**Hewlett Packard Enterprise**
Operated by Selectium

Global and regional leader
# HPE's solutions for HPC and AI

Đorđe Ristić, Milan Avramović
May 2025

# The World's Only Three Verified Exascale Supercomputers

All exascale supercomputers are built on HPE Cray Supercomputing EX – 100% Liquid cooled

## #1

LLNL's El Capitan supercomputer is #1 on the TOP500

**1.742 exaflops** of FP64 performance[1]

**10 952** liquid-cooled APUs

## #2

ORNL's Frontier supercomputer is #2 on the TOP500

**1.353 exaflops** of FP64 performance[1]
**10.2 exaflops** of HPL-MxP

**9 604** liquid-cooled CPUs
**38 416** liquid-cooled GPUs

## #3

ANL's Aurora supercomputer is #3 on the TOP500

**1.012 exaflops** of FP64 performance[1]
**10.6 exaflops** of HPL-MxP

**21 248** liquid-cooled CPUs
**63 744** liquid-cooled GPUs

## HPE Cray Supercomputing EX broke the exascale barrier for the THIRD time!

· Source:Top500.org, November 2024 list

# The European fastest Supercomputers
## First three supercomputers in Europe are built on HPE Cray Supercomputing EX

**#1**

ENI's HPC6 supercomputer is #5 on the TOP500

**477 petaflops** of FP64 performance[1]

**3 330** liquid-cooled CPUs
**13 320** liquid-cooled GPUs

**#2**

CSCS's Alps supercomputer is #7 on the TOP500

**434 petaflops** of FP64 performance[1]

**10 400** liquid-cooled CPUs
**10 400** liquid-cooled GPUs

**#3**

CSC's LUMI supercomputer is #8 on the TOP500

**379 petaflops** of FP64 performance[1]
**2.35 exaflops** of HPL-MxP

**2 919** liquid-cooled CPUs
**11 664** liquid-cooled GPUs

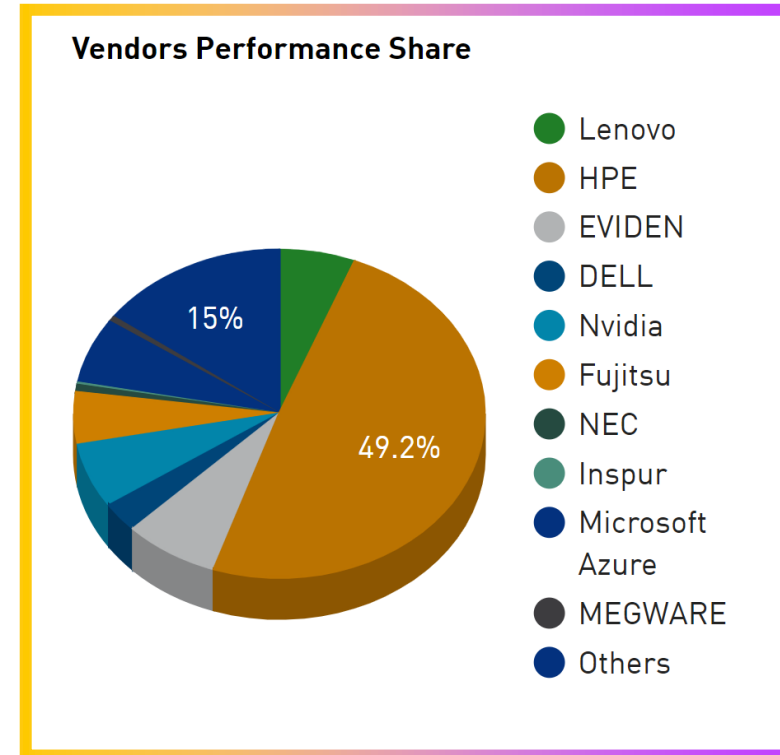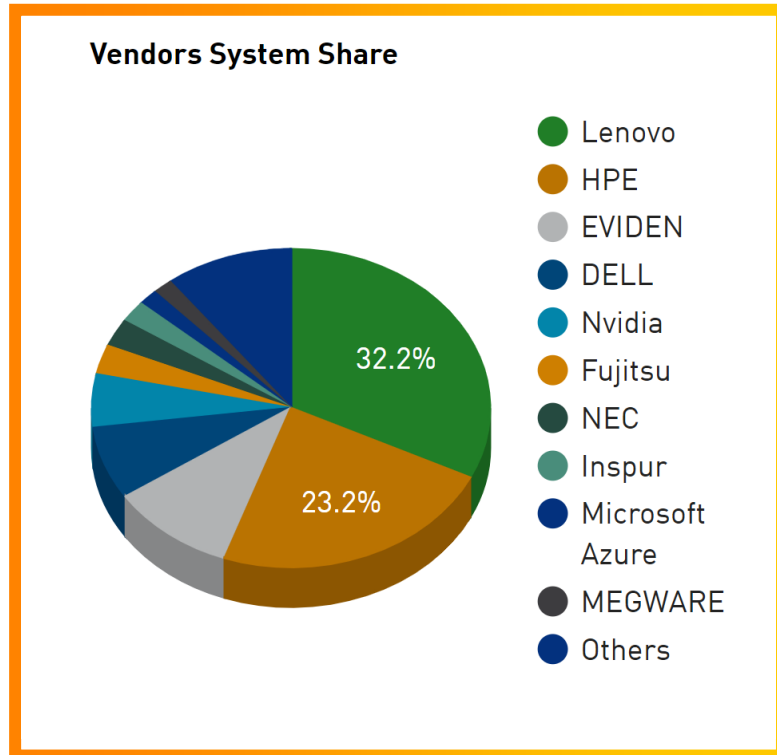**HPE Cray Supercomputing EX deliver 1.5 exaflops of aggregate performance in Europe**

