# SUPERMICRO AND HABANA® HIGH-PERFORMANCE, HIGH-EFFICIENCY AI TRAINING SYSTEM

*Enables up to 40% better price/performance for Deep Learning training than traditional AI solutions*

**Supermicro SYS-420GH-TNGR**

**Habana Gaudi Processor**

## TABLE OF CONTENTS

## SUPERMICRO

As a global leader in high performance, high efficiency server technology and innovation, we develop and provide end-to-end green computing solutions to the data center, cloud computing, enterprise IT, big data, HPC, and embedded markets. Our Building Block Solutions® approach allows us to provide a broad range of SKUs, and enables us to build and deliver application-optimized solutions based upon your requirements.

## Executive Summary

Demand for high-performance AI/Deep Learning (DL) training compute has doubled in size every 3.5 months since 2013 (according to OpenAI) and is accelerating with the growing number of applications and services based on computer vision, natural language processing, recommendation systems, and more. With the increased demand for greater training performance, throughput, and capacity, the industry needs training systems that offer increased efficiency, lower cost, flexibility to enable customization and ease of implementation, and scaling of training systems. AI is becoming an essential technology for diverse areas such as virtual assistants, manufacturing operations, autonomous vehicle operations, and medical imaging, to name a few. Supermicro has partnered with Habana Labs to address these growing requirements. The Habana® Gaudi® AI processor is designed to maximize training throughput and efficiency while providing developers with optimized software and tools that scale to many workloads and systems.

Supporting Gaudi processor deployment is the Habana® SynapseAI® Software Platform, created with developers and data scientists in mind, providing versatility and ease of programming to address end-users unique needs while allowing for simple and seamless transition of their existing models over to Gaudi.

This solution brief provides an overview of the Supermicro X12 Gaudi AI Training system to enable high-performance computing for deep learning workloads.

## Supermicro X12 Gaudi AI Training System Overview

The Supermicro X12 Gaudi AI Training System (SYS-420GH-TNGR), powered by Habana Gaudi Deep Learning Processors, pushes the boundaries of deep learning training and can scale up to hundreds of Gaudi processors in one AI cluster. Gaudi is the first DL training processor with integrated RDMA over Converged Ethernet (RoCE v2) engines on-chip. With bi-directional throughput of up to 2 TB/s, these engines play a critical role in the inter-processor communication needed during the training process. This native integration of RoCE allows customers to use the same scaling technology, both inside the server and rack (scale-up) and across racks (scale-out). These can be connected directly between Gaudi processors or through any number of standard Ethernet switches.

With high compute utilization for GEMM and convolutions, low-power system design, and Bfloat16 support enabling FP32 accuracy with 16bit training speed, the Supermicro X12 Gaudi AI Training System is built to prioritize two key real-world considerations: training an AI model as fast as possible and reducing the cost of training. The system enables high-efficiency AI model training for vision applications such as manufacturing defects, resulting in better products with fewer warranty issues and fraud detection, saving billions of dollars annually. Inventory management is another area that benefits from AI technologies by allowing enterprises to become more efficient. Using AI technologies, medical imaging becomes more accurate and faster at detecting abnormalities, and identification from photos or videos can enhance security where needed. AI also enables language applications, including question answering, subject matter query, chatbots, translations, and sentiment analysis for recommendation systems, enhancing customer service organizations with more accurate and consistent knowledge.

## System Specifications

The 420GH-TNGR system contains eight Gaudi HL-205 mezzanine cards, dual 3rd Gen Intel® Xeon® Scalable processors, two PCIe Gen 4 switches, four hot swappable NVMe/SATA hybrid hard drives, fully redundant power supplies, and 24 x 100GbE RDMA (6 QSFP-DDs) for unprecedented scale-out system bandwidth. This system contains up to 8TB of DDR4-3200MHz memory, unlocking the Gaudi processors' full potential. The HL-205 is OCP-OAM (Open Compute Project Accelerator Module) specification compliant. Each card incorporates the Gaudi HL-2000 processors with 32GB HBM2 memory and ten natively integrated ports of 100GbE RoCE v2 RDMA.

