to set world records for high-performance networking, delivering 400Gb/s per port—2X higher bandwidth compared to the previous generation, 3X higher switch silicon port density, 5X higher switch system capacity, 32X higher AI acceleration power per switch, and the ability to connect over one million nodes at 400Gb/s in a 3-hop Dragonfly+ topology.

Featuring the third generation of NVIDIA SHARP technology, the platform allows virtually unlimited scalability for large data aggregation for small and large data aggregations through the network — 32X higher AI acceleration power compared to the previous generation. Moreover, the third generation of SHARP technology enables multiple tenants or multiple parallel applications to share the infrastructure without any performance degradation. MPI_Alltoall acceleration and MPI tag matching hardware engines, and other features like advanced congestion control, adaptive routing, and self-healing networking provide critical enhancements to high-performance computing (HPC) and AI clusters, enabling them to reach even higher levels of performance than before.

NVIDIA Quantum-2 Portfolio

The NVIDIA Quantum-2 switch application-specific integrated circuit (ASIC) delivers 64 400Gb/s InfiniBand ports or 128 200Gb/s InfiniBand ports, the third generation of NVIDIA SHARP technology, and other InfiniBand features. The NVIDIA Quantum-2 platform includes fixed-configuration and modular switch systems based on the NVIDIA Quantum-2 switch ASIC.

The NVIDIA ConnectX®-7 InfiniBand host channel adapter (HCA) ASIC delivers 400Gb/s data throughput and supports 32 lanes of PCIe Gen5 or Gen4 for host connectivity. Built on the most advanced 100Gb/s-per-lane serializer/deserializer (SerDes) technology, 400Gb/s InfiniBand physical connectivity is based on octal small form-factor pluggable (OSFP) connectors on both the switches and HCA endpoints. Each switch OSFP connector holds two 400Gb/s InfiniBand ports or four 200Gb/s InfiniBand ports. Each HCA OSFP connector carries a single 400Gb/s InfiniBand port. The 400Gb/s cabling offering includes active and passive copper cables, transceivers, and multi-fiber push on (MPO) optical cables.

Fixed-Configuration Switches

The NVIDIA Quantum-2 family of fixed configuration switches consists of switches with 64 400Gb/s ports on 32 physical OSFP connectors, which can be split to deliver up to 128 200Gb/s ports. The compact 1U switches offering includes internally managed and externally managed (aka unmanaged) versions. The switches carry an aggregated bidirectional throughput of 51.2 terabits per second (Tb/s), with a landmark capacity of more than 66.5 billion packets per second. An ideal rack-mounted InfiniBand solution, the fixed configuration switch allows maximum flexibility as it enables a variety of topologies, including Fat Tree, DragonFly+, multi-dimensional Torus, and more.

Modular Switches

The NVIDIA Quantum-2 modular switch offering includes several flavors: 2,048 400Gb/s InfiniBand ports (that can be split into 4,096 200Gb/s InfiniBand ports); 1,024 400Gb/s ports (that can be split into 2,048 200Gb/s ports); or 512 400Gb/s ports (that can be split into 1024 200Gb/s ports). The large modular switch is based on a non-blocking, two-level Fat Tree “fabric-in-a-rack” occupying a three rack width. It carries a total unidirectional throughput of 819Tb/s or bidirectional throughput of 1.64 petabits per second (Pb/s), which is 5X over the previous generation of InfiniBand modular switches.

Host Channel Adapters

ConnectX-7 HCAs are offered in various form factors, delivering single or dual ports at 400Gb/s speeds with OSFP connectors or 200Gb/s speeds with quad small form-factor pluggable (QSFP112) connectors. The popular form factors meet the Card Electromechanical (CEM) Specification with 16 lanes of PCIe Gen5 or Gen4. Some versions come with the option to connect an additional 16-lane auxiliary card, which leverages NVIDIA Socket Direct® technology, to achieve 32 lanes of PCIe Gen4. Other form factors include Open Compute Project (OCP) 3.0 with OSFP connectors, OCP 3.0 with QSFP112 connectors, and CEM PCIe x16 with QSFP112 connectors.

The ConnectX-7 HCAs deliver advanced In-Network Computing with MPI_Alltoall and MPI tag matching hardware engines, and other fabric enhancement features such as quality of service (QoS), congestion control, and more.

Transceivers and Cables

400Gb/s InfiniBand connectivity provides the maximum flexibility to build a topology of choice, using single-mode and multimode transceivers, MPO fiber cables, active copper cables (ACCs), and direct attached cables (DACs).

