# Characterizing I/O in Machine Learning Workloads
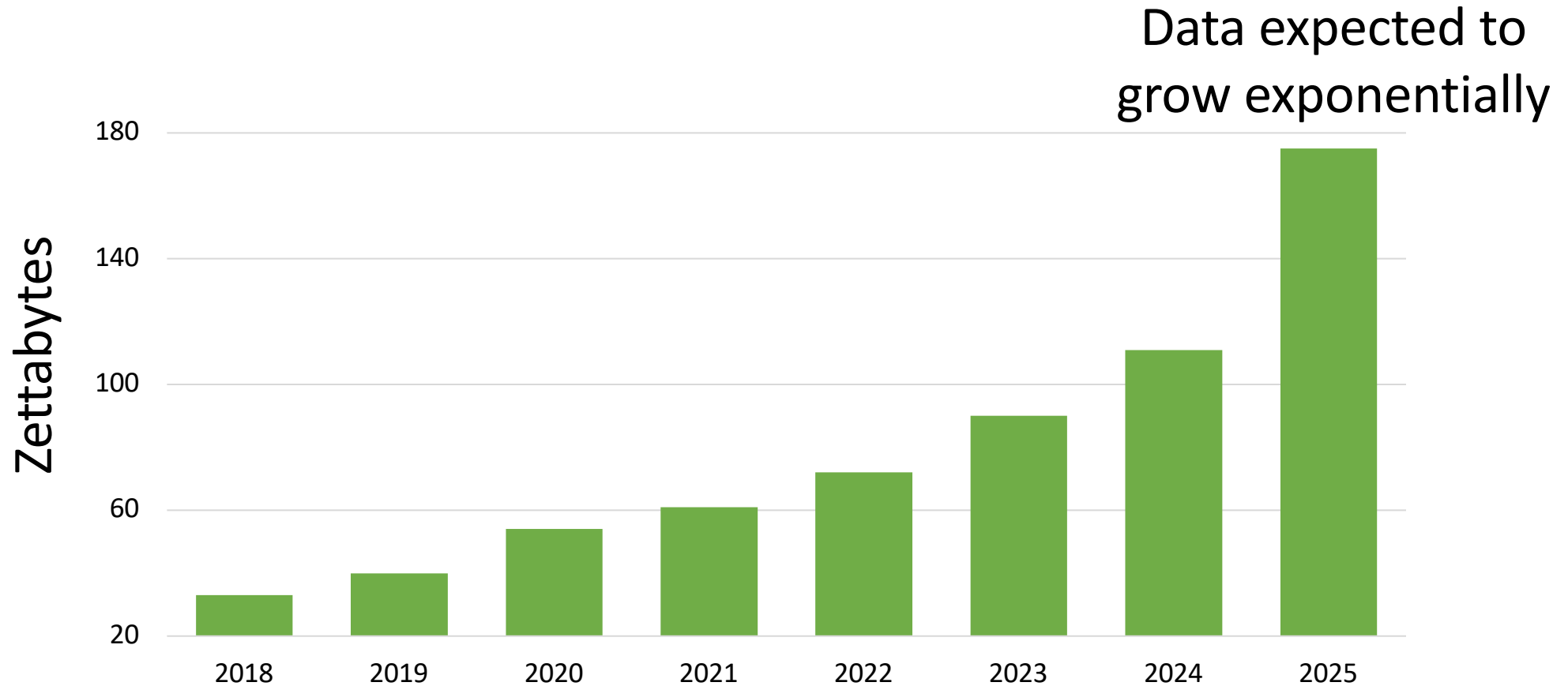
*Oana Balmau*

*Resource-Aware ML Day @ ITU, February 13th, 2023*
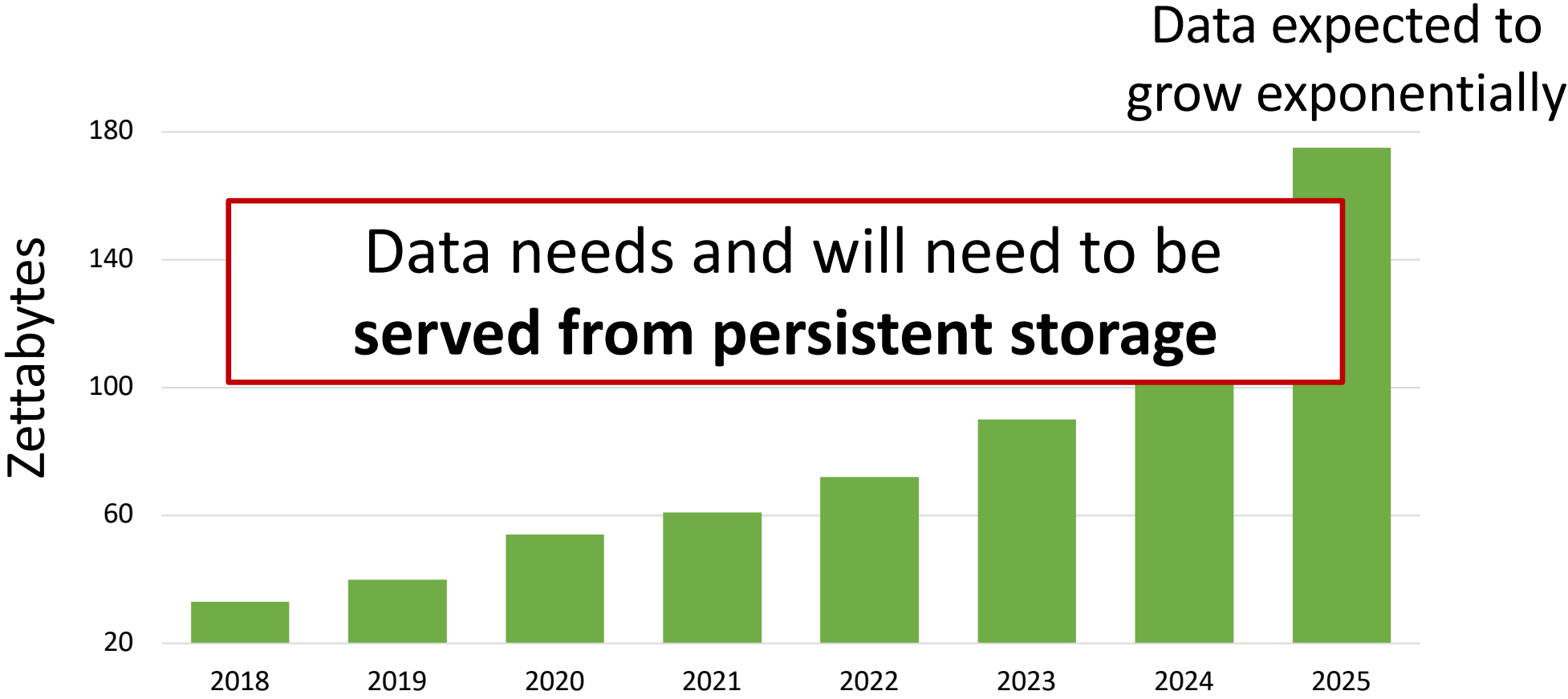
Oana Balmau

