Open. Together.



John Hu

Vice President Inspur Group, Inspur





inspur

Open Innovation Driven by AI. Dynamic, accelerated solutions for the open data center

AI for Open Rack and OCP Hardware

John Hu VP/CTO, Inspur Group



Open Platform, Open Eco-system



Embrace Open Sourcing

Innovate Open Technology

JDM for Open Design

SR-AI Rack: 1st Rack Server with Resource Pooling for AI Compute



Flexible GPU pooling solution with 1-4 GPU nodes; designed to run various AI applications such as AI cloud, deep learning for training, and image recognition.



PCIe Fabric Interconnected Infra GPU/FPGA/MIC Compatible 54 x GPU/FPGA/MIC Expandability 512TFlops Peak Performance 2-4X Comm Bandwidth Boost 50% Reduced comm latency

Resource Pooling Infrastructure to Enable Equipment's Physical Layer Elasticity and Optimize Cost **\$** 31%

Decoupling Compute and Storage Pooling
Flexible Configuration



Efficiency Improvement



High-Density Cloud-Optimized Platform by Inspur



inspur



- A lead cloud-optimized platform design partner with Intel
- Cloud workload optimization system solutions
- Joint GTM and 1st High-density Cloud-optimized Platform Contributor to open source community



High-Density Cloud-Optimized Platform

NF8260M5 Flexible High Density 2U 4-Socket Server

Validated for the future generation Intel® Xeon[™] Scalable Processor (Codename Cascade Lake), optimized with Intel® Optane[™] DC persistent memory

Designed and optimized for cloud Infrastructure-aaS, Function-aaS, and Bare-Metal-aaS solutions

Operational Efficiency

- High density 4-socket design offers double digit % TCO savings compares to 2-socket in same generation
- Offset processor placement for efficient cooling

23

 Front hot-swap accessibility U.2 drives, PCIe and OCP 2.0/3.0 module choice

Power Saving

	Fan Power Saving	System Total Saving	Power Saving over 3 Years
<u>2 * 2S</u> vs <u>1 * 4S</u>	30watt	87watt	87watt * \$6 = \$522

inspur

Increased VM Density, Cloud-Optimized Performance

2U2S in TWO 42U Rack

2U4S in ONE 42U Rack



Additional saving on Network Fabric & Boot Device will further lower the TCO

New JBOG Contribution to OCP

- Contribution and OCP Accepted recognition In Process
- New contribution to OCP for 2nd choice of 8 GPU box solution (XM2)
- Two new topology contributed to OCP for different AI applications



Advanced New Architecture for AI / First 16 GPU Box to Market

4-Socket Olympus: Contribution and OCP Accepted recognition In Process



• First 16 GPU Box launch to market

- Most powerful in computing and GPU combined for deep learning

 80 CPU cores : 16 GPUs
- Capability to maximize throughput with multi workload

 Performance increase 20%-60% under various workload with 16 GPU

> Inspur 4-Socket + GPU Box (16 GPUs)

Most Efficient 32 GPU Topology

4-Socket Olympus and 32 GPU

Inspur 4-Socket + 2*16GPU Box (32 GPUs) UPI CPU 0 CPU 1 **4-Socket Olympus** JUPI UPI **Contribution and OCP** CPU 2 CPU 3 UPI Accepted recognition In Process 1.1 PCIe x8 PCIe x16 Slot 100G NIC PCIe x16 00G NIC PCIe x16 .00G NIC PEX 9797 Base Mode PEX8713 PEX8713 1111 HGX-2 Base board 0 HGX-2 Base board 1 HGX-2 Base board 0 HGX-2 Base board NF5888M5 with 16*GPUs NF5888M5 with 16*GPUs

Most efficient 32 GPU topology vs. InfiniBand

OCP GLOBAL SUMMIT 2019

Ø,

Open Together

OCP Community

OCP OAM Project

Embrace open collaboration Commit to OCP community Collaborate on the module Plan in place to support baseboard and chassis design

SPEC_ML

Establish unified & credible benchmarks for ML

Chair the committee with 12 other members







As a top 3 worldwide server vendor, Inspur is fully committed to "Open Together".



Moving You Forward