

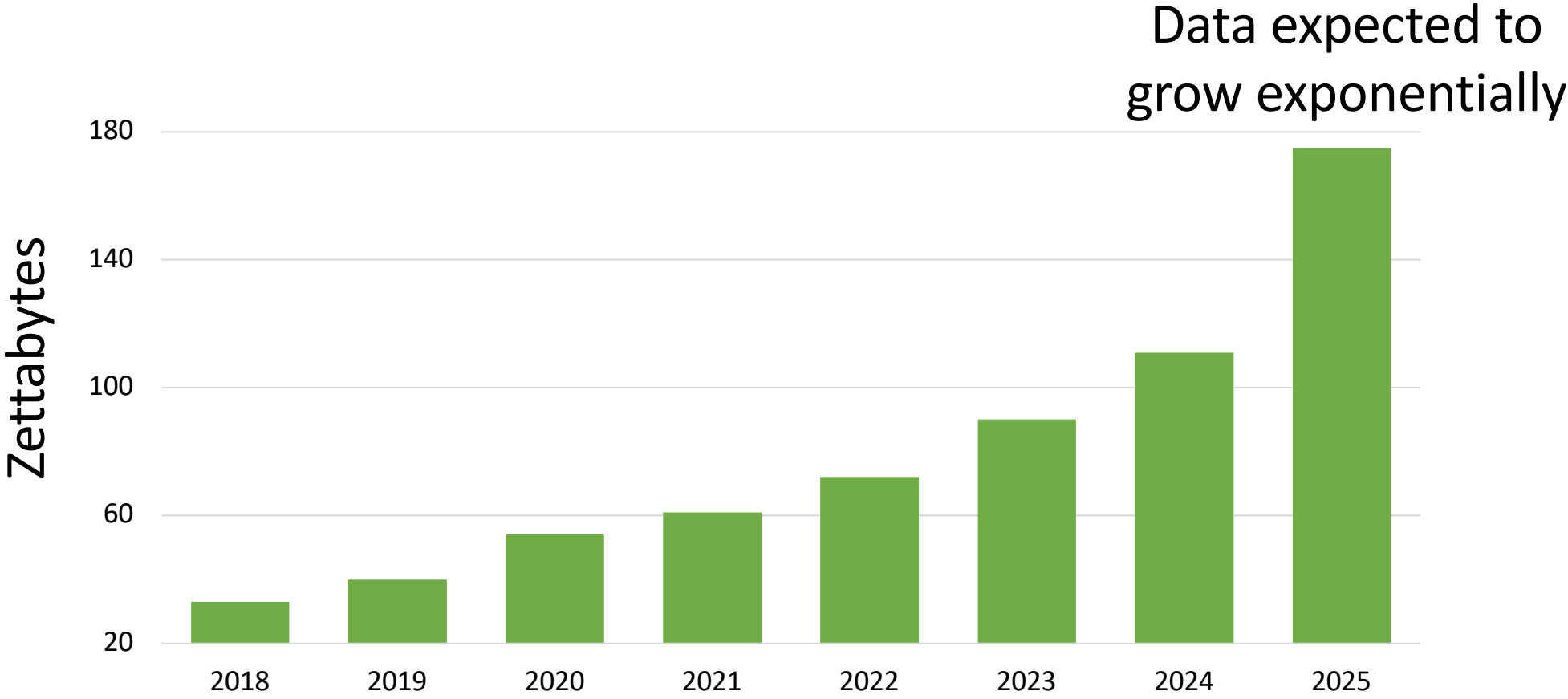
Characterizing Machine Learning I/O with MLPerf Storage

Oana Balmau

CHEOPS @ EuroSys, May 8th, 2023

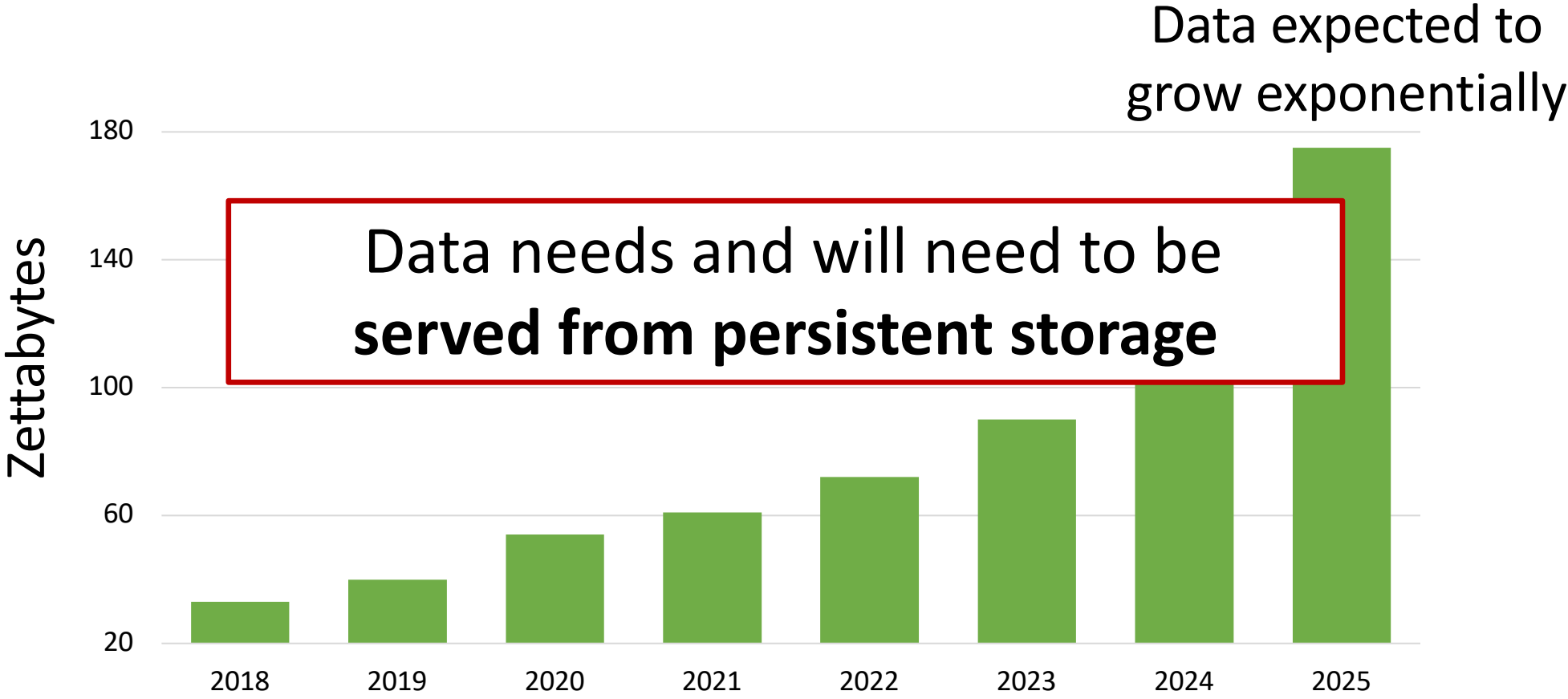


Humanity produces a lot of data



Source: IDC 2022

Humanity produces a lot of data



Data needs and will need to be served from persistent storage

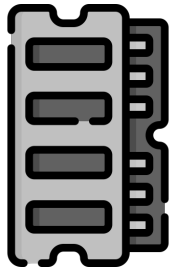
Source: IDC 2022

Data is the moving force of ML algorithms

... but in many projects the **storage decision is an afterthought**

Inefficient I/O can slow down ML Workloads

Dataset fits in system memory



Dataset = 2x system memory



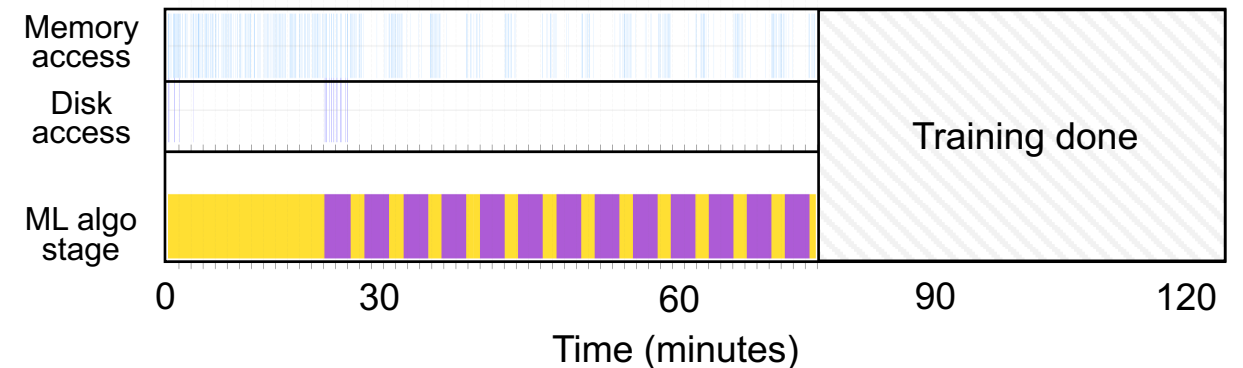
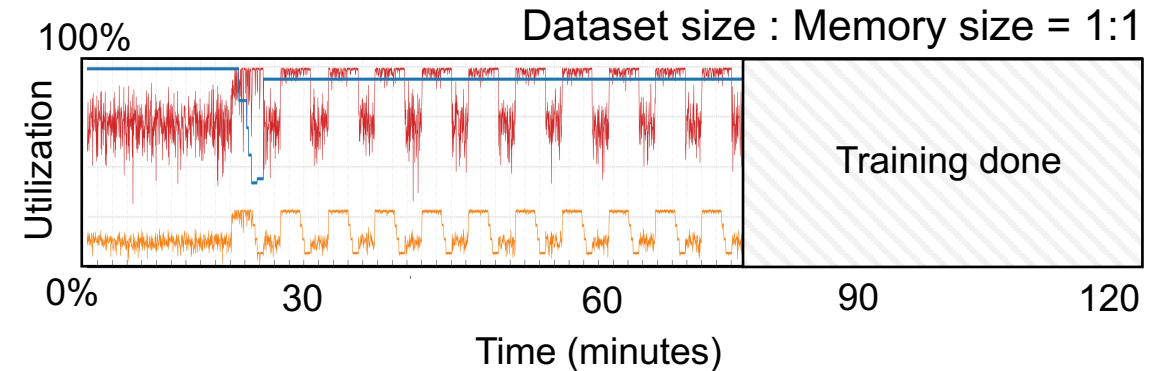
Training time increased by 3x

Inefficient I/O can slow down ML Workloads

Experiment setup

- DGX-1 server
 - 8 x V100 GPUs, 32GB GPU memory
 - 512GB DRAM
- Image segmentation workload:
 - Unet3D, Pytorch
 - MLPerf Training implementation
 - KiTS19 dataset