Hewlett Packard Enterprise
Operated by Selectium

· Source:Top500.org, November 2024 list

4

# Global Leader in Supercomputing

Providing our customers with performance they require



**Vendors System Share**

- Lenovo — 32.2%
- HPE — 23.2%
- EVIDEN
- DELL
- Nvidia
- Fujitsu
- NEC
- Inspur
- Microsoft Azure
- MEGWARE
- Others

**Vendors Performance Share**

- Lenovo
- HPE — 49.2%
- EVIDEN
- DELL — 15%
- Nvidia
- Fujitsu
- NEC
- Inspur
- Microsoft Azure
- MEGWARE
- Others

**HPE Cray Supercomputing EX systems delivers outstanding performance**

Source:Top500.org, November 2024 list

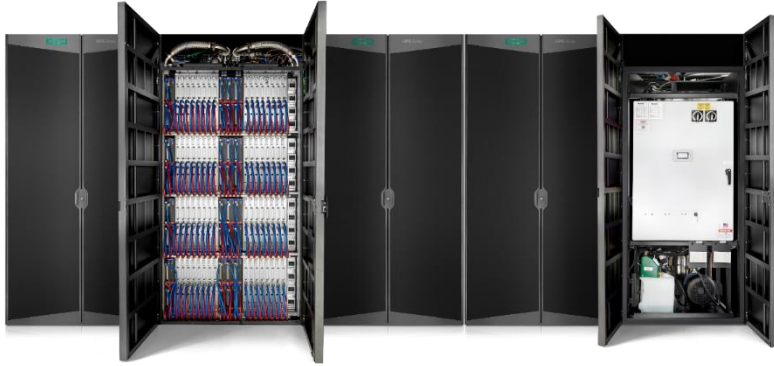**Hewlett Packard Enterprise**
Operated by Selectium

5

# Why not just go to the
# Public Cloud
## for training, fine-tuning and serving generative AI?

Because supercomputing-type workloads like generative AI that run 24x7 can be **4 to 5 times more expensive** to run in the public cloud than in a private environment.*

Because the public clouds were built on a cloud-native architecture. They are **not built on AI-native** infrastructure that has proven its scalability and cost-effectiveness for many years.

