# The AI-Driven Storage Revolution

## Solution Brief

**Ultra-High Performance**
- 100GB/s throughput per storage node

**Data Protection**
- Software RAID 0/10/5/6/N+2
- Async Replication

**Block/File Presentation**
- NFSv3 RDMA/TCP
- NFSv4.x RDMA/TCP
- NVMe-oF RDMA/TCP

**Networking**
- NICS: CX-5,6 & 7
- Ethernet (RoCE)
- InfiniBand
- Multipath support

**Storage:**
- Capacity 80TB –1.3PB usable (24 drive chassis)
- Drives supported: 7.69/15.3/30/61TB

**HW Vendors:**
- Kaytus, Dell, HPE, Lenovo, ASUS, Supermicro Gigabyte

**Sustainability:**
- 1.1KWatts
- 2 Rack Units

**SW Supported:**
- NVIDIA OS Native (no drivers needed)
- GPUDirect
- VMware for vGPUs
- Kubernetes CSI

## All Flash AI Data Server designed for AI

Groundbreaking AI Storage providing HPC level performance with the simplicity of a traditional NAS and reducing costs by up to 75%. The AI Data Server revolutionises the way data is stored and processed for AI workloads, delivering unparalleled performance, price effectiveness, scalability and sustainability, empowering organizations to accelerate their AI initiatives and achieve ground-breaking results.
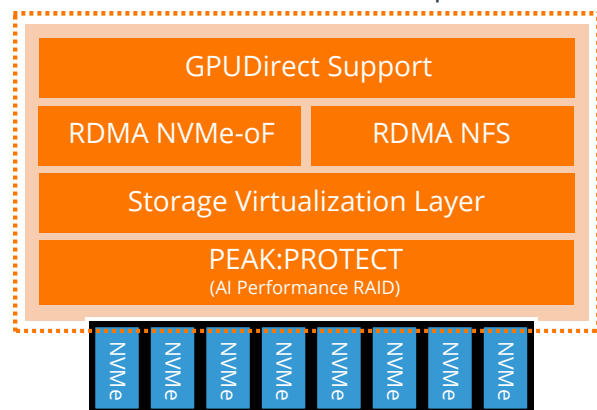
## AI Storage simplified

PEAK:AIO software converts an off-the-shelf NVMe server into an ultra-low latency, protected, shareable filesystem with plug-n-play simplicity. Built from the ground up around NVIDIA's AI ecosystem and fully compliant with modern linux kernels, PEAK:AIO is designed and tuned for AI performance with data shareability. For shared project data exceeding 1.3PB, simply add another PEAK:AIO AI Data Server, or for archive data, the AI Data Server automatically replicates data to PEAK:ARCHIVE.

### PEAK:AIO AI DATA SERVER

**NVIDIA Mellanox ConnectX-6/7  |  RoCE/IB**

Deployed on off-the-shelf NVMe Server

Up to 1.3PB usable per node



GPUDirect Support

RDMA NVMe-oF | RDMA NFS

Storage Virtualization Layer

PEAK:PROTECT
(AI Performance RAID)

NVMe NVMe NVMe NVMe NVMe NVMe NVMe NVMe
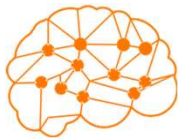
## What sets PEAK:AIO aside:

- **Fast**: Accelerated performance that keeps pace with AI processing.
- **Affordable**: Releases more funds for the GPU budget.
- **Purpose-Built for AI**: Designed to meet AI's unique storage needs.
- **Simplicity**: No storage administration, maintains focus on innovation.
- **Scalability:** Each active storage and archive node scales up to an ultra dense 1.3PB usable
- **Proven**: PEAK:AIO is at the core of a long list of world leading AI projects.

# How AI has changed Storage

- GPU servers are the primary commodity
  - But not storage

- GPUs demand high-end performance storage
  - Clichés: Data hungry GPUs
  - Feed the beast

- Storage vendors were highly focused on AI
  - Come from Enterprise / HPC solutions background
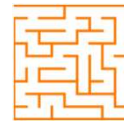
- **Simply not selling**

# Storage for AI Data?

**AI demands super fast compute and storage**

In addition, AI projects start small and scale in smaller than expected stages

**Traditional Storage vendors have not designed for this scale**

High performance storage is too expensive and large, while cost effective storage provides inadequate performance

