



Accelerating AI with Storage Scale

Storage Scale User Group

May 13th, 2024, ISC, Hamburg Germany

Ted Hoover

Product Manager, Storage for Data and AI

Disclaimer



IBM's statements regarding its plans, directions, and intent are subject to change or withdrawal without notice at IBM's sole discretion. Information regarding potential future products is intended to outline our general product direction and it should not be relied on in making a purchasing decision. The information mentioned regarding potential future products is not a commitment, promise, or legal obligation to deliver any material, code, or functionality. The development, release, and timing of any future features or functionality described for our products remains at our sole discretion.

IBM reserves the right to change product specifications and offerings at any time without notice. This publication could include technical inaccuracies or typographical errors. References herein to IBM products and services do not imply that IBM intends to make them available in all countries.

To unlock the full potential of AI we must overcome the challenges of enterprise infrastructure



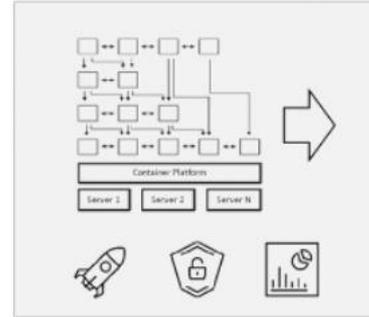
Infrastructure limitations & platforms to scale AI

AI is the fastest growing workload driving spending on compute and storage infrastructure².



Growing resource demands & silos

82% of organizations cite siloed data as a key obstacle to more effective AI development¹.



Operational and physical resource efficiencies

Increasing operational overhead with new AI apps challenge IT budgets and energy efficiencies



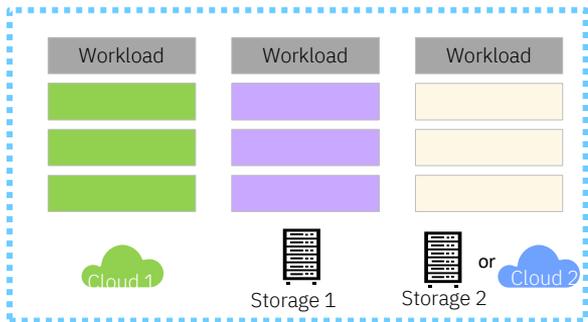
Security and data resiliency

Data must be trusted and security of sensitive information from cyberthreats, loss or downtime is high priority.

What if your organization could accelerate AI workloads with a storage infrastructure designed to accelerate business growth?

Isolated with Silos

Difficult for AI

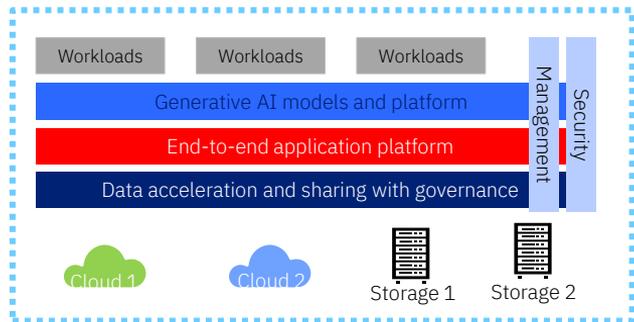


Client IT environment

- Siloed and slow-to-adopt innovation
- Sub-optimal use of resources
- Hard to align across business
- AI constrained

