



Exploring Supermicro AI Solutions Inside Out

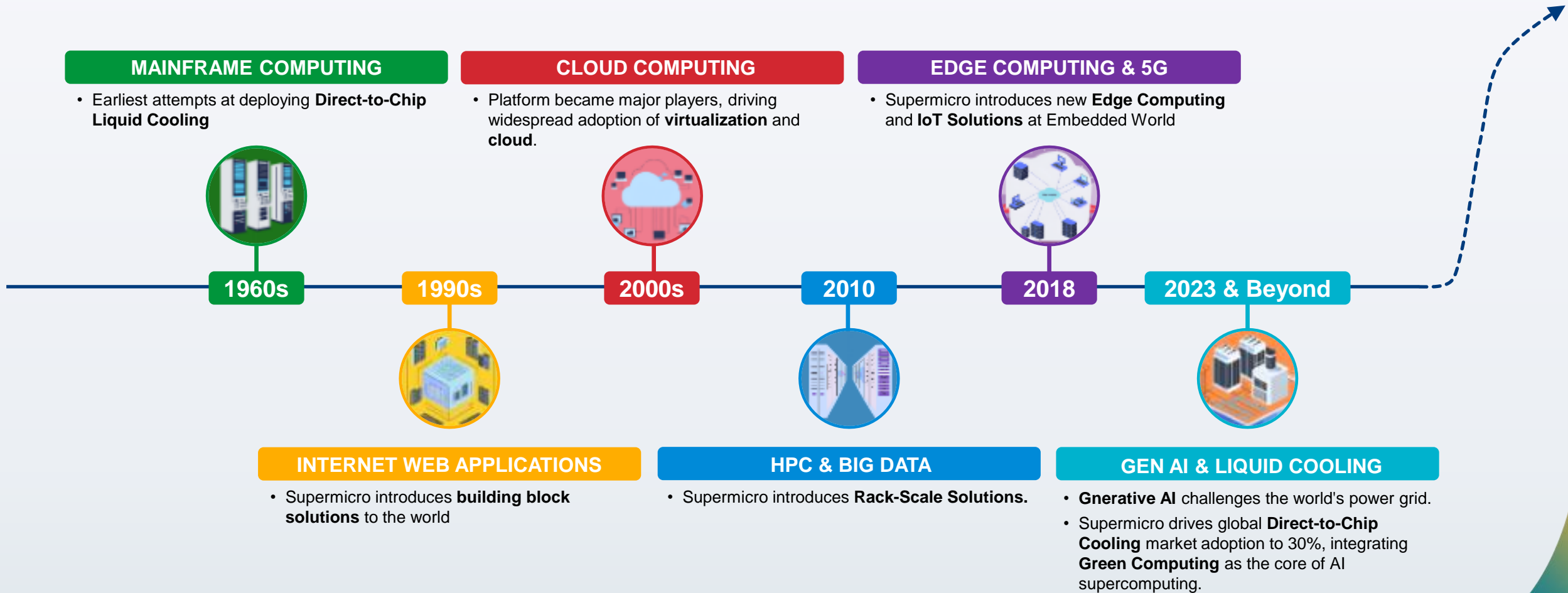
Speaker: **Hendry Chang**, Senior Product Manager



Agenda

- 01** AMD Product Portfolio Lineup
- 02** Liquid Cooling Total Solutions
- 03** Software Management Solutions (SCC)
- 04** GPUaaS & AI Cloud

Industry Shifts: How do we design our Data Centers?



