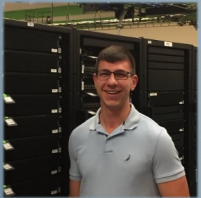
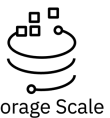


IBM Storage Scale

*High Performance Global Data
Platform for AI/ML*



Matthew Klos
Senior Solutions Architect
Americas SWAT Team
matthew.klos@ibm.com

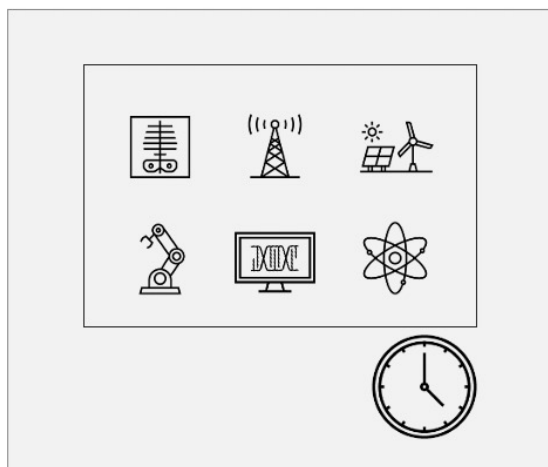


IBM Storage Scale



Data Service Requirements for Enterprise AI Infrastructure

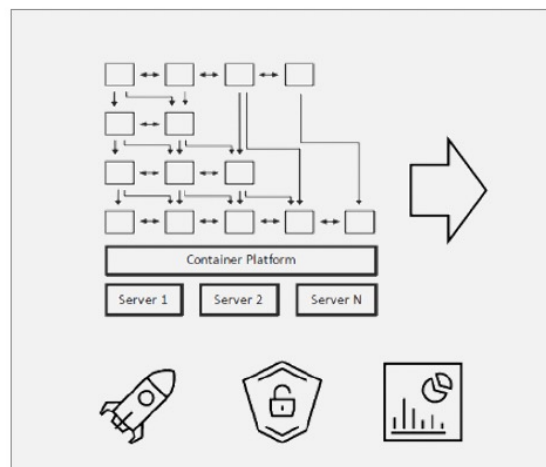
Performance and Resiliency



Growing Resource Demands & Platform Reliability

32% of IT professionals rate increasing data complexity & data silos as a top barrier to AI adoption

Data Orchestration



Lack of Tools & Platforms to Scale AI

28% of IT professionals rate lack of infrastructure a top barrier to AI adoption

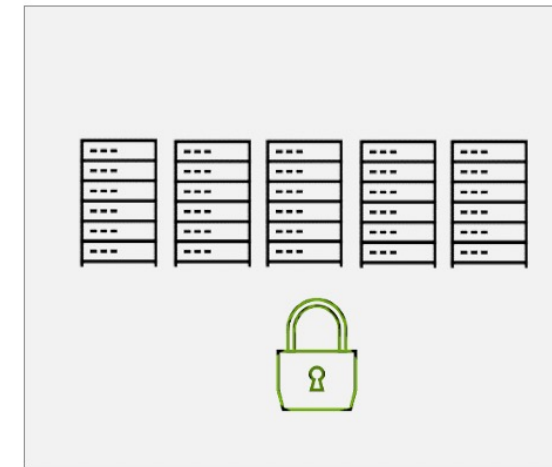
Unified Analytics



New Apps, New Silos

Increasing operational overhead with new AI apps challenge IT standards and practices compliance