Accelerated Innovation

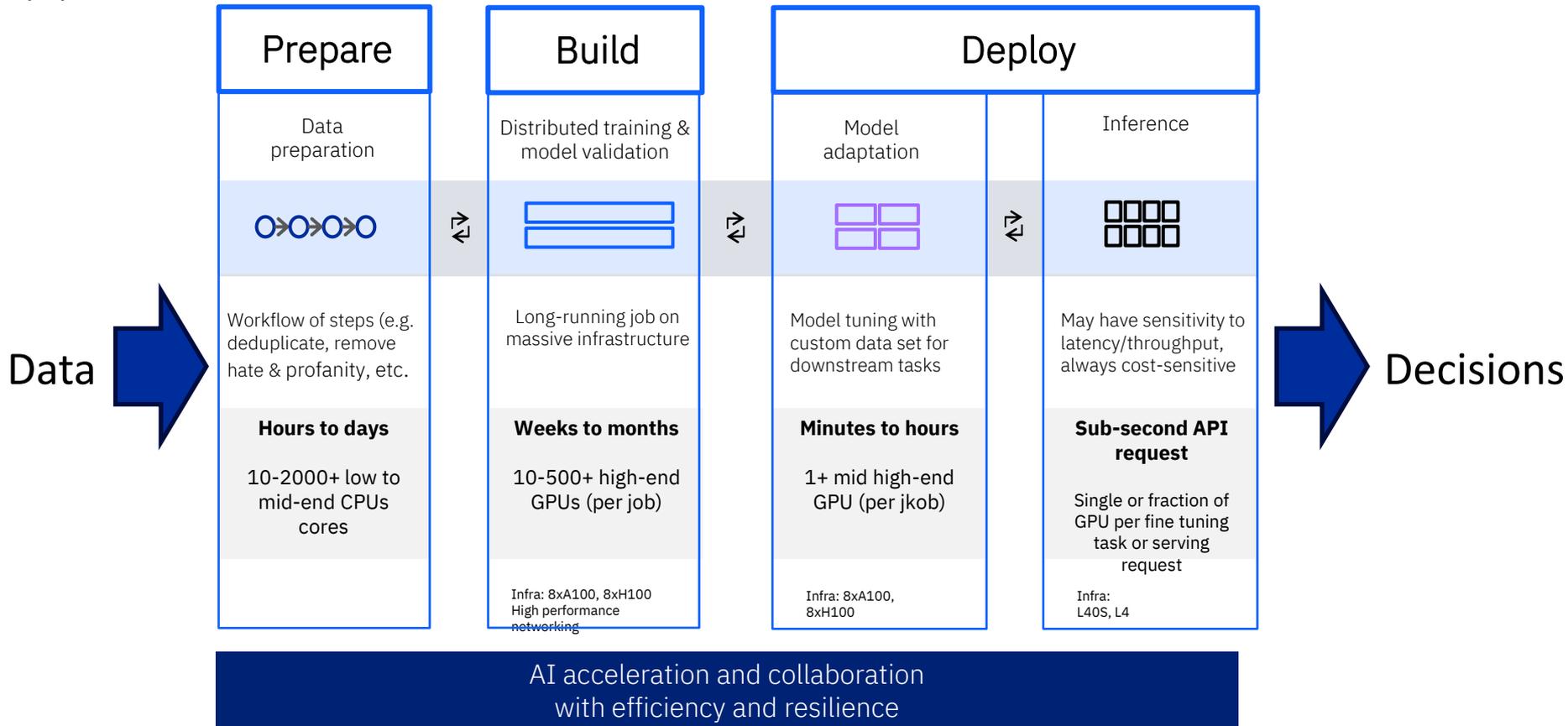
Optimized for AI



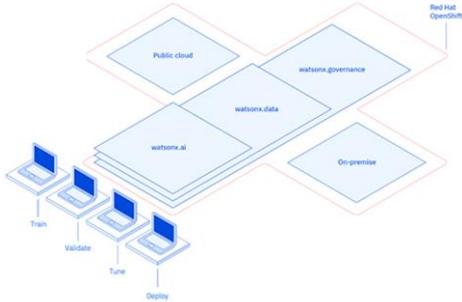
Client IT environment

- Continuous and speedy innovation
- Integrated and automated operating model
- Accelerated value of investments
- Generative AI at scale

Customers need an end-to-end data strategy to bring accelerated results for the AI pipeline

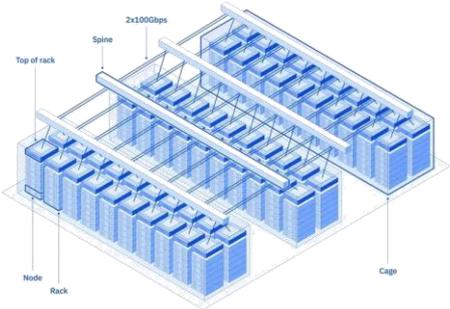


Storage Requirements for AI



AI Tuning/Inferencing

AI Workloads	Storage Acceleration	Efficient GPU support	Rapid deployment
AI Platform		Storage Abstraction	HA/DR/Backup
Optimized for AI		Metadata catalog integration	Simplified Day-2 operations



