

SCSI Trade Association A SNIA O Community

Storage Trends in Al



Patrick Kennedy Principal Analyst, ServeTheHome



J Metz Chair, SNIA AMD, Technical Director



Jeremiah Tussey

At-Large Board Member, STA Microchip Technology Inc. Alliances Manager, Product Marketing, Data Center Solutions

SNIA Legal Notice

- The material contained in this presentation is copyrighted by SNIA unless otherwise noted.
- Member companies and individual members may use this material in presentations and literature under the following conditions:
 - Any slide or slides used must be reproduced in their entirety without modification
 - SNIA must be acknowledged as the source of any material used in the body of any document containing material from these presentations.
- This presentation is a project of SNIA.
- Neither the author nor the presenter is an attorney and nothing in this presentation is intended to be, or should be construed as legal advice or an opinion of counsel. If you need legal advice or a legal opinion please contact your attorney.
- The information presented herein represents the author's personal opinion and current understanding of the relevant issues involved. The author, the presenter, and the SNIA do not assume any responsibility or liability for damages arising out of any reliance on or use of this information.

NO WARRANTIES, EXPRESS OR IMPLIED. USE AT YOUR OWN RISK.



Housekeeping

- Ask questions in the dialogue box on your screen
- Questions will be answered at the end of the webinar
 - If we don't get to your question today, we will post a follow up blog to answer it
- Slides are available for download via the SNIA Educational Library
- Polls will be presented and participation is encouraged



Polling question: What is the biggest challenge when it comes to AI storage?

- Managing data volume
- Data security
- Impact on workload performance
- Cost & ROI

(answer in the polling window)





Storage Architecture



SAS Technology Roadmap



24G+ SAS Feature Overview

	Feature	Benefit
Already Defined:		
Command Duration Limits (CDLs)	Gives user control of HDD latencies	Better supports SLAs
Format with presets	Ability to format as SMR or CMR HDD technology.	Reduction of SKUs
Online & Offline Media Depopulation (Depop)	Ability to isolate media components and logically remove them from data set.	Less service calls in the data center; allows drives to degrade in place, extending the life of the device, thus improving sustainability.
Rebuild Assist for SSDs	Device reports failed LBAs to reduce read time of stripes;	Improves the rebuild time for RAID systems
Persistent Connections	Ability to sustain a connection	Less overhead & higher performance (overhead is associated with opens & closes)
SlimSAS-HD connector (SFF-TA-1016)	Adding a higher density connector	More connectivity out of a low profile board and adding compatibility with CopperLink internal cables.
Under Consideration:		
SPDM Capability	Attestation and authentication of devices	Advanced systems security
Grow the PI (Protection Information) field	Greater capacity, support 64-bit CRC	Increase data reliability
Fairness enhancements for Large	Expanders assess quality of service across the topology	Improved quality of service for all devices.

Designed to enhance SAS architectures that include 24G SAS, 12Gb/s SAS, and 6Gb/s SATA products 7 | ©SNIA. All Rights Reserved.



Topologies

Phases of AI and Storage Workload Requirements





Al Storage Problems To Be Solved



The Al Monster

- AI workloads need:
 - Ever-increasing Memory Bandwidth
 - Ever-increasing Memory Capacity
 - (Near) Instantaneous Data Access (Exabytes)
- Intermittent data surges
- "Straggler" data (tail latency) significantly impacts completion time
- Extended operation duration (hours, days)





What workloads are really like*

- Storage needs vary at each stage
 - Ingest -
 - High capacity and sequential write performance
 - High capacity, heavy read, moderate write
 - Data prep -
 - Sequential read and write performance
 - Moderate capacity, heavy read/write
 - Training -
 - Random read performance
 - Moderate capacity, high read/write
 - Checkpointing -
 - Sequential write performance
 - Moderate-to-high capacity, high read, moderateto-high write
 - Inference -
 - Random read/write performance
 - Archive -
 - Very High Capacity
 - Heavy capacity, heavy read, moderate write



Phases of AI and Storage Workload Requirements **SNIA Initiatives**



• DNA (future)

Storage Al Needs

- Scalability and Performance
- Data Diversity and Edge Computing
- Cloud Integration
- Cost-Effectiveness
- Al-Specific Features





Memory Infrastructures



High-Concept Futures

- Computational Fabric-Attached Memory
- Hierarchical memory pooling
- Intra- and Inter-processor network fabric end-points
- Disaggregated multi-access Ethernet-based storage/data
 Low-Level Efficiency Improvements
- Kernel-Bypass for memory access
- In-process data mutation
- Processing-near-data



SNIA Addresses Al Storage Needs



- Standards for Al-Driven Data Storage
 - Computational Storage Architecture 1.0 & API 1.0
 - SNIA Emerald[™] Power Efficiency Measurement Specification v4
 - Native NVMe-oF Drive Specification v1.0.1
 - Persistent Memory (PM) Performance Test Specification (PTS) v1.0
- Best Practices for AI Data Management
 - Swordfish[™] Scalable Storage Management API Specification v1.2.6
 - Flexible Data Placement; Zoned Storage Models v1.0
- R&D Initiatives
- Advocacy for AI-Friendly Policies



Industry Ecosystem - It's a Big Problem!

• Everyone has a stake in the game:



Polling question: Which aspect of storage for AI workloads is most critical?

- Speed
- Capacity
- Reliability
- Cost

(answer in the polling window)



17 | ©SNIA. All Rights Reserved.



Real World Insights









The Fix Fix Fix Fix

Shared Cluster Storage

Two Camps in AI Clusters



Where is Data Today?

Enterprise Arrays

Cloud Storage

Local Storage

Back-up Storage

Key Questions

- Are the arrays fast enough to feed Al systems?
- Does the data proximity to Al accelerators meet the enormous speed and latency demands of an Al cluster?
- Folks focus on the GPUs and Al accelerators, but how do you solve for interconnects?



What is the storage for AI Agents?







Where Does AI Generated Data Land in 2027?

Enterprise Arrays

Cloud Storage

Local Storage

Back-up Storage

Key Questions

- Compliance needs
- Cost
- Performance
- When to tier



Opportunities & Challenges



Thank you!



Patrick Kennedy Principal Analyst



https://www.servethehome.com/ https://www.youtube.com/ServeTheHomeVideo



Dr. J Metz Chair, SNIA

SNIA[®]

https://www.snia.org/ https://www.linkedin.com/in/jmetz/



25 | ©SNIA. All Rights Reserved.

Follow STA



https://www.snia.org/sta-forum



https://www.linkedin.com/company/snia/



Serial Attached SCSI Playlist on SNIAVideo https://www.youtube.com/@SNIAVideo