- > Twin-port transceivers with finned OSFP connectors should be plugged into the (air-cooled) fixed configuration switch, whereas twin-port transceivers with flat OSFP connectors can be plugged into the liquid-cooled modular switches and HCAs.
- > When optical connectivity between the switch and HCA is required, the twin port transceivers are used at the switch side. On the HCA side, a single-port transceiver is used—either an 400Gb/s or 200Gb/s single-port transceiver. For the latter case, two split MPO fibers are used to achieve connectivity to four 200Gb/s adapters.
- > MPO fibers are offered at 3–150 meters (m) and one-to-two splitter fibers are offered at 3–50m.
- > Switch-to-HCA connectivity offers both DACs (up to 1.5m), and ACCs (up to 3m). The one-to-two split cable connects an OSFP switch port (that holds two 400 Gb/s InfiniBand ports) to two individual 400Gb/s HCAs. The one-to-four split cable is used to connect an OSFP switch port to four 200Gb/s HCAs.
- > OSFP to OSFP DAC (2m) and ACC (up to 5m) are used to connect two parallel 400Gb/s connections between two neighboring switches. This option enables the deployment of cost-effective topologies, such as DragonFly+, that leverage the colocation of spine switches in a single rack.
- > Backward compatibility AOC and DAC are also available. With these cables, it's possible to connect the new 400Gb/s clusters to existing 200Gb/s or 100Gb/s infrastructures.

UFM Cyber-AI

The NVIDIA Unified Fabric Manager (UFM) Cyber-AI platform offers enhanced and real-time network telemetry, together with AI-powered intelligence and advanced analytics. Enabling IT managers to discover operational anomalies and predict network failures, it improves security and data center uptime and decreases overall operating expenses.

Ordering Information

NVIDIA Quantum-2 Modular Switches

Orderable Part Number (OPN)	Description
MCS9500	NVIDIA 1,600Tb/s, 2,048-port 400Gb/s InfiniBand chassis
MCS9510	NVIDIA 800Tb/s, 1,024-port 400Gb/s InfiniBand chassis
MCS9520	NVIDIA 400Tb/s, 512-port 400Gb/s InfiniBand chassis
MCS9505-CDU	NVIDIA MCS95XX director systems liquid-to-liquid CDU
MCS9505-AHX	NVIDIA MCS95XX director systems liquid-to-air heat exchanger
MCS9505-KIT	NVIDIA MCS95XX chassis equipment kit
MCS9505-FWL	NVIDIA MCS95XX firewall appliance
MQM9510-N	NVIDIA Quantum-2 400Gb/s InfiniBand, 2U leaf blade, 128 400Gb/s ports, 64 OSFP ports
MQM9520-N	NVIDIA Quantum-2 400Gb/s InfiniBand, 2U spine blade, 128 400Gb/s ports
MTDF-LIQ-D	Modular chassis 19 liters PG25 coolant
MMB9500	NVIDIA MCS95XX management switch
MTDF-PDU-A	NVIDIA MCS95XX modular system PDU

NVIDIA Quantum-2 Fixed-Configuration Switches

Orderable Part Number (OPN)	Description
MQM9790-NS2F	NVIDIA Quantum-2-based 400Gb/s InfiniBand switch, 64 400Gb/s ports, 32 OSFP ports, non-blocking switching capacity of 51.2Tb/s, two power supplies (AC), standard depth, unmanaged, power-to-connector (P2C) airflow, rail kit
MQM9790-NS2R	NVIDIA Quantum-2-based 400Gb/s InfiniBand switch, 64 400Gb/s ports, 32 OSFP ports, non-blocking switching capacity of 51.2Tb/s, two power supplies (AC), standard depth, unmanaged, connector-to-power (C2P) airflow, rail kit
MQM9700-NS2F	NVIDIA Quantum-2-based 400Gb/s InfiniBand switch, 64 400Gb/s ports, 32 OSFP ports, non-blocking switching capacity of 51.2Tb/s, two power supplies (AC), standard depth, managed, P2C airflow, rail kit
MQM9700-NS2R	NVIDIA Quantum-2-based 400Gb/s InfiniBand switch, 64 400Gb/s ports, 32 OSFP ports, non-blocking switching capacity of 51.2Tb/s, two power supplies (AC), standard depth, managed, C2P airflow, rail kit

ConnectX-7 Host Channel Adapter (HCA)