**High Purchase and Support Costs**

Users are forced to choose between overly expensive or poor performing solutions
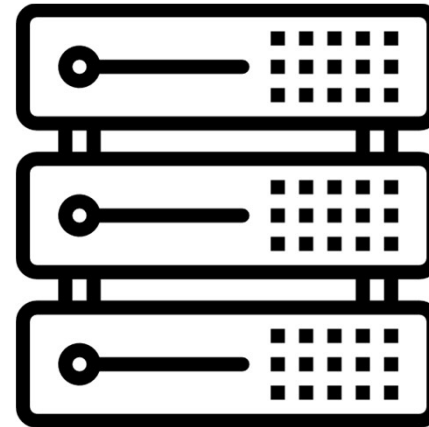
# Enterprise Storage

## Ebay

- Data is created by users
- Lose it, never see it again
- Not only a backup, could not afford loss of access for a moment
- Risk of interruption mean $HIGH
- Meaning the features protecting the data and access are valuable

## Local Retailer

- Data could be:
    - Stock
    - Deliveries
    - Wages
    - Creditors

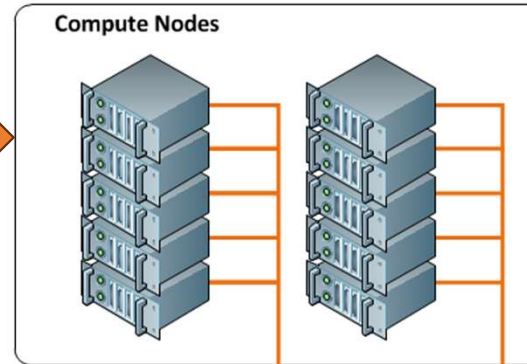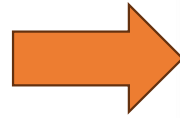- May not be as time critical as Ebay, but still important data

**DATA HAS VALUE**

Storage

Features:
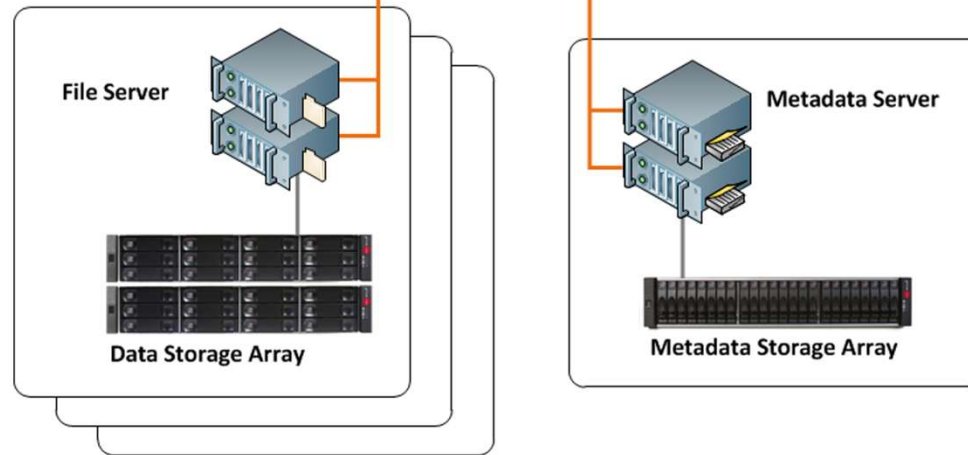- Snapshots
- Versioning
- Replication
- DR
- Integration with Apps
    - SAP
    - Vmware
    - Oracle …
- Cloud Backup
- High Cost $$$

7

# HPC Storage

## 100's / 1000's Compute Nodes



Data Sources:
- Sensors
- Particle Accelerators
- Microscopes
"High Value Data sets"

Features:
- Very high performance
- High Multi Stream performance
- Difficult to Configure/Maintain
- High cost $$$

# AI Storage

Data Sources:
- Corporate systems
- Document Repository
- Standard Training Data
  - Wikipedia

**Central NAS**

Model/Decision/..

GPU Server   GPU Server

**Storage**

**Storage**

Features:
- Very high performance
- Single Stream performance
- GPUDirect
- Low cost $

PEAK AIO

8

PEAK
AIO

# TECHNICAL BREAKDOWN

# PEAK:AIO AI Data Server

## HPC level performance for AI scale



Single storage node in 1U & 2U form factors with up to 24 NVMe SSDs (2U).