Global Data Platform

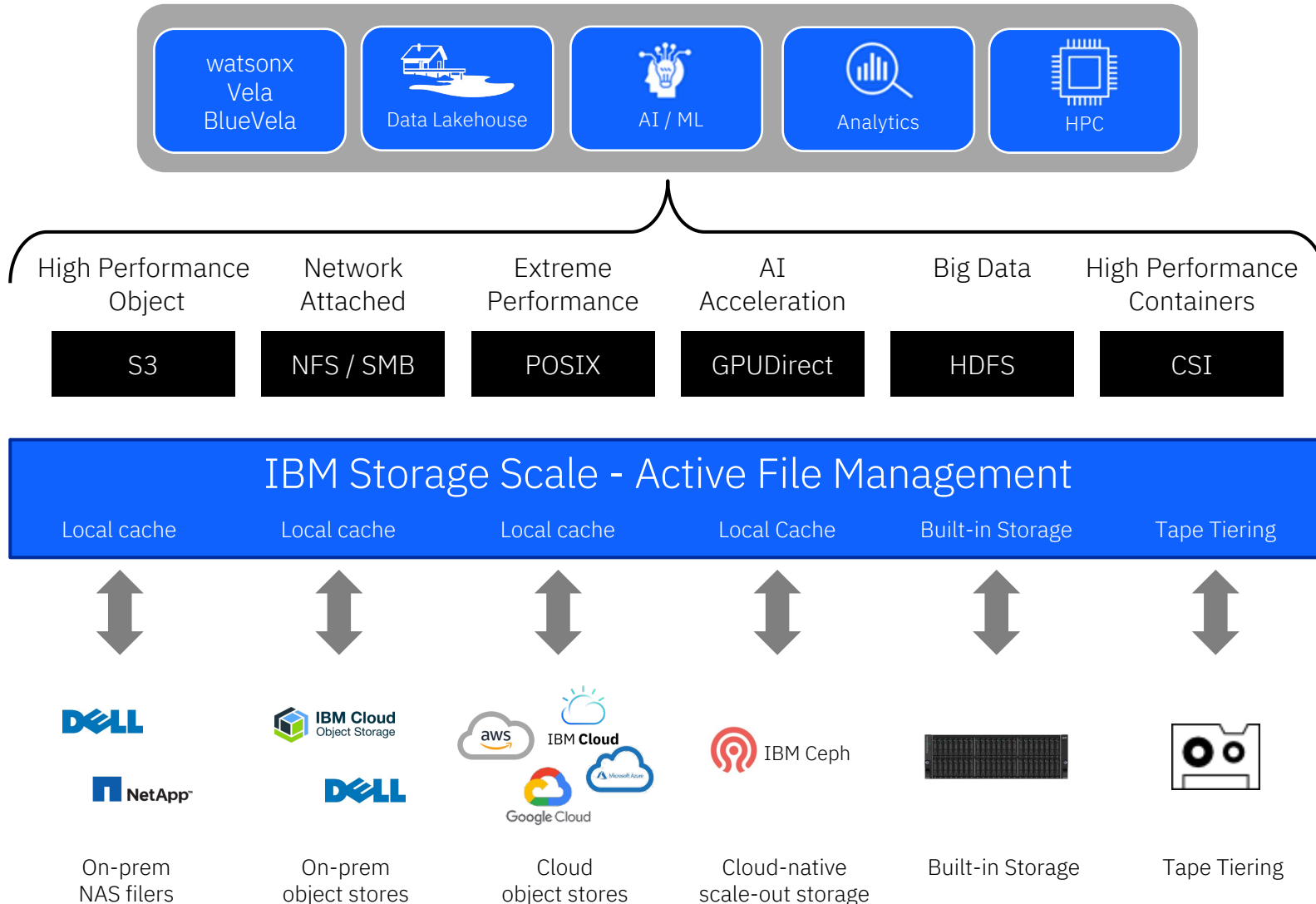


Security

Security of sensitive information from cyberthreats, especially with the widespread adoption of cloud computing

IBM Storage Scale – Global Data Platform for AI

Storage Access, Abstraction and Acceleration to maximize CapEx investment



Multi-Protocol Support Access

Simultaneous multi-protocol access including GPUDirect support

Outcome: Enable globally dispersed teams to collaborate on data regardless of protocol, location or format

Storage Acceleration

Automatic, transparent caching of back-end storage systems

Outcome: Accelerates data queries and improves economics by fronting lower performance storage