AI Training

Maximum Performance	Efficient GPU support High bandwidth Low latency
Scalability	Linear scaling of performance and capacity High density

Storage: Why Matters

High Performance Storage, Demanded by Checkpoint

- While LLM cases often do not require as much read performance for training, peak performance for reads and writes are needed for creating and reading checkpoint files.
- This is a synchronous operation and training stops during this phase.
- HPS must provide:
 - High-performance, resilient, POSIX-style file system optimized for multi-threaded read and write operations across multiple nodes.
 - Native RDMA support.
 - Local system RAM for transparent caching of data.
 - Leverage local disk transparently for caching of larger datasets.
- Checkpoint file size: 40GB (estimated) x tensor parallelism x pipeline parallelism.
- Read size: Data parallelism x checkpoint file size simultaneously.
- Example in theory: For 1K GPUs, 8 x 8 x 16, store: 41s – 366s, load: 82s – 683s.

<https://docs.nvidia.com/https://docs.nvidia.com/dgx-superpod-reference-architecture-dgx-h100.pdf>

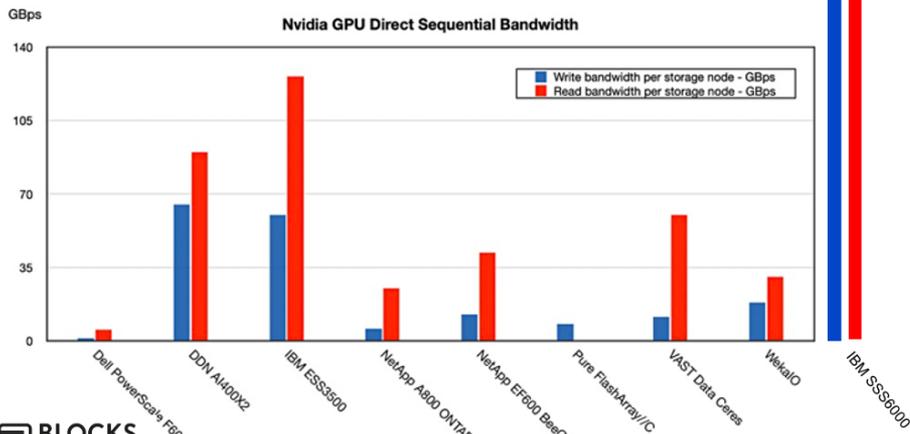
Performance Characteristic	Good (GBps)	Better (GBps)	Best (GBps)
Single-node read	4	8	40
Single-node write	2	4	20
Single SU aggregate system read	15	40	125
Single SU aggregate system write	7	20	62
4 SU aggregate system read	60	160	500
4 SU aggregate system write	30	80	250

- Lot of I/O (Yes, customers will downplay it)
- 2:1 Read: Write
- Most important are Read & Re-Read
- Writes are massive with large parameters models with 175B+ Scalable Performance really matters
- High Performance Parallel File Storage (PFS) is a scratch space, not long-term storage
- Tiering to Object/NL-SAS/Tape is common practice

IBM Tops Nvidia GPU data delivery charts

IBM Storage Scale System 6000 is over 2x more performant than ESS 3500

IBM Tops Nvidia GPU data delivery charts



<https://blocksandfiles.com/2023/08/15/ibm-nvidia-gpu-data-delivery/>

<https://blocksandfiles.com/2023/08/15/ibm-nvidia-gpu-data-delivery/>

Why IBM Storage for NVIDIA GPUs?