- 1x 2U 24 bay server chassis
- 2x 32 Core AMD Genoa / Intel SPR
- 512GB RAM
- 2x NVIDIA ConnectX-6 200GBe / ConnectX-7 400GBe
- Minimum 7 NVMe SSDs

| | Starter: 200TB usable | Large Tier: 1.3PB usable | |
|---|---|---|---|
| Capacity / Drives / Protection | 30TB - 1.34PB Usable | 7.69TB / 15.3TB / 30TB / 61TB Drives | PEAK-PROTECT: RAID 0, 1, 10, 5, 6 |
| PEAK:PROTECT Performance | RDMA NFS 40GB/s (2x CX-6) 80GB/s (2x CX-7) | NVMe-oF (Read): 10M IOPS (Write): 1M IOPS | Performance achieved with single host |
| | NFS3/4 (RDMA / TCP) / NVMe-oF | Kubernetes CSI | NVIDIA GPUDirect |

# PEAK:AIO AI Data Server

**NVIDIA Mellanox ConnectX-6/7** ETH/IB

## Planned Plug-in Modules

**Medical Based**

**HPC / Large Scale**

## AI Plug-in Modules

(GPUDIRECT)

(GPUDIRECT)

**Medical Near Edge Inference**

| MONAI | AIDE | DICOM |
|-------|------|-------|

| PEAK:FS | pNFS | S3 Object | Lustre |
|---------|------|-----------|--------|

NVMe-oF
RDMA / TCP

NFS
RDMA / TCP

**GPU for Edge**

Storage Virtualization Layer

PEAK:PROTECT
(AI performance RAID)

NVMe NVMe NVMe NVMe NVMe NVMe NVMe NVMe

**Generation 4/5 NVMe**

# PEAK:PROTECT

**Benefit from the stability of the world's most mature RAID (MD), boosted to modern-day ultra-fast Performance**

### LINUX MD

- **Decades of Stability**
- **Legacy Limitations**
- **Advanced Parallelization**
- **Stability & Performance**

# PEAK:PROTECT v DEFAULT MD

Gen 4 NVMe

Matches RAID0
All 24 drives at 7GB/sec

Linux MD Random Write RAID6 - 3GB

Randon Read Perfromance
PEAK:PROTECT RAID6 - **171GB/sec**

Linux MD Random Read RAID6 - 38GB

PEAK:PROTECT Randon Write RAID6 - 18GB

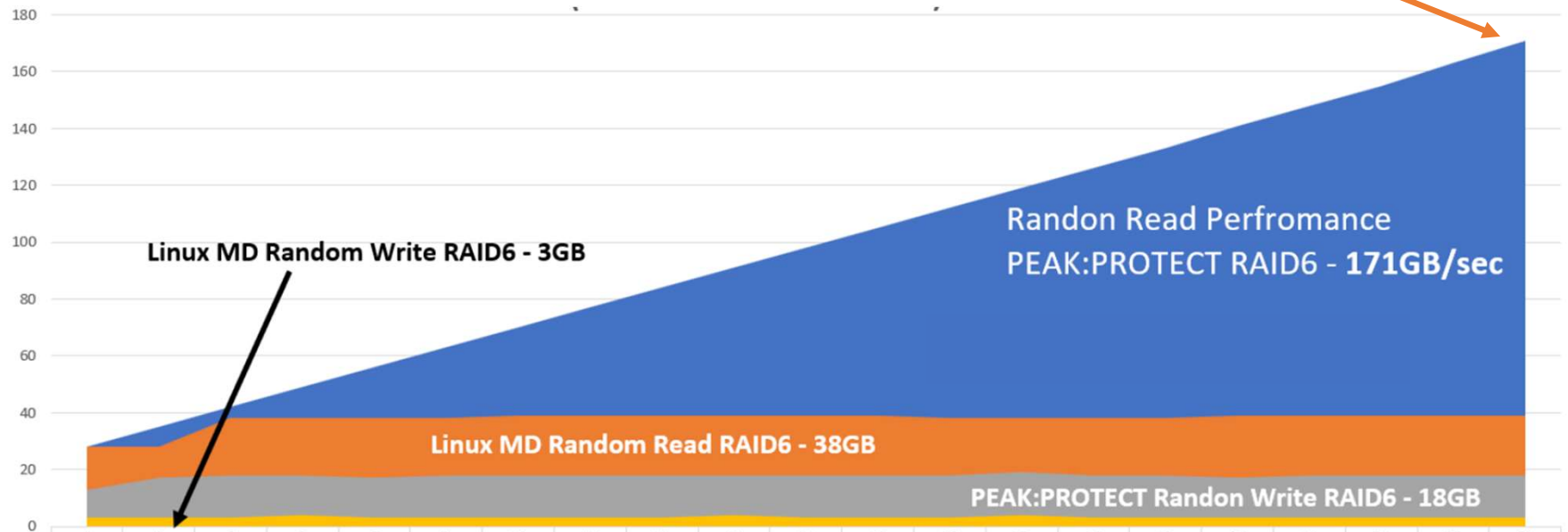| | | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PEAK:PROTECT 6 | Reads | 28 | 35 | 42 | 49 | 56 | 63 | 70 | 77 | 84 | 91 | 98 | 105 | 112 | 119 | 126 | 133 | 141 | 148 | 155 | 163 | 171 |
| MD RAID 6 | Reads | 28 | 28 | 38 | 38 | 38 | 38 | 39 | 39 | 39 | 39 | 39 | 39 | 38 | 38 | 38 | 38 | 39 | 39 | 39 | 39 | 39 |
| PEAK:PEOTECT 6 | Writes | 13 | 17 | 18 | 18 | 17 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 19 | 18 | 18 | 17 | 18 | 18 | 18 | 18 |
| MD RAID 6 | Writes | 3 | 3 | 3 | 4 | 3 | 3 | 3 | 3 | 3 | 4 | 3 | 3 | 3 | 4 | 3 | 3 | 3 | 3 | 3 | 3 | 3 |