April 2021

**SYS-420GH-TNGR**

System Rear View

System Front View

**System Features**

- High Density, 4U System for 8x Habana Gaudi HL-205 AI Processors
- Purpose Built for AI/Deep Learning Training
- Lower system cost with build-in 100GbE Ethernet ports
- 24 x 100Gb E RDMA (6 QSFP-DDs) for scale-out

**Processor Support**
Dual 3rd Gen Intel Xeon Scalable Processors, 3 UPI designed for up to 11.2 GT/s, Up to 300W

**Memory Capacity**
32 DIMM ECC DDR4 for up to 8TB memory

**System Management**
Two motherboard BMC and OAM BMC for system and module monitoring and debugging

**Mezz Card**
8 Habana Gaudi HL-205 OAM Mezz cards

**Expansion slots & I/O ports**
1 PCIe x16 LP from CPU, AIOMs Support, 1 x VGA, 2x USB 3.0

**Drive Bay**
4 hot-swap 2.5" drive bays (NVMe/SATA/Hybrid)

**System Cooling**
5 removable heavy duty fans

**Power Supply**
4 x 3KW (3+1) redundant power supply

**Figure1: SYS-420GH-TNGR Specifications**

Each of the Gaudi processors dedicates seven of its ten 100GbE RoCE ports to an all-to-all connectivity within the system, with three ports are available for scaling out for a total of 24 x100GbE RoCE ports per 8-card system. This allows end customers to scale their deployment using standard 100GbE switches. The high throughput of RoCE bandwidth inside and outside the box and the unified standard protocol used for scale-out make the solution easily scalable and cost-effective. The diagram below shows the Gaudi HL-205 processors and the communication paths between processors and the server CPUs.
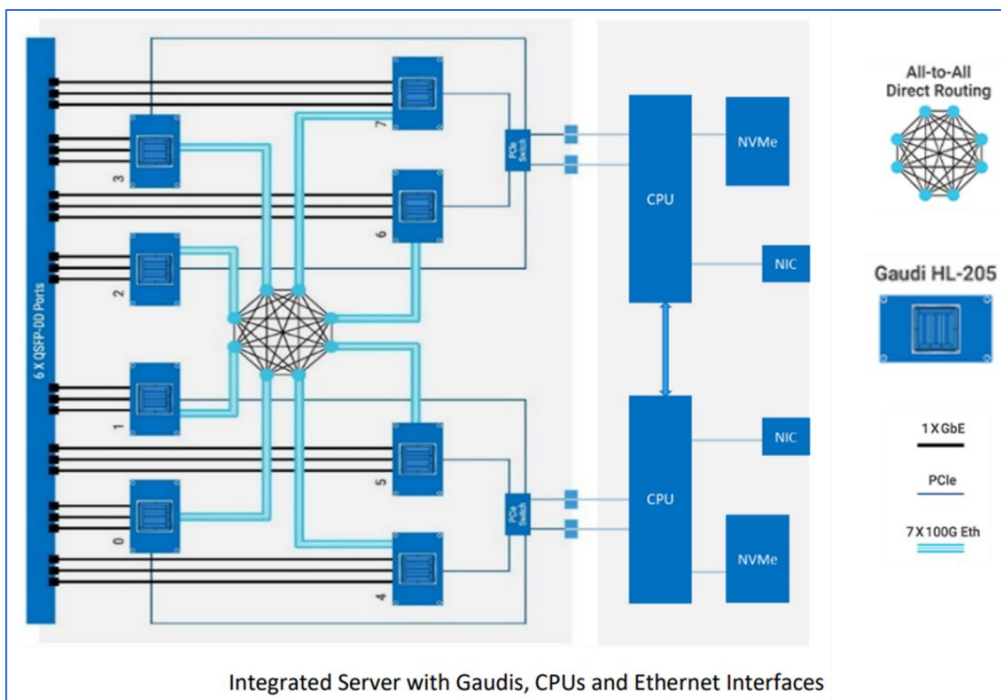