# Humanity produces a lot of data

Data expected to grow exponentially



*Source:* IDC 2022

# Humanity produces a lot of data

Data expected to grow exponentially

Data needs and will need to be **served from persistent storage**

Zettabytes

180

140

100

60

20

2018 2019 2020 2021 2022 2023 2024 2025

*Source:* *IDC 2022*

Data is the moving force of ML algorithms

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

Data is the moving force of ML algorithms

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

Why create an ML Storage benchmark?

# Current ML/AI benchmarks

**Many existing ML/AI benchmarks**
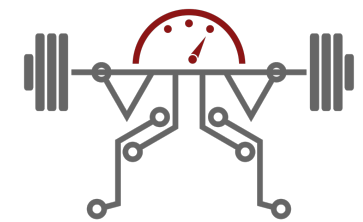
DeepMind Lab

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

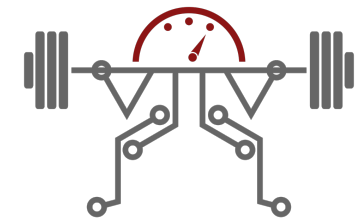- Require **extensive entry knowledge**

DeepMind Lab

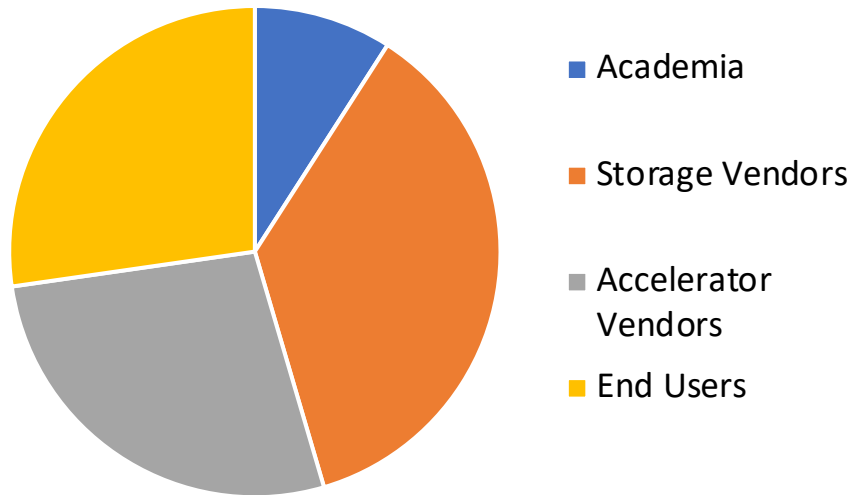MLPerf     OpenAI

DLBT

PMLDB     DAWNBench

# Why create an ML Storage benchmark?

- **Understand <u>storage</u> bottlenecks in ML workloads**
**and propose optimizations**

- **Help AI/ML researchers and practitioners**
***make an informed <u>storage</u> decision***

# MLPerf Storage working group
# Who are we?

Mix of industry and academia



- ■ Academia
- ■ Storage Vendors
- ■ Accelerator Vendors
- ■ End Users

McGill

NVIDIA

intel

NUTANIX

WEKA

Red Hat

PANASAS

Micron

Argonne
NATIONAL LABORATORY

tenstorrent

# 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**

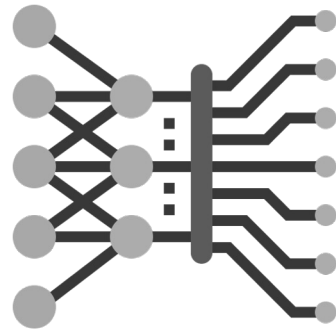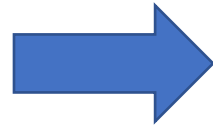Realistic **storage & pre-processing** settings

**No accelerator required** to run
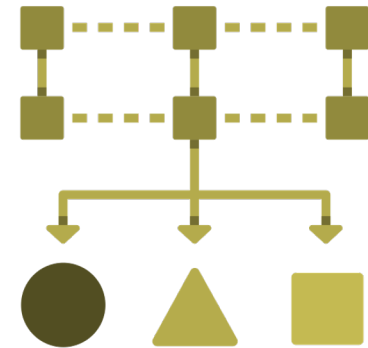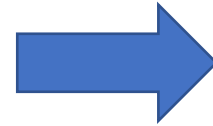
**Minimal AI/ML knowledge** required

