# IO Characteristics of AI models and Workloads
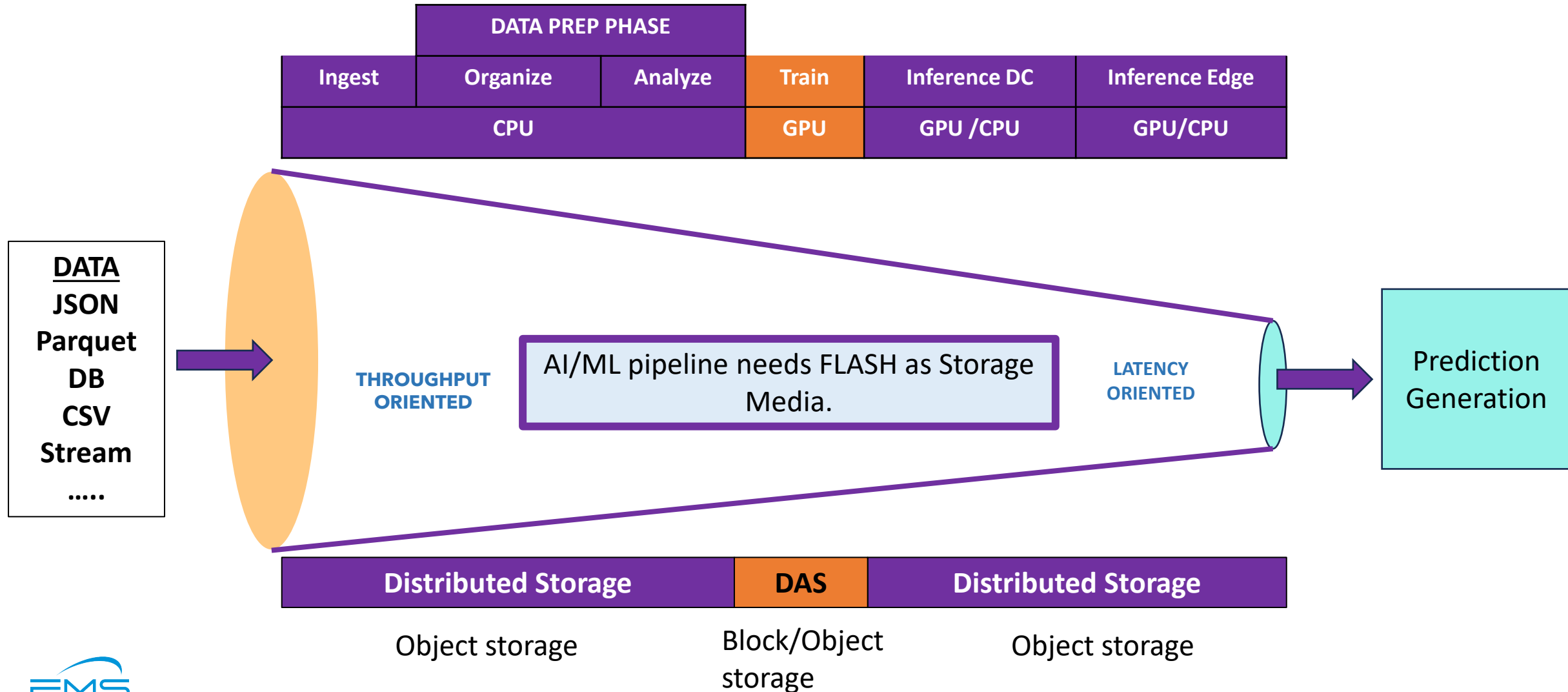
Aug 7th 2025

Kiran Bhat, Product Manager, Solidigm

Alessandro Goncalves – Solution Architect, Solidigm

FMS
*the Future of Memory and Storage*

# Data pipeline and Storage

| DATA PREP PHASE | | | | | |
|---|---|---|---|---|---|
| Ingest | Organize | Analyze | Train | Inference DC | Inference Edge |
| CPU | | | GPU | GPU /CPU | GPU/CPU |

**DATA**
**JSON**
**Parquet**
**DB**
**CSV**
**Stream**
**.....**

THROUGHPUT ORIENTED

AI/ML pipeline needs FLASH as Storage Media.