**Hewlett Packard Enterprise**
Operated by Selectium

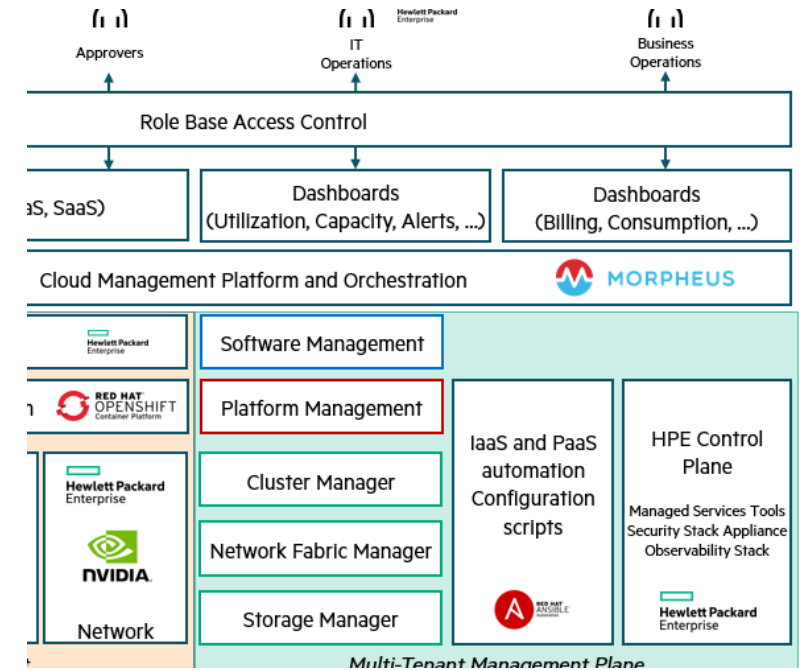# Three solutions to cover key use-cases for AI



## Supercomputing Solution for National Level AI Training

- Highly efficient solution for national AI. Best for LLM training. 1000s of projects.
- The fastest and the most efficient AI infrastructure in the world.

## Private Cloud AI for Enterprise

- Turn-key solution, with full stack of the infrastructure hardware and software for enterprise
- Solution accelerators, full stack of a validated software, RAG & finetuning
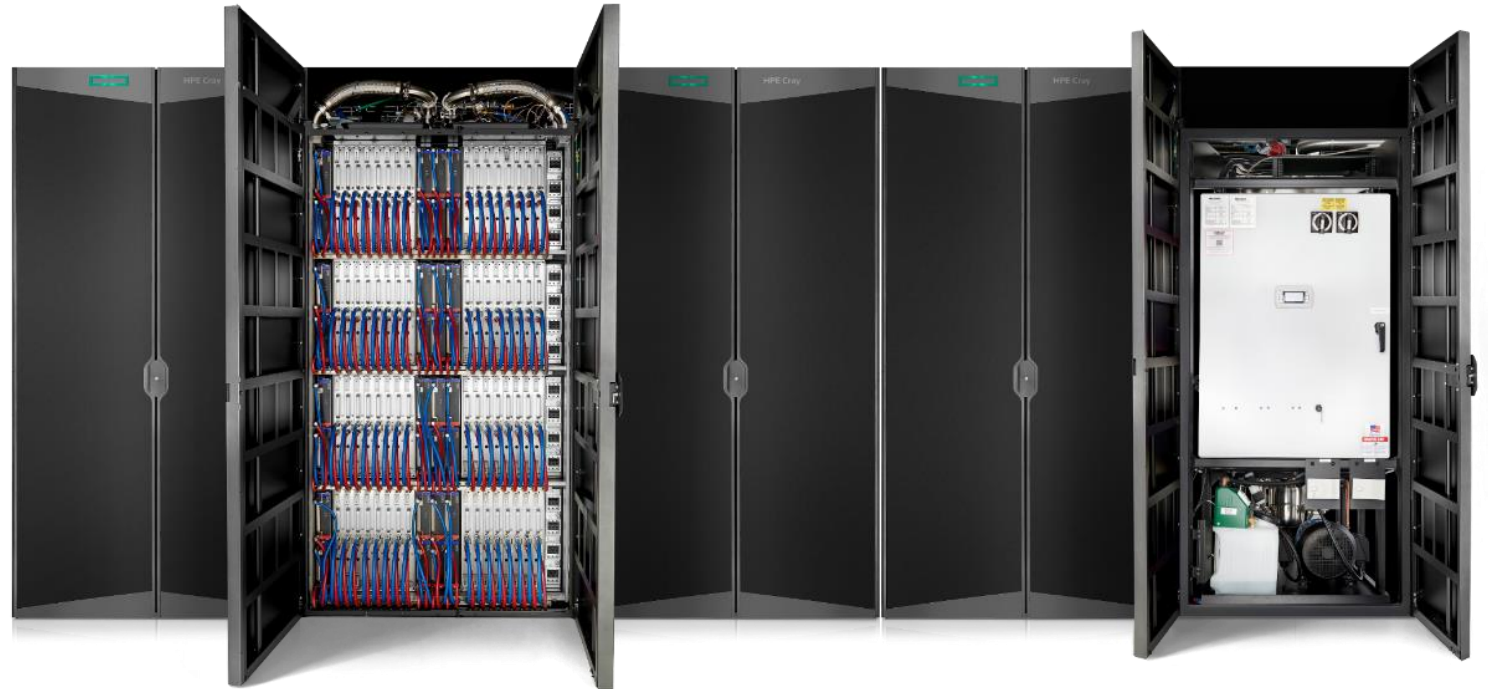
## AI factory for Cloud Solution Providers