# Stages of the ML Pipeline

**Data cleaning & pre-processing**

**Training**

**Inference**

# Stages of the ML Pipeline

**I/O intensive [1,2] – Our focus**

==*As much as **50% of the Watts** can go into storage and data cleaning [2]*==
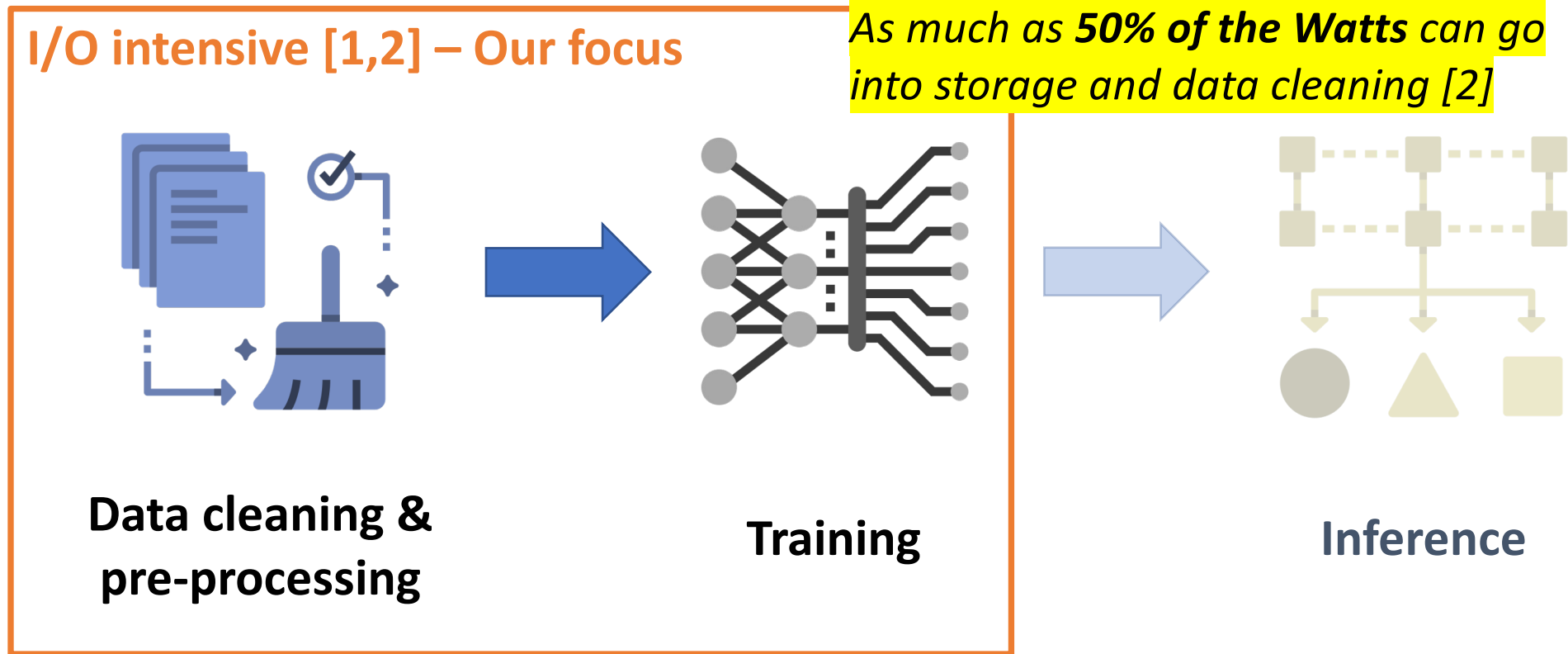


**Data cleaning & pre-processing**

**Training**

**Inference**

[1] Murray et al. **tf.data: A Machine Learning Data Processing Framework**, VLDB 21.
[2] Zhao et a. **Understanding Data Storage and Ingestion for Large-Scale Deep Recommendation Model Training** ISCA 22.