Supermicro H14 Systems with AMD EPYC™ 9005/9004 Series Processors



Hyper DP

Flagship Power and Flexibility
Rackmount Server



CloudDC with DC-SCM

Ultimate Scalability and Flexibility for
Cloud Data Center



GrandTwin®

High Density Combined with I/O
Flexibility



FlexTwin™

Max Density, High-Performance
Multi-Node System with Liquid
Cooling



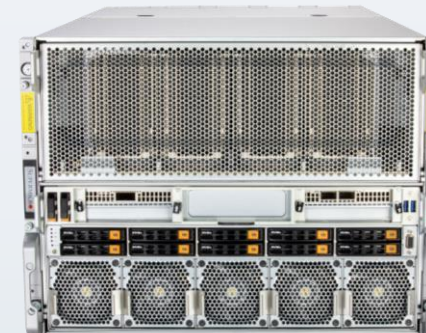
5U PCIe GPU System

GPU-Dense Servers Optimized for
Compute Intensive Workloads



4U Liquid cooling Universal GPU System

For AI/ Deep Learning and HPC MI325X



8U Universal GPU System

Next Generation MI325 Machine
Learning Platform



10U GPU System

Superior AI Applications and
Large Language Models

Advantage of Supermicro & Instinct™ MI300X

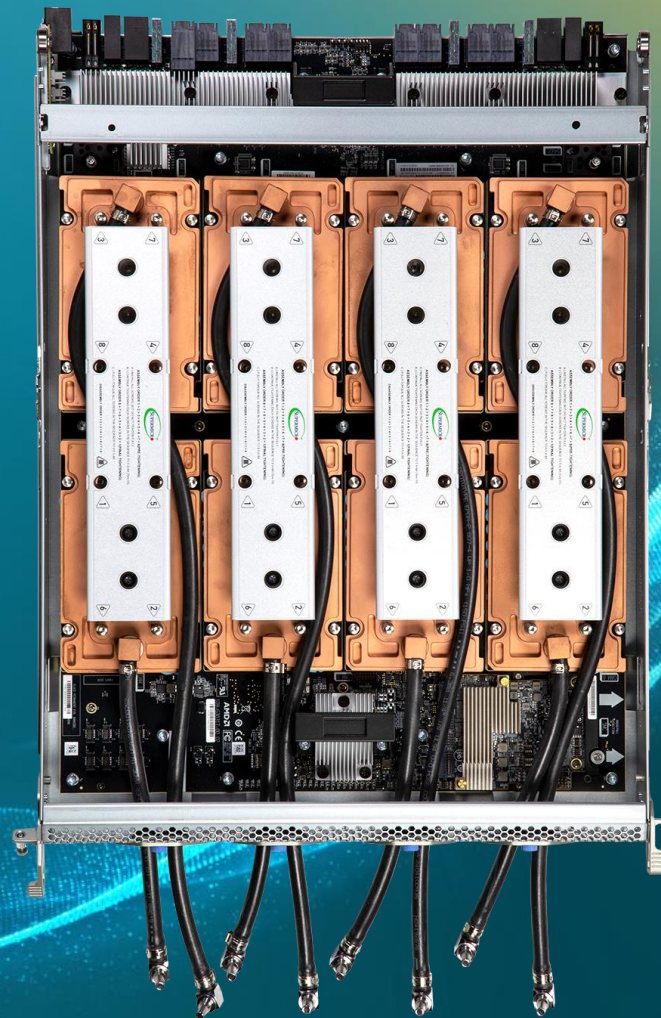
Supermicro Universal GPU Platform

- AI Training & Inferencing
- Front Access Design for Easy Maintenance
- Maximum Density for Space Saving
- Direct-to-Chip End-to-End Single Vendor Solution