**Figure 2: SYS-420GH-TNGR Specifications**

The Figure below shows how a large scale, distributed training solution is built using the Gaudi AI System as a basic component with standard Ethernet connectivity. It offers three reduction levels – one within the system, another between 11 Gaudi AI Systems, and another between 12 islands. Altogether, this system hosts 8*11*12 = 1056 Gaudi cards. Larger systems can be built with an additional aggregation layer or with less bandwidth per Gaudi.
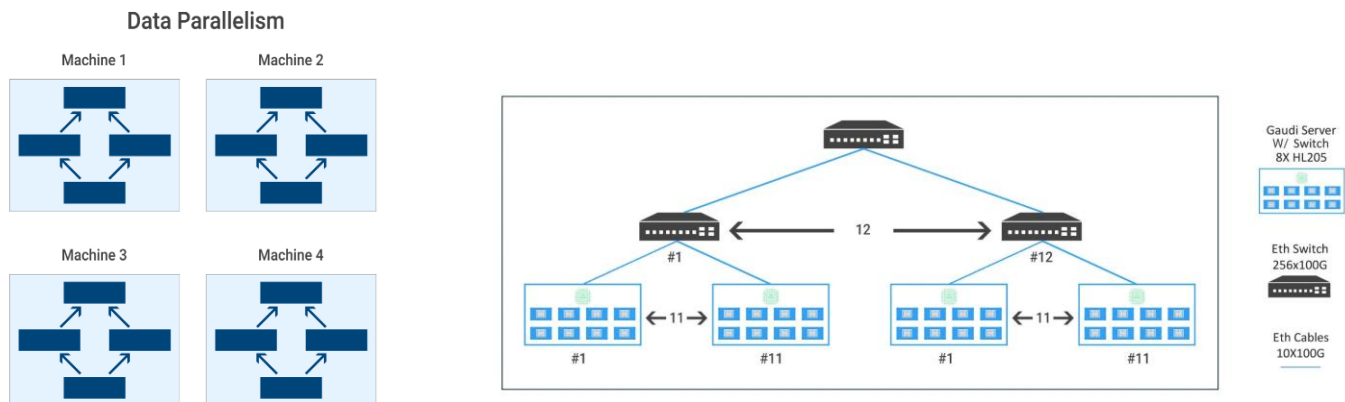


**Figure 3: Topologies for scaling Different Training Models**

## SUPERMICRO 420GH-TNGR SHOWCASES THE BENEFITS OF THE HABANA GAUDI AI TRAINING PROCESSOR

- High-performance, high efficiency
  - Price performance that enables more access to more end-customers to AI training
  - Performance at scale: high throughput at low batch size
  - Low system power
- First-of-its-kind scalability with native Ethernet scale-out
  - Avoids proprietary interfaces with industry standard Ethernet
  - Eliminates bottlenecks with integrated  NIC with built-in RDMA over Converged Ethernet (RoCE v2)
- Reduces system complexity, cost, and power with component integration
  - Leverages wide availability of standard Ethernet switches
- Solution flexibility and support for customization and ease of implementation
  - Habana SynapseAI Software Platform featuring
    - Programmable TPC and rich TPC kernel libraries
    - Software Infrastructure and Tools
    - Graph compiler and runtime
    - Support for popular frameworks and models
- Open Compute Project (OCP) Accelerator Module (OAM) compliance