# Data Pipeline in ML: Pre-processing
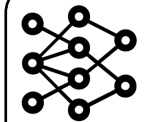
Storage resources

Compute resources

**Disk**                                    **Memory**
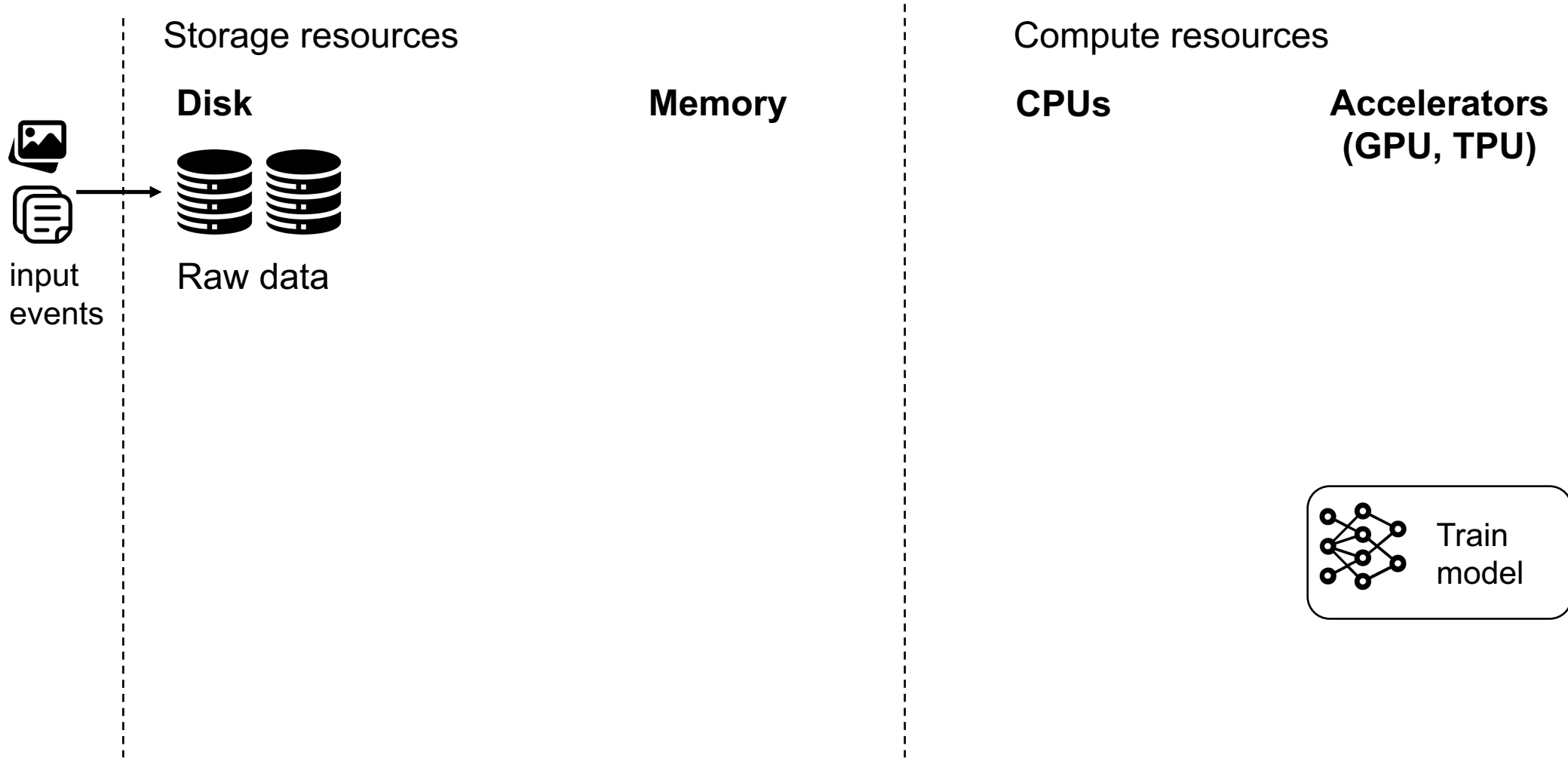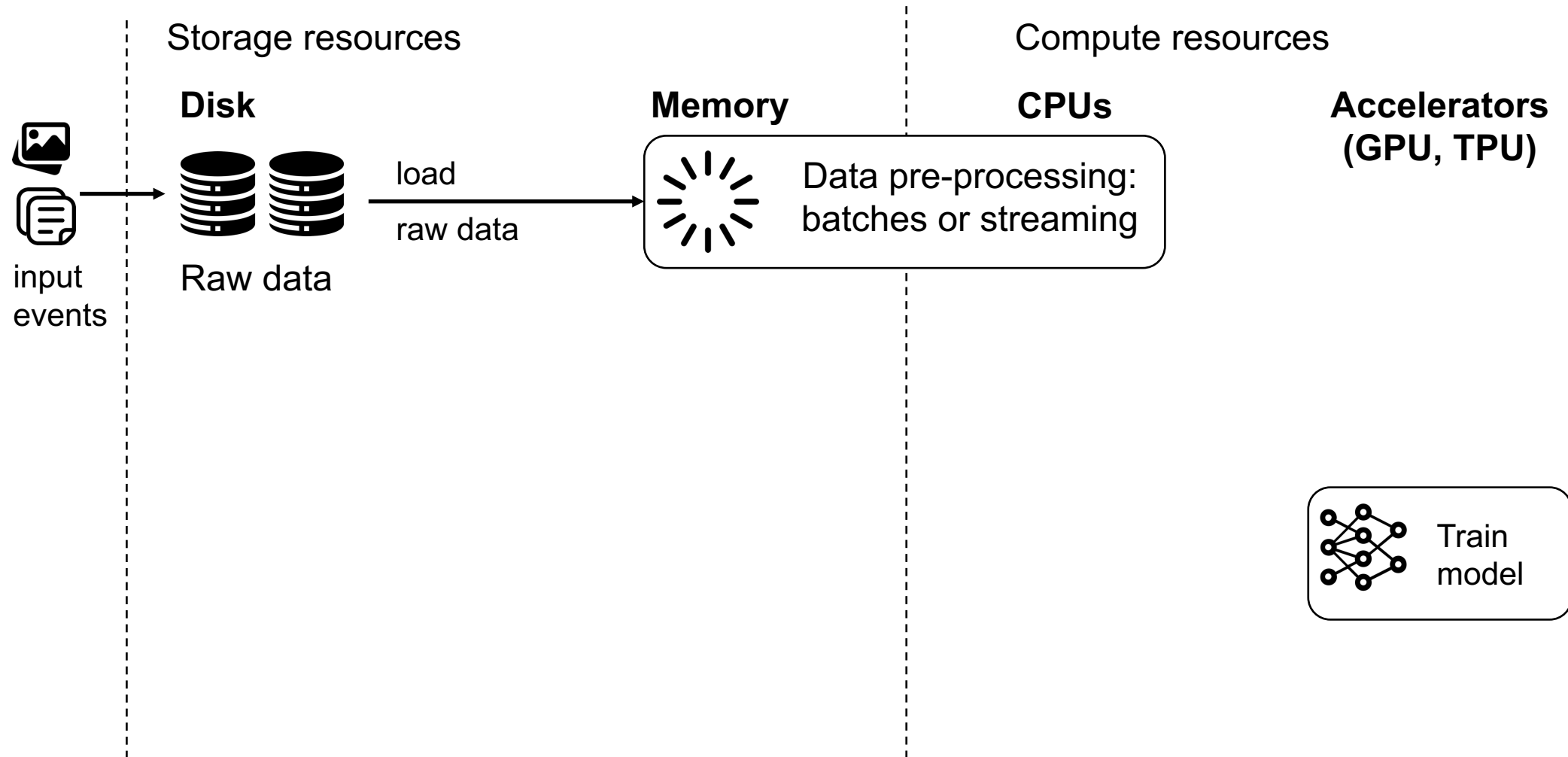
**CPUs**                          **Accelerators (GPU, TPU)**