The world's fastest systems need the world's best storage. IBM has the best storage for NVIDIA GPUs

Highest Performance Platform

- Fastest performance for reads, writes, and density
- Linearly scalability for future growth

A Robust Enterprise Platform

- Six 9's for all apps: AI, Analytics, HPC, Back-up, Archive, Cloud
- Cyber-resilient, encryption, WORM, and immutability

Collapse Layers & Simplify Data Integration

- Eliminate extra copies and share data globally with all protocols
- Data cataloging and tiering for economics and data flexibility

Why IBM Storage and **NVIDIA** are better together to accelerate AI innovation

IBM Storage Scale accelerates your infrastructure with a hybrid cloud by design for AI platform



Servers with NVIDIA GPUs



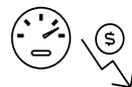
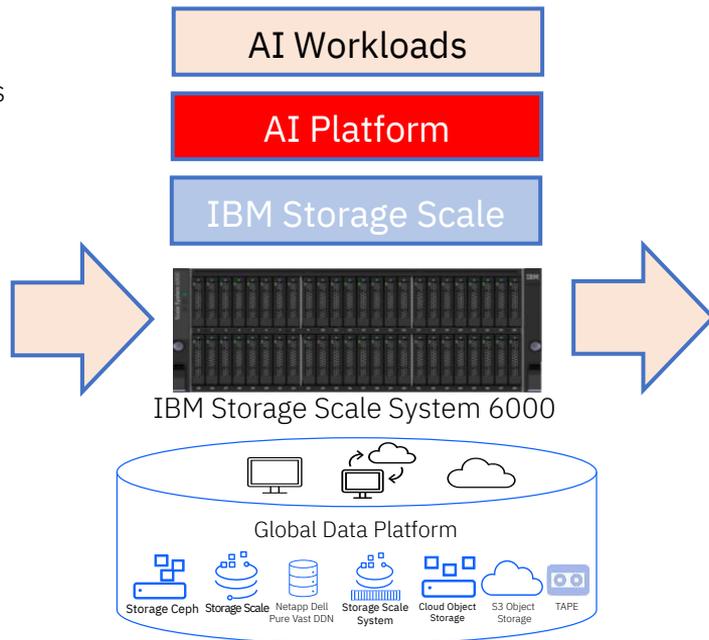
NVIDIA DGX BasePOD



NVIDIA DGX SuperPOD



NVIDIA DGX Grace Hopper



Accelerate discovery

Multi-protocol parallel data access w/ up to 310GB/s, 13M IOPs and NVIDIA GPUDirect® support



Increase collaboration

Data abstraction with remote data, non-IBM storage and cloud data directly to NVIDIA Systems



Support lower cost and green initiatives

New QLC computational storage with transparent archive optimization



Safeguard data from the unknown

Cyber enhanced 99.9999% availability w/data catalog/namespaces to enhance trust

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

AI Entrant



1 x DGX A100/H100
or
1x NVIDIA Certified Server



- 12 NVMe Half Populated 3500
- Up to 60 GB/s Read
- **1 x 6000 w/ 12 NVMe**
- **Up to 80+ GB/s read**

AI Medium



4 x DGX A100/H100
or
4 x NVIDIA Certified Servers



- 1 x 3500
- Up to 125 GB/s read
- **1 x 6000**
- **Up to 310 GB/s read**
- **Up to 155 GB/s write**

AI Master



8 x DGX A100/H100
or
8x NVIDIA Certified Servers



- 2 x 3500
- Up to 250 GB/s read
- **1 x 6000**
- **Up to 310 GB/s read**
- **Up to 155 GB/s write**

AI Scaler



NVIDIA SuperPOD
32 x DGX H100
or
32 x NVIDIA Certified Servers



- 2 x 3500
- Up to 250 GB/s read
- **2 x 6000**
- **Up to 620 GB/s read**
- **Up to 310 GB/s write**

IBM Storage:

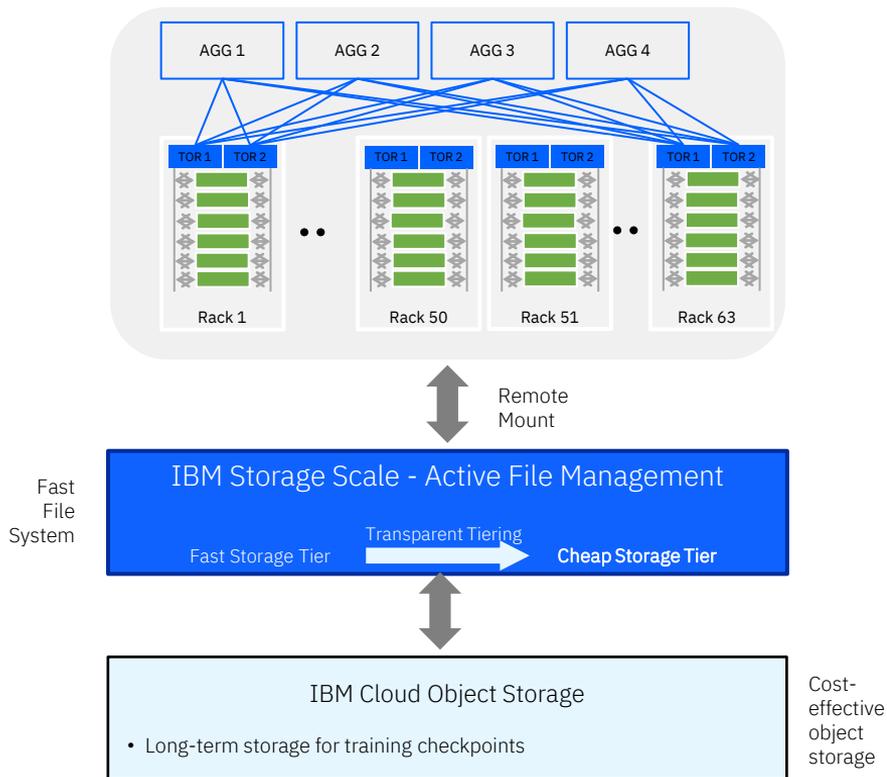
- *Simple building block – simple, scalable seamless upgrade path*
- *Enterprise features– performance, scalability, data protection and security*
- *Global Data Platform Services – Integrate with current storage, multi-site active-active, edge to cloud to core, single namespace across multiple installations*
- *IBM expertise and services*
- *Successful deployments across the globe –*

Telco, Automobile, Banking and Finance, Healthcare, Retail, Academic/ Research and Public Sector

A simple, scalable upgrade path

IBM Storage Scale

An integral part of the Vela architecture



- Built completely on **IBM Cloud** infrastructure
- Dedicated **IBM Storage Scale cluster** on IBM Cloud instances
 - Cloud-Native Scale Access (CNSA) on GPU compute cluster
 - 200 nodes, 1600 GPUs
 - Shared POSIX file system semantics
 - One volume for training data
 - Fit complete training dataset
 - One volume for checkpointing
 - Can accumulate ~10 days of checkpointing
- Large cost-effective data repository using **IBM Cloud Object Storage**
 - Two-tier architecture where AFM transparently moves data between the object storage and file system

Raw performance improvements:

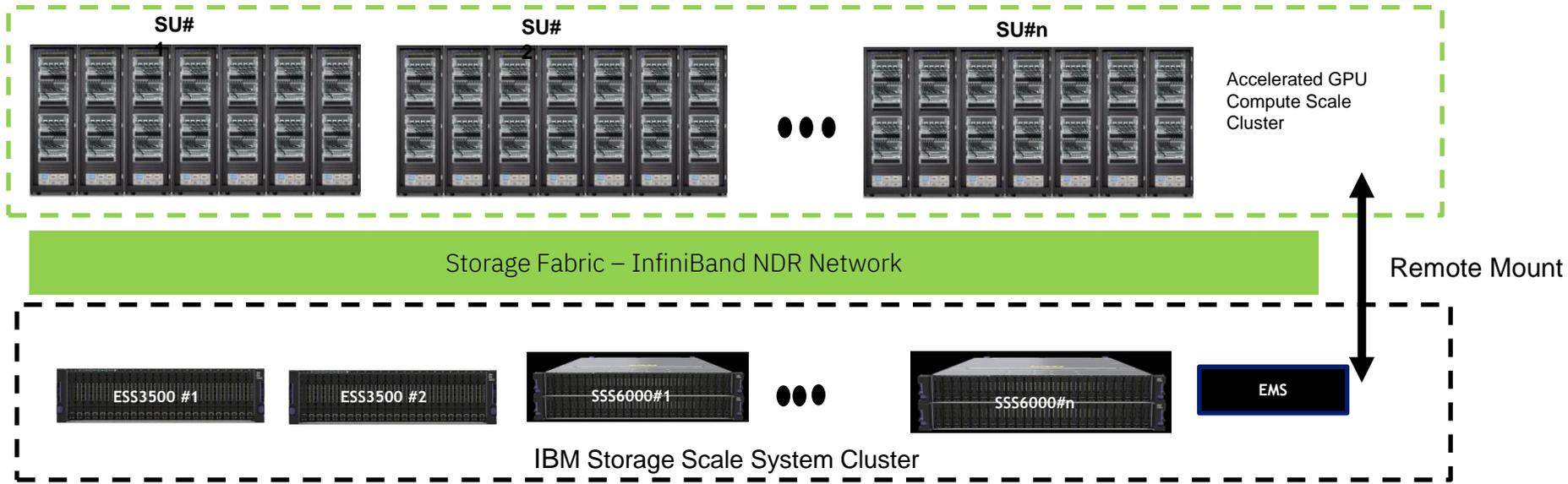
- 3x write bandwidth compared to COS-only (15GB/s vs 5GB/s)
- 40x read bandwidth over NFS (40GB/s vs 1GB/s)