PCIe Standup Adapters

OPN	Description
MCX75510AAS-NEAT	NVIDIA ConnectX-7 adapter card, 400Gb/s InfiniBand, single-port OSFP, Socket Direct ready, PCIe 5.0 x16 with IPEX connectors for extension, no crypto, tall bracket
MCX75310AAS-NEAT	NVIDIA ConnectX-7 adapter card, 400Gb/s InfiniBand, single-port OSFP, PCIe 5.0 x16, no crypto, tall bracket
MCX75510AAS-HEAT	NVIDIA ConnectX-7 adapter card, 200Gb/s, single-port OSFP, Socket Direct ready, PCIe 5.0 x16 with IPEX connectors for extension, no crypto, tall bracket
MCX75310AAS-HEAT	NVIDIA ConnectX-7 adapter card, 200Gb/s, single-port OSFP, PCIe 5.0 x16, no crypto, tall bracket
MCX75210AAS-NEAT	NVIDIA ConnectX-7 adapter card, 400Gb/s InfiniBand, single-port OSFP, PCIe 5.0 2x8 in a row, no crypto, tall bracket
MCX75210AAS-HEAT	NVIDIA ConnectX-7 adapter card, 200Gb/s, single-port OSFP, PCIe 5.0 2x8 in a row, no crypto, tall bracket
MCX755105AS-HEAT	NVIDIA 200Gb/s single-port InfiniBand, QSFP, PCIe Gen5 x16, half-height, half-length (HHHL), extension option
MCX755106AS-HEAT ¹	NVIDIA 200Gb/s dual-port virtual protocol interconnect (VPI), QSFP, PCIe Gen5 x16 HHHL, extension option

¹ This card supports one port of InfiniBand, and a second port as either InfiniBand or Ethernet.

Transceivers and Cables

Transceivers

OPN	Description
MMA4Z00-NS	NVIDIA twin port transceiver, 800Gb/s, 2x 400Gb/s, OSFP, 2xMPO, 850nm MMF, SR8, up to 30m, finned
MMA4Z00-NS400	NVIDIA single port transceiver, 400Gb/s, OSFP, MPO, 850nm MMF, SR4, up to 30m, flat top
MMA4Z00-NS200	NVIDIA single port transceiver, 200Gb/s, OSFP, MPO, 850nm MMF, SR2, up to 30m, flat top
MMA1Z00-NS400	NVIDIA single port transceiver, 400Gb/s, QSFP112, MPO, 850nm MMF, SR4, up to 30m, flat top
MMA1Z00-NS200	NVIDIA single port transceiver, 200Gb/s, QSFP112, MPO, 850nm MMF, SR2, up to 30m, flat top
MMS4X00-NL	NVIDIA twin-port transceiver, 800Gb/s, 2x 400Gb/s, OSFP, 2x MPO, 1,310 nanometers (nm) single-mode fiber (SMF), up to 30m
MMS4X00-NL400	NVIDIA single-port 400Gb/s transceiver (to be used at the HCA side)
MMS4X00-NL200	NVIDIA single-port 200Gb/s transceiver (to be used at the HCA side)
MMS4X00-NS	NVIDIA twin-port transceiver, 800Gb/s, 2x 400Gb/s, OSFP, 2x MPO, 1,310nm SMF, DR8, up to 150m
MMS4X00-NS400	NVIDIA single-port transceiver, 400Gb/s, OSFP, MPO, 1,310nm SMF, DR4, up to 150m
MMS4X00-NS200	NVIDIA single-port transceiver, 200Gb/s, OSFP, MPO, 1,310nm SMF, DR2, up to 150m

Auxiliary Cards

OPN	Description
MTMK9100-T15	NVIDIA auxiliary kit for additional PCIe Gen4 x16 connection, PCIe Gen4 x16 passive auxiliary card, two 150 millimeters (mm) IPEX cables
MTMK9100-T25	NVIDIA auxiliary kit for additional PCIe Gen4 x16 connection, PCIe Gen4 x16 passive auxiliary card, two 250 millimeters (mm) IPEX cables
MTMK9100-T35	NVIDIA auxiliary kit for additional PCIe Gen4 x16 connection, PCIe Gen4 x16 passive auxiliary card, two 350mm IPEX cables
MTMK9100-T55	NVIDIA auxiliary kit for additional PCIe Gen4 x16 connection, PCIe Gen4 x16 passive auxiliary card, two 550mm IPEX cables