LATENCY ORIENTED

Prediction Generation

| Distributed Storage | DAS | Distributed Storage |
|---|---|---|
| Object storage | Block/Object storage | Object storage |

# Video Streamer - Pipeline



Fig: Video streamer Application Pipeline

Source: AI Pipeline Optimization on Xeon® Processors | Intel®

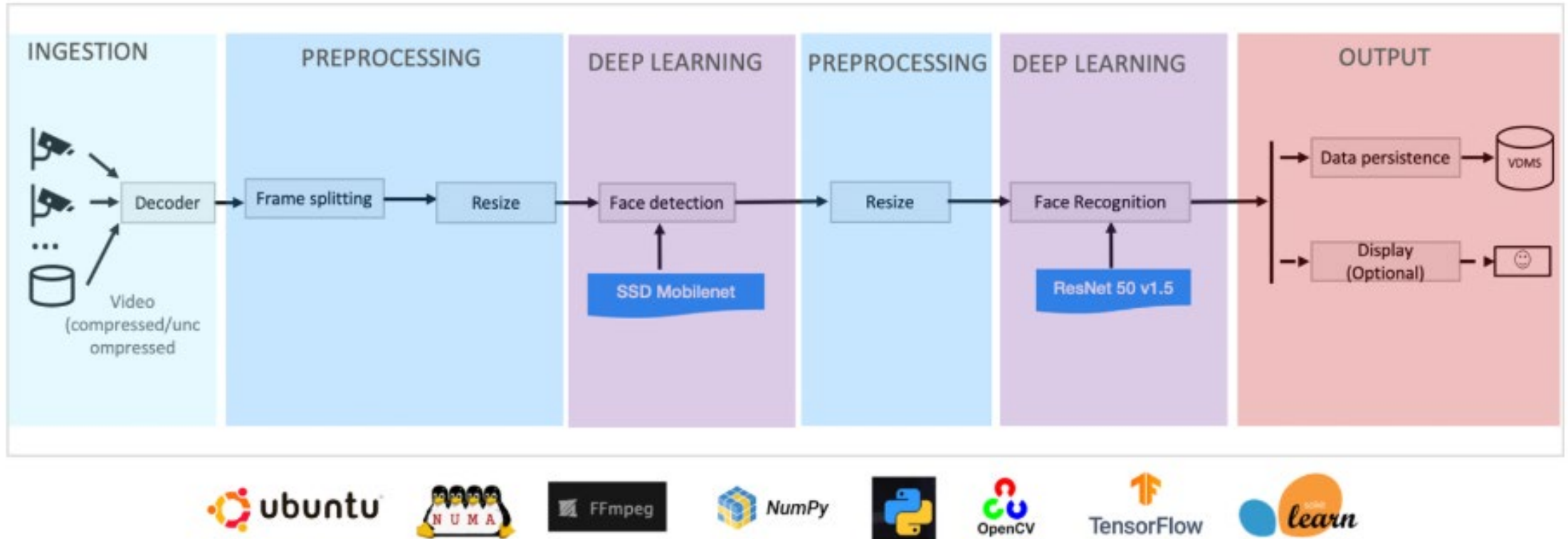# Face Recognition Pipeline



Fig: Face Recognition Application Pipeline

# Census - Pipeline



Fig: Census application pipeline

%time in pre/post processing vs AI cycles
3rd Generation Intel® Xeon® Scalable processor