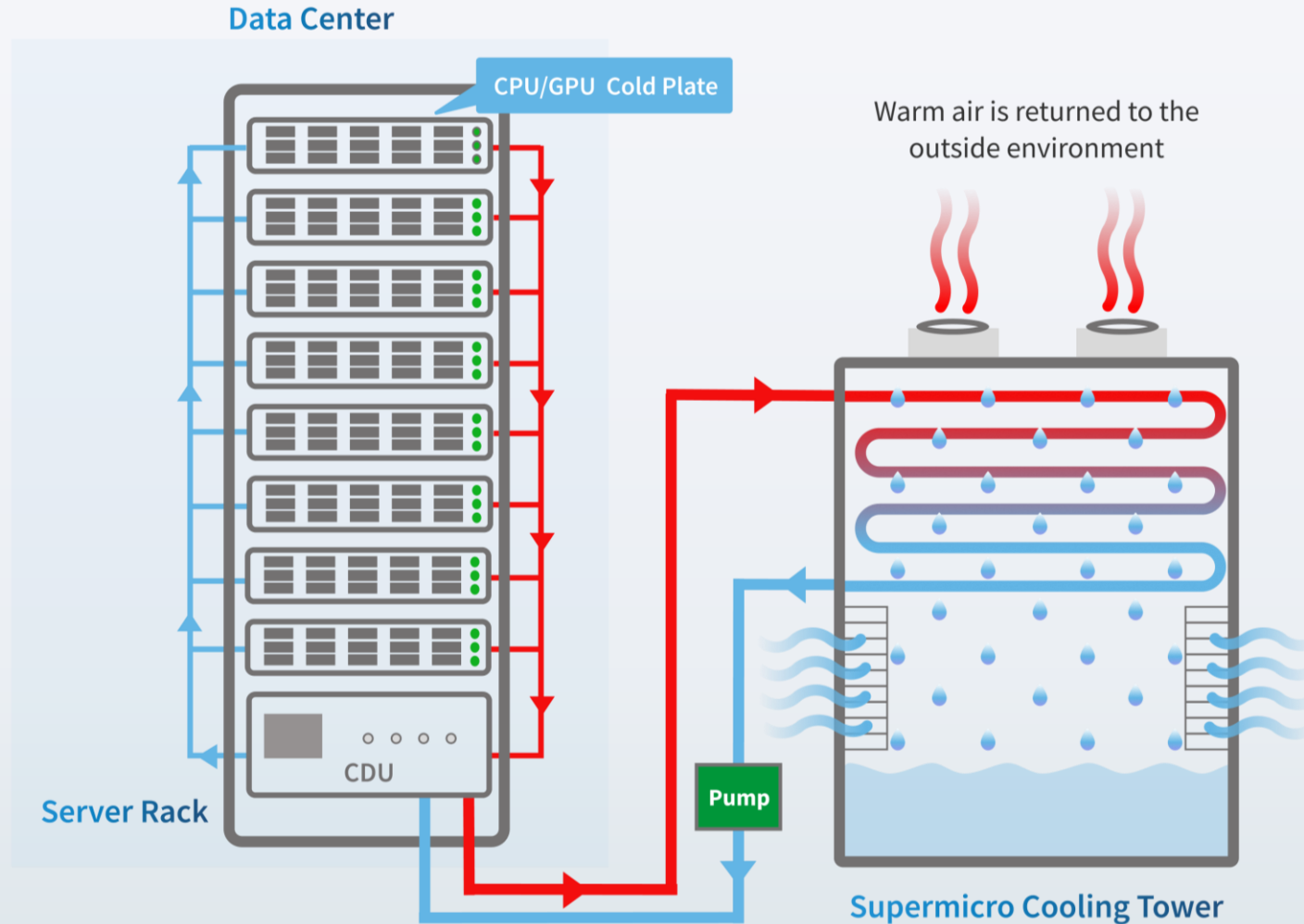


Scan to learn more system detail!

AS-4125GS-TNMR2-LCC with 8-MI300X OAM Baseboard Shelf



Liquid Cooling Overview



How Direct-to-Chip Can Help

DLC Blooming AI Infra to Peak

- Allows GPUs to Reach Maximum Performance
- Lower Electricity Costs
- Lower Acoustic Noise
- Buy more, Save more!

