# Dnevi SLING

Supermicro HPC Solutions

3.-5. december 2024

# Supermicro HPC Solutions

Dec 5th, 2024

Petr Karbus

Senior Sales Manager

# ABOUT SUPERMICRO

| | |
|---|---|
| **Revenue** | **$14.7B+** (FY2024 guidance)<br>$7.1B (FY2023)<br>$5.2B (FY2022) |
| **Worldwide Presence** | **6M+ Sq ft. Facilities Worldwide**<br>1. Silicon Valley (HQ),<br>2. Taiwan,<br>3. The Netherlands,<br>4. Malaysia and others |
| **Production** | **$25B/yr Production Capacity (CY24)**<br>Top 5 Largest Server System Provider Worldwide (IDC & Gartner 2022), ~1.3M units annually |
| **Human Resource in 4 Campuses** | ~6000 headcount Worldwide,<br>~50% Technical / R&D |
| **Key Growth Matrix** | **#1** in Generative AI and LLM Platforms<br>500%+ YoY Growth in Accel. Computing |

# Industry's Most Comprehensive Portfolio

## Rack Mounted

Multi Processor
8 x CPU Sockets

Multi Processor
4 x CPU Sockets

Hyper
2 or 1 x CPU Sockets

CloudDC
2 or 1 x CPU Sockets

WIO
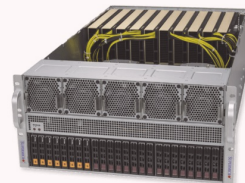1 x CPU Socket

## GPU

HGX
8 x GPU SXM

HGX DLC
8 x GPU SXM

PCIe GPU

MGX
Grace Hooper

## Multi Node

Superblade

MicroCloud

FlexTwin

GrandTwin

BigTwin

## Storage

All Flash

Front Load

Simply Double

Top Load

# OPTIMIZED RACK-SCALE TOTAL SOLUTIONS



**5,000+ Racks** per month global capacity

**2,000+ DLC Racks** per month

**100KW** (150KW) Racks ready to deploy

## One-Stop Total IT Solutions

### Hardware + Software + Services

- Rack-scale plug-and-play with optimized building block architecture

- Scalable compute, storage, network, infrastructure, cooling, software & service

- Free-air, liquid cooling (DLC) & liquid immersion technologies enable flexible deployment

- Lowest TCO & leading energy efficiency

Server software management solutions drive optimization and higher infrastructure security

# Supermicro Rack Integration Services provides a "one-stop-shop" for your data center needs

## Optimized and Lab Tested Components for Superior Performance

**Turn-Key Data Center**

**Accelerate Your Deployment**

**Professional Rack Level Design**

**Validation and Benchmarking**

Hi Throughput Switch Solutions for High Performance Connectivity

Value Optimized, Twin Technology or GPU Enabled 1U Server Solutions

Value Optimized, Twin Technology or GPU Enabled or High Capacity Storage 2U, 3U and 4U Server Solutions

Software Tested and Certified

High Quality, Well Organized Cabling for High Reliability and Maintenance

High Density Storage, High Performance Computing and GPU Blade Solutions

- Server
- Storage
- Network
- Software
- Cabling
- Power and Cooling
- Testing
- Benchmarking
- Full Rack Burn-in

# Supermicro Networking - Ethernet

## SSE-T7132S/SR (32 ports)

**Key Features**
- 32x 400Gbps Ethernet ports (QSFP-DD)
- SONiC Networking Operating System
- Fully shared packet buffering
- Redundant hot-pluggable power supplies
- 1U form factor ideal for spine/super-spine
- Regular and reverse airflow models

## SSE-C4632SB/SRB (32 ports)

**Key Features**
- 32 x100Gbps Ethernet ports (QSFP28)
- 1:1 Non-blocking connectivity
- 1U form factor for flexible installation
- Data-Center friendly - regular and reverse airflow models
- Hot-pluggable power supplies
- Broadcom Advanced Enterprise SONiC Switch Software pre-installed

## SSE-T8032S

**Key Features**
- 64x 400Gbps Ethernet ports
- Broadcom Advanced Enterprise SONiC Networking Operating System
- Fully shared packet buffering
- Redundant hot-pluggable power supplies
- 1U form factor ideal for leaf/spine/super-spine
- Regular airflow model

