

Regional SDC Denver April 30, 2025

# Storage for Al 101 An overview of Al Workloads from a Storage Perspective

Curtis Ballard SNIA Technical Council, Co-Chair Al Task Force













**SNIA Technical Council** 



**Hewlett Packard** Enterprise

Bio available at: SNIA 2024/2025 Technical Council



### What this presentation IS

High level introduction to storage for Al

Foundation for other Al presentations Thinking Food for planning storage for Al



## What Kind of Storage?

- File storage?
- Object storage?
- Block storage?

- Direct attached storage?
- External (SAN/NAS) storage?

### It Depends

- This presentation will cover general characteristics of AI workloads, from the perspective of the "storage".
- The specific storage implementation for any given workload has lots of choice points that need deeper analysis.



### Storage, A Key Part of a Solid AI Foundation



77% of companies with AI in production plan to upgrade their storage!\*



81% of companies evaluating AI plan to upgrade their storage.\*

<sup>\*</sup> IDC, Future Proofing Hardware for Artificial Intelligence, August 2022

5 | ©2025 SNIA. All Rights Reserved.



## Why is storage for AI different?

Requires huge quantities of data	<ul> <li>With different quantities in different phases</li> </ul>
Extensive Data Manipulation	<ul> <li>With different types of data manipulation in different phases</li> </ul>
It is a multi-phase workload	<ul> <li>Most traditional workloads, like databases, have predictable access patterns</li> <li>AI has widely different workload patterns for different phases</li> </ul>
Performance and capacity varies widely for different AI tasks	<ul> <li>Overview of these variations follows in this presentation</li> </ul>
Highly Parallel Operations	<ul> <li>Usually multiple parallel operations with their own storage workloads</li> </ul>

## Food for AI

AI

 It runs on GPUs and CPUs

but

• Eats Data! and that data

• Eats Storage!



### The Two Sides of AI: Training and Inferencing



### Several phases behind these



### Storage Phases of Al one perspective



Regional SDC Denver April 30, 2025

## **AI Affects All of Your Storage**

Think about your needs for today and tomorrow How does using AI change your storage requirements?



10 | ©2025 SNIA. All Rights Reserved.

### Example: Data Ingest

- Your business processes generate data today
- You already have storage for data ingest
   or do you?
- Business data is already being captured, But:
  - How does AI affect what you capture
  - How does AI affect how you store your business data
  - How does AI affect how you access your business data



### One Example Company

### What they had



#### Al unlocked value in data that they weren't saving!

- Logs of data like mfg sensors, customer interactions, etc.
- Valuable business insights were hiding in sea of data
- AI PoC's demonstrated that they needed to save more data
- Their existing storage capacity was far too small!

#### What they need







## So what do the storage requirements look like for these Storage for AI phases?



## Data Aggregation



- Raw source data has to be prepared for use in AI
  - Logs, pictures, video, documents, etc.
- Data needs organized before becoming training data
  - Clean out noise
  - De-duplicate
  - Normalize

Random

- Privacy and Ethical processing, (anonymizing PII, removing bias, etc)
- Data is read from the ingest storage
  - Cleaned data needs written to storage for data preparation
  - Process may be able to be partially automated applying Al

#### REGIONAL **SD**



Capacity

Sequential



- Data Scientists serve as translators
  - Raw data  $\rightarrow$  Food for AI (Numbers)



- Exploring the data identifying patterns, outliers, relationships, etc.
- Splitting data for training and testing
- Feature extraction converting key features into consumable nuggets
- Data transformation converting data types (Vectorizing)
- Often highly parallel



💻 Random

Sequential

## Model Training - Storage Interactions

- Loading memory for training
- Checkpointing
  - A point-in-time backup of the model copied off for recovery
  - GPU paused while model state is copied out of GPU Memory
  - Checkpointing may be synchronous or asynchronous





 $\Pi\Pi\Pi$ 

## Fast Object for Al

Object storage is convenient but often has higher latency than file

REGIONAL

- File has been more popular than object for AI workloads
- Fast object has been gaining in popularity
  - SSD based object in most cases
  - RDMA accelerated object in some cases



17 | ©2025 SNIA. All Rights Reserved.

## Model Training



- Checkpointing saving model weights and other state
  - Model weights are expensive when training takes a long time

Capacity



- Checkpointing saves state to allow restart after an error
- Checkpoint files are written sequentially
  - May be multiple sequential writes in parallel
- Checkpoint restoration is reversed
  - high sequential read, parallel reads to restore to multiple GPUs
- Training is paused, part of checkpoint, and all of restore time
- Storage performance determined be save/restore time goals

### Sequential Random 18 | ©2025 SNIA. All Rights Reserved.

## Model Training - General Storage Planning

- GPUs drive the cost maximizing GPU utilization optimizes investment
- Design for a balanced architecture
  - Balance storage performance with GPU requirements
- Consider data sources
  - May require both file and object access
- If known training workloads match storage performance to workload \*
  - AI GPU benchmarks can show peak performance for various models
  - MLCommons MLPerf Training benchmarks is a good source
  - Determine size of training examples
  - Multiply throughput and size to estimate required read bandwidth
- For general purpose training may need to support GPU max read speed
  - Can be up to 1GB/s per GPU for high end GPUs today, increasing regularly

REGIONA

19 | ©2025 SNIA. All Rights Reserved. \* <u>https://www.snia.org/educational-library/storage-requirements-ai-2024</u>



- Evaluation measuring how well the results of the model match expectations
  - Accuracy how often is it correct?
    - Precision/Recall roughly a measure of how often wrong vs right
  - Measures such as F1 Score and AUC-ROC (area under the curve/receiver operating characteristics)
- Tuning Adjusting hyperparameters to improve evaluation
- Produces a dataset containing the Model Parameters
  - Internal representation of the neural network
- Model Parameters size is constant, based on # of weights

#### REGIONAL **SDC**

Capacity





💻 Random

Sequential

Inference



- Running production data through the finished model to generate business value
- Inference = Inferring information from the data
- Multiple types of Inference
  - Retrieval Augmented Generation from LLMs
  - Predictive analytics
  - Computer Vision

**Read Perf** 

Capacity

Write Perf

R

S R

- Anomaly detection (e.g., malware, fraud)
- Access pattern can vary some depending on type of inference
  - RAG can produce a random workload similar to databases



21 | ©2025 SNIA. All Rights Reserved.

Sequential

Random

### **RAG:** Retrieval Augmented Generation



Agentic Al





The next wave of AI digital transformation

An application of AI to perform tasks on behalf of users

Learn and act autonomously in complex environments

- Still in the early phases but generating a lot of excitement and activity
- May use RAG (Agentic RAG)
- May use LLMs or small domain specific models



### And Don't Forget! - Archive





- Often overlooked, not core AI, but important AI storage
  - Mandated by regulations for some AI applications
  - Similar, but not traditional "archive"
    - Archived data may be brought back for training or new insights

Write Perf



- Performance needs vary but "just fast enough"
- No accepted terminology, maybe "Cold Storage"
- Continually growing data set
- Requires low cost and low carbon footprint storage
  - Opportunity for zero power storage such as DNA and Optical



Random

Sequential