Dataset fits in system memory



Inefficient I/O can slow down ML Workloads

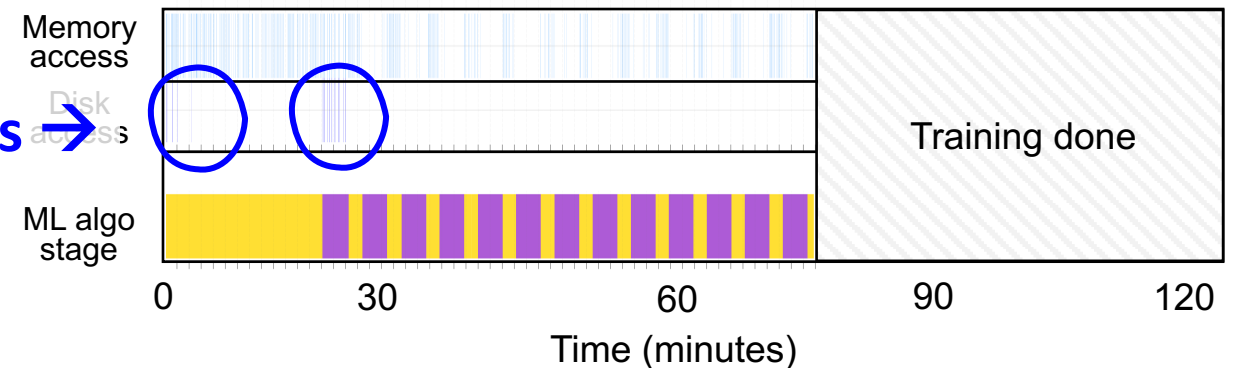
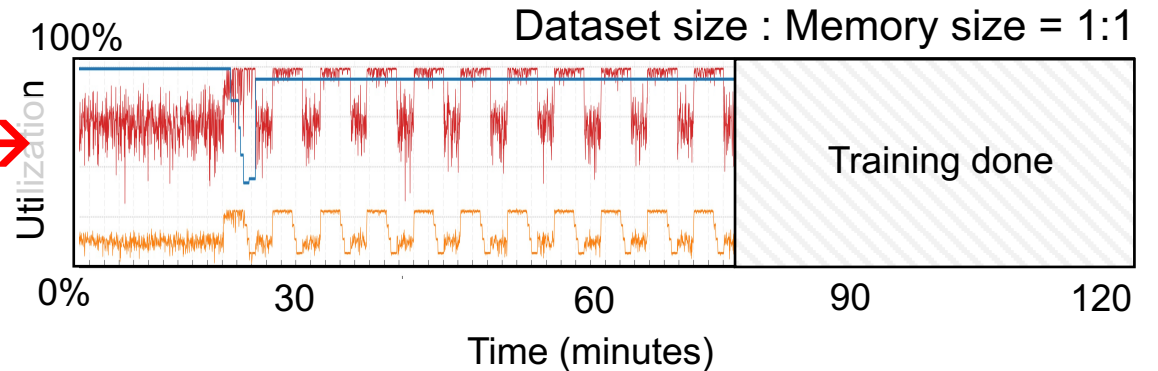
Experiment setup

- DGX-1 server
 - 8 x V100 GPUs, 32GB GPU memory
 - 512GB DRAM
- Image segmentation workload:
 - Unet3D, Pytorch
 - MLPerf Training implementation
 - KiTS19 dataset

High GPU utilization →

Little disk access →

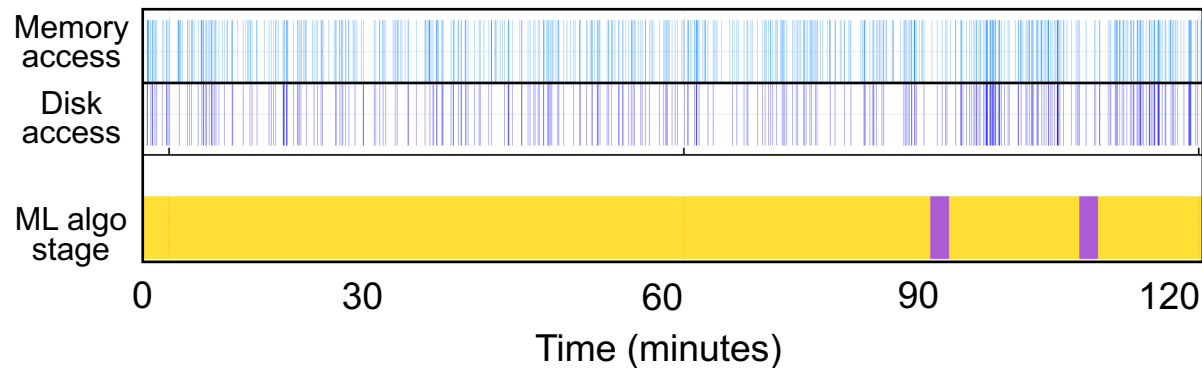
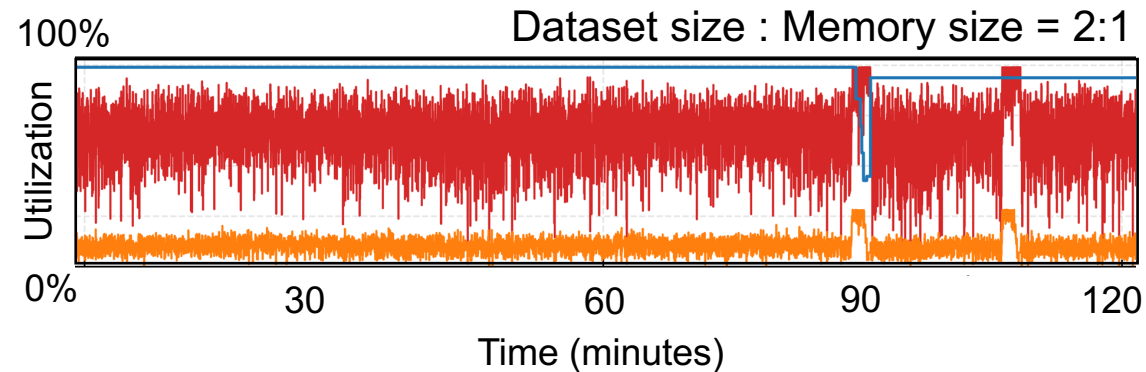
Dataset fits in system memory



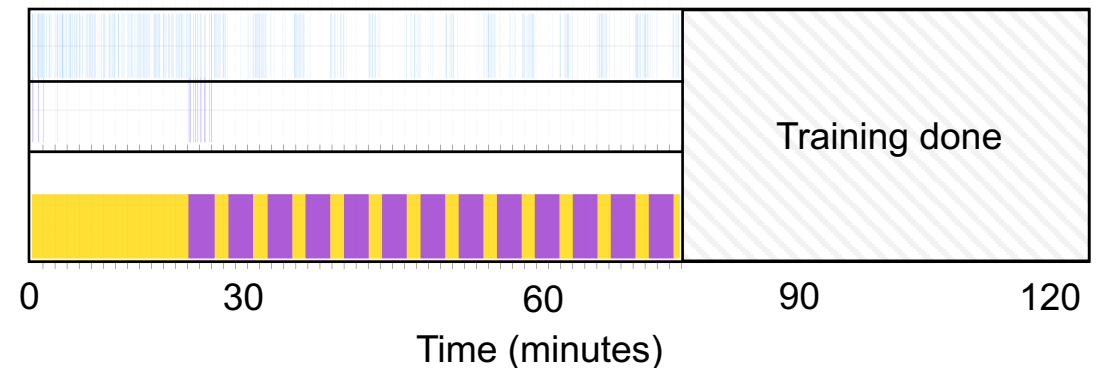
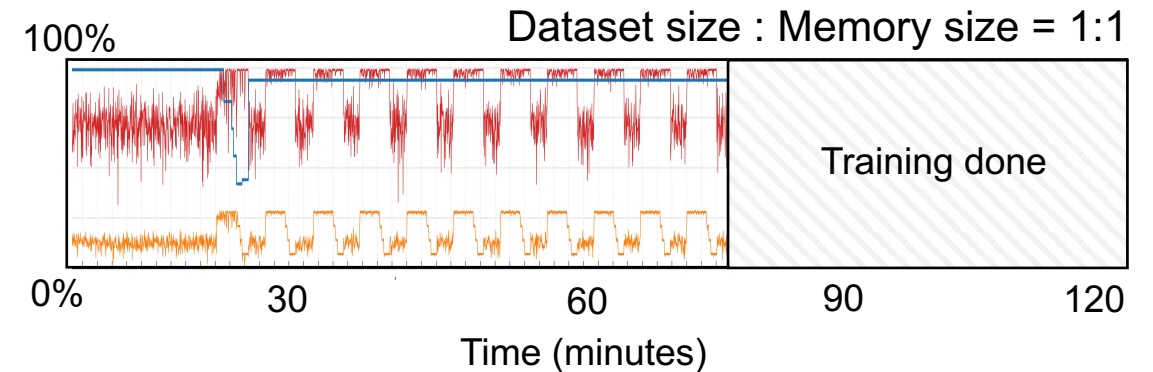
■ ML Training ■ ML Evaluation ■ Disk I/O Read ■ In-memory Read ■ GPU ■ CPU ■ GPU Memory