- Cloud Solution Provider class multitenant platform for multi-national level deployments
- Easy integration with multiple clouds, unlimited customization capabilities, groundbreaking security

**Hewlett Packard Enterprise**
Operated by Selectium

# Supercomputing Solution for National Level AI Training

- Highly efficient solution for national AI. Best for LLM training. 1000s of projects.
- The fastest and the most efficient AI infrastructure in the world.



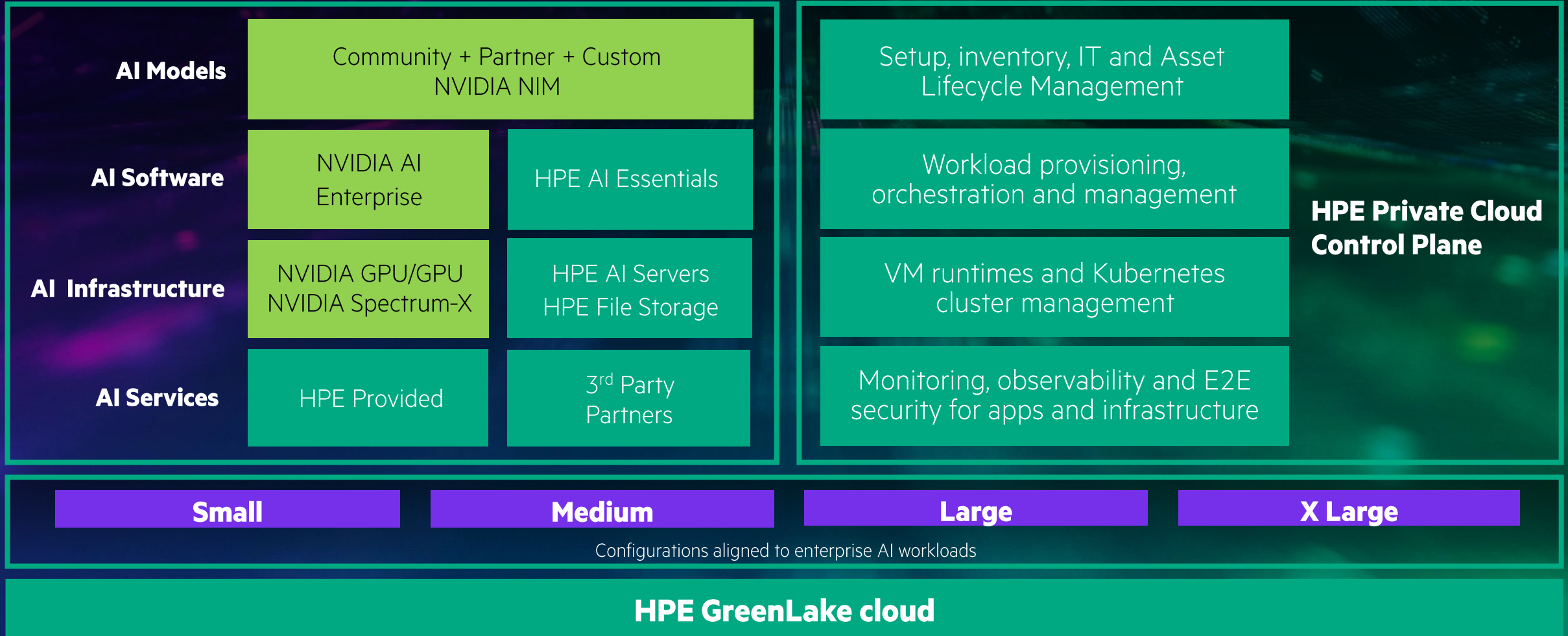**Hewlett Packard Enterprise**
Operated by Selectium