Training performance improvements:

- Storage Scale improved training step time variation by 5X

IBM Blue Vela

HGX “SuperPOD” Storage Fabric (IBM Cloud/ IBM Research)



- A leading global AI & Hybrid Cloud company
- AI Supercomputer Scalable up to 5000 H100 HGX Systems
- 1st Phase 1 SU with 32 HGX node
- 2nd phase will have 20 Scalable Units; 384 HGX nodes
- ESS3500 for initial Phase 1 deployment; 32 SSS6000 for Phase 2
- NDR is the Network Fabric for both compute & storage

- AI and Data platform to deliver enterprise AI service
- Training LLM models with 100B+ parameters
- Faster results – quality & speed of the training models.



IBM Storage Scale on ARM

GA with IBM Storage Scale 5.2.0

On April 26, 2024

Storage Scale User Group

May 13th, 2024, ISC, Hamburg Germany

Ingo Meents

IT Architect

Storage Scale Development

Why ARM? Increasing demand in AI & HPC



- Advanced RISC Machine
- Processor design licensed from ARM limited
- Simple RISC architecture 64 bit (and 32)
- Efficiency: embedded, mobile devices
- Growing into HPC, AI, ML

<https://www.arm.com/markets/computing-infrastructure/high-performance-computing>



TOP 500 list
Fugaku super computer

<https://www.top500.org/system/179807/>



European
Processor Initiative

<https://www.european-processor-initiative.eu/>



Grace-CPU
DPU

<https://www.nvidia.com/de-de/data-center/grace-cpu/>



AWS
Graviton 2 and 3

<https://aws.amazon.com/de/ec2/graviton/>

ARM Neoverse Family

Group of 64-bit ARM processor cores

Neoverse Series	Intended Usage	Level	Instruction Set	Examples
Neoverse N-series (scale out performance)	Data center usage	N1	ARMv8.2-A	Ampere Altra (2-socket 80 cores) AWS Graviton2 (64 cores) Huawei Kunpeng 920
		N2	ARMv9.0-A	Alibaba Yitian 710
Neoverse E-series (efficient throughput)	Edge computing	E1	ARMv8.2-A	
		E2	ARMv9.0-A	
Neoverse V-series (max performance)	High performance computing	V1	ARMv8.4-A	AWS Graviton3 (64 cores) Center for Dev of Advanced Computing (C-DAC) AUM
		V2	ARMv9.0-A	Nvidia Grace (144 cores) Nvidia Blue Field 3 AWS Graviton 4 Google Axion
A64FX, Fujitsu	HPC		Armv8.2-A + SVE	Supercomputer Fugaku

N3 and V3 have been presented in Feb 2024

This is a general list of where ARM can be found, how it can be categorized and some examples. This is not a Scale support list.