Inefficient I/O can slow down ML Workloads

Dataset does not fit in memory



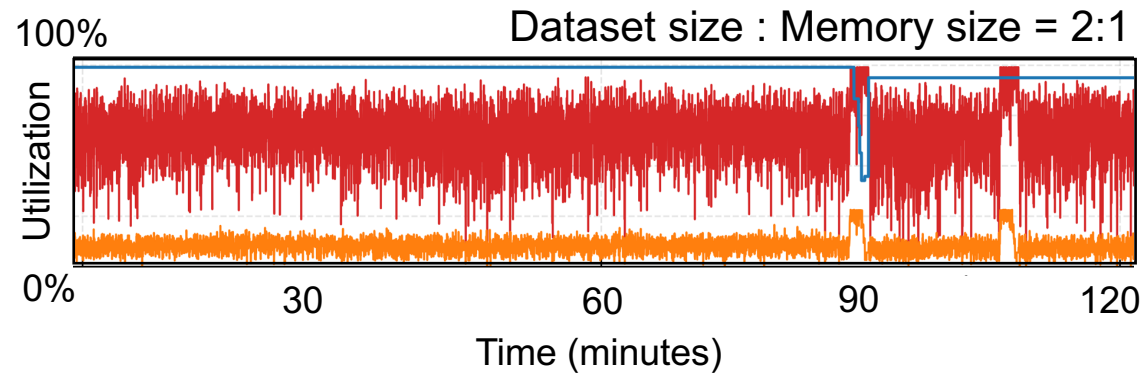
Dataset fits in system memory



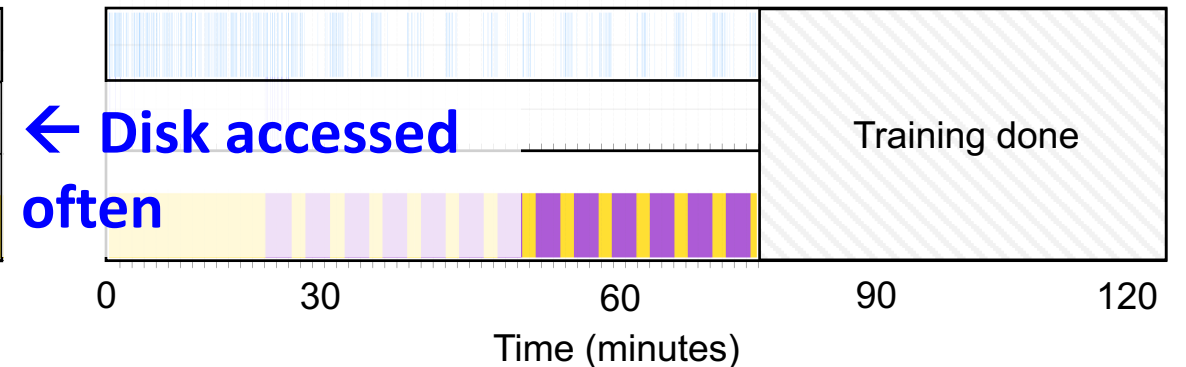
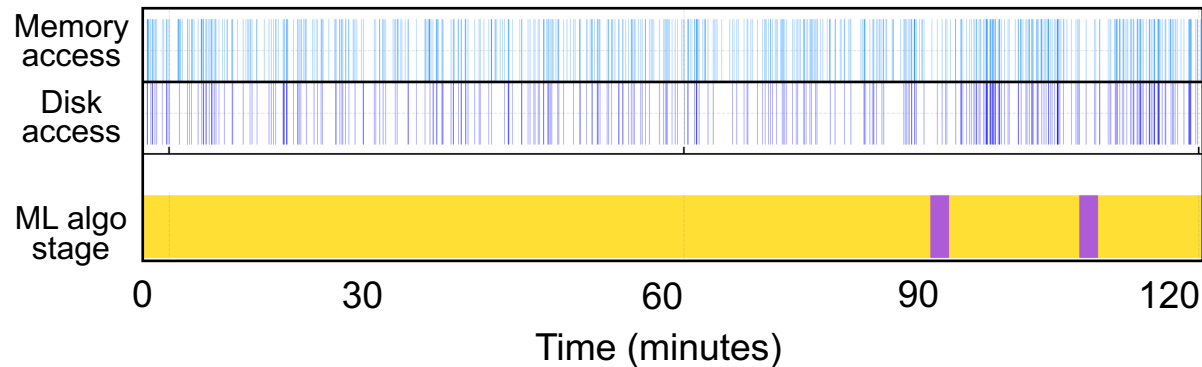
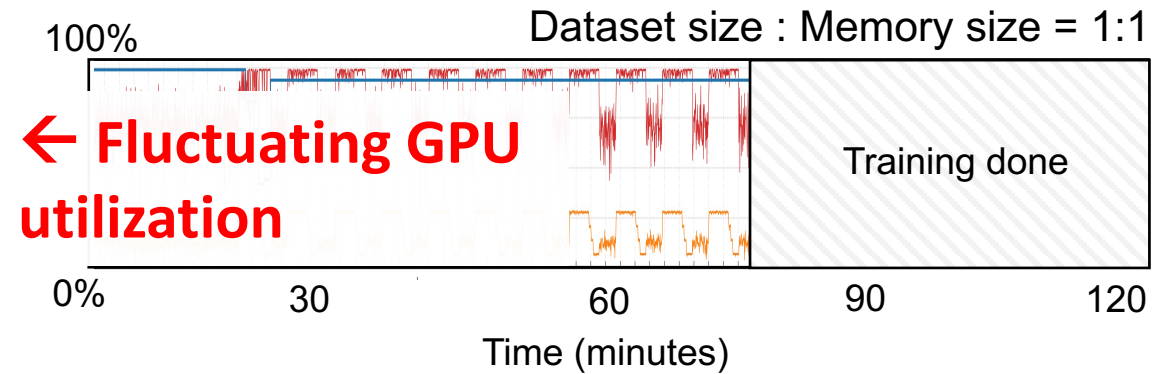
ML Training ML Evaluation Disk I/O Read In-memory Read GPU CPU GPU Memory

Inefficient I/O can slow down ML Workloads

Dataset does not fit in memory



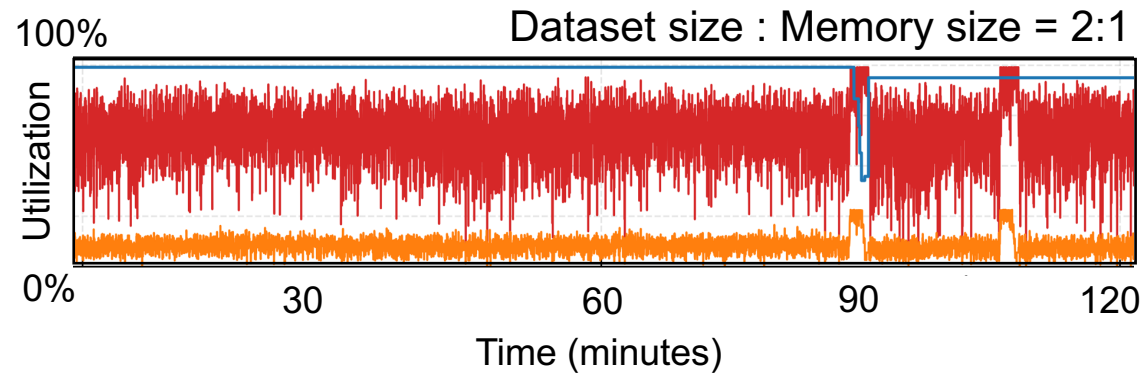
Dataset fits in system memory



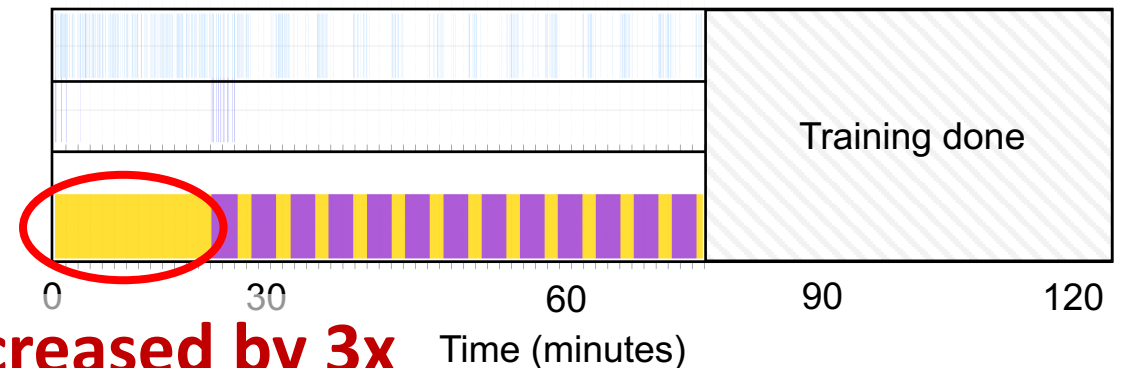
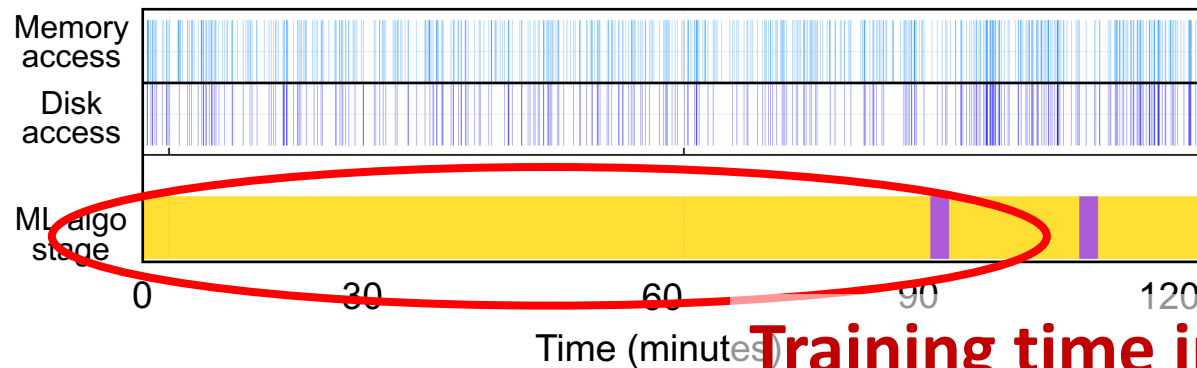
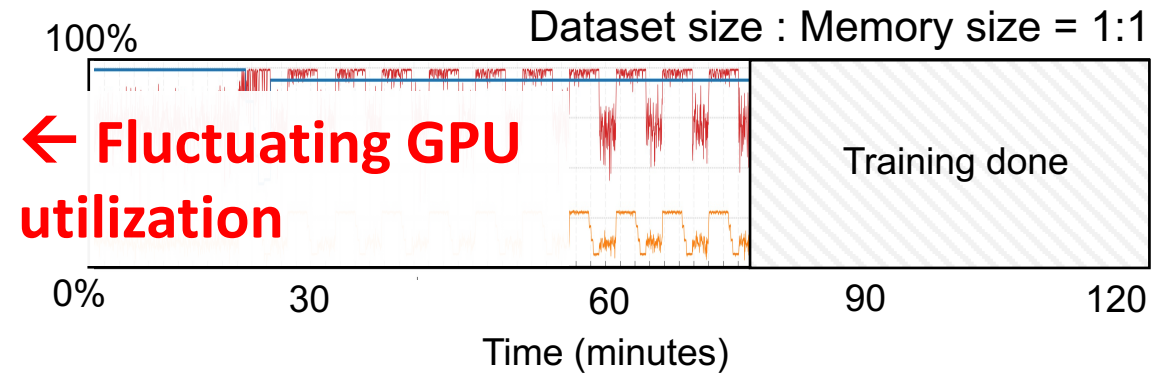
■ ML Training ■ ML Evaluation ■ Disk I/O Read ■ In-memory Read ■ GPU ■ CPU ■ GPU Memory

Inefficient I/O can slow down ML Workloads

Dataset does not fit in memory



Dataset fits in system memory



Training time increased by 3x

ML Training ML Evaluation Disk I/O Read In-memory Read GPU CPU GPU Memory

Data is the moving force of ML algorithms

... but in many projects the **storage decision is an afterthought**

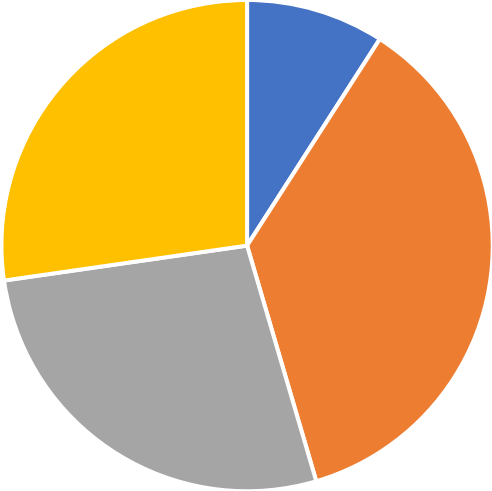
Why create an ML Storage benchmark?

Why create an ML Storage benchmark?

- Understand storage bottlenecks in ML workloads and propose optimizations
- Help AI/ML researchers and practitioners make an informed storage decision

MLPerf Storage working group

Mix of industry and academia



- Academia
- Storage Vendors
- Accelerator Vendors
- End Users



tenstorrent

Current ML/AI benchmarks

Many existing ML/AI benchmarks



DeepMind Lab



MLPerf

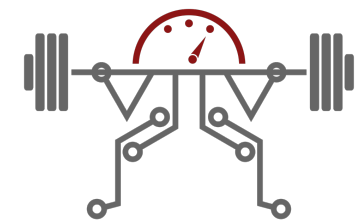


OpenAI

DLBT



PMLDB



DAWNBench

Current ML/AI benchmarks

- Focus on **end-to-end testing**
 - hard to isolate value of each component
- Insist on **training and inference** speed
 - tend to simplify storage
 - ignore pre-processing
- **Expensive accelerators** needed to run
- Require **extensive entry knowledge**



DeepMind Lab



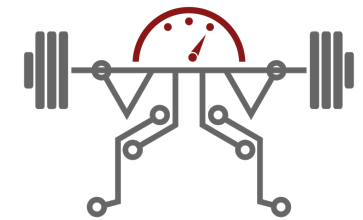
MLPerf



OpenAI

DLBT 

PMLDB



DAWNBench

Benchmark Vision

Existing benchmarks

Focus on **end-to-end testing**

Simplified storage setup

Expensive accelerators needed to run

Require **extensive entry knowledge**

Our work

Focus on **storage impact in ML/AI**

Realistic **storage & pre-processing** settings

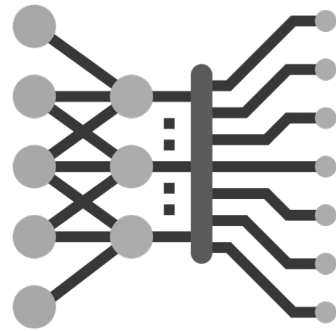
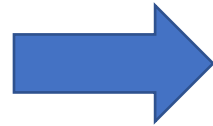
No accelerator required to run

Minimal AI/ML knowledge required

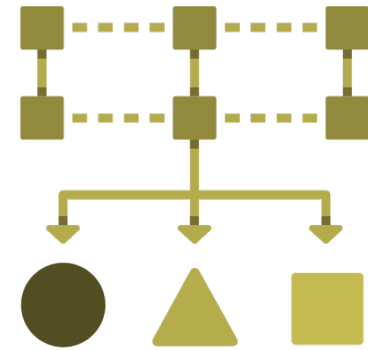
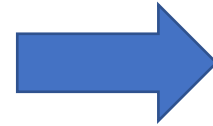
Stages of the ML Pipeline