Legend: Pre/post processing, AI

Source: AI Pipeline Optimization on Xeon® Processors | Intel®
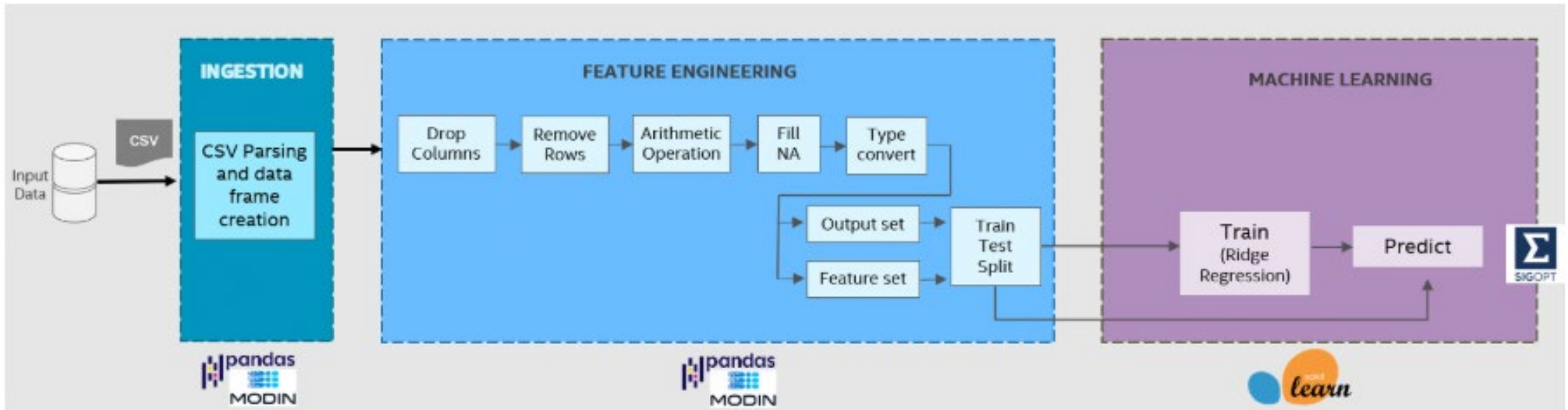