<https://www.nextplatform.com/2023/09/13/other-than-nvidia-who-will-use-arms-neoverse-v2-core/>

Hardware Examples

NVIDIA ARM HPC Developer Kit Server - Ampere® Altra® Max ARM Server

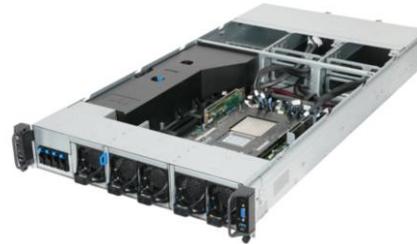
- Single socket Ampere® Altra® Max or Altra® Processor
- Up to 2 x NVIDIA® A100 PCIe Gen4 GPU cards
- Up to 2 x NVIDIA® BlueField-2 DPUs
- 8-Channel RDIMM/LRDIMM DDR4, 16 x DIMMs



Our development & test platform

QuantaGrid S74G-2U

- NVIDIA GH200 Grace™ Hopper™ Superchip
- NVIDIA Grace™ 72 Arm® Neoverse V2 cores
- 1 Processor
- NVIDIA® NVLink®-C2C 900GB/s
- 3 PCIe 5.0 x16 FHFL Dual Width slots



Grace Hopper

Blue Field3 DPU

- Up to 16 Armv8.2+ A78 Hercules cores (64-bit)
- 16GB on-board DDR5



DPU

CPU Fujitsu A64FA

<https://www.fujitsu.com/global/products/computing/servers/supercomputer/a64fx/>

Positive Feedback



AWS Graviton-Processor in Amazon EC2

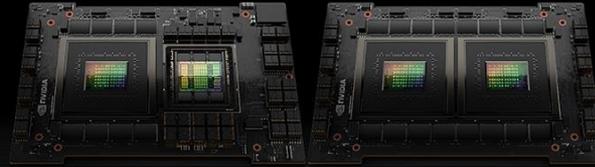
<https://aws.amazon.com/de/ec2/graviton/>

Basic tests successful



NVIDIA Grace CPU

Purpose-built to solve the world's largest computing problems.



Grace Grace Super Chip
Grace Hopper Super Chip

NVIDIA GB200 NVL72

The NVIDIA GB200 Grace Blackwell Superchip combines two NVIDIA Blackwell Tensor Core GPUs and a Grace CPU and can scale up to the GB200 NVL72, a massive 72-GPU system connected by NVIDIA® NVLink®, to deliver 30X faster real-time inference for large language models.

[Learn More >](#)

[Read the Press Release >](#)



Grace Blackwell Super Chip
Just announced in March
at GTC24



Grace Blackwell & Grace Hopper

Grace = ARM CPU where our clients runs
Hopper or Blackwell = GPU where we can put data with **GDS**

ARM support with Storage Scale 5.2.0

- **Included**
 - SE package / install toolkit / rpm based install
 - NSD client
 - Scale base functionality (IO, policies, remote mounts, snapshots, quotas, etc.)
 - Manager roles: file system manager / token manager / cluster manager
 - RDMA (IB or RoCE) including GDS
 - Health Monitoring
 - Target OS: RHEL 9.3 and Ubuntu 22.04 (ask to open RFE for customers assign for RHEL 8)
 - File audit logging, watch folders folders
 - Call home
 - GUI (can display ARM node, but cannot run on ARM)
- **Excluded, but planned for future releases**
 - NSD servers
 - GNR/ECE
- **Excluded**
 - SNC
 - Protocols
 - BDA / HDFS
 - CNSA
 - TCT

Where to get the SE package

- <https://www.ibm.com/support/fixcentral>
- Data Access and Data Management editions

Data Management

Details zum Fix filtern:

▲ Beschreibung

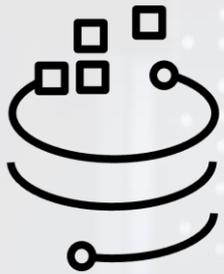
Releasedatum

<input type="checkbox"/>	1	Fixpack: → Storage_Scale_Data_Management-5.2.0.0-aarch64-Linux	2024/04/26
		Product Information	Readme

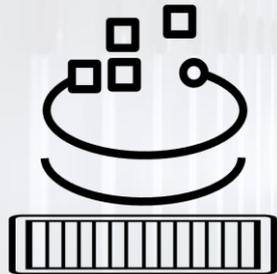
Supported Operating Systems

- RHEL 9.3
 - `gpfs.base-5.2.0-0.aarch64.rpm`
- Ubuntu 22.04
 - `gpfs.base_5.2.0-0_arm64.deb`

Thank you for using



Storage Scale



Storage Scale
System