Data cleaning & pre-processing

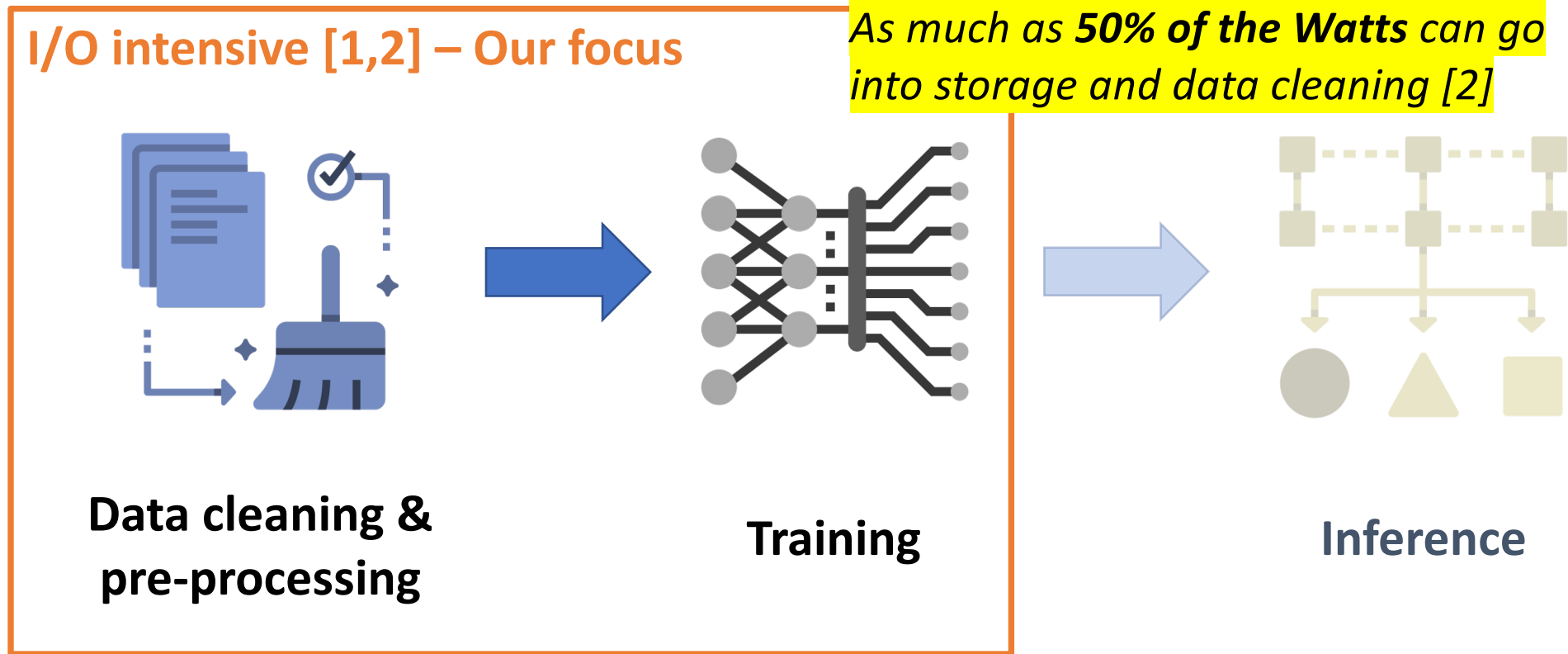


Training



Inference

Stages of the ML Pipeline



[1] Murray et al. *tf.data: A Machine Learning Data Processing Framework*, VLDB 21.

[2] Zhao et al. *Understanding Data Storage and Ingestion for Large-Scale Deep Recommendation Model Training* ISCA 22.