Open Compute Project (OCP) Adapters

OPN	Description
MCX75343AAS-NEAC	NVIDIA 400Gb/s single port, OSFP, PCIe Gen5 x16 OCP3.0 ² small form factor (TSFF)
MCX753436AS-HEAB	NVIDIA 200Gb/s VPI dual port, QSFP, PCIe Gen5 x16 OCP3.0 (SFF)

² Pre OCP3.2 Specification

MPO Fiber

OPN	Description
MFP7E30-N003	NVIDIA passive fiber cable, SMF, MPO to MPO, 3m
MFP7E30-N005	NVIDIA passive fiber cable, SMF, MPO to MPO, 5m
MFP7E30-N007	NVIDIA passive fiber cable, SMF, MPO to MPO, 7m
MFP7E30-N010	NVIDIA passive fiber cable, SMF, MPO to MPO, 10m
MFP7E30-N015	NVIDIA passive fiber cable, SMF, MPO to MPO, 15m
MFP7E30-N020	NVIDIA passive fiber cable, SMF, MPO to MPO, 20m
MFP7E30-N030	NVIDIA passive fiber cable, SMF, MPO to MPO, 30m
MFP7E30-N050	NVIDIA passive fiber cable, SMF, MPO to MPO, 50m
MFP7E30-N100	NVIDIA passive fiber cable, SMF, MPO to MPO, 100m
MFP7E30-N150	NVIDIA passive fiber cable, SMF, MPO to MPO, 150m

Transceivers and Cables (continued)

MPO Split Fiber

OPN	Description
MFP7E40-N003	NVIDIA passive fiber cable, SMF, MPO to 2x MPO, 3m
MFP7E40-N005	NVIDIA passive fiber cable, SMF, MPO to 2x MPO, 5m
MFP7E40-N007	NVIDIA passive fiber cable, SMF, MPO to 2x MPO, 7m
MFP7E40-N010	NVIDIA passive fiber cable, SMF, MPO to 2x MPO, 10m
MFP7E40-N015	NVIDIA passive fiber cable, SMF, MPO to 2x MPO, 15m
MFP7E40-N020	NVIDIA passive fiber cable, SMF, MPO to 2x MPO, 20m
MFP7E40-N030	NVIDIA passive fiber cable, SMF, MPO to 2x MPO, 30m
MFP7E40-N050	NVIDIA passive fiber cable, SMF, MPO to 2x MPO, 50m

Direct Attached Copper (DAC) Switch to HCA

OPN	Description
MCP7Y00-N001	NVIDIA passive copper splitter cable, InfiniBand 800Gb/s to 2x 400Gb/s, OSFP to 2x OSFP, 1m
MCP7Y00-N01A	NVIDIA passive copper splitter cable, InfiniBand 800Gb/s to 2x 400Gb/s, OSFP to 2x OSFP, 1.5m
MCP7Y50-N001	NVIDIA passive copper splitter cable, InfiniBand 800Gb/s to 4x 200Gb/s, OSFP to 4x OSFP, 1m
MCP7Y50-N01A	NVIDIA passive copper splitter cable, InfiniBand 800Gb/s to 4x 200Gb/s, OSFP to 4x OSFP, 1.5m
MCP7Y10-N001	NVIDIA passive copper splitter cable, InfiniBand 800Gb/s to 2x 400Gb/s, OSFP to 2x QSFP112, 1m
MCP7Y10-N01A	NVIDIA passive copper splitter cable, InfiniBand 800Gb/s to 2x 400Gb/s, OSFP to 2x QSFP112, 1.5m
MCP7Y40-N001	NVIDIA passive copper splitter cable, InfiniBand 800Gb/s to 4x 200Gb/s, OSFP to 4x QSFP112, 1m
MCP7Y40-N01A	NVIDIA passive copper splitter cable, InfiniBand 800Gb/s to 4x 200Gb/s, OSFP to 4x QSFP112, 1.5m

Direct Attached Copper (DAC) Switch to Switch

OPN	Description
MCA4J80-N001	NVIDIA active copper cable, 400Gb/s InfiniBand, up to 800Gb/s, OSFP, 1m
MCA4J80-N002	NVIDIA active copper cable, 400Gb/s InfiniBand, up to 800Gb/s, OSFP, 2m
MCA4J80-N003	NVIDIA active copper cable, InfiniBand, up to 800Gb/s, OSFP, 3m
MCP4Y10-N00A	NVIDIA passive copper cable, InfiniBand 400Gb/s, up to 800Gb/s, OSFP, 0.5m
MCP4Y10-N001	NVIDIA passive copper cable, InfiniBand 400Gb/s, up to 800Gb/s, OSFP, 1m
MCP4Y10-N01A	NVIDIA passive copper cable, InfiniBand 400Gb/s, up to 800Gb/s, OSFP, 1.5m
MCP4Y10-N002	NVIDIA passive copper cable, InfiniBand 400Gb/s, up to 800Gb/s, OSFP, 2m