Storage Abstraction

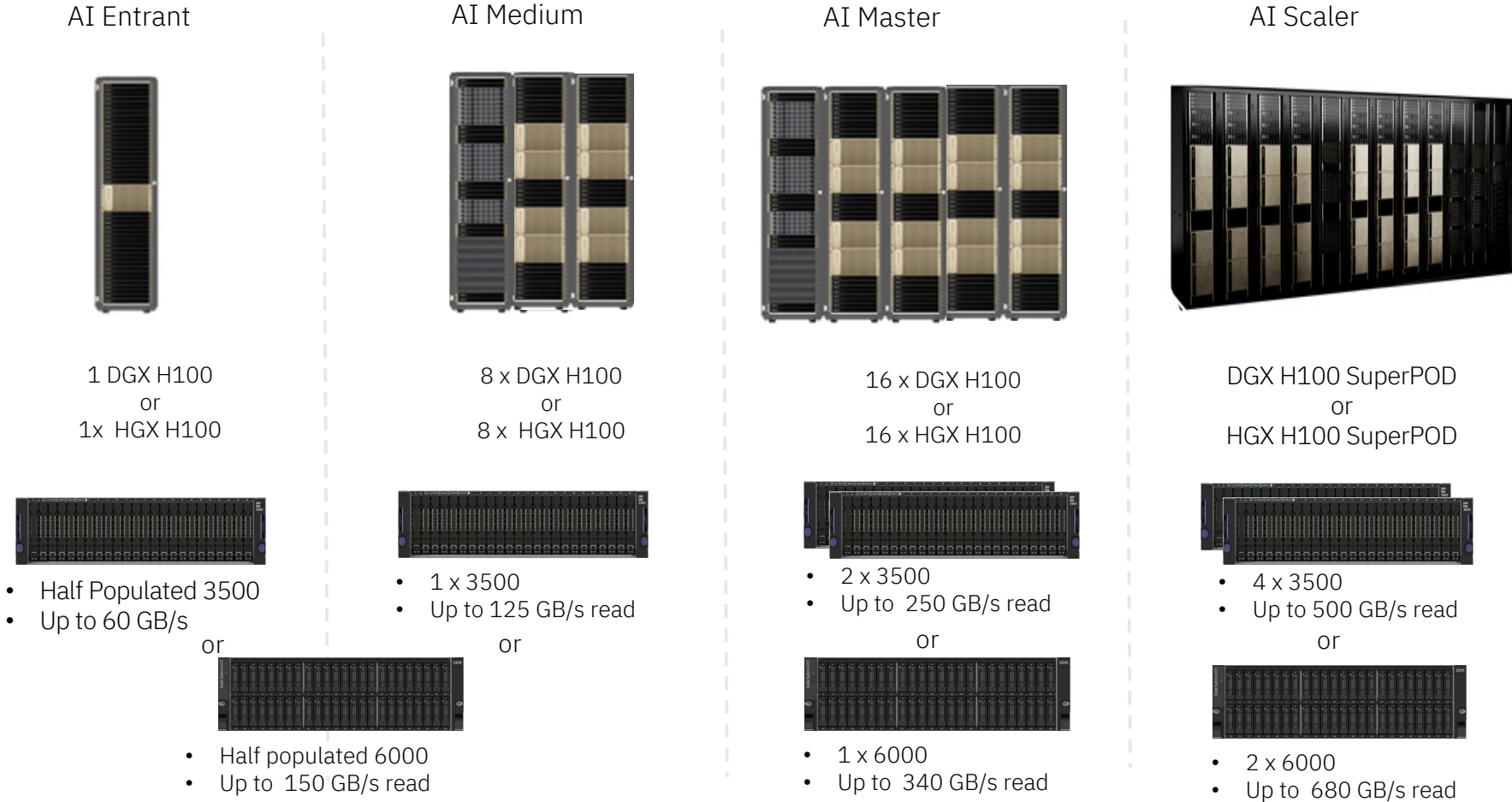
Single global namespace delivers a consistent, seamless experience for new or existing storage

Outcome: Reduce unnecessary data copies and improve efficiency, security and governance

IBM Storage for Data and AI & NVIDIA GPU Solutions

A full spectrum of scalable AI solutions

Start small and scale predictably in response to business demand with the same IBM Storage Software



A simple, scalable upgrade path

IBM Storage Scale

- Simple building blocks – scalable seamless storage upgrade path as needs grow from 1st DGX to AI CoE DGX SuperPOD
- Global Data Platform – Data fidelity capabilities to automate AI workflows.
- Data Economics – Eliminate copies and transparently tier
- Trusted, global enterprise level support and services.
- Successful deployments across the globe

IBM Scale System 6000



up to 1PB Usable Flash
up to 340GB/S Read & 175GB/s Write

IBM Storage Scale: HyperStore

Accelerates AI and Analytics by allowing data to be stored in two pools at once, enabling current use of:

(1) a GNR reliable copy

AND a performance copy, which could be:

(2a) a copy of data that is as close as possible to the computing client (where it's operated on) OR

(2b) a copy of the data that leverages the small I/O performance benefits of NVMeoF.

The first release will support just (case 2b) SSS 6000 shared storage accessed via NVMeoF, and later will exploit (case 2a) local storage (e.g. local NVMe, Persistent Memory, or DRAM).

This first use of Asymmetric Replication on Storage Scale System is called the Hyper Store NVMeoF Performance Pool, providing one erasure-encoded copy of the data and one performance copy that can dramatically accelerate the read performance of smaller I/Os.

Can create a Shared (in model 2a) Co-operative Cache across all compute nodes. Any node can access all cached data, regardless of physical location.

NVIDIA DGX SuperPOD with IBM Scale System



Specifications

GPU	8x NVIDIA H100 Tensor Core GPUs
GPU memory	640GB total
Performance	32 petaFLOPS FP8
NVIDIA® NVSwitch™	4x
System power usage	10.2kW max
CPU	Dual Intel® Xeon® Platinum 8480C Processors 112 Cores total, 2.00 GHz (Base), 3.80 GHz (Max Boost)
System memory	2TB
Networking	4x OSFP ports serving 8x single-port NVIDIA ConnectX-7 VPI ➤ Up to 400Gb/s InfiniBand/Ethernet 2x dual-port QSFP112 NVIDIA ConnectX-7 VPI ➤ Up to 400Gb/s InfiniBand/Ethernet
Management network	10Gb/s onboard NIC with RJ45 100Gb/s Ethernet NIC Host baseboard management controller (BMC) with RJ45
Storage	OS: 2x 1.92TB NVMe M.2
Internal storage:	8x 3.84TB NVMe U.2