April 2021

## ResNet-50 Performance

The Supermicro 420GH-TNGR with dual 3$^{rd}$ Gen Intel Xeon Scalable Processors and eight Habana Gaudi AI training processors has shown excellent performance and scalability when running ResNet-50 in the TensorFlow framework. Below is a chart showing how, as the number of Habana Gaudi AI training cards is increased, the throughput, as measured in the number of images per second, increases with near linear scale. When compared to other systems on a price and performance basis, this solution shines.
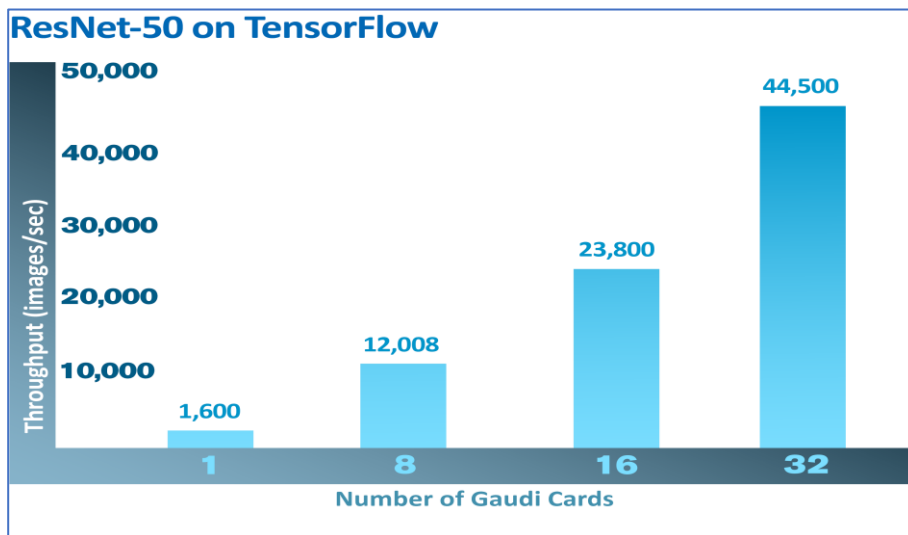


**Figure 4 - Performance as Gaudi Cards are Increased**

## Case Study: Large scale, distributed AI training at San Diego Supercomputer Center

The San Diego Supercomputer Center (SDSC) is a leader and pioneer in high-performance and data-intensive computing, providing cyberinfrastructure resources, services, and expertise to the national research community, academia, and industry. Located on the UC San Diego campus, SDSC supports hundreds of multidisciplinary programs spanning a wide variety of domains, from astrophysics and earth sciences to disease research and drug discovery.

The National Science Foundation (NSF) has awarded the San Diego Supercomputer Center (SDSC) at UC San Diego a grant to develop a high-performance resource for conducting artificial intelligence (AI) research across a wide swath of science and engineering domains. Called Voyager, the system will be the first-of-its-kind available in the NSF resource portfolio. As part of their mission to develop algorithms for these domains, SDSC required a cost effective yet powerful system to accelerate AI Training algorithms. SDSC chose the combination of Supermicro Intel Xeon-based CPU servers and the Habana Gaudi AI training systems. When complete, the Voyager system will contain over 42 Supermicro X12 Gaudi AI Training Systems, 336 Habana Gaudi processors, and 16 Habana Goya processors for inference.

Supermicro has worked closely with Habana and SDSC to create a large-scale AI training system. An efficient and performance-oriented system can be implemented based on the users' needs by selecting the right components, as shown below. Expertise in many areas, combined with a wide range of server and storage hardware, enables Supermicro to deliver custom solutions

based on industry standards. Supermicro builds the systems from the ground up with a complete manufacturing facility, from board testing up through multi-rack configuration and testing. This results in more confidence and fewer issues when delivered to innovative customers.