## SSE-SN3700-VS2

**Key Features**
- 32 x 200 Gbps Ethernet ports (QSFP56)
- Connectivity at different speeds with throughput of 12.8Tb/s
- Cumulus Linux Networking Operating System
- Fully shared packet buffering
- Best-in-class VXLAN scale
- Redundant hot-pluggable power supplies
- 1U form factor ideal for ToR super spine

# Supermicro Networking - Infiniband

## NVIDIA Quantum-2 QM9700 Series

Scaling out data centers with 400G InfiniBand smart switches.

### System Specifications

| | |
|---|---|
| **Performance** | 400Gb/s per port |
| **Switch radix** | 64 400Gb/s non-blocking ports with aggregate data throughput up to 51.2Tb/s |
| **Connectors and cabling** | 32 octal small form-factor pluggable (OSFP) connectors; passive or active copper or active fiber cable; optical module |

## NVIDIA Quantum-X800 InfiniBand Switches

Accelerate AI workloads with 800G InfiniBand.

### System Specifications

| | Q3200-RA | Q3400-LD | Q3400-RA |
|---|---|---|---|
| **Performance** | Two switches, each of 28.8Tb/s throughput | 115.2Tb/s throughput | 115.2Tb/s throughput |
| **Switch radix** | Two switches, each of 36 800Gb/s non-blocking ports | 144 800Gb/s non-blocking ports | 144 800Gb/s non-blocking ports |
| **Connectors and cabling** | Two groups of 18 OSFP connectors | 72 OSFP connectors | 72 OSFP connectors |

# Supermicro Solutions

## GPU Acceleration (AI/ML, HPC, Omniverse)



**Workload Sizes**

Extra Large — Large — Medium — Storage

HGX H100/H200, H100 NVL & H200 NVL

Grace Hopper Superchip          L40S

## Data Base & ERP



ORACLE

SAP

## Cloud & Virtualization



Canonical Ubuntu — Management Switches

kubernetes — Data Switches

Infrastructure Nodes

ceph

Red Hat — Cloud Nodes

Azure Stack HCI

# Supermicro Solutions - Storage

## All Flash



### SSG-122B-NE316R

1U front-loading all-flash storage server with 16 E3.S NVMe drives and PCIe 5.0


MINIO


VAST


WEKA

## Hybrid

### Storage SuperServer SSG-620P-E1CR24H



**Key Features**
- Dual socket 3rd Gen Intel® Xeon® Scalable processors, up to 72 Cores Per Node;
- 16 ECC DDR4-3200: LRDIMM/RDIMM;
- Dedicated PCIe 4.0 AIOM slot; 3 x PCIe 4.0 x16 Slots;
- Server remote management: IPMI 2.0 / KVM over LAN / Media over LAN per node;
- 24 3.5" Hot-swap SAS3/SATA3 drives, 4x Rear SATA/NVMe Slots, 2x SATA/NVMe M.2 (form factor: 2280);
- 5x 8cm hot-swap counter-rotate redundant PWM cooling fans;
- 1600W Redundant Power Supplies Titanium Level (96%);
- HW RAID support via Broadcom® 3908;

### Storage SuperServer SSG-640SP-E1CR90



**Key Features**
- 16 ECC DDR4-3200: LRDIMM/RDIMM;
- 3 x PCIe 4.0 x16 HHHL PCIe slots;
- 90 3.5"/2.5" Hot-swap SAS3/SATA3 drives, 2x Fixed slim SATA SSD, 2x NVMe M.2 (form factor: 2280 and 22110);
- 6 x 8cm hot-swap counter-rotate redundant PWM cooling fans;
- 2600W Redundant Power Supplies Titanium Level (96%);
- Drive Controller support via Broadcom® 3916 or 3616; Server remote management: IPMI 2.0 / KVM over LAN / Media over LAN;


SCALITY    OSNEXUS    Quantum.