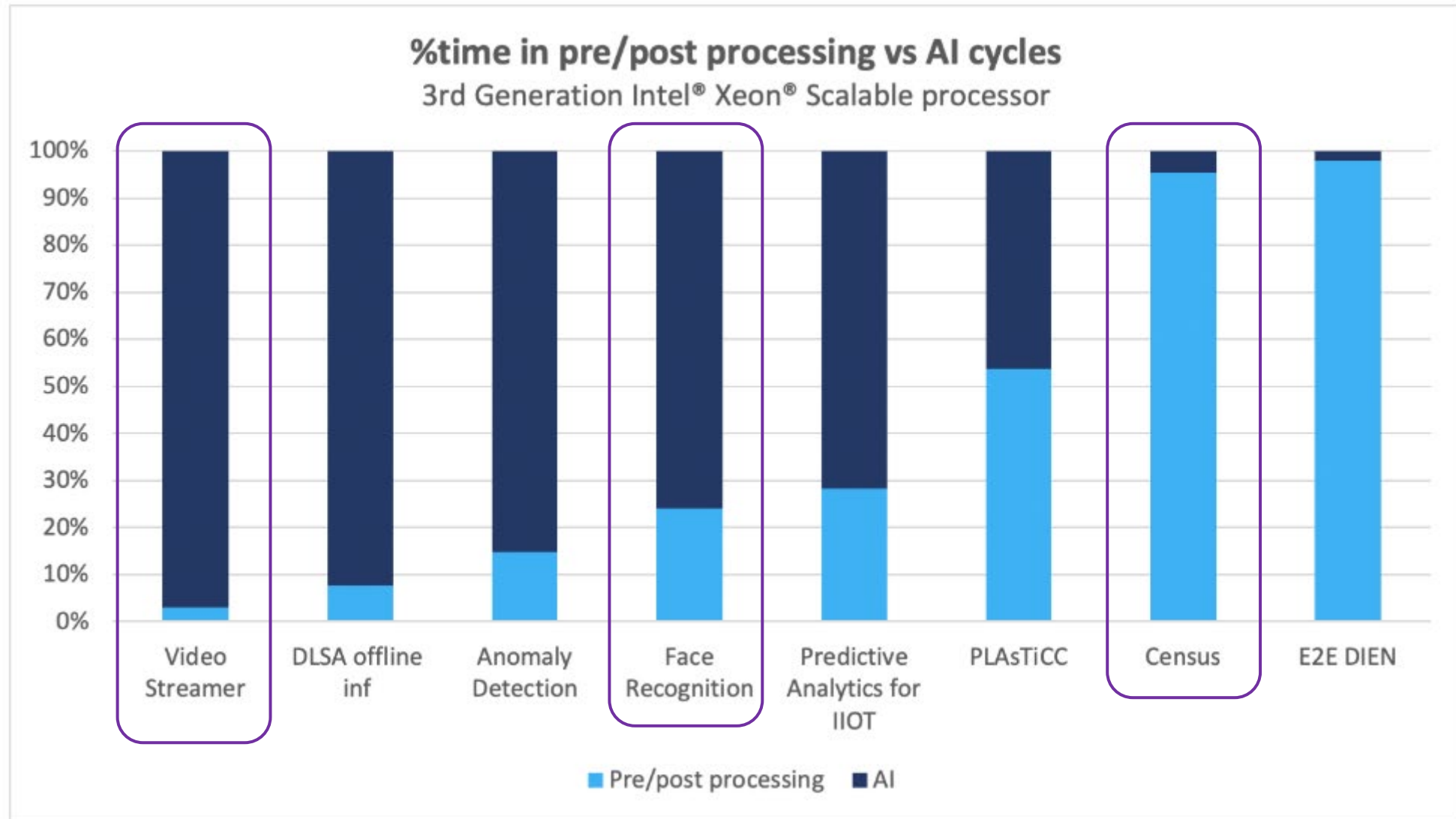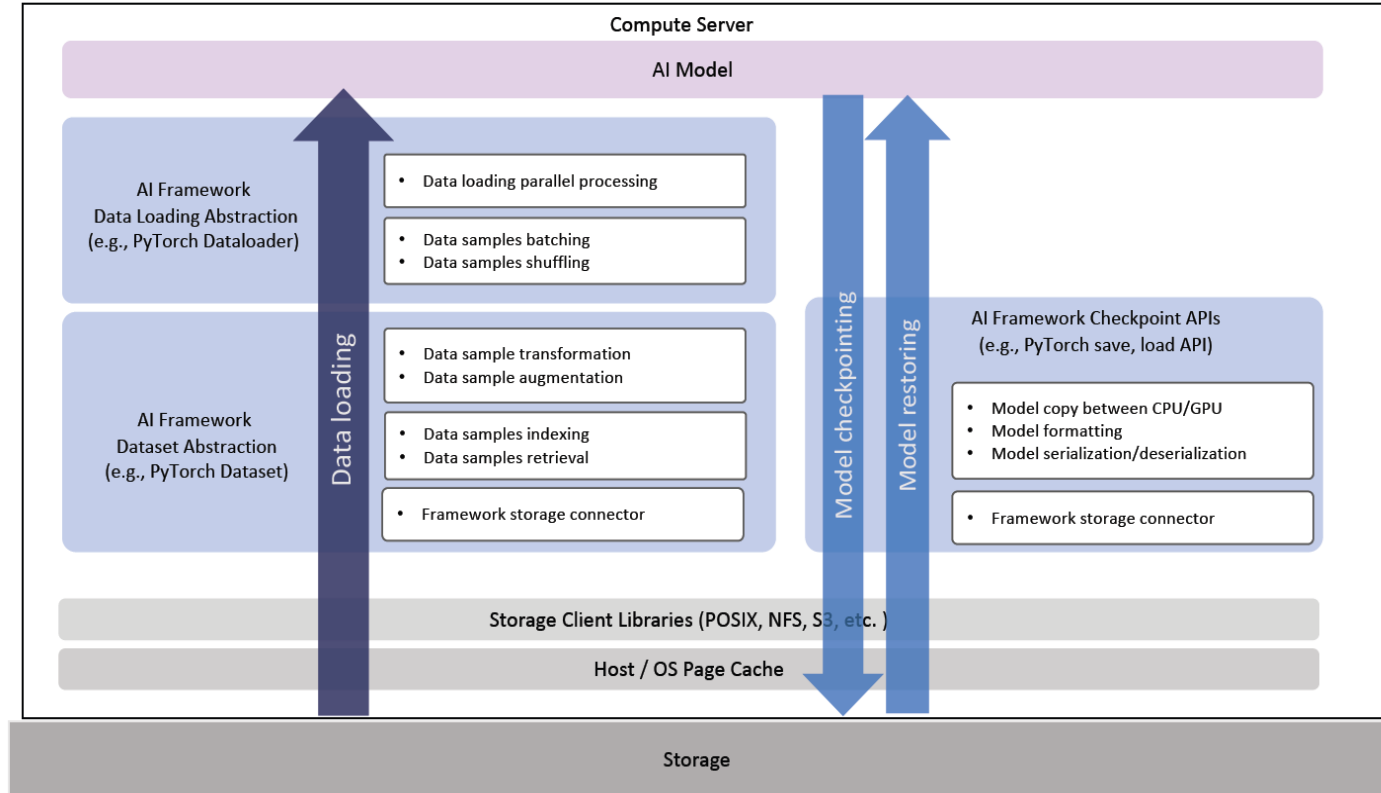
# Complex Solution impacts Storage Performance