8

# Private Cloud AI for Enterprise

- Turn-key solution, with full stack of the infrastructure hardware and software for enterprise
- Solution accelerators, full stack of a validated software, RAG & finetuning

# NVIDIA AI Computing by HPE

## HPE Private Cloud AI

| | | |
|---|---|---|
| **AI Models** | Community + Partner + Custom NVIDIA NIM | |
| **AI Software** | NVIDIA AI Enterprise | HPE AI Essentials |
| **AI Infrastructure** | NVIDIA GPU/GPU NVIDIA Spectrum-X | HPE AI Servers HPE File Storage |
| **AI Services** | HPE Provided | 3rd Party Partners |

Setup, inventory, IT and Asset Lifecycle Management

Workload provisioning, orchestration and management

VM runtimes and Kubernetes cluster management

Monitoring, observability and E2E security for apps and infrastructure

**HPE Private Cloud Control Plane**

| Small | Medium | Large | X Large |
|---|---|---|---|

Configurations aligned to enterprise AI workloads

**HPE GreenLake cloud**

# Curated and supported set of AI tools
End-to-end software platform

# AI optimized and scalable hardware

## For your Inferencing, RAG, and Fine-tuning workloads

**Unified experience through HPE GreenLake cloud**

|  | **Small** | **Medium** | **Large** | **Extra Large** |
|---|---|---|---|---|
| **Compute** | 4 or 8 NVIDIA L40S GPUs | 8 or 16 NVIDIA L40S GPUs | 16 or 32 NVIDIA H100 NVL GPUs | 16 or 32 NVIDIA GH200 NVL2 |
| ***Storage** | 109 TB to 500TB | 217 TB to 500 TB | 670 TB to 1088 TB | 670 TB to 1088 TB |
| **Networking** | 100GbE NVIDIA Networking | 200GbE NVIDIA Networking | 400GbE NVIDIA Networking | 800GbE NVIDIA Networking |
| **Power** | up to 8 kW rack | up to 17.7 kW | up to 16.5 kW x 2 | up to 16.5 kW x 2 |

**Hewlett Packard Enterprise**

**NVIDIA**

\* Storage Includes 10+ TB Data Fabric Lakehouse capacity in addition to File storage.
Upgraded configurations shown

**Hewlett Packard Enterprise**

# Fueling the future with improved seismic imaging

ExxonMobil improves decision-making and doubles its chances of discovering oil and gas with advanced imaging technology

Sustainable cultivation that's good for biodiversity and the economy

University of Zagreb researchers leverage supercomputing to speed up crop gene analysis by more than 30x

# Multi-purpose HPC cluster

Institute of Physics Belgrade hosts a supercomputer for University-wide projects

Dozens of projects for researchers from all over the country and region

Regional pioneer, in HPC since 1997, moving to HPE in 2013, expanded in 2024

With assistance from local HPE team

**Hewlett Packard Enterprise**
Operated by Selectium

# IPB HPC clusters from HPE

## 2013

105 nodes
Compute HPE SL250
Management HPE DL380 Gen8
Infiniband interconnect
Lustre cluster-wide filesystem
Storage HPE MSA 2000 G3
IO performance ~100MB/s per client
Room neutral with HPE ARCS

## 2024

24 nodes
Compute HPE XL220
Management HPE DL380 Gen10
Infiniband EDR interconnect
Lustre cluster-wide filesystem
Storage HPE MSA 2060
IO performance ~200MB/s per client
Room neutral with HPE ARCS

## Future

Considering merging clusters
Considering expansion
Software updates for interoperability
Lustre PFS updates

**Hewlett Packard Enterprise**
Operated by Selectium

# Including a worry-free operational experience

Maintain the value of your investment with HPE Tech Care and Complete Care Services

## Assigned HPE team

Account Support Manager
Technical Account Manager
Assigned Customer Engineer

### Experts on your environment

**Know-me guidance** provides contextualized recommendations and expertise while i**nventory management** improves environmental awareness and control.