Liquid Cooling Advantages



REDUCTION

Electricity Costs of Cooling Infrastructure in Server



REDUCTION

in Electricity Costs for Entire Data Center

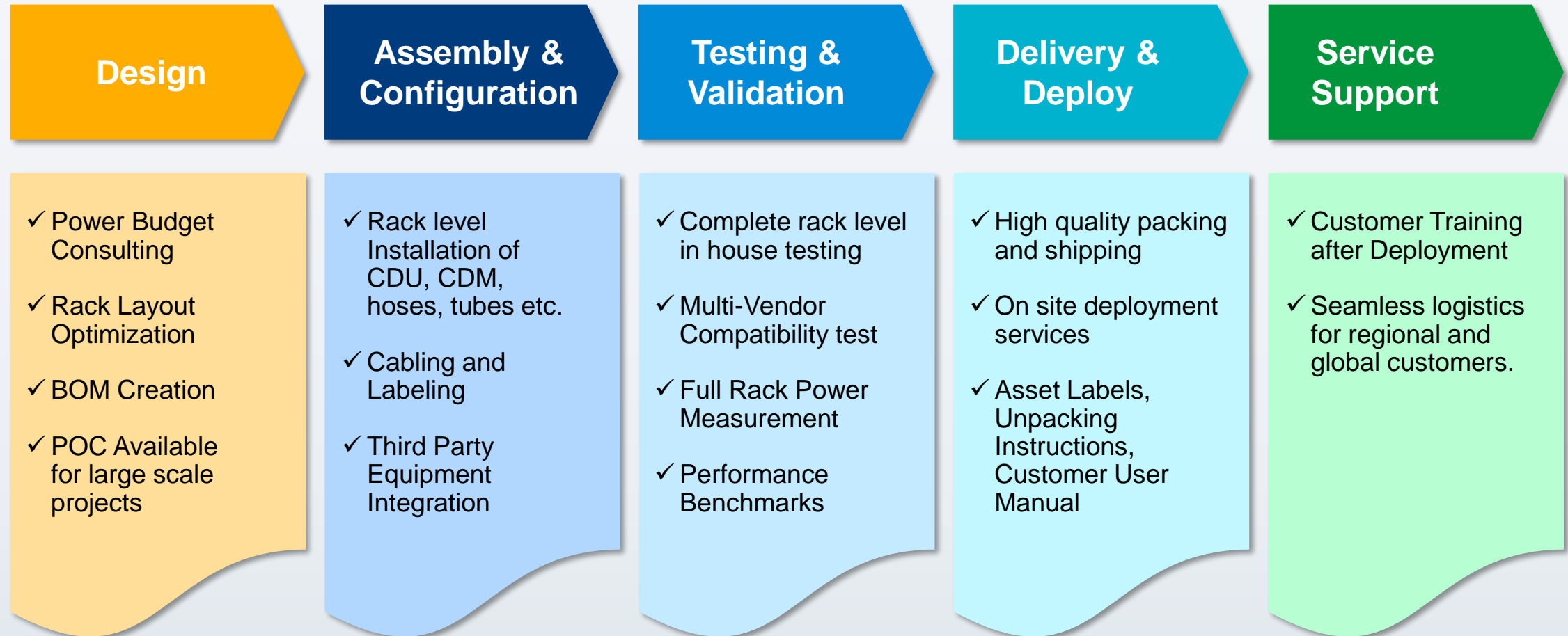


REDUCTION

in Data Center Server Noise



Liquid Cooling End-to-End Single Vendor



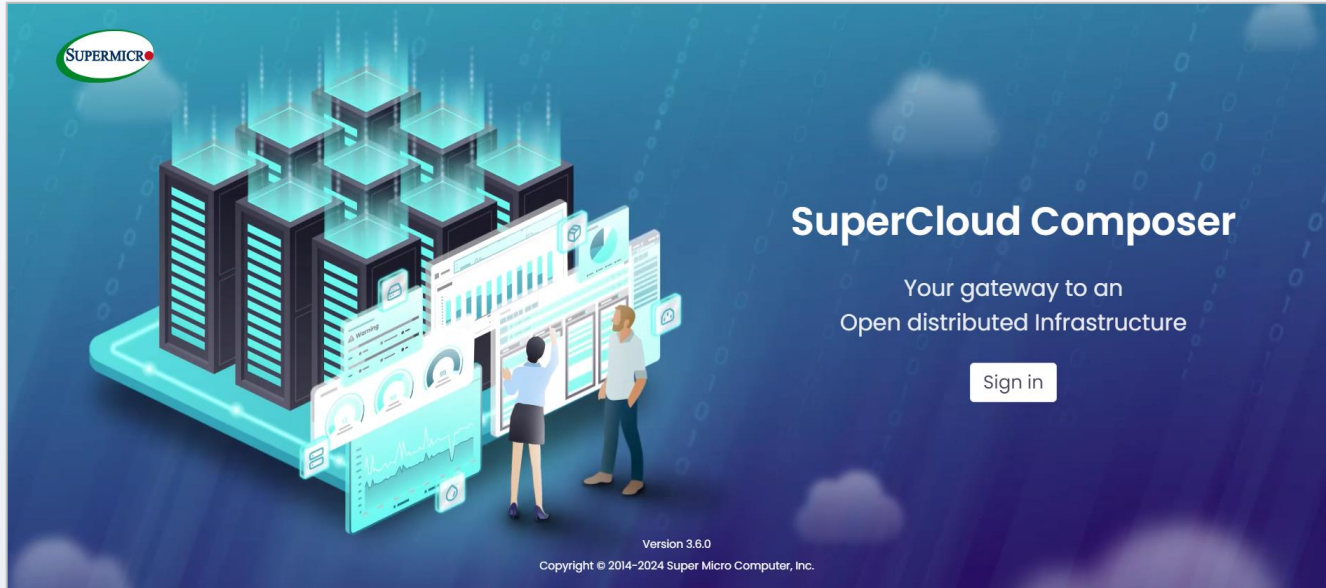
Too Many Systems to Manage?