AI Framework Stack and Data Flow

Source: https://snia.org/sites/default/files/ESF/AI-Storage-The-Critical-Role-of-Storage-in-Optimizing-AI-Training-Workloads.pdf

# It's All About the Use Case

- AI/ML storage systems encompass a wide range of requirements and architectural complexities.

- Data Pipeline varies by each use case and model used

- Each use case has different overall architecture.
  - Software Stack – framework, application and library
  - Data set
  - Model
  - Training Parallelism

# Ingest Phase

- Source data is different for each use case.
  - File / Block / Object Storage.
  - Parquet, CSV, JSON,image files, etc...

- Source data can be static, based in very slow tiers.
  - LLM , Image Search, etc...
  - Data is prepared Offline and put into the pipeline.

- Source also can be real time devices.
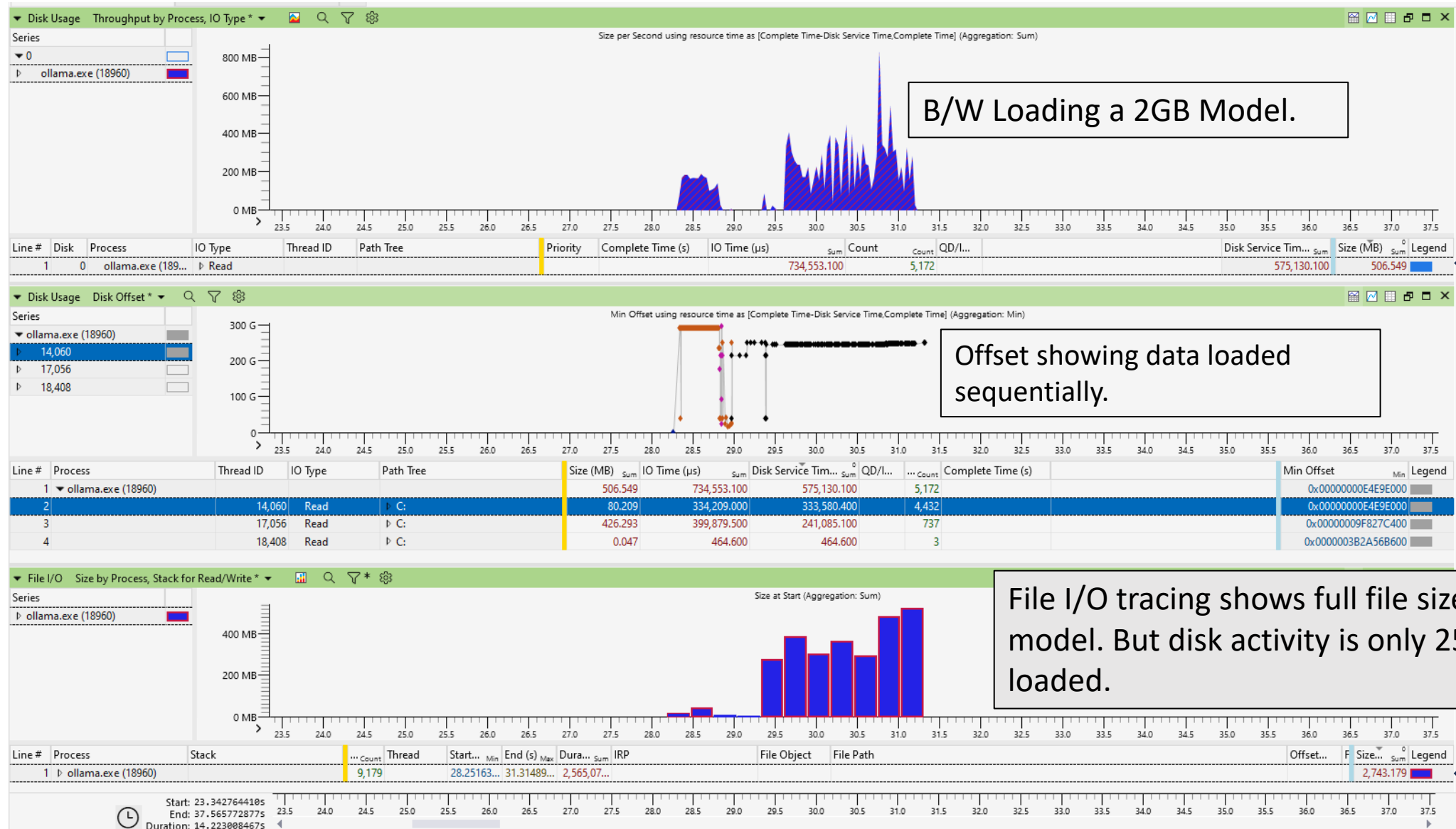  - Facial / Defect Recognition