Data Pipeline in ML: Pre-processing

Storage resources

Disk

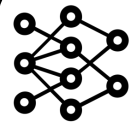


Memory

Compute resources

CPUs

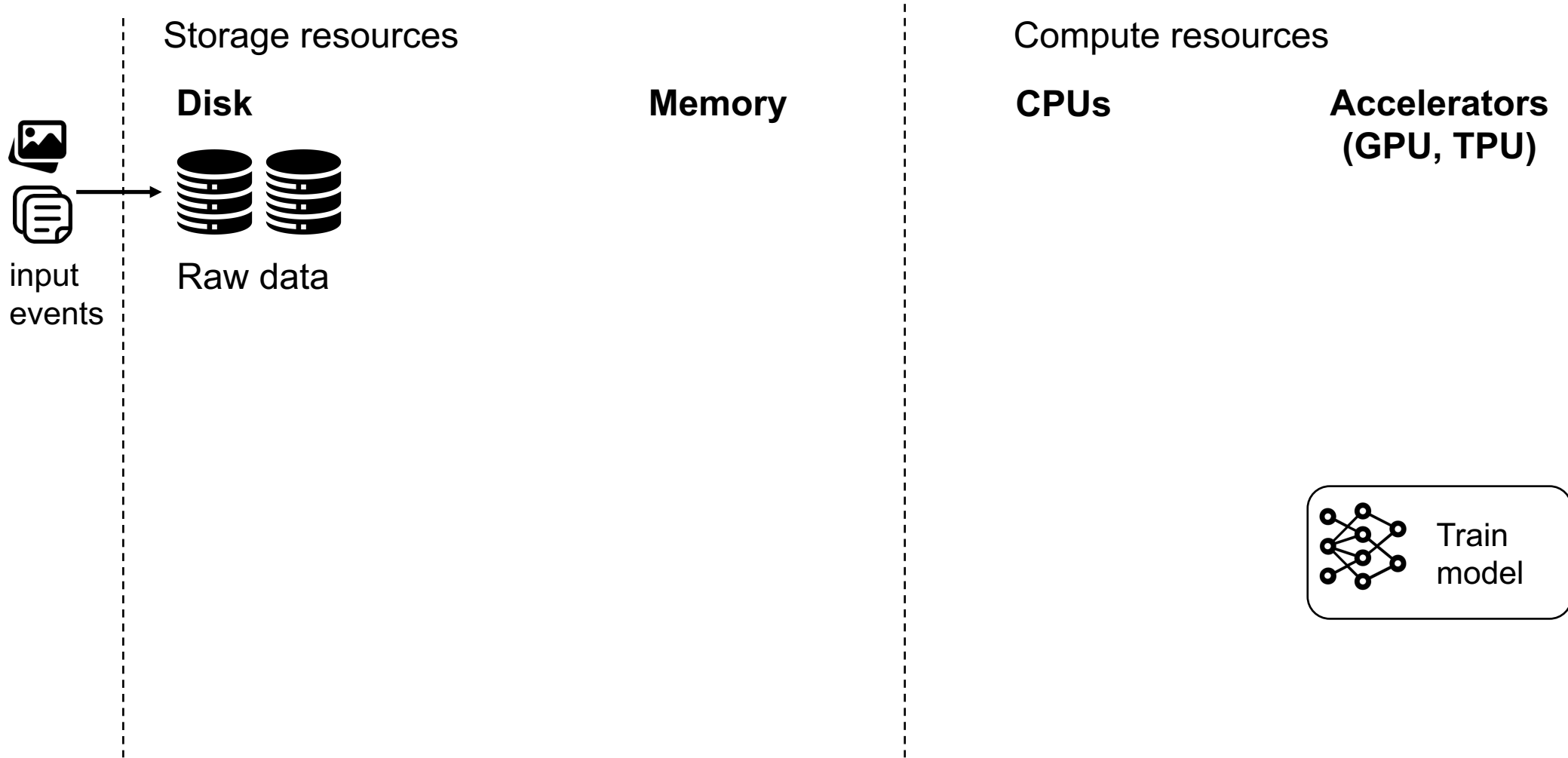
**Accelerators
(GPU, TPU)**

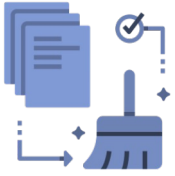


Train
model

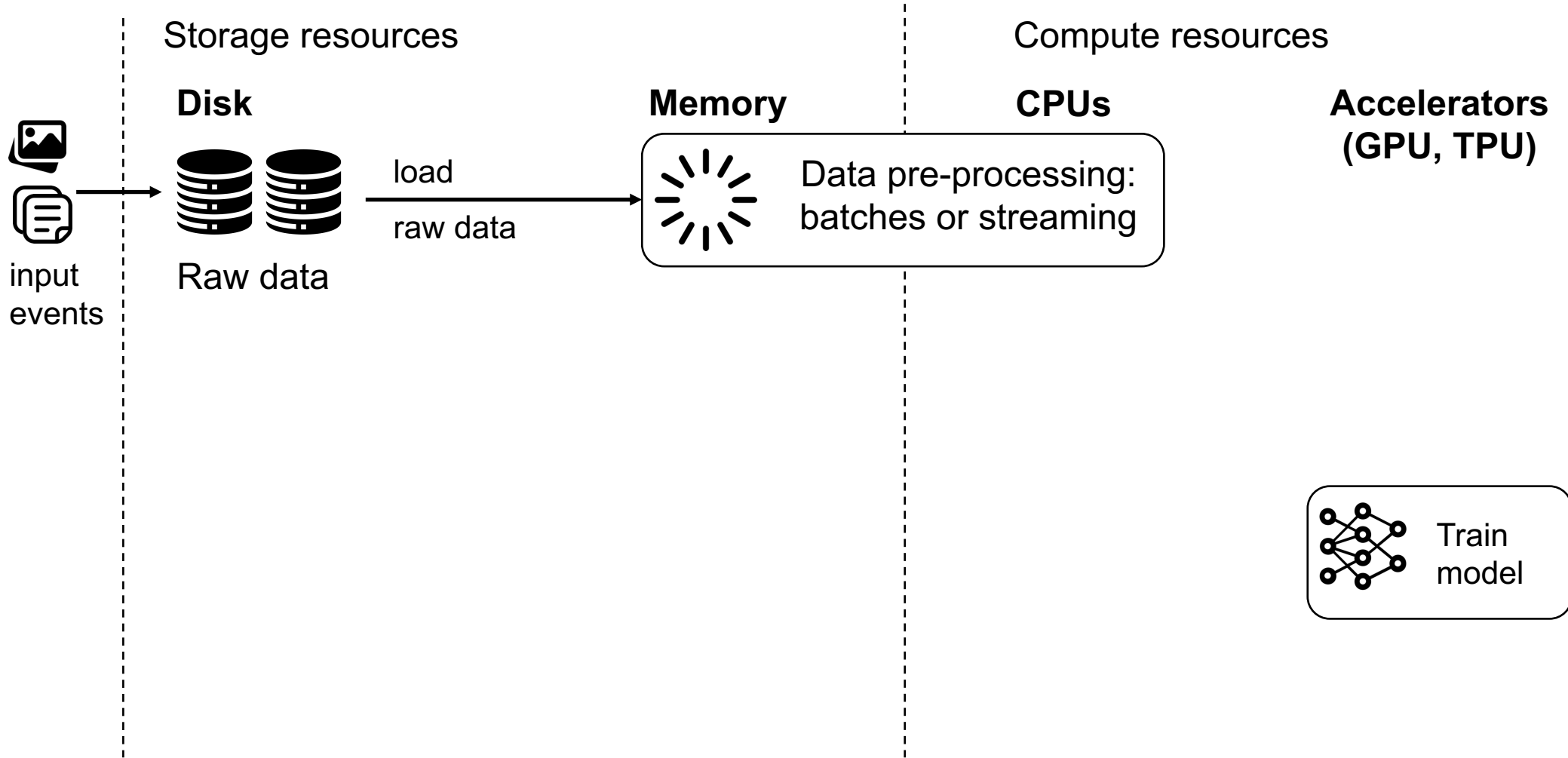


Data Pipeline in ML: Pre-processing



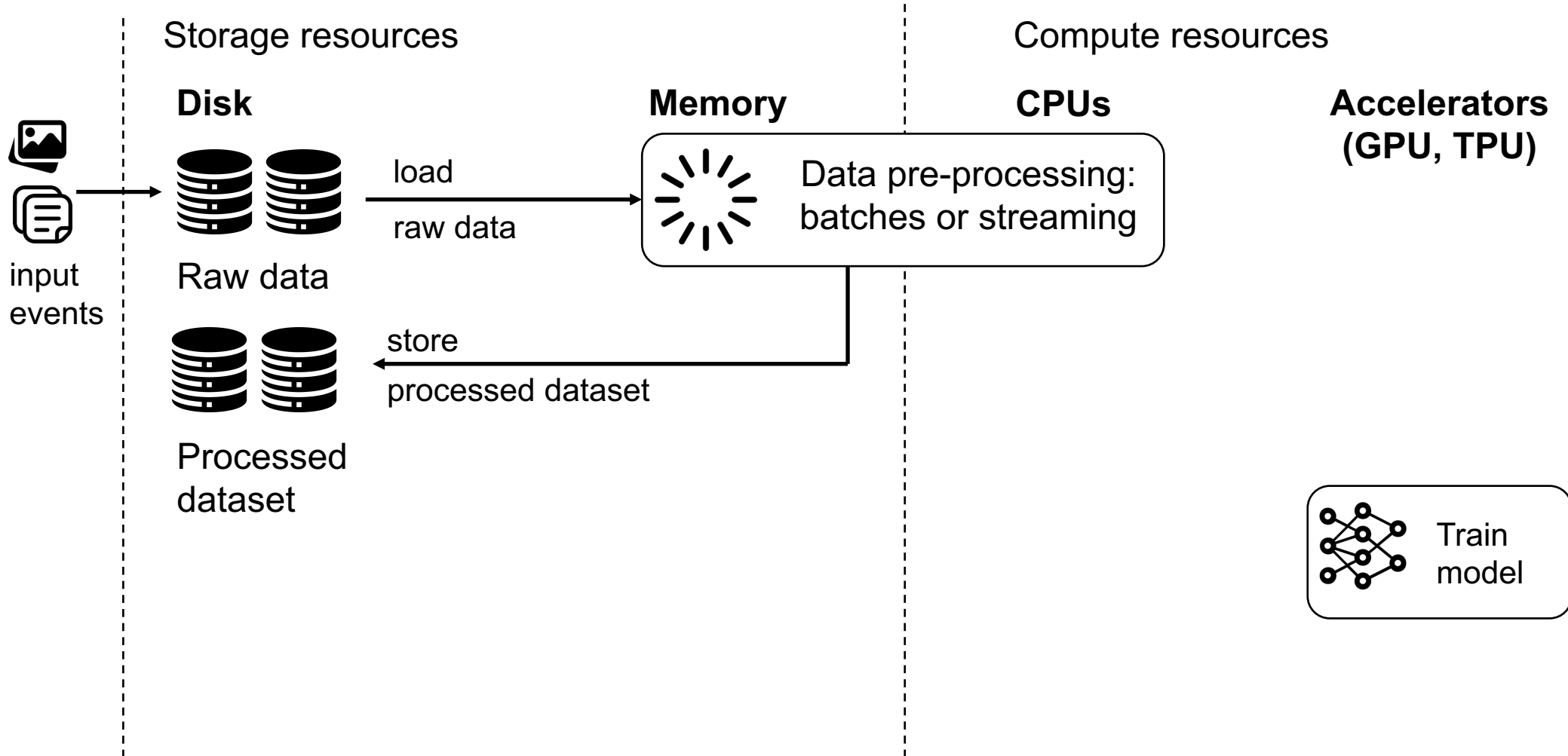


Data Pipeline in ML: Pre-processing

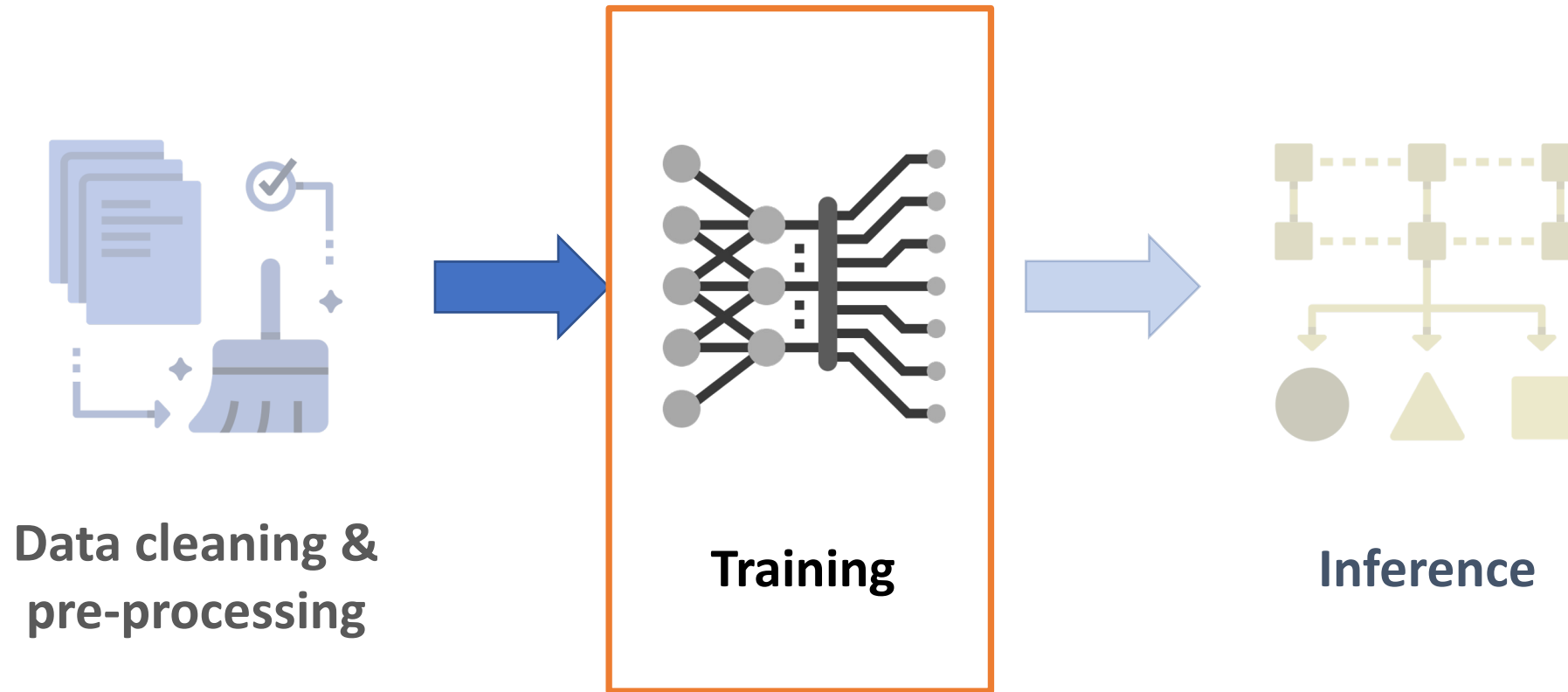




Data Pipeline in ML: Pre-processing

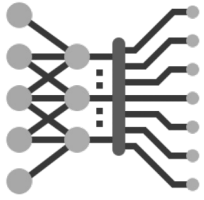


Stages of the ML Pipeline



[1] Murray et al. *tf.data: A Machine Learning Data Processing Framework*, VLDB 21.

[2] Zhao et al. *Understanding Data Storage and Ingestion for Large-Scale Deep Recommendation Model Training* ISCA 22.



Data pipeline in ML: Training

Storage resources

Disk



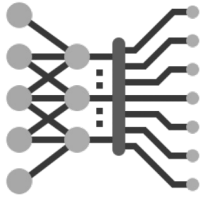
Cleaned
dataset

**System
Memory (DRAM)**

Compute resources

CPUs

**Accelerators
(GPU, ASIC)**



Data pipeline in ML: Training

Storage resources

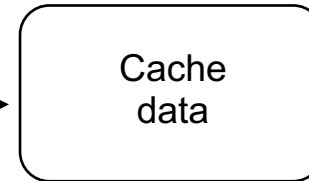
Disk



Cleaned dataset

 TensorFlow
PYTORCH

load
data

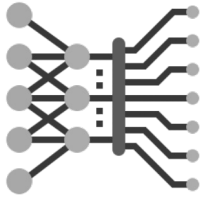


System
Memory (DRAM)

Compute resources

CPUs

Accelerators
(GPU, ASIC)



Data pipeline in ML: Training

Storage resources

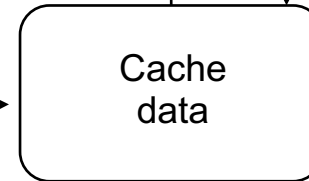
Disk



Cleaned dataset

 TensorFlow
PYTORCH

load
data



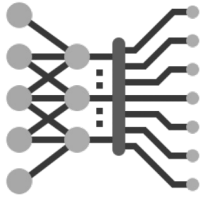
System
Memory (DRAM)