Qumulo

# GPU Accelerated Workloads

ML/DL/Inference/Generative AI



HPC



Enterprise AI



Virtualization & Design



Content Delivery (GPT, Copilot)



Edge AI

# Supported By Supermicro

**SUPERMICRO**
**Confidential**

**NVIDIA**                    **AMD**                    **intel**

## Multi socket

**HGX**
H100
H200
*Coming* → B100
*Soon* B200

**CDNA3**
MI300X

**Gaudi3 UBB**

## PCIe

H100 NVL
L40S
L4

MI210

**Gaudi3 PCIe**

## CPU+GPU

**Grace Hooper**
*Coming* GH200
*Soon* → GB200

**CDNA3**
MI300A

# What GPU Fits The Best for Your Workload?

| Manufacturer | GPU Model | Architecture | DL Training & DA | DL Inference | HPC / AI | Omniverse / Render Farms | AI Video | Far Edge Acceleration |
|---|---|---|---|---|---|---|---|---|
| NVIDIA | H200 | Multi Socket | Best | Better | Best | | | |
| NVIDIA | H100 | Multi Socket | Best | Better | Best | | | |
| AMD | MI300X | Multi Socket | Best | Better | | | | |
| intel | GAUDI3 | Multi Socket | Good | Best | Best | | | |
| NVIDIA | H100NVL | PCIe | Best | Best | | | | |
| NVIDIA | L40S | PCIe | Good | Good | Good | Best | Best | |
| AMD | MI300A | CPU+GPU | | | Best | | | |
| intel | GAUDI3 | PCIe | | Better | Good | | | |
| NVIDIA | L4 | PCIe | | Good | | Good | Best | Best |
| NVIDIA | GH200 | CPU+GPU | Best | Best | Best | | | |

**Price-performance** comparison relative across each entire workload column. This chart should be used in conjunction with measured data for targeted workloads.

Best · Better · Good

# Why Use GPU for AI Workloads?

**1.Parallel Processing Power**: GPUs are designed to handle multiple tasks simultaneously, making them highly efficient for parallel computations. In deep learning, many operations (like matrix multiplications) can be parallelized, which GPUs excel at due to their architecture with numerous cores.

**2.High Performance**: GPUs are optimized for handling large amounts of data and performing complex calculations quickly. They can process thousands of arithmetic operations in parallel, significantly speeding up model training compared to CPUs.

**3.Deep Learning Framework Support**: Most deep learning frameworks (like TensorFlow, PyTorch, and MXNet) are designed to leverage GPU acceleration. They have libraries that automatically distribute computations across multiple GPU cores, maximizing performance.

**4.Memory Bandwidth**: GPUs have high memory bandwidth, allowing them to efficiently handle the large amounts of data involved in deep learning tasks. This helps prevent bottlenecks that can slow down training on CPUs.

**5.Specialized Architectures**: Modern GPUs often include specialized cores and features specifically tailored for deep learning tasks, such as Tensor Cores for accelerated matrix operations (e.g., in NVIDIA GPUs).

**6.Cost-Effectiveness**: GPUs can offer significant speedups in model training time compared to CPUs. This means that training large models or processing extensive datasets can be done more quickly, potentially reducing overall training costs in terms of time and resources.

# Building AI Infrastructure

*What is AI Infrastructure?*

- *Compute GPU nodes*
- *Fast Interconnect (Network)*
- *Supporting Sys (Storage, MNGM)*
- *Orchestration Tools (Cluster Management, Cloud & Virtualization)*

*What do we need for AI Infrastructure?*

- *Concept*
- *Planning*
- *Data Center*

# Building AI Infrastructure

**Concept. Why do you need it?**

- Tasks you solve
- Workloads you accelerate
- Monetization
- AI Factory? Cloud/Multitenancy? Hybrid?

**Data Center.**

- How much power do we have?
- Liquid or Air?
- Racks Layout
- HW Layout inside racks.

