

# GigaIO™ SuperNODE™ - MI300X

## IMMENSE INFERENCE PERFORMANCE

Breaks traditional server limits by consolidating up to 32 GPUs, simplifying scalability for demanding processing tasks

## INNOVATIVE AI MEMORY FABRIC

Optimizes AI applications with a memory-centric infrastructure, aggregating resources for faster, low-latency data transfer

## EMPOWERING LARGE MODEL PROCESSING

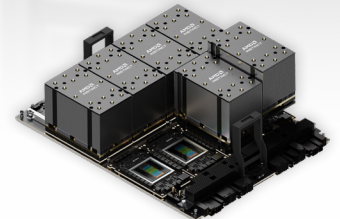
Vast GPU memory pool handles large model inference, reducing data transfers and boosting performance

## SIMPLIFIED AI DEPLOYMENT

By minimizing the complexity of AI and HPC infrastructure, organizations can quickly get their LLMs and applications up and running, without the usual delays associated with InfiniBand infrastructure setup

The GigaIO SuperNODE seamlessly presents all accelerators to a single server in an easy, scale-up configuration, enabling expansion without infrastructure complexity.

This scale-up architecture minimizes latency, enhances performance, and reduces power consumption by allowing models to efficiently utilize all GPUs as if locally connected, ensuring rapid scaling of AI workloads with minimal operational overhead.



## Specifications

<b>Accelerators</b>	32x AMD Instinct™ MI300X OAM GPUs
<b>Data rate (each direction)</b>	512Gb/s (Accelerator-to-Accelerator)
<b>FP64</b>	2.61 PFLOPS
<b>FP8</b>	83.52 PFLOPS
<b>GPU memory</b>	6.14TB HBM3
<b>Failover<sup>1</sup></b>	Primary and secondary servers
<b>CPU cores<sup>2</sup></b>	64 AMD EPYC™ “Genoa” 9534
<b>System memory<sup>2</sup></b>	3.0 TB
<b>Storage<sup>2</sup></b>	2x 15.3TB NVMe-U.2 (30TB total)
<b>Boot drive<sup>2</sup></b>	2x 960GB NVMe-M.2
<b>Network<sup>2</sup></b>	1x 400G QSFP112-DD 2x 25G/10G SFP28
<b>Rack Units<sup>3</sup></b>	42U
<b>Rack Power<sup>3</sup></b>	36.5kW
<b>Weight<sup>3</sup></b>	1,131 lbs (513.2 kgs)
<b>Cooling</b>	Air cooled, airflow front-to-rear (rack handle to power side)
<b>Environmental</b>	Operating Temperature: 10°C to 35°C (50°F to 95°F)



## READY TO GET STARTED?

Contact a GigaIO authorized representative today.

[info@gigaio.com](mailto:info@gigaio.com)

<sup>1</sup> Dual head node servers configured identically provide failover standby to maintain up-time

<sup>2</sup> Each server, based on max spec head nodes with 3.0 TB memory, 64 cores, 30 TB storage

<sup>3</sup> Does not include available top of the rack Ethernet switches, rack and PDU(s)