Compute resources

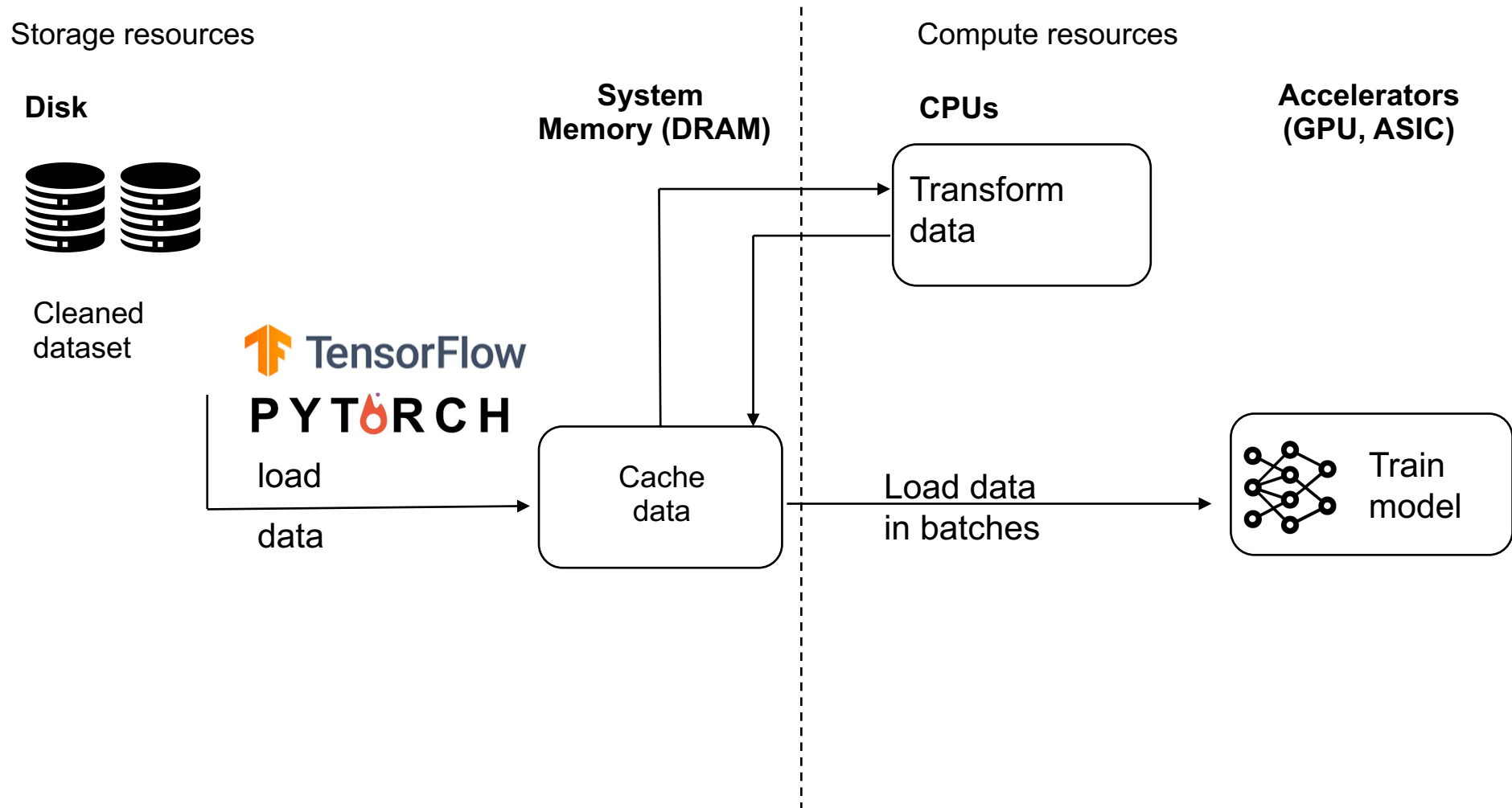
CPUs



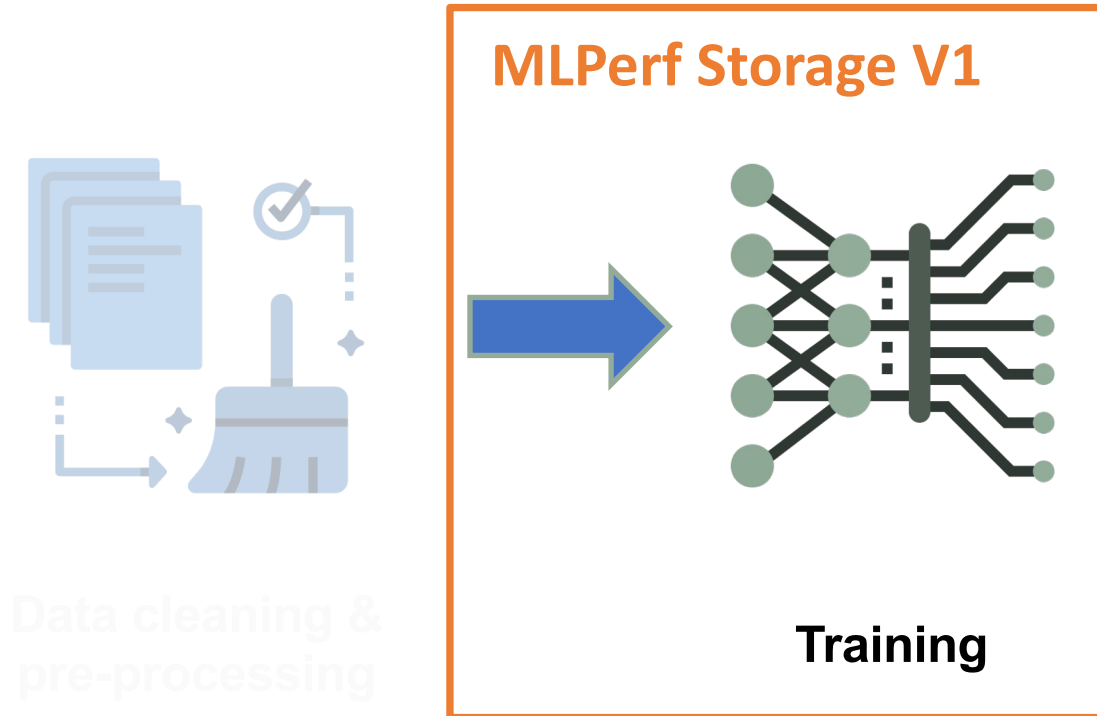
Accelerators
(GPU, ASIC)



Data pipeline in ML: Training



MLPerf Storage



Focus on **storage impact in ML/AI**

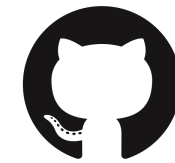
Realistic **storage** settings in
training phase

No accelerator required to run

Minimal AI/ML knowledge

MLPerf Storage – workloads

Workload	Image segmentation	Natural language processing	Recommender Systems
Model	Unet3D	BERT	DLRM
Seed data	KiTS19 Set of images	Wikipedia 2020 Text	Criteo Terabyte Click logs
Framework	Pytorch	Tensorflow	Pytorch
I/O behavior	Random access inside many small files	Sequential access of small subset of files, streamed.	Random access inside one large file

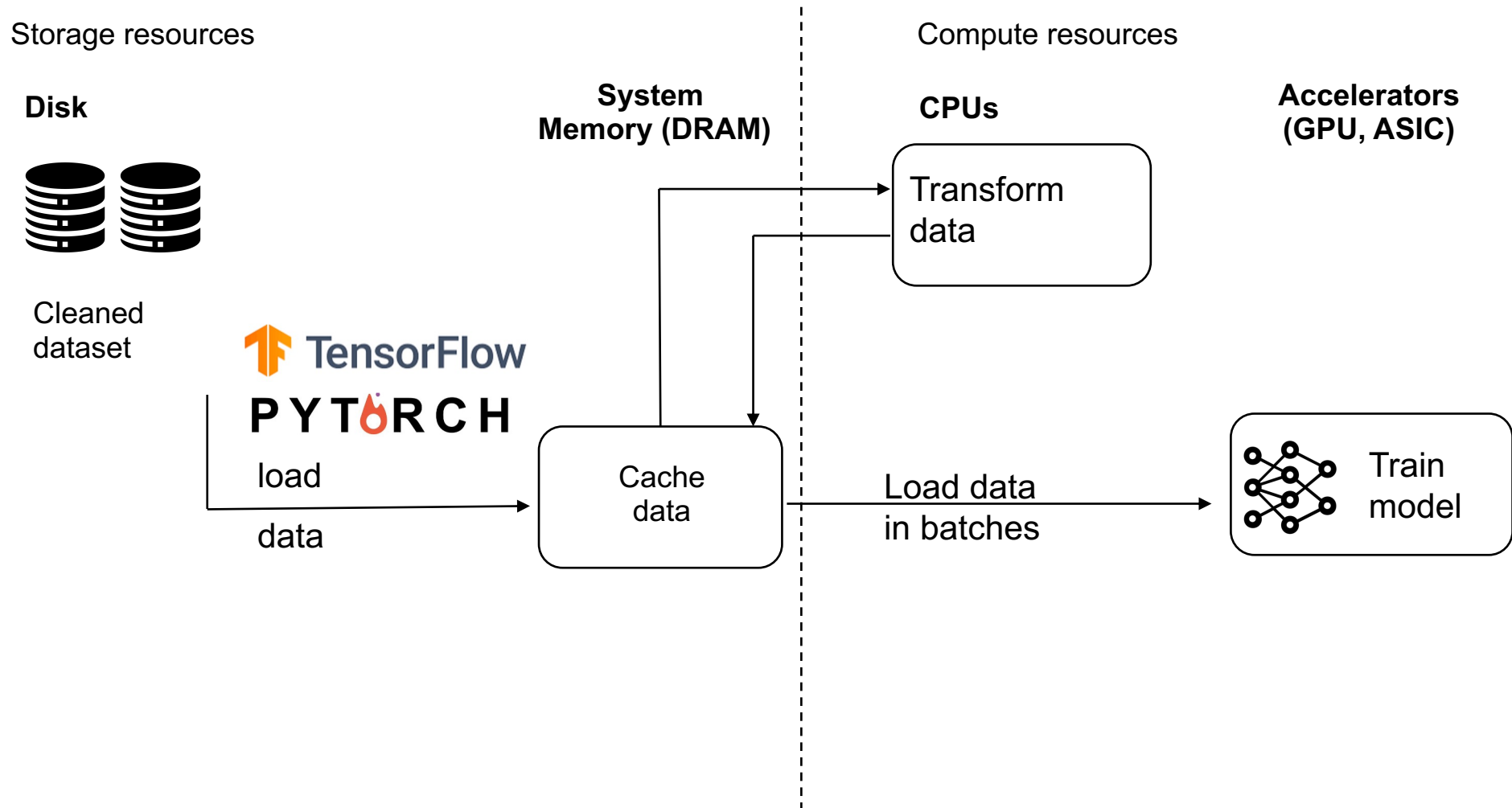


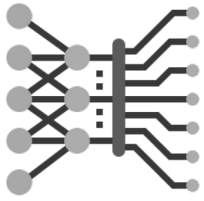
<https://github.com/mlcommons/storage>

- **Single node**
- **Synthetic datasets** generated from real dataset seed.
- Many **simulated accelerators.**
- **Local storage**

