## GigaIO Introduces the First Ever 32 GPU Single-Node Supercomputer for Generative AI and Accelerated Computing

GigaIO SuperNODE™ with FabreX™ dynamic memory fabric delivers unprecedented compute capabilities with AMD Instinct™ accelerators for low-power, accelerated computing.

**Carlsbad, California, July 13, 2023** – GigaIO, the leading provider of workload-defined infrastructure for AI and technical computing, recently announced that it successfully configured 32 AMD Instinct MI210 accelerators to a single-node server utilizing the company's transformative FabreX ultra-low latency PCIe memory fabric. Available today, the 32-GPU engineered solution, called S[...] system capable of scaling multiple accelerator technologies such as GPUs [...] latency, cost, and power overhead required for multi-CPU systems.

As large language model applications demand even more GPU performan[...] the number of required node-to-accelerator data communications are cr[...] compute power at improved infrastructure TCO.

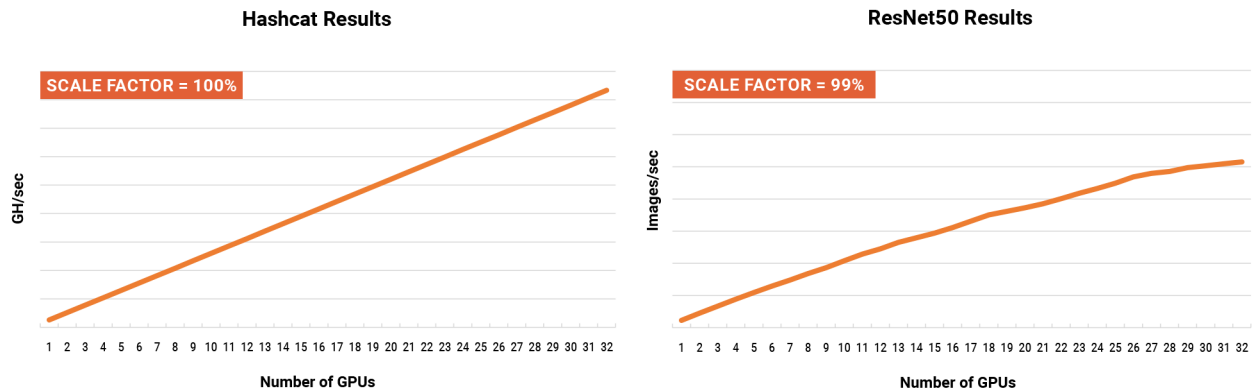"As AI workloads become more broadly adopted, systems that offer the a[...] power of multiple GPUs and better manage data saturation at ultra-low l[...] Nossokoff, Research Director, Hyperion Research. "And as large language[...] demand for more GPU performance, technologies that work to minimize[...] better positioned to provide the necessary performance for a robust AI in[...]

"AMD collaborates with startup innovators like GigaIO in order to bring u[...] workload demands of AI and HPC," said Andrew Dieckmann, corporate vi[...] manager, Data Center and Accelerated Processing, AMD. "The SuperNOD[...] and powered by AMD Instinct accelerators offers compelling TCO for bot[...] generative AI workloads."

GigaIO's SuperNODE system was tested with 32 AMD Instinct MI210 acce[...] server powered by dual 3rd Gen AMD EPYC™ processors.

- *Hashcat:* Workloads that utilize GPUs independently, such as HashCat, scale perfectly linearly all the way to the 32 GPUs tested.

- *Resnet50:* For workloads that utilize GPU Direct RDMA or peer-to-peer, such as Resnet50, the scale factor is slightly reduced as the GPU count rises. There is a one percent degradation per GPU, and at 32 GPUs, the overall scale factor is 70 percent.

These results demonstrate significantly improved scalability compared to the legacy alternative of scaling the number of GPUs using MPI to communicate between multiple nodes. When testing a multi-node model, GPU scalability is reduced to 50 percent or less. The following charts show two real-world examples of these two use cases:

**Hashcat Results**

SCALE FACTOR = 100%

GH/sec

1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32

Number of GPUs

**ResNet50 Results**

SCALE FACTOR = 99%

Images/sec

1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32

Number of GPUs

More testing results can be found here.

"This testing shows the enormous value of using GigaIO's SuperNODE to get all the benefits of composability, without any of the hassles," said Alan Benjamin, CEO & President, GigaIO. AMD and GigaIO engineered the entire hardware and software stack of the SuperNODE up to and including the TensorFlow and PyTorch libraries so that applications "just run" without any software changes. "Customers can scale GPU performance without the overhead of multiple servers using our FabreX software, and get unprecedented flexibility. When a large job needs results fast, 32 GPUs can be deployed on a single compute node simply and efficiently, with leadership low latency and power usage. Those same accelerators can then be easily and quickly reallocated to other servers, thus optimizing their utilization. Let the job define your system, and not the other way around," added Benjamin.

**About GigaIO**
GigaIO provides workload-defined infrastructure through its dynamic memory fabric, FabreX, which seamlessly composes rack-scale resources and integrates natively into industry-standard tools. FabreX lets customers build impossible servers for AI and technical computing — from storage to accelerators to memory — at a fraction of cloud TCO, by optimizing the utilization and efficiency of their existing hardware, allowing them to run more workloads faster at lower cost through higher utilization of resources and more agile deployment. Visit www.gigaio.com, or follow on Twitter and LinkedIn.

**Contact**
Danica Yatko
760-487-8395
danica@xandmarketing.com

*AMD, the AMD Arrow logo, EPYC, AMD Instinct, and combinations thereof, are trademarks of Advanced Micro Devices, Inc.*