Supermicro's Got Your Back!



SuperCloud Composer (SCC)



Advantages:

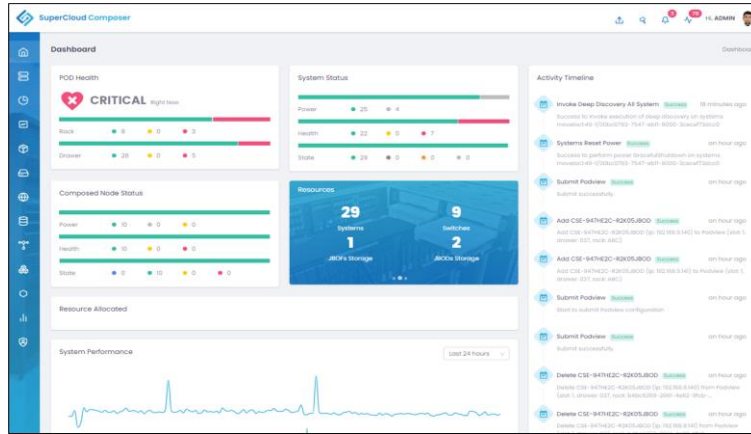
- GUI (Graphic User Interface) for easy server configuration
- Monitors Supermicro servers & 3rd party systems
- Server Monitoring with Historical Data
- Components to Cooling Tower Level Health Alert
- Reliable Scalability for servers of all sizes