Train model

# Data Pipeline in ML: Pre-processing

Storage resources

Compute resources

**Disk**

**Memory**

**CPUs**

**Accelerators (GPU, TPU)**

input events

Raw data

Train model

# Data Pipeline in ML: Pre-processing

Storage resources

Compute resources

**Disk**

**Memory**

**CPUs**

**Accelerators (GPU, TPU)**

input events

Raw data

load raw data

Data pre-processing: batches or streaming

Train model

# Data Pipeline in ML: Pre-processing

Storage resources

Compute resources

**Disk**

**Memory**

**CPUs**

**Accelerators (GPU, TPU)**

input events

Raw data

load
raw data

Data pre-processing: batches or streaming

store
processed dataset

Processed dataset

Train model

# Data Pipeline in ML: Pre-processing



Storage resources

Compute resources

**Disk**

**Memory**

**CPUs**

**Accelerators (GPU, TPU)**

input events

Raw data

load raw data

Data pre-processing: batches or streaming

store processed dataset

Processed dataset

Train model

# MLPerf Storage V0.5

**MLPerf Storage V0.5**

Training

Data cleaning & pre-processing

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

Realistic **storage** settings in

**training phase**

**No accelerator required** to run

**Minimal AI/ML knowledge**

required

# Data pipeline in ML: Training

Storage resources

Compute resources

**Disk**

**System Memory (DRAM)**

**CPUs**

**Accelerators (GPU, ASIC)**

Cleaned dataset

# Data pipeline in ML: Training

Storage resources

Compute resources

**Disk**

**System Memory (DRAM)**

**CPUs**

**Accelerators (GPU, ASIC)**

Cleaned dataset

TensorFlow

PYTORCH

load data

Cache data

# Data pipeline in ML: Training

Storage resources

Compute resources

**Disk**

**System Memory (DRAM)**

**CPUs**

**Accelerators (GPU, ASIC)**

Cleaned dataset

TensorFlow

PYTORCH

load

data

Cache data

Transform data

# Data pipeline in ML: Training

Storage resources

Compute resources

**Disk**

**System Memory (DRAM)**

**CPUs**

**Accelerators (GPU, ASIC)**

Cleaned dataset

TensorFlow

PYTORCH

load data

Transform data

Cache data

Load data in batches

Train model

# MLPerf Storage V0.5 – workloads

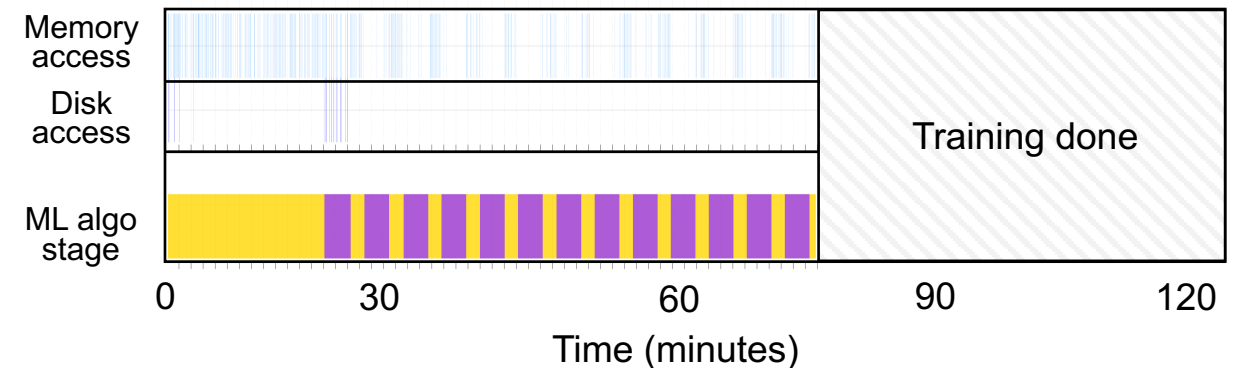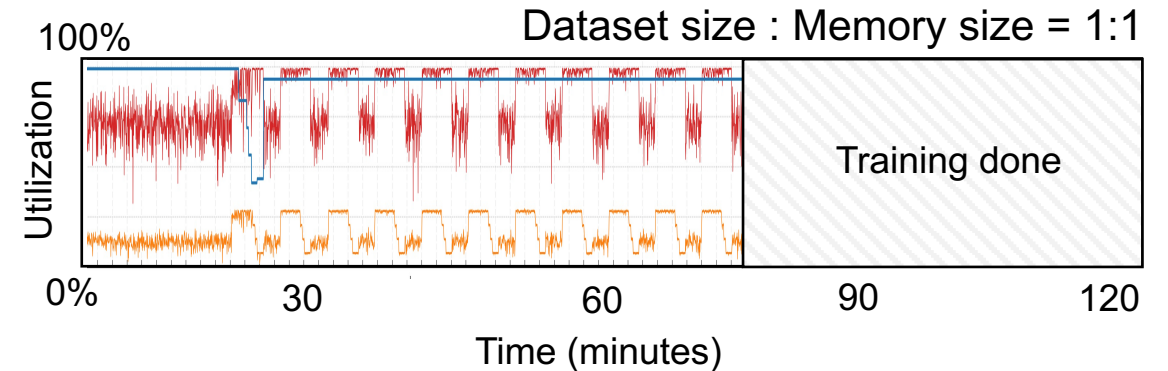| Workload | Image segmentation | Natural language processing | Recommender Systems |
|---|---|---|---|
| Model | Unet3D | BERT | DLRM |
| Seed data | KiTS19 Set of images | Wikipedia 2020 Text | Criteo Terabyte Click logs |
| Framework | Pytorch | Tensorflow | Pytorch |
| I/O behavior | Random access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

https://github.com/mlcommons/storage

**Preview package**

- **Single node**

- Many **simulated accelerators.**

- **Synthetic datasets** generated from real dataset seed.