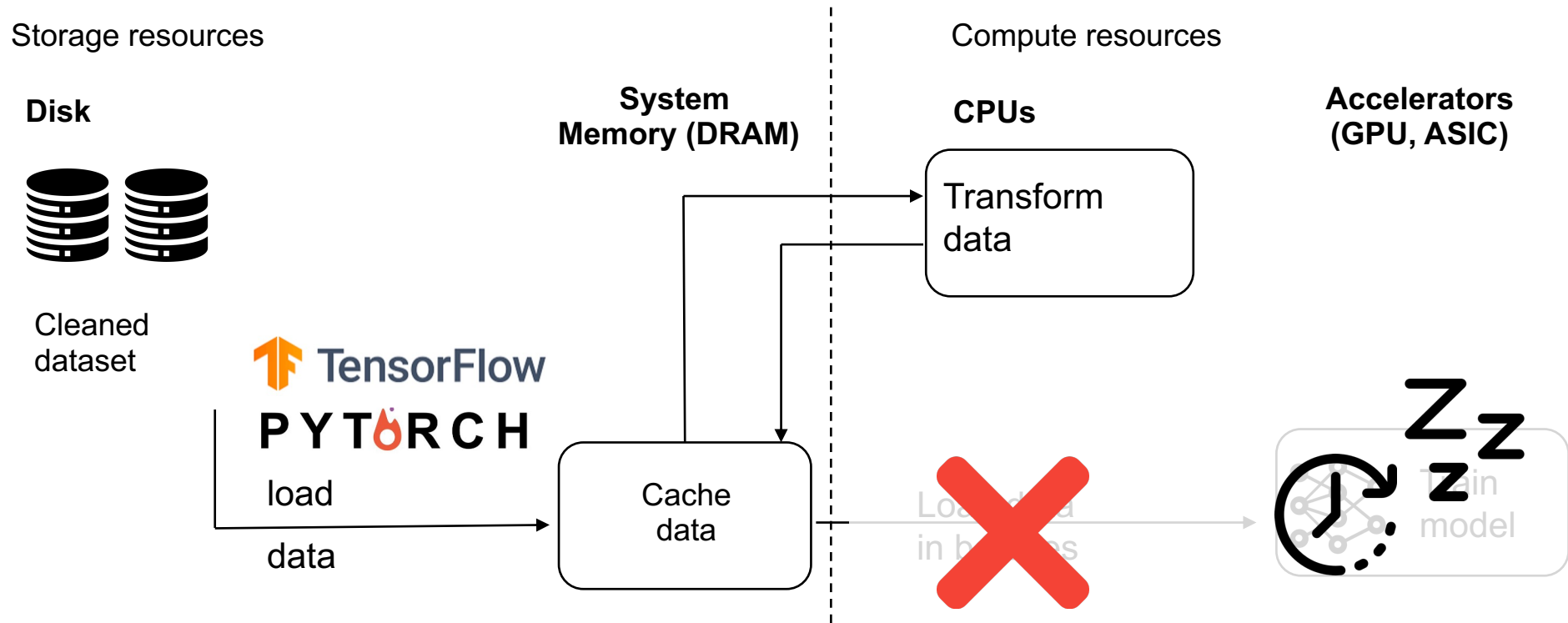


Data pipeline in ML: Training



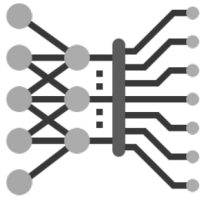


Data pipeline in MLPerf Storage benchmark



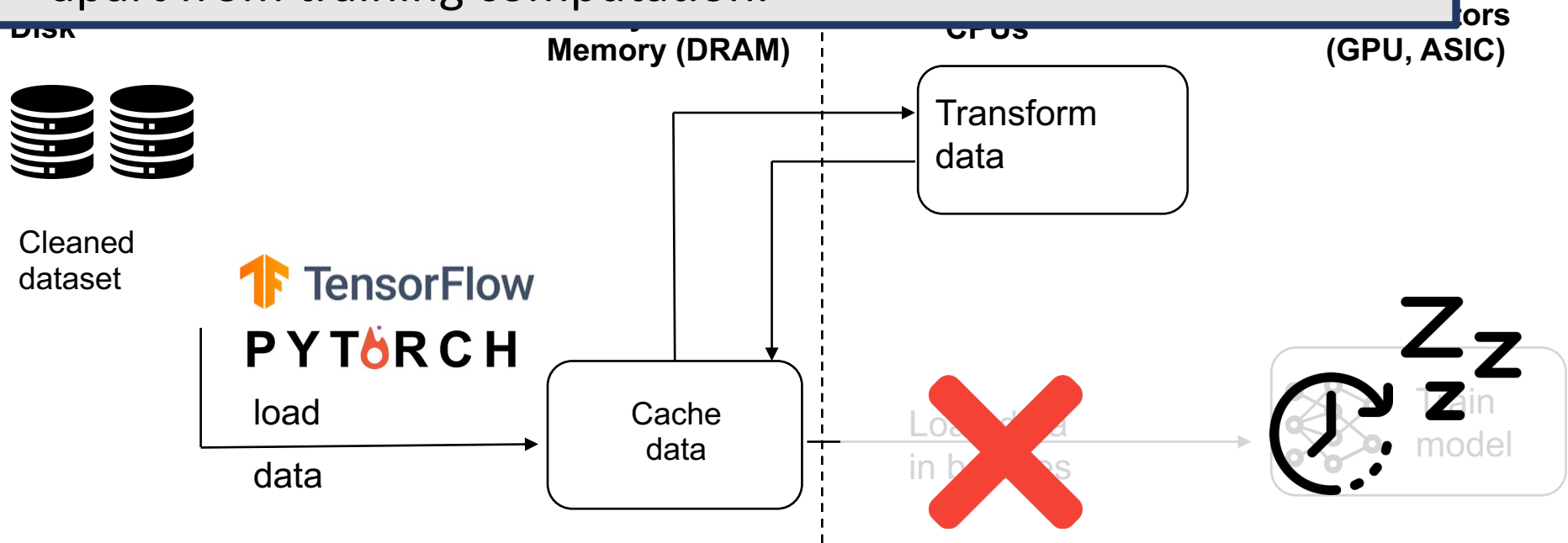
Benchmark is built as an extension of DLIO [1]

[1] H. Devarajan, H. Zheng, et al. DLIO: A Data-Centric Benchmark for Scientific Deep Learning Applications, CCGrid '21.



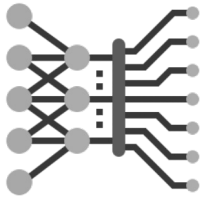
Data pipeline in MLPerf Storage benchmark

✓ Realistic storage settings: nothing changes in data pipeline, apart from training computation.

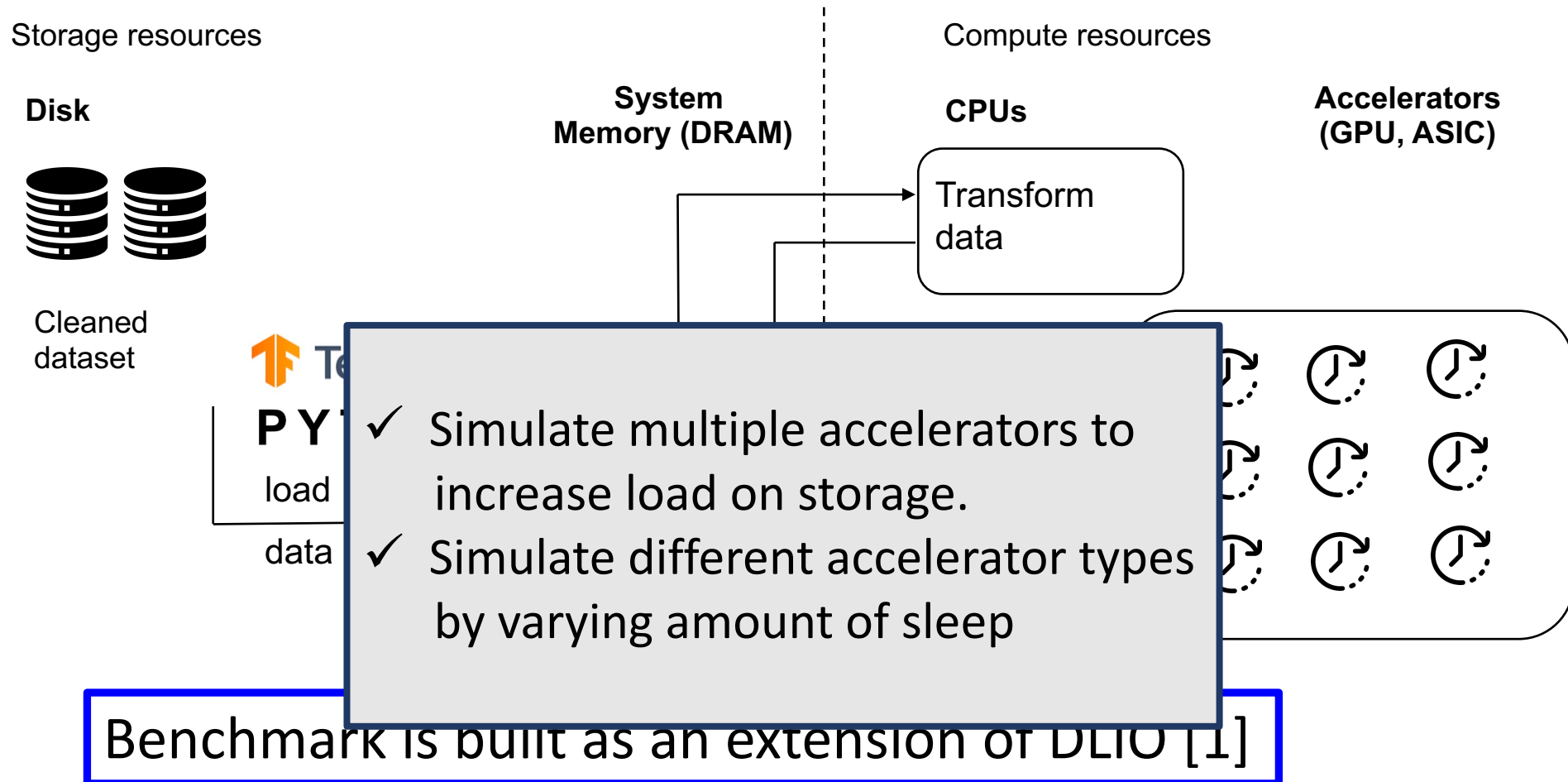


Benchmark is built as an extension of DLIO [1]

[1] H. Devarajan, H. Zheng, et al. DLIO: A Data-Centric Benchmark for Scientific Deep Learning Applications, CCGrid '21.



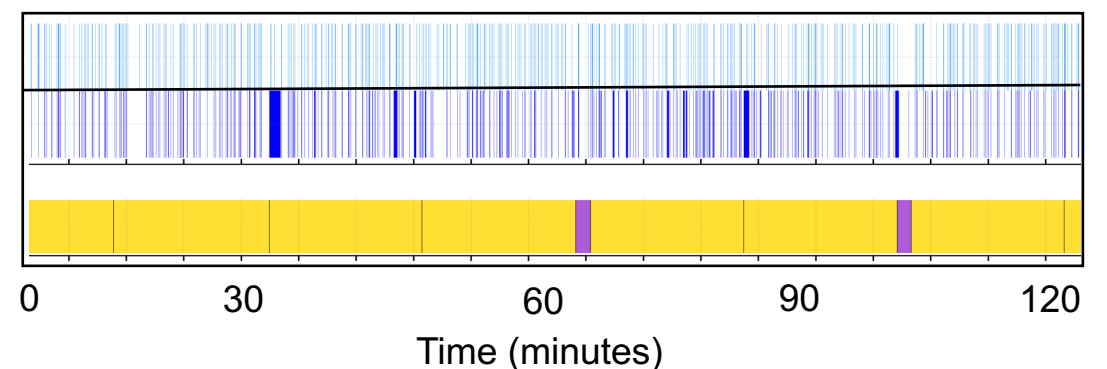
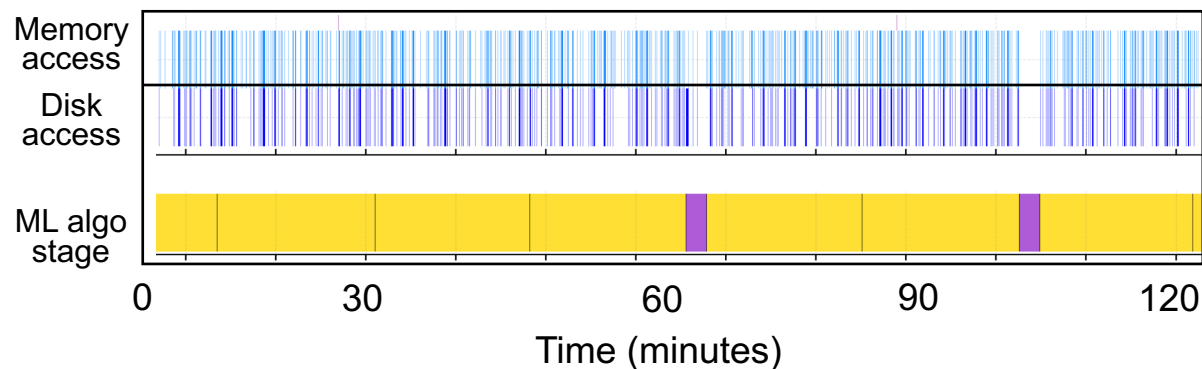
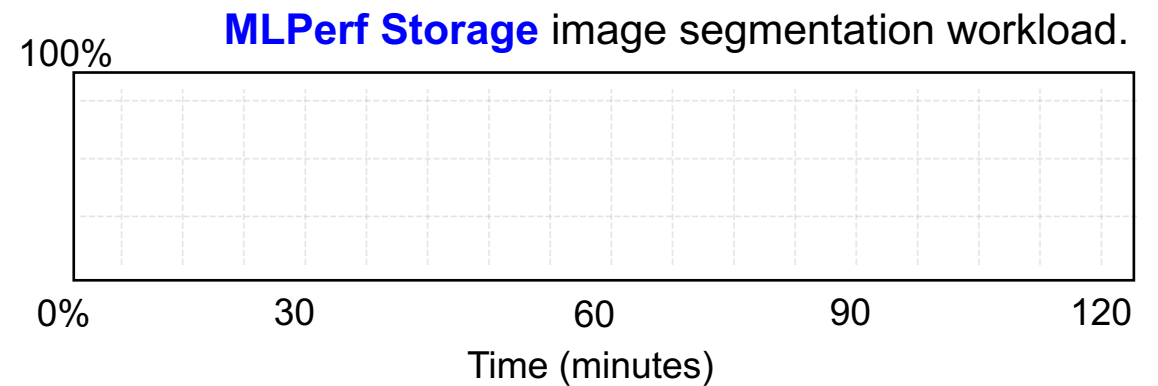
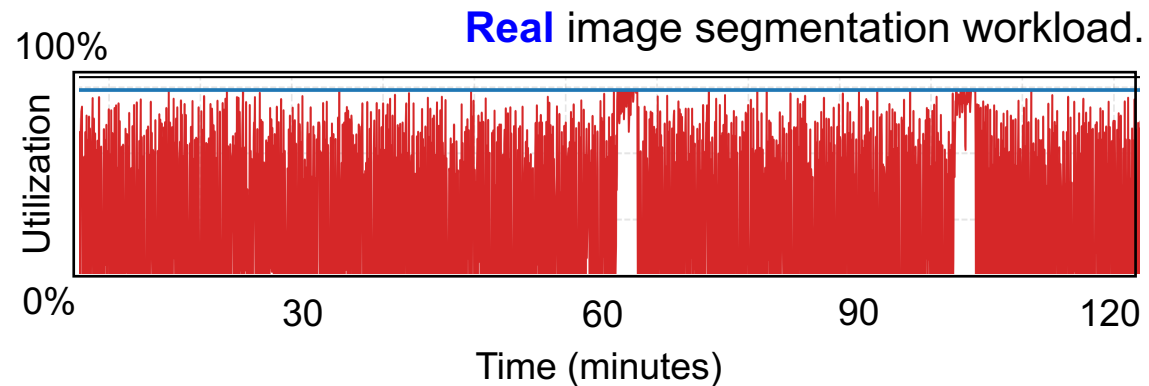
Data pipeline in MLPerf Storage benchmark



[1] H. Devarajan, H. Zheng, et al. DLIO: A Data-Centric Benchmark for Scientific Deep Learning Applications, CCGrid '21.