The large and efficient Voyager installation for AI training and inference consists of multiple racks of the Habana Gaudi AI Training systems, a Goya Inference rack, storage nodes, and an appropriate rack of networking equipment to support the high speeds necessary for the very fast training scale-out system and inference, as this diagram depicts.
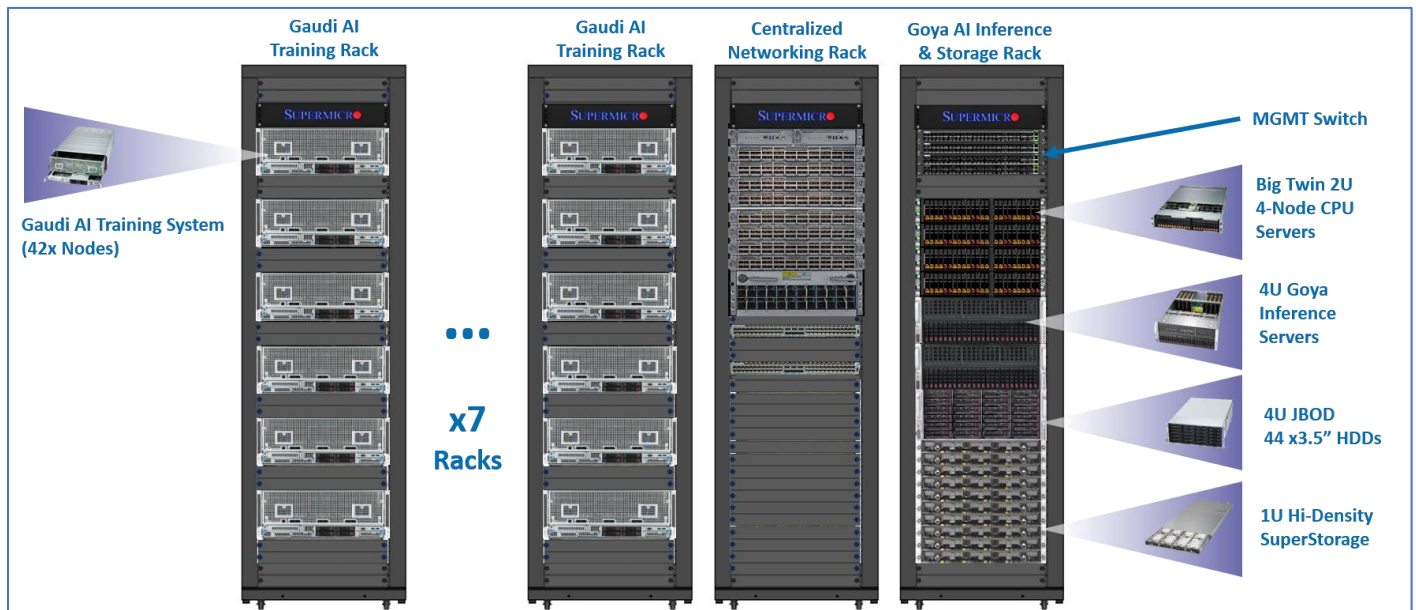


**Figure 5 – Supermicro Large Scale Gaudi Training System example**

## Summary

Deep Learning revolutionizes computing, impacting enterprises across multiple industrial and research sectors. The computational complexity of deep neural networks is becoming exponentially greater, driving massive demand for compute power. The challenge of deep neural network training is to improve upon multiple criteria at once: first, to complete the job faster with reduced training time; second, to achieve improved price/performance, thus lowering total cost of ownership to enable access to more AI training to more end-users; third, to reduce overall system energy consumption; fourth, to provide flexible scalability with standard interfaces that eliminate vendor lock-in; and lastly, to enable end-customers to customize workloads to address their unique needs.

Supermicro has partnered with Habana Labs and Intel to deliver a high performance and cost-effective system for AI training, coupled with the 3rd Gen Intel Xeon Scalable Processor. The combination of the expertise of Supermicro in systems design and Habana Labs with AI processor design will enable high-performance training at reasonable prices and make AI training more accessible to a wide range of industries. The Supermicro X12 Gaudi AI Training System provides superior scalability and has been substantiated by demanding AI customers. By designing, configuring, and testing, customers can be confident that the clusters are operational and optimized for their intended use.

The 420GH-TNGR presents several key advantages over traditional AI training solutions:

- Performance leadership that results in significantly lower training time, improved price/performance efficiency, and lower system size
- High throughput at small batch size and near linear scaling
- Low overall system power
- Native integration of Ethernet for scaling
- Lower system cost through wide availability of Ethernet switches of any size from multiple vendors
- Large scale distributed training with near linear scaling
- Support for popular frameworks and models, with the Habana Synapse AI Software Platform
- Open Compute Project (OCP) Accelerator Module (OAM) compliant

The 420GH-TNGR provides organizations with an opportunity to lower the total cost of ownership (TCO) and a flexible and easy path to scale their systems as they grow and evolve.

**Additional Resources:**

Supermicro X12 Gaudi AI Training System
Habana Training Site
Habana Developer Platform