17

# A storage solution designed for the AI era

**Performance**

- Full performance to a single compute node

- Single storage node drives BasePod of GPU servers at full performance

- Full performance from 30TBs

- RAID6 equivalent perf. to RAID0, no overheads

- Scales Linearly with Drives

**Nvidia GPU Direct Sequential Bandwidth per Node**



Legend:
- Write bandwidth per storage node - GBps
- Read bandwidth per storage node - GBps

X-axis categories: Dell PowerScale F600, DDN AI400X2, IBM ESS3500, Huawei A310 (8 nodes per chassis), NetApp A800 ONTAP, NetApp EF600 BeeGFS, Pure FlashArray//C, VAST Data Ceres, WekaIO, PEAK:AIO HPE White Paper, PEAK:AIO Gen5 RDMA NFS, PEAK:AIO Gen5 NVME-oF

# A storage solution designed for the AI era

**PEAK AIO≡**

## Simplified Data Shareability with end-to-end Nvidia compatibility

**NVIDIA GPUDirect® Storage Compatibility**
- RDMA Protocols including RDMA over Converged Ethernet (RoCE) with GPUDirect compatibility.

**NVIDIA Supported File Storage**
- RDMA NFS, GPUDirect supported, for shared data.

**NVIDIA Supported Block Storage**
- RDMA NVMe-oF, GPUDirect supported for analytical IO intensive workloads.

**NVIDIA OS Native**
- Out of the box compatibility with NVIDIA OS with no propriety storage drivers required. NVIDIA Kernel support for all features and performance.
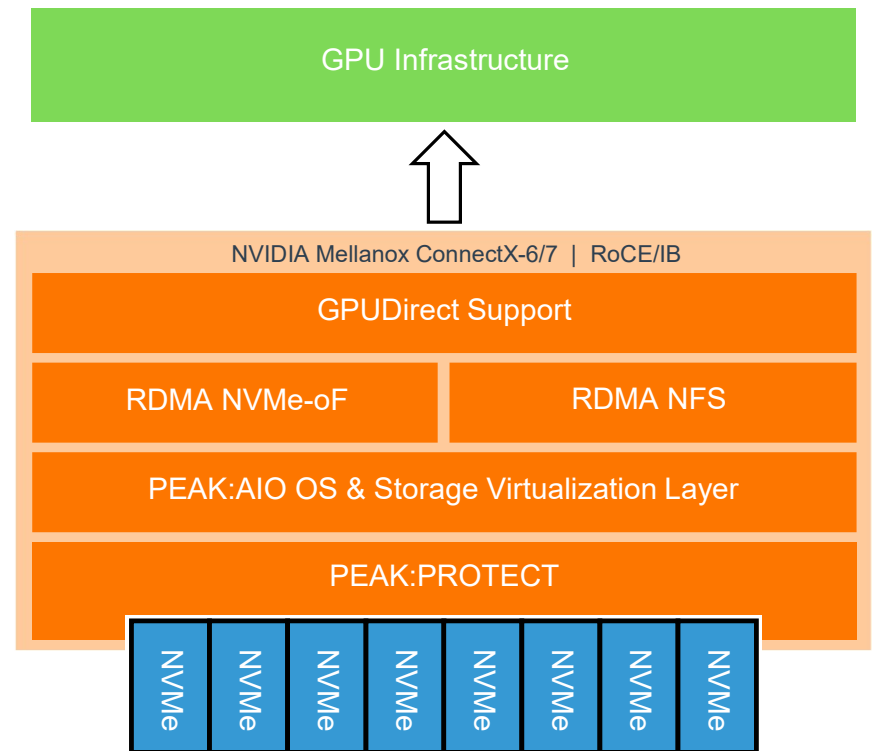
**NVIDIA Performance**
- Saturates CX-6 links & sustains min. of 200Gb/sec of bandwidth per ConnectX-6 to a single GPU server for both file and block values.