Active Copper (ACC) Switch to HCA

OPN	Description
MCA7J60-N002	NVIDIA active copper splitter cable, split port InfiniBand 800Gb/s to 2x 400Gb/s, OSFP to 2x OSFP, 2m
MCA7J60-N003	NVIDIA active copper splitter cable, 400Gb/s InfiniBand, 800Gb/s to 2x 400Gb/s, OSFP to 2x OSFP, 3m
MCA7J70-N002	NVIDIA active copper splitter cable, 400Gb/s InfiniBand, 800Gb/s to 4x 200Gb/s, OSFP to 4x OSFP, 2m
MCA7J70	NVIDIA active copper splitter cable, 400 Gb/s InfiniBand, 800Gb/s to 4x 200Gb/s, OSFP to 4x OSFP, 3m
MCA7J65-N002	NVIDIA active copper splitter cable, split port InfiniBand 800Gb/s to 2x 400Gb/s, OSFP to 2x QSFP112, 2m
MCA7J65-N003	NVIDIA active copper splitter cable, 400Gb/s InfiniBand, 800Gb/s to 2x 400Gb/s, OSFP to 2x QSFP112, 3m
MCA7J75-N002	NVIDIA active copper splitter cable, 400Gb/s InfiniBand, 800Gb/s to 4x 200Gb/s, OSFP to 4x QSFP112, 2m
MCA7J75-N003	NVIDIA active copper splitter cable, 400 Gb/s InfiniBand, 800Gb/s to 4x 200Gb/s, OSFP to 4x QSFP112, 3m

Active Copper (ACC) Switch to Switch

OPN	Description
MCA4J80-N001	NVIDIA active copper cable, 400Gb/s InfiniBand, up to 800Gb/s, OSFP, 1m
MCA4J80-N002	NVIDIA active copper cable, 400Gb/s InfiniBand, up to 800Gb/s, OSFP, 2m
MCA4J80-N003	NVIDIA active copper cable, InfiniBand, up to 800Gb/s, OSFP, 3m
MCA4J80-N004	NVIDIA active copper cable, InfiniBand, up to 800Gb/s, OSFP, 4m
MCA4J80-N005	NVIDIA active copper cable, InfiniBand, up to 800Gb/s, OSFP, 5m

Transceivers and Cables (continued)

Backward-Compatible AOC, Connecting NVIDIA Quantum-2 to ConnectX-6 200Gb/s HCA

OPN	Description
MFA7U10-H003	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 2x 200Gb/s, OSFP to 2x QSFP56, 3m
MFA7U10-H005	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 2x 200Gb/s, OSFP to 2x QSFP56, 5m
MFA7U10-H010	NVIDIA AOC splitter cable, InfiniBand twin-port 200Gb/s, 400Gb/s to 2x 200Gb/s, OSFP to 2x QSFP56, 10m
MFA7U10-H015	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 2x 200Gb/s, OSFP to 2x QSFP56, 15m
MFA7U10-H020	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 2x 200Gb/s, OSFP to 2x QSFP56, 20m
MFA7U10-H030	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 2x 200Gb/s, OSFP to 2x QSFP56, 30m

Backward-Compatible Split AOC, Connecting NVIDIA Quantum-2 to ConnectX-6 100Gb/s

OPN	Description
MFA7U20-H003	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 4x 100Gb/s, OSFP to 4x QSFP56, 3m
MFA7U20-H005	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 4x 100Gb/s, OSFP to 4x QSFP56, 5m
MFA7U20-H010	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 4x 100Gb/s, OSFP to 4x QSFP56, 10m
MFA7U20-H015	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 4x 100Gb/s, OSFP to 4x QSFP56, 15m
MFA7U20-H020	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 4x 100Gb/s, OSFP to 4x QSFP56, 20m
MFA7U20-H030	NVIDIA AOC splitter, InfiniBand twin-port 200Gb/s, 400Gb/s to 4x 100Gb/s, OSFP to 4x QSFP56, 30m

[Learn more](#)

To learn more about NVIDIA Quantum-2 InfiniBand platform, visit [NVIDIA.com/InfiniBand](https://www.nvidia.com/InfiniBand).