**Assigned Technology Specialists** provide solution-specific expertise when you need it.

### Focused on your outcomes

**HPE experts** provide personalized expertise and deliverables based on individual priorities.

**Environment profile, support plan, and activity reviews** ensure your unique business needs are being served.

Hewlett Packard Enterprise

NVIDIA

# HPC / AI Products Portfolio 2025

**Hewlett Packard Enterprise**
Operated by Selectium

# HPC and AI systems portfolio

## Leadership-class supercomputing

## Accelerated AI

## Mainstream HPC/AI

## Purpose-built storage

**100% fanless direct liquid cooling**

**70% direct liquid cooling, liquid to air cooling**

The next frontier of supercomputing systems redesigned for HPC, AI, and converged workloads

**HPE Cray Supercomputing EX4000**

**HPE Cray Supercomputing EX2500**

Purpose-built, 8-way AI servers for AI model training, tuning and inference

**HPE ProLiant Compute XD680**

**HPE ProLiant Compute XD685**

**HPE Cray XD670**

**HPE Cray XD675**

HPE ProLiant Compute accelerating AI applications for Enterprises

**HPE ProLiant Compute DL380a Gen12**

**HPE ProLiant Compute DL384 Gen12**

Density-optimized, scale-out compute for HPC and AI workloads

**HPE Cray XD2000**

**HPE Cray XD665**

Unprecedented data storage price/performance for HPC, AI, and converged workloads

**HPE Cray Supercomputing Storage Systems E2000**

**Cray ClusterStor E1000 Storage Systems**

**HPE Cray Storage Systems C500**

**HPE Slingshot** combines the performance of a supercomputing interconnect with the cost-effectiveness of Ethernet

Integrated HPC and AI software portfolio, including application and software development ecosystem, system management suite, orchestration tools, enhanced compute environment & more

HPE Services Experts available globally to accelerate your strategic HPC and AI initiatives

# HPE – Nvidia GPU Servers - Subway MAP

DL 380a

x8 B300A NVL

x8 B40

XD 295v

XD 685

EX 154n

"Mandalore"

GB200 NVL4

DL 384

GB300
NVL72

GB300A
NVL36

x8 B300

x8 H200 NVL

GB200 NVL2

EX 254n

Gen12

XD 670

V2

GB200 NVL72

x8 B200

x8 H200

x4 H200 NVL

XD 665

2025

2024

x2 H200

x8 H200

x2 GH200

x4 L40S

x8 H100

x4 GH200

x4 H100

x2 H100

x4 H100

Scale-out

Scale-Up

2023

HGX

NVL options

PCIe

x86

ARM

NVIDIA

**HPE ProLiant DL portfolio**

**HPE ProLiant/Cray XD portfolio**

**HPE Cray EX portfolio**

Hewlett Packard
Enterprise
Operated by Selectium

20

# HPE ProLiant Compute XD
Optimized for natural language processing, large language model training and multimodal training

**HPE ProLiant
Compute XD680**



## New
## HPE-designed
## modular 5U chassis

- Faster go-to-market

- Choice of GPUs, CPUs, boards
  and cooling methods

- Optimized pod deployment
  leveraging a compact
  8-nodes-per-rack arrangement

## DLC Capable

- Energy-efficiency, enablement of
  >1KW GPUs. Optional HPE ARCS
  for room-neutrality