- **Local storage**

# Dataset Size to Memory Ratio is Important

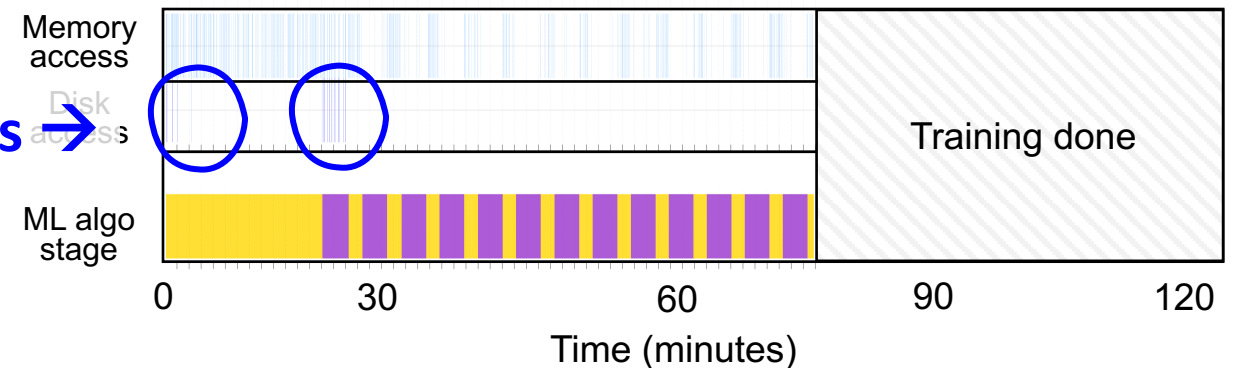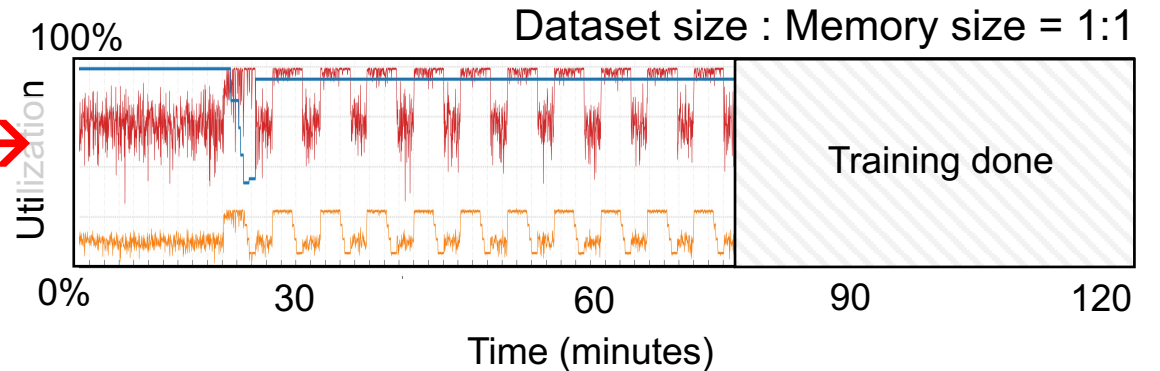**Dataset fits in system memory**

**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 size : Memory size = 1:1

Training done

Training done

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

# Dataset Size to Memory Ratio is Important

**Dataset fits in system memory**

Dataset size : Memory size = 1:1

**Experiment setup**
- DGX-1 server
  - 8 x V100 GPUs, 32GB GPU memory
  - 512GB DRAM

- Image segmentation workload:
  - Unet3D, Pytorch
  - MLPerf Training implementation
  - KiTS19 dataset
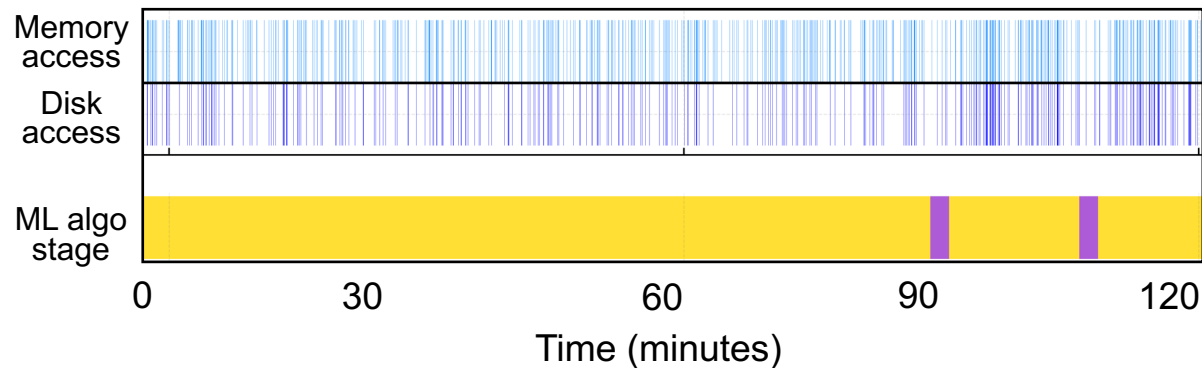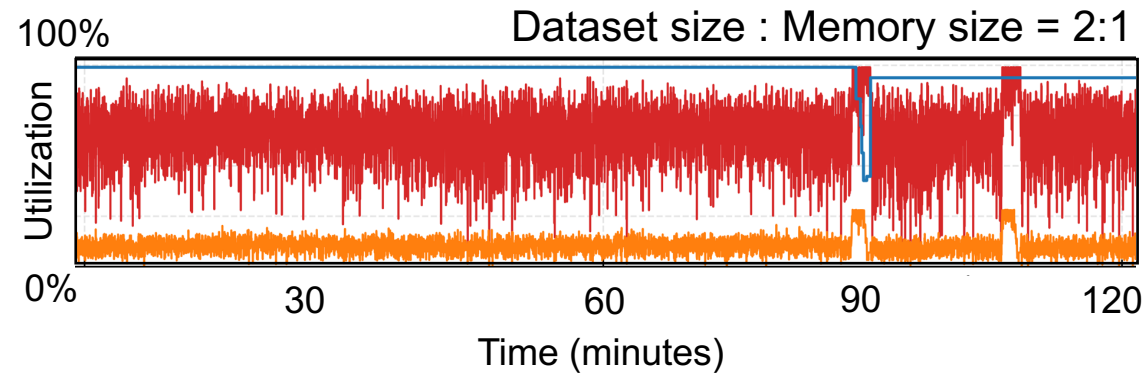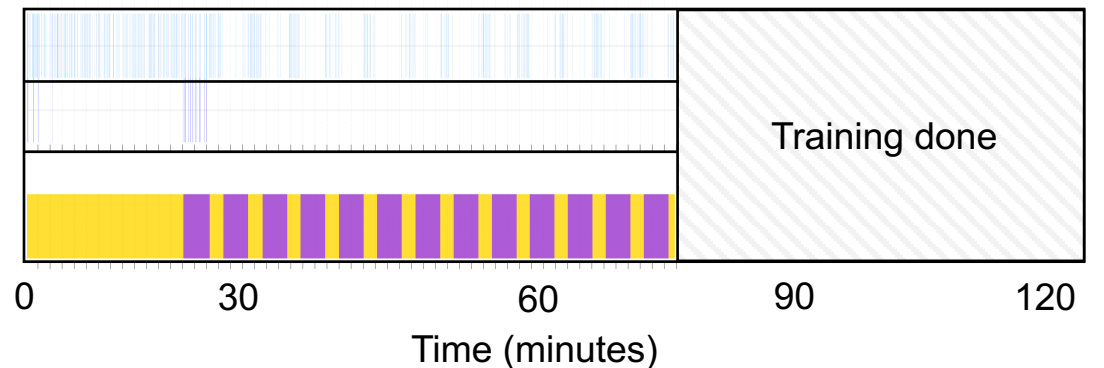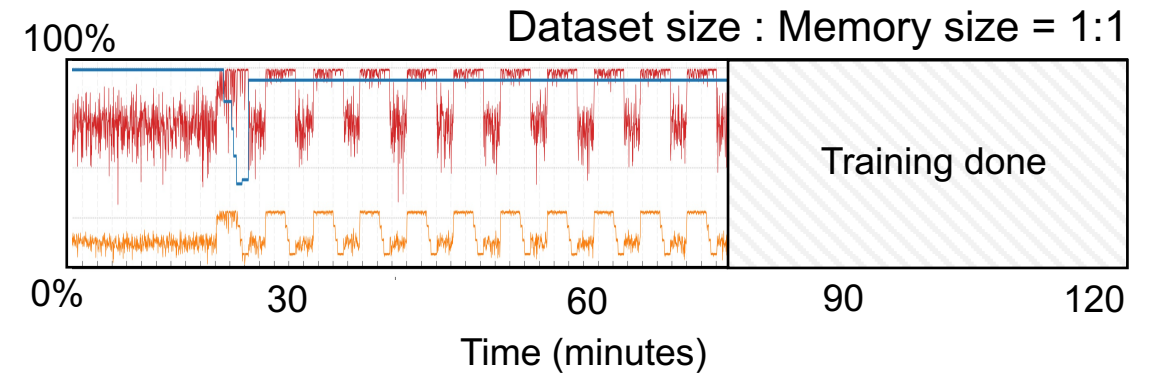
**High GPU utilization →**