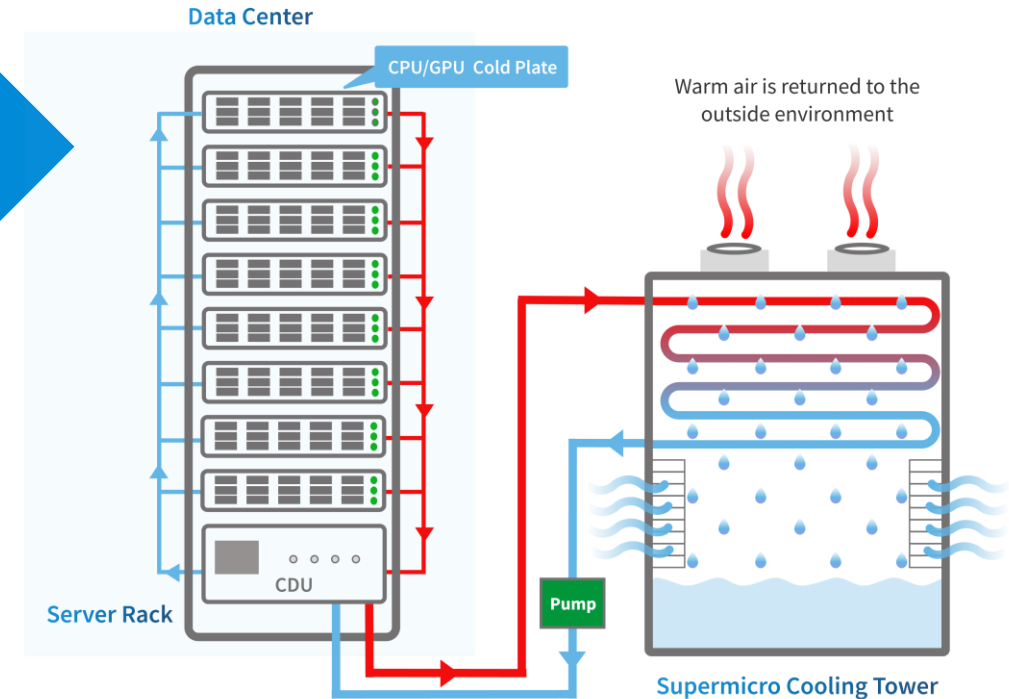
SCC Liquid Cooling Management

SCC

IT Admin



Supermicro Total Liquid Cooling Solution



SCC Features:

- Liquid Cooling Rack Space Management
- CDU and Cooling Tower Physical Assets Management
- CDU and Cooling Tower Live Sensor Reading with Historical Data



Supermicro AI Cloud Platform Solution

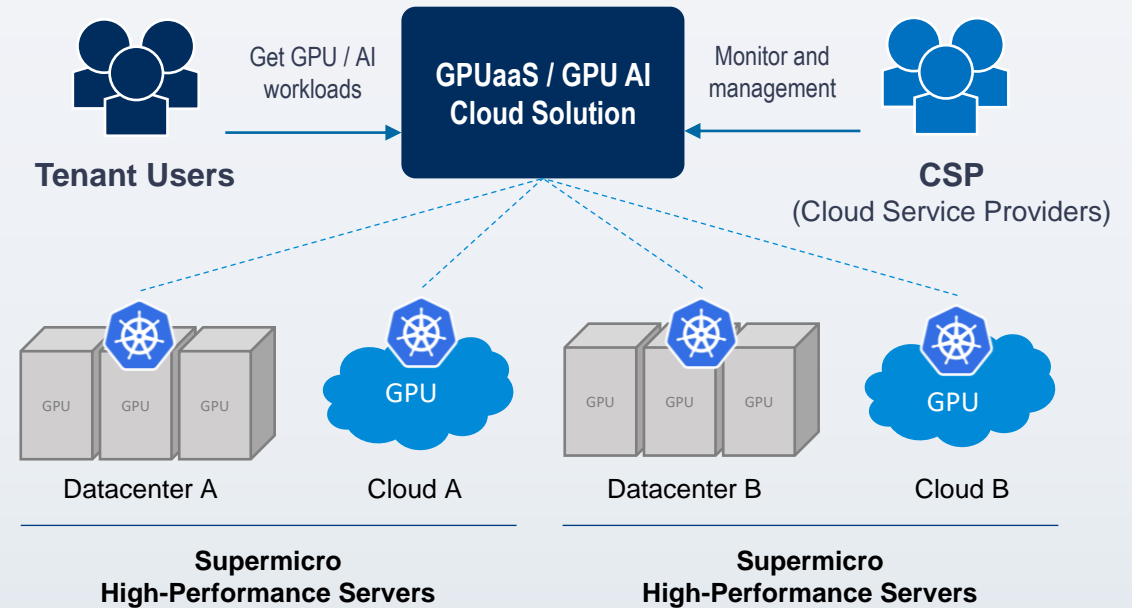
Speed-up your GPU / AI Cloud Business

Enables Cloud Service Provider(CSPs) to Sell Compute Resources on Pay-as-You-Go Basis