**NVIDIA® Port Compatibility**
- NVIDIA ConnectX®-6 200GB Ports for full NVIDIA to NIVIDIA network compatibility and Ethernet / InfiniBand connectivity.

**NVIDIA® AI Enterprise compatibility**
- VMware for vGPUs

GPU Infrastructure

NVIDIA Mellanox ConnectX-6/7 | RoCE/IB

GPUDirect Support

RDMA NVMe-oF | RDMA NFS

PEAK:AIO OS & Storage Virtualization Layer

PEAK:PROTECT

NVMe NVMe NVMe NVMe NVMe NVMe NVMe NVMe

# A storage solution designed for the AI era

PEAK
AIO

**Performance**

- Full perf. to a single compute node
- Single storage node drives BasePod of GPU servers at full performance
- Full performance from 30TBs
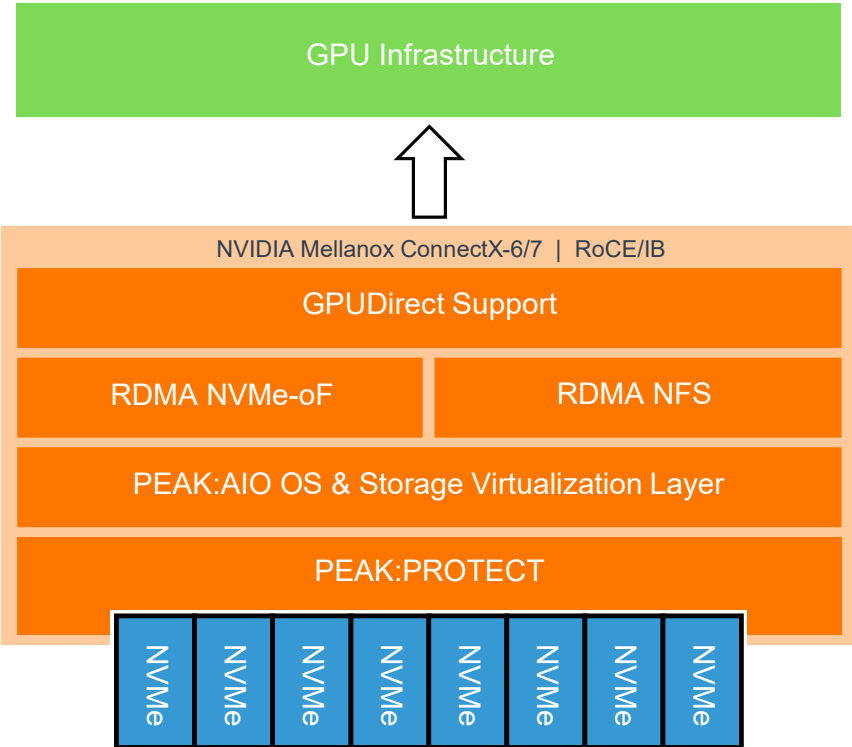- RAID6 equivalent perf. to RAID0, no overheads
- Scales Linearly with Drives

**Simplified Data Shareability**

- Native compatibility with NVIDIA's AI ecosystem
- Uses GPUDirect, MOFED, RDMA
- NVidia Supported file (RDMA NFS) & block storage (RDMA NVMe-oF)
- Both InfiniBand & ROcE enabled
- VMware enabled for vGPU environments

**Zero Maintenance**

- NVIDIA OS Native
- No proprietary drivers
- Plug-and-play design for simplified user management; boots off USB, self-installs in 5 minutes.

---

**GPU Infrastructure**

⇧

NVIDIA Mellanox ConnectX-6/7 | RoCE/IB

**GPUDirect Support**

| RDMA NVMe-oF | RDMA NFS |

**PEAK:AIO OS & Storage Virtualization Layer**

**PEAK:PROTECT**

| NVMe | NVMe | NVMe | NVMe | NVMe | NVMe | NVMe | NVMe |

# PEAK:ARCHIVE Server

The Ultimate AI Data Archiving Solution

- **High Performance:** 1.4PB all-flash storage in just 2U,

- **Immutable Archiving:** Ensures data integrity / protection against tampering and ransomware.

- **Seamless Integration:** Integrates with PEAK:AIO Data Server for automated archiving.

- **Regulatory Compliance:** Meets stringent requirements for healthcare, finance, and legal sectors.

- **Rapid Recovery:** Quick access to archived data for retraining and updating AI models.