Simulating training time does not impact I/O patterns

■ ML Training ■ ML Evaluation ■ Disk I/O Read ■ In-memory Read ■ GPU

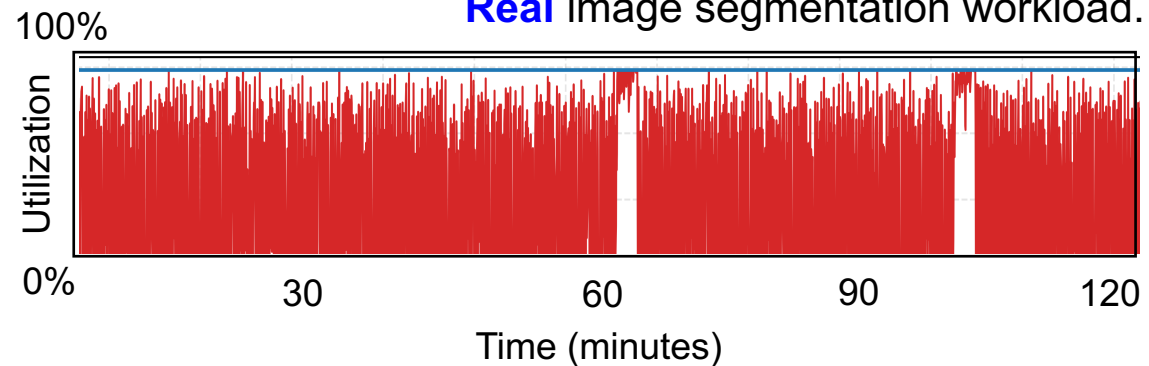


Experiment setup: DGX-1 with 8xV100 GPUs, 512GB DRAM. Dataset : KiTS19, Dataset size:Memory size ratio 2:1

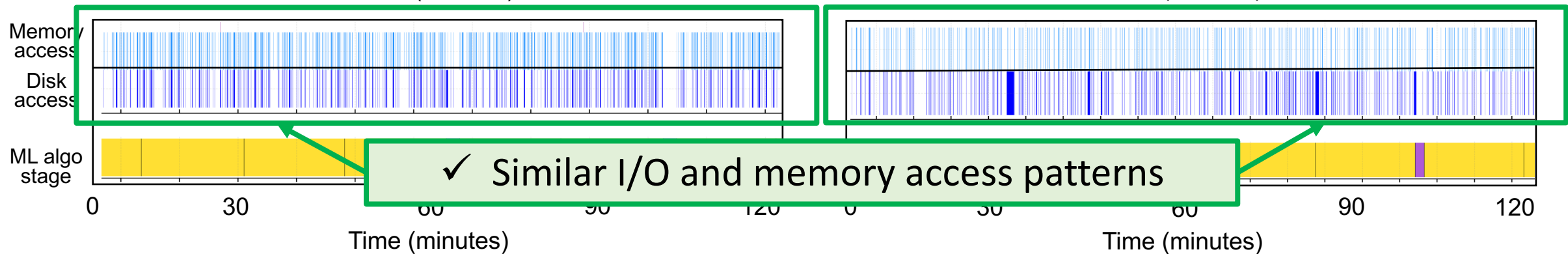
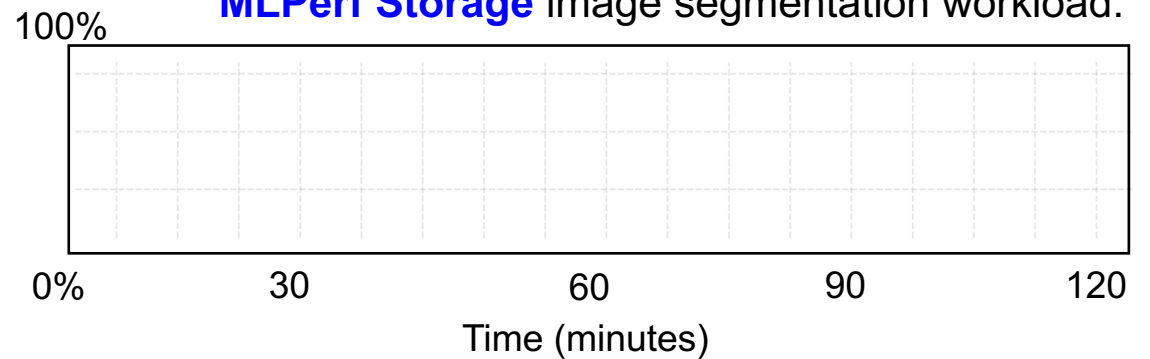
Simulating training time does not impact I/O patterns

■ ML Training ■ ML Evaluation ■ Disk I/O Read ■ In-memory Read ■ GPU

Real image segmentation workload.

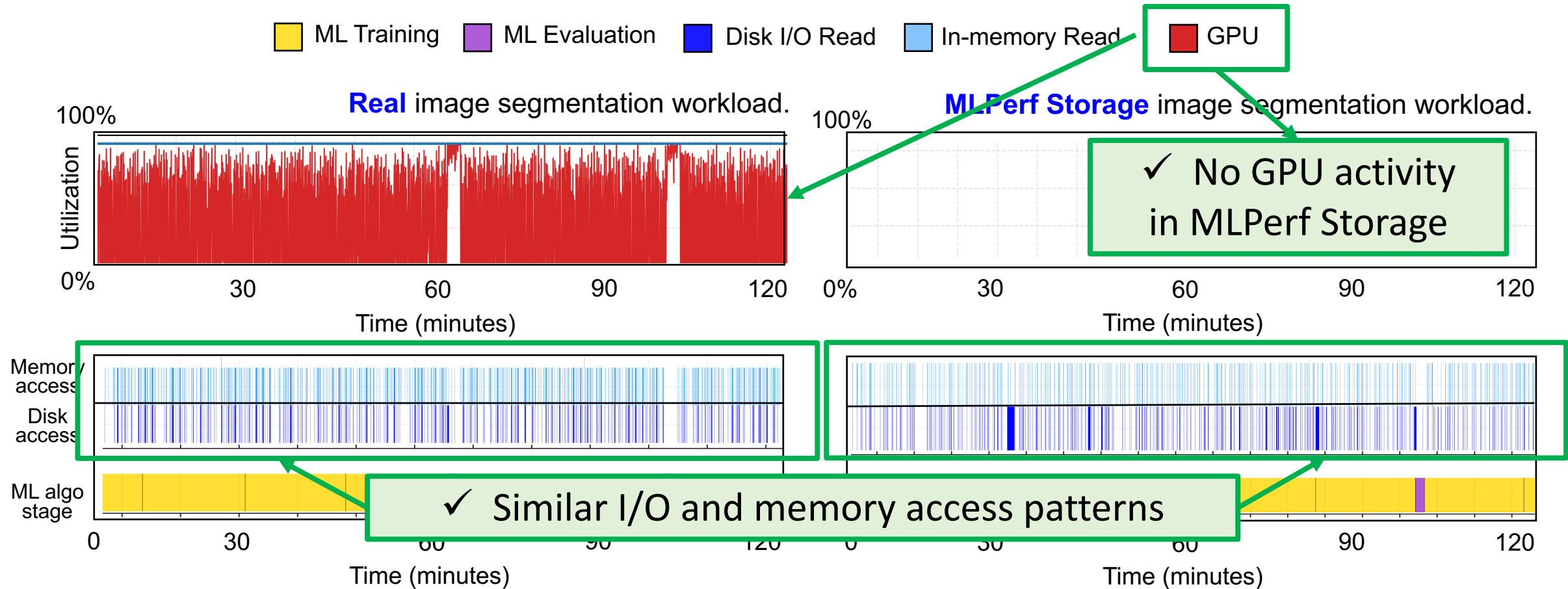


MLPerf Storage image segmentation workload.



Experiment setup: DGX-1 with 8xV100 GPUs, 512GB DRAM. Dataset : KiTS19, Dataset size:Memory size ratio 2:1

Simulating training time does not impact I/O patterns



Experiment setup: DGX-1 with 8xV100 GPUs, 512GB DRAM. Dataset : KiTS19, Dataset size:Memory size ratio 2:1

Next Steps

Collect **processing times** for different accelerator types.

Open benchmark for submissions.

→ <https://github.com/mlcommons/storage>

I/O in distributed training

Trace and benchmark **ML pre-processing phase.**

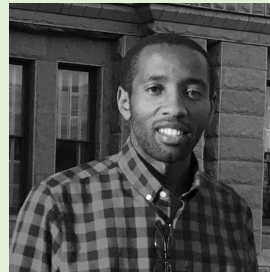
McGill DISCS Lab

Postdoctoral
Researcher



Dr. Stella Bitchebe

PhD
Candidates:



Nelson Bore



Jiaxuan Chen



discslab.cs.mcgill.ca
gitlab.cs.mcgill.ca/discs-lab

Masters
Students



Sebastian Rolon



Loïc Ho-Von



Aayush Kapur



Aidan Goldfarb

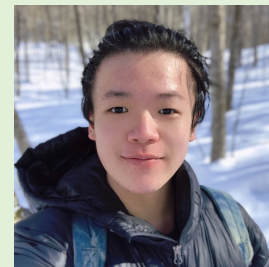


Rahma Nouaji

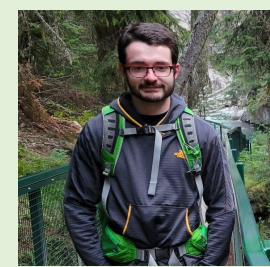
Undergraduate
Students



Zachary Doucet



Zhongjie Wu



Olivier Michaud

Key Takeaways – MLPerf Storage

MLPerf Storage is a new benchmark

Realistic **storage** settings

No accelerators required to run

Follow MLPerf Storage repository for updates:

<https://github.com/mlcommons/storage>

Get involved

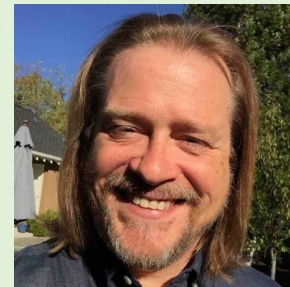
mlcommons.org/en/get-involved/

We appreciate your feedback

Share your thoughts

Email oana.balmau@mcgill.ca

Thanks to all working group co-chairs!



Curtis Anderson
Panasas



Huihuo Zheng
Argonne National Labs



Johnu George,
Nutanix