- 👍 Efficient GPU allocation
- 👍 Unified containerized resource pooling
- 👍 Multi-tenant flexible shared/dedicated policy
- 👍 Streamlined GPU cloud operations

AI DEVELOPMENT

GenAI INFERENCE



What is GPUaaS and AI Cloud?

GPUaaS

Provides on-demand access to GPU resources over the cloud.

Use Cases:

- AI / ML infrastructure provision
- High-performance data analysis

Customer Key Requirement:

- Scalability and Flexibility
- Cost Efficiency
- Multi-Tenant Architecture

AI Cloud

Specifically designed to support AI workloads using GPUs.

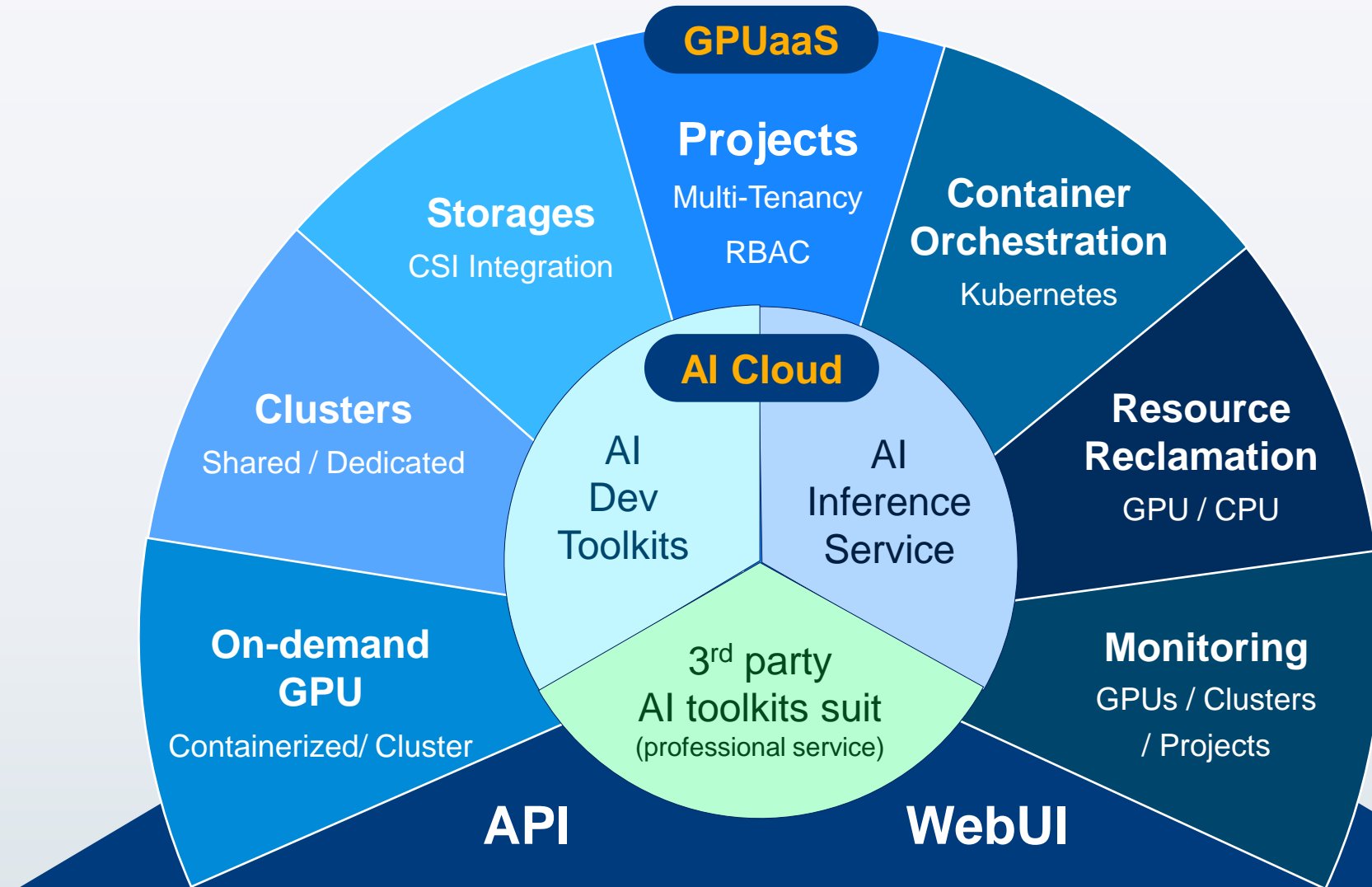
Use Cases:

- AI model training
- AI inferencing service

Key Benefits to End-user:

- End-to-End AI Environment
- Integrated AI Tools
- Generative AI & LLM Scalability
- All the GPUaaS requirements.

Supermicro GPUaaS / AI Cloud Platform Solution



Success Story

Taiwan's Largest CSP for AI: NCHC

TWCC TAIWAN COMPUTING CLOUD

TAIWAN AI ENABLER
Taiwan Computing Cloud - Where the future is

- Rapid & Lightweight Deployment**
- Effective Command & Control**
- Smart & Speedy Computations in One**
- Secure Data Repository**

News

Architecture Diagram:

- VM Container, GPU Container, AI HPC
- Admin Portal, ASUS Portal, API Gateway, CLI
- GOC Pass, Slurm
- TensorFlow, Caffe, Torch, DIGITS
- VM, K8S, Docker / Nvidia-Docker
- Nova Scheduler, OpenStack, Kubernetes, Singularity
- CentOS 7, SLES 12 SP3, CentOS 7
- CPU Node, GPU Node
- SES, RBD, NFS-Ganesha, RGW, GPFS
- IBM TSM, Mellanox Cumulus / NSX-T, Eth, IB

Well-Known Medical Institution Private GPU Cloud

北榮一號雲 (VGHTPE No.1 Cloud)

iThome 新聞 產品與技術 專題 AI · Cloud · 醫療IT 資安 · 研討會 · 社群 · IT EXPLAINED 搜尋

臺北榮總自建超級電腦北榮一號雲, 採32張A100 GPU要加速發展智慧醫院

臺北榮總昨日揭曉自建超級電腦「北榮一號雲」, 由4節點共32張A100 GPU組成一層無邊際雲平台, 要加速建構醫療和醫學發展, 如基因序列資料分析、數位病理和智能藥物開發等, 未來將擴展應用程度來補充超級電腦。同時, 北榮也展開新設立的智慧科技加裝房, 配置了AI防護交際層與解鎖新設備。

文/王郁傑 | 2023-12-07 9:58

與您一起邁進 智慧醫院的共榮路

臺北榮總昨日(12/6)亮相自家超級電腦「北榮一號雲」, 由4節點共32張A100 GPU組成, 要來加速北榮的精準醫療和醫學發展, 如基因序列資料分析、數位病理和智能藥物開發等, 同時, 北榮一號雲應用來訓練醫院專屬的大型語言模型(LLM), 要在安全環境中, 提供生成式AI應用給醫院使用者。接下來, 他們將規畫超級電腦算力用量, 來持續擴充。另一方面, 臺北榮總也揭曉其他智慧醫療成果, 如啟用智慧科技加裝房, 重粒子中心完成100例患者治療, 以及遠端西手術業務完成5,000例。

Enterprise Private GPU Cloud for Smart Manufactory

Vanguard International Semiconductor Corporation (VIS)

