

Regional SDC Denver April 30, 2025

# Storage Devices for the AI Data Center

Erich F. Haratsch, Marvell



# Outline

- Evolution of Al
- Al Model Complexity
- AI Data Processing Pipeline
- Storage in the AI Data Center
- Performance Considerations for SSDs
- Conclusion



# Recent Evolution of Al

2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026
Deep Learning				Transformers				Generative AI			Multimodel			
												Reaso	ning	
													Chain	of
													Inoug	int
													World Found	lation
Image, Speech Recognition			<b>Recommendation Models</b>				Generative Al			Agentic Al Physical Al				
AlexN	et					GPT-2				ChatG	РТ		Cosmo	<b>DS</b>

Compute, memory and storage requirements continue to increase



# Al Model Complexity

Model	Release Year	Parameter Count	Model Size	Training Tokens	Raw training data
AlexNet	2012	60 million	240 MB	Not applicable	~1.2 TB
GPT-3	2020	175 billion	700 GB	300 billion	~45 TB
GPT-4	2023	1.76 trillion (*)	7 TB	13 trillion (*)	1 PB (*)
Llama2	2023	70 billion	280 GB	2 trillion	N/A
Llama3	2024	405 billion	1.6 TB	15 trillion	N/A

(\*) estimated Assumption: 4 bytes per parameter



# AI Processing Phases and Storage Workloads



REGIONAL **SDC** 

5 | ©2025 SNIA. All Rights Reserved.

# Training vs Inference

### Training

- One large job on supercomputer with 10000s or 100000s of GPUs
- Bandwidth is important
- Frequent checkpointing to save model state in storage

### Inference

- Many threads in parallel
- Time to answer (latency) is important
- RAG drives additional need for storage



## Exemplary AI Compute Tray and Rack Configuration: **GB200 NVL72**



# Al Storage Tiers



# SSD Performance Considerations

- In the past, SSDs adopted next generation PCIe interfaces later than CPUs
- Al is now driving adoption of next generation PCIe SSDs
- SSDs typically designed to saturate sequential and 4KB RR performance
- However, AI workloads (especially Inference) have random accesses smaller than 4KB.



# Saturating Read Performance For PCIe Gen6



- Nvidia's Storage-Next initiative targets ~200 MIOPS (512B) for 16 lanes
  - See Nvidia presentations at SC 2024, OCP 2024 and GTC 2025
- 16 PCIe lanes seen by GPU
  - 128 GB/s Raw
  - ~110 GB/s Effective
  - ~26.8 MIOPS (4KB)
  - ~215 MIOPS (512B)
- 4 PCIe lanes seen by SSD
  - 32 GB/s Raw
  - ~27.5 GB/s Effective
  - ~6.7 MIOPS (4KB)
  - ~53.7 MIOPS (512B)



# Optimizations for High IOPS SSDs



### NAND

- SLC, MLC vs TLC NAND
- Read Time
- Page and Plane Architecture
- SSD Controller
  - CPUs
  - HW acceleration for FW Offload
- Host interfaces
- Form factors
  - 8x random read performance will increase power





 New AI models continue to be released with increased capabilities and complexity

Storage is an essential component in AI data centers

• Al drives adoption of next generation PCIe interfaces for storage

 Increasing Random Read IOPS by 8x requires optimizations in NAND media, SSD controller, host interfaces and potentially new form factors