# Train - Model Sharding / Parallelism

- Tensor
  - Splits individual weight tensors into multiple chunks on different devices.
- Pipeline
  - Partitions the model vertically into stage by layers. Different devices can process in parallel different stages of the full model.
- Context
  - Divides the input context into segments, reducing memory bottleneck for very long sequence length inputs.
- Data parallelism
  - Full model is inserted in each device HBM, and data is processed in parallel in multiple devices. Synchronization happens after each train step.

# TRAIN - Checkpoint

- Checkpoint size is the model size plus model and Optimizer state.
- Model size is set by parameter and precision.
  - FP32 = 4Bytes
  - FP16 = 2Bytes
- Model State is added to checkpoint size.
  - Normally 2 states are saved.
  - Each state has the same number of parameters and normally with the same precision.
- Checkpointing saved as one or more files is based on the model parallelism and implementation
- The higher the GPU count utilized, the higher the checkpoint frequency.

# Inference IO Profile – LLM –Llama 3



B/W Loading a 2GB Model.

Offset showing data loaded sequentially.

File I/O tracing shows full file size for the model. But disk activity is only 25% of data loaded.

*the Future of Memory and Storage*

# Why Storage Matters in Inference

| Mainstream | | Emerging | |
|---|---|---|---|
| Data Ingest and Archive | RAG Database Scaling | KV Cache (Context) Offload | Model Weight Offload |

**Network Storage**
GB / W | GB / in³ | GB / s

**Direct Attach Storage**
IOPS | IOPS / W | IOPS / $

Modern storage overcomes memory constraints to enable **larger models, longer interactions, and better outputs**

FMS
*the* **Future** *of* **Memory** *and* **Storage**

# Summary

- AI/ML storage systems encompass a wide range of requirements and architectural complexities

- IO characteristics (IO size, Read/write ratio, Queue Depth, threads) depends on the model, framework, Data set and Training parallelism, type of storage

- There is no one size fits all; storage needs to be selected based on the usage case scenario

the **Future** of **Memory** and **Storage**

The Solidigm Advantage in the AI Data Pipeline

| Stage: | Data Ingest | Data Prep | Training | Fine-Tuning | Inference | Archive |
|---|---|---|---|---|---|---|
| Storage Requirements: | High Capacity and Sequential Write Performance | Sequential Read and Write Performance | Random Read Performance | Sequential Write Performance | Random Read and Write Performance | High Capacity |

**Class-Leading PCIe 5.0 Performance**

**Recommended Solution:**

**Solidigm D5-P5336** PCIe 4.0 QLC SSD

Capacity · Read · Write

**Solidigm D7-PS1010 and D7-PS1030** PCIe 5.0 TLC SSD

Capacity · Read · Write

**Solidigm D5-P5336** PCIe 4.0 QLC SSD

Capacity · Read · Write

**Up to 122.88TB Capacity**

FMS — the Future of Memory and Storage

## Performance Leaderboard

**Farm GPU**

**#1** Single Storage Node (2U) · 24x NVMe    **116** GB/s

116 GB/s from 1 host

**Multi-Node Competitor**

**#2** 15+ Storage Nodes · Distributed System    **100** GB/s

~6.7 GB/s per host

| | |
|---|---|
| **1**<br>Farm GPU Storage Nodes | **15+**<br>Competitor Storage Nodes |
| **17.3x**<br>Per-Host Efficiency | **93%**<br>Rack Space Reduction |

**Solidigm D7-PS1010 offers** the highest throughput based on the latest MLPerf 2.0 measurement results submitted, for 3D-UNET model training carried out by farm GPU.

[Redefining AI Storage Economics: A Deep Dive into Single-Node Performance and System-Level Optimization](#)

the *Future of Memory and Storage*

# Backup

the **Future** of **Memory** and **Storage**

# References

- [AI Pipeline Optimization on Xeon® Processors | Intel®](#)
- [2108.09373] Understanding Data Storage and Ingestion for Large-Scale Deep Recommendation Model Training

the **Future** of **Memory** and **Storage**