**Little disk access →**



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

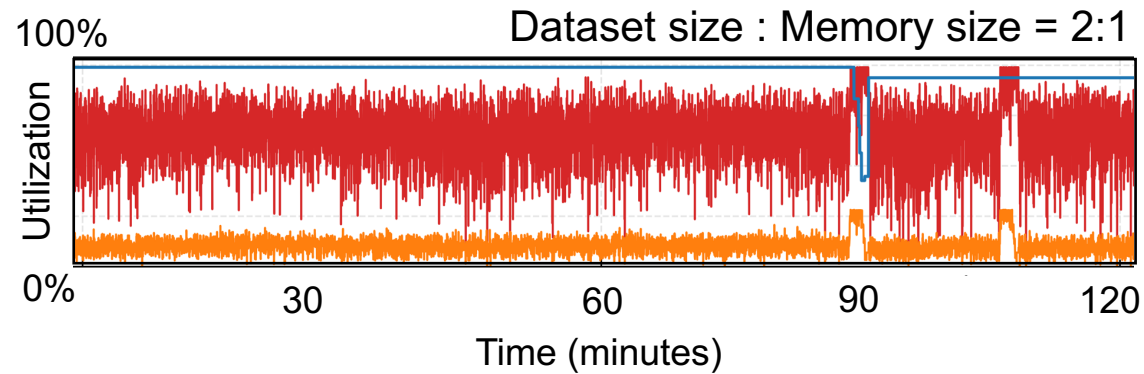# Dataset Size to Memory Ratio is Important

