H100 CNX CONVERGED ACCELERATOR PRODUCT OVERVIEW **NVIDIA**.

GPU-ACCELERATED WORKLOADS NEEDING DEEPER NETWORK INTEGRATION



NVIDIA H100 CNX

Converged Accelerator

H100 Tensor Core GPU

- 80GB HBM2e memory, > 2.0 TB/s
- Up to 7 MIG instances







Al on 5G

Multi-node Training

5G vRAN



ConnectX-7

- 400 Gb/s
- Network/security accelerators

PCI Gen5, 128 GB/s

2 slot FHFL

350 W



FASTER DATA SPEEDS WITH CONVERGED ARCHITECTURE

Ideal for Mainstream Servers



H100 CNX USE CASES



Any GPU-accelerated, I/O intensive application can benefit

Availability

- Mass Production in Q3
- NVIDIA-Certified Servers in Q4

HIGHER PERFORMANCE FOR GPU-ACCELERATED I/O INTENSIVE WORKLOADS

Up to 1.5 faster on large-scale training*



Challenges with discrete accelerators

- Bound by speed of host PCIe backplane
- Contention due to bottleneck at CPU
- Number of PCIe lanes limits additional devices

*assumes a cluster of 16 hosts, 4 GPUs each



Benefits of Converged Accelerator

- PCIe Gen5 speed regardless of host
- Tight coupling of GPU and network avoids contention
- Scale up capabilities with fewer devices

🧆 NVIDIA.

H100CX KEY BENEFITS



NVIDIA H100 CNX - H100 AND CONNECTX-7 CONVERGED ACCELERATOR

Unprecedented performance for GPU-powered, IO-intensive workloads