**Planning.**

- Deployment timeline.
- Selection of GPU models.
- Selection of Interconnect
- Selection of orchestration tools.
- Design

# Inside the xAI Colossus, Powered by Supermicro

## World's Largest Liquid-Cooled AI Cluster

- xAI Colossus Supercomputer features 6,144 Supermicro NVIDIA HGX 8-GPU 4U Liquid-Cooled Systems
- A multi-billion-dollar cluster, deployed in 122 days
- The basic building block for Colossus is the Supermicro liquid-cooled rack
- 8 4U servers each with 8 GPUs, for a total of 64 GPUs per rack, plus a CDU
- Supermicro's design is from the ground up to be liquid-cooled, and all from one vendor
- Runs on Ethernet, 3,6 Terabit per second each server

Featuring

# 6144

Supermicro 8-GPU 4U Liquid-Cooled Systems

**Supermicro NVIDIA HGX Systems**

**NVIDIA Spectrum-X Ethernet networking platform**

**Supermicro Liquid Cooling Total Solution**

*Coverage courtesy of @ServeTheHome* STH
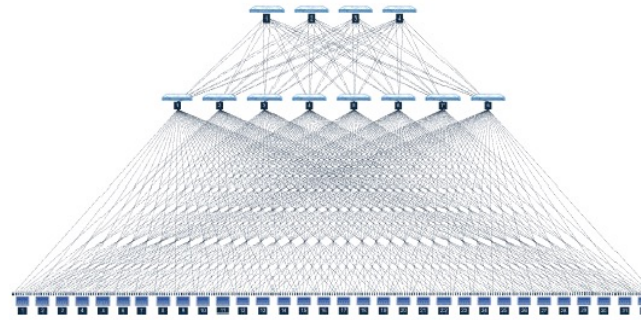
Storage Network Diagram

# AI Infrastructure Complete Solution from Supermicro

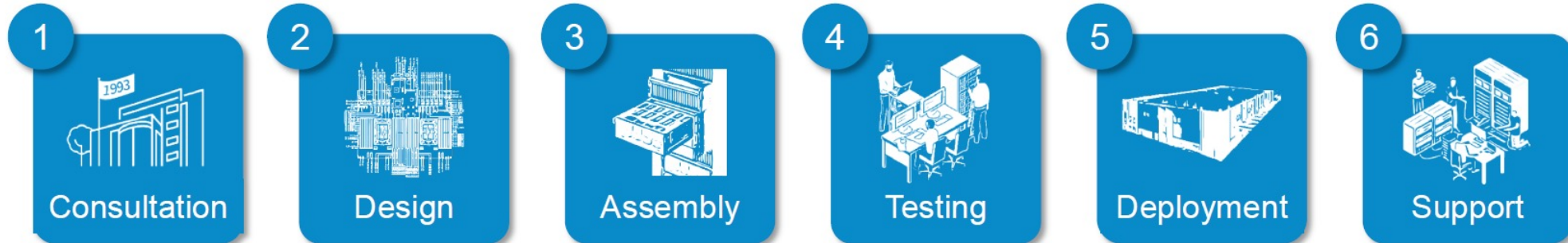## Supermicro Rack-Scale Advantage:

**Leverage Proven Building Blocks**

**Simplify Cluster-Scale Architecture**

**Deploy Plug & Play Racks**

## Rack Solution Design & Deployment Steps:

1. Consultation
2. Design
3. Assembly
4. Testing
5. Deployment
6. Support

**DISCLAIMER**

Super Micro Computer, Inc. may make changes to specifications and product descriptions at any time, without notice. The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions and typographical errors. Any performance tests and ratings are measured using systems that reflect the approximate performance of Super Micro Computer, Inc. products as measured by those tests. Any differences in software or hardware configuration may affect actual performance, and Super Micro Computer, Inc. does not control the design or implementation of third party benchmarks or websites referenced in this document. The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to any changes in product and/or roadmap, component and hardware revision changes, new model and/or product releases, software changes, firmware changes, or the like. Super Micro Computer, Inc. assumes no obligation to update or otherwise correct or revise this information.

SUPER MICRO COMPUTER, INC. MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION.

SUPER MICRO COMPUTER, INC. SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL SUPER MICRO COMPUTER, INC. BE LIABLE TO ANY PERSON FOR ANY DIRECT, INDIRECT, SPECIAL OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF SUPER MICRO COMPUTER, Inc. IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

ATTRIBUTION

# Thank You

www.supermicro.com

# Hvala za vašo pozornost

REPUBLIKA SLOVENIJA
MINISTRSTVO ZA VISOKO ŠOLSTVO,
ZNANOST IN INOVACIJE

EuroHPC
Joint Undertaking