**Dataset does not fit in memory**
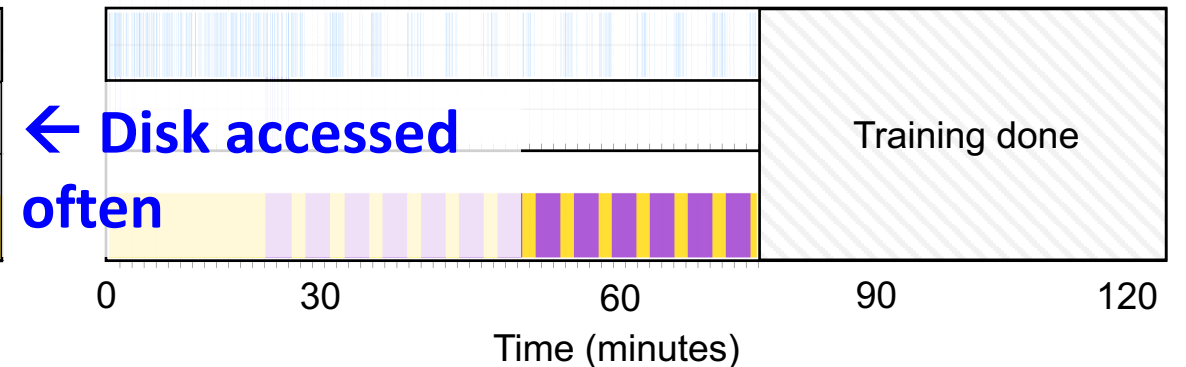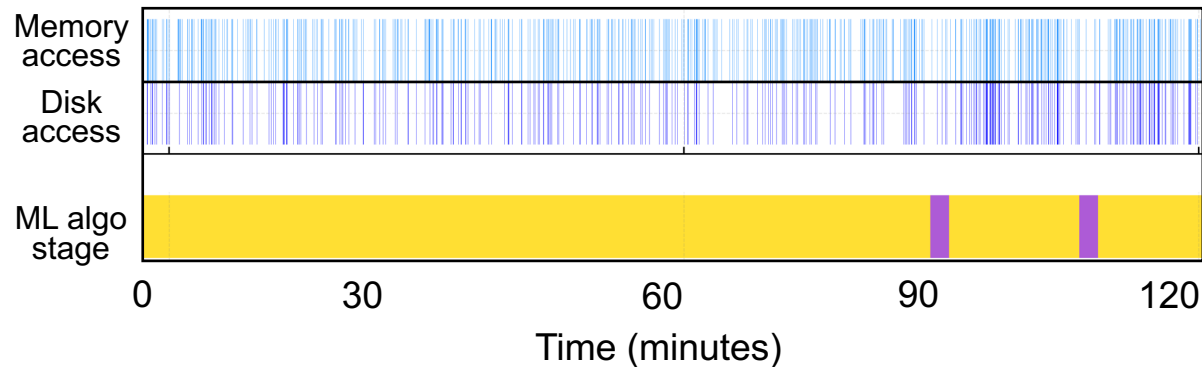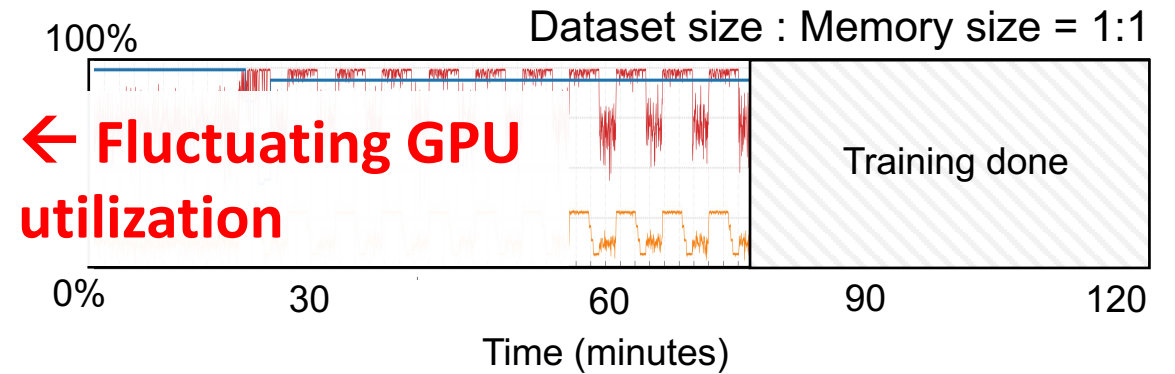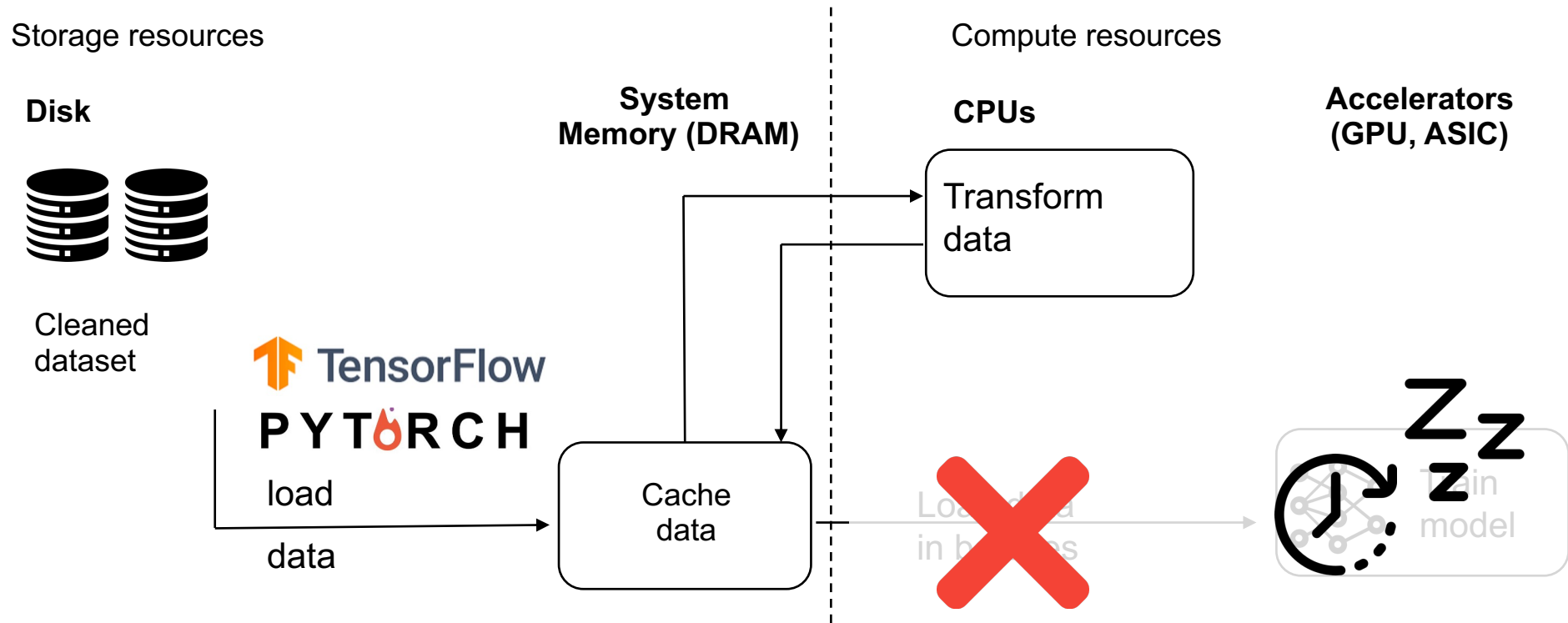
**Dataset fits in system memory**

# Dataset Size to Memory Ratio is Important

**Dataset does not fit in memory**

**Dataset fits in system memory**



← **Fluctuating GPU utilization**

← **Disk accessed often**

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

# Data pipeline in MLPerf Storage benchmark

Storage resources

Compute resources

**Disk**

**System
Memory (DRAM)**

**CPUs**

**Accelerators
(GPU, ASIC)**

Cleaned
dataset

TensorFlow

PYTORCH

load

data

Transform
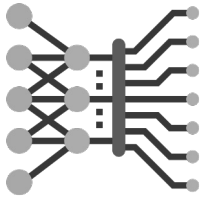