## Adaptable GPUs

- Nvidia B200 (Air | DLC)
- AMD MI300x (Air | DLC)
- AMD MI325x (Air | DLC)
- NVIDIA H200 (Air | DLC)
- Intel Gaudi 3 (Air | DLC)
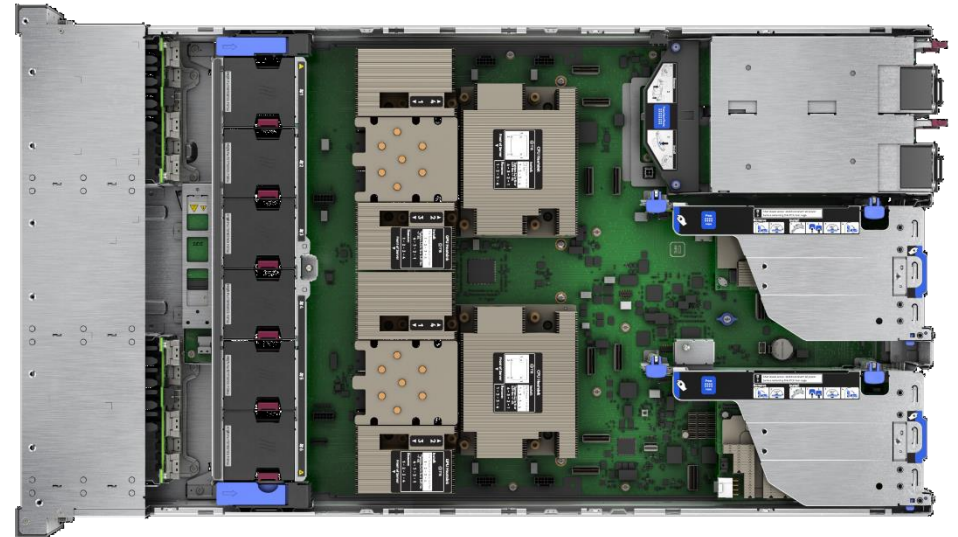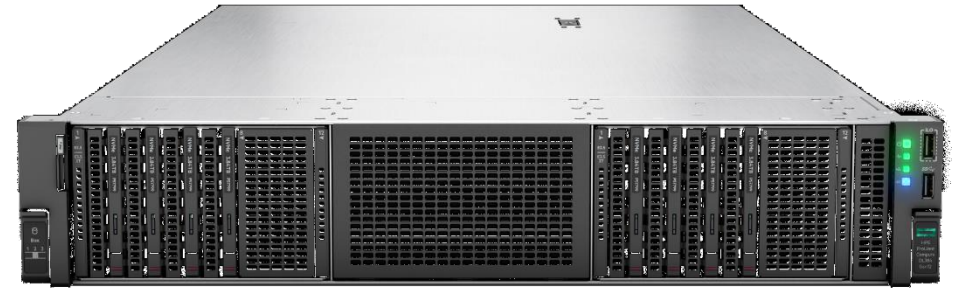- Future GPUs

**HPE ProLiant
Compute XD685**



**Hewlett Packard
Enterprise**
Operated by Selectium

# Grace Hopper in the HPE ProLiant Compute DL384 Gen12

- Cutting-edge NVIDIA GH200 Superchip
  - NVIDIA GH200 NVL2 for up to 1248GB coherent memory for the largest models
  - Enterprise high-end AI Inference workloads including generative AI models
  - AI Training and HPC

- Enterprise Features with the Latest Air-cooled Superchip:
  - Maximum inferences/watt for generative AI inference
  - Maximum performance/GPU and performance/$ in our air-cooled portfolio
  - Enterprise management (iLO-based) for enterprise AI
  - Broad operating system, options, and services available globally

# Rack-level comparison between parallel and NFS storage

For a job with 5 hours run time checkpointing every 15 minutes 3 TB of checkpoint data

Every 15 minutes your compute nodes wait for
## 4 seconds
for the checkpoint files to be written

Every 15 minutes your compute nodes wait for
## One minute and 7 seconds
for the checkpoint files to be written

Rack of Cray ClusterStor E1000 with
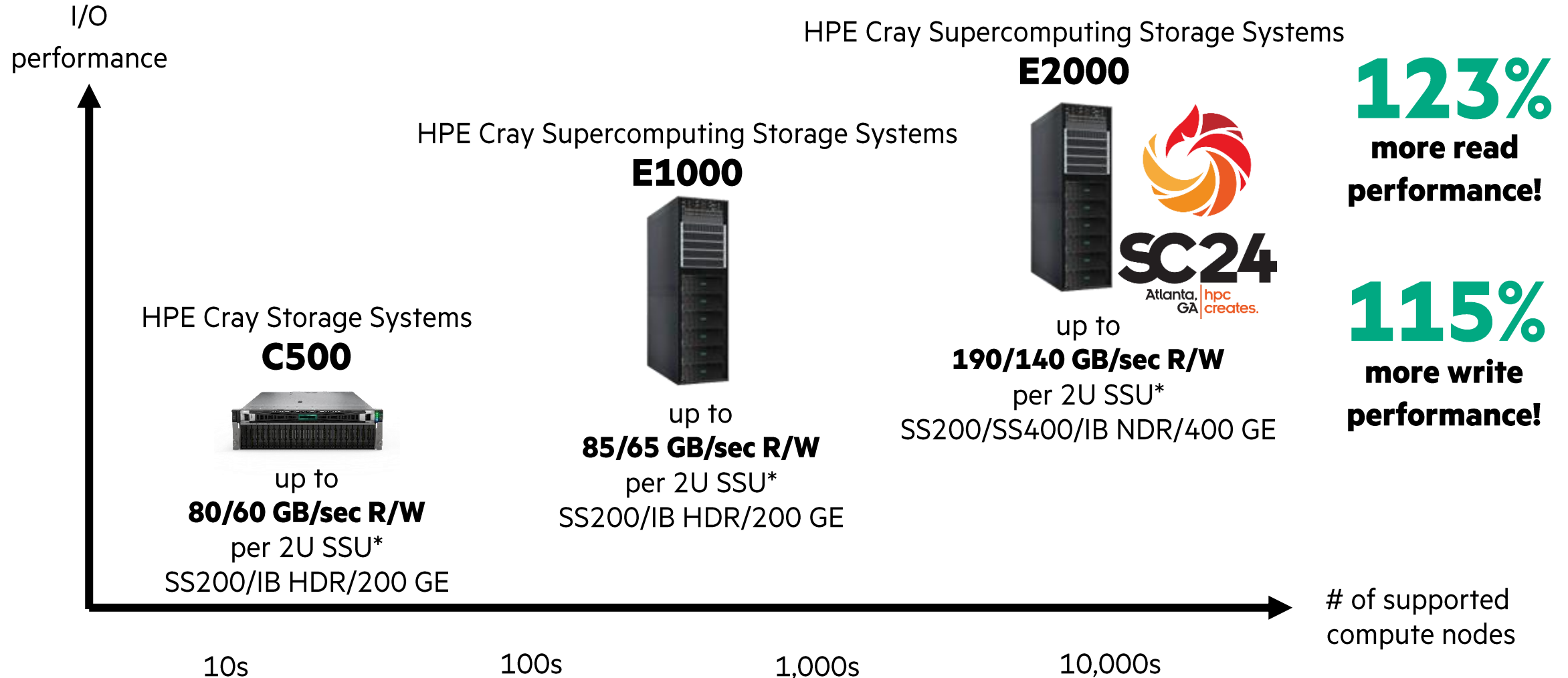960 GB/sec write throughput

Rack of NFS storage with
45 GB/sec write throughput

Hewlett Packard
Enterprise
Operated by Selectium

# HPC Storage portfolio

Storage solutions for any performance requirements and compute cluster sizes

I/O performance

HPE Cray Supercomputing Storage Systems
## E2000

HPE Cray Supercomputing Storage Systems
## E1000

HPE Cray Storage Systems
## C500

up to
**80/60 GB/sec R/W**
per 2U SSU*
SS200/IB HDR/200 GE

up to
**85/65 GB/sec R/W**
per 2U SSU*
SS200/IB HDR/200 GE

up to
**190/140 GB/sec R/W**
per 2U SSU*
SS200/SS400/IB NDR/400 GE

**123%** more read performance!

**115%** more write performance!

SC24
Atlanta, GA | hpc creates.

# of supported compute nodes

10s          100s          1,000s          10,000s

*Scalable Storage Unit

# Thank you!

d.ristic@selectium.com
m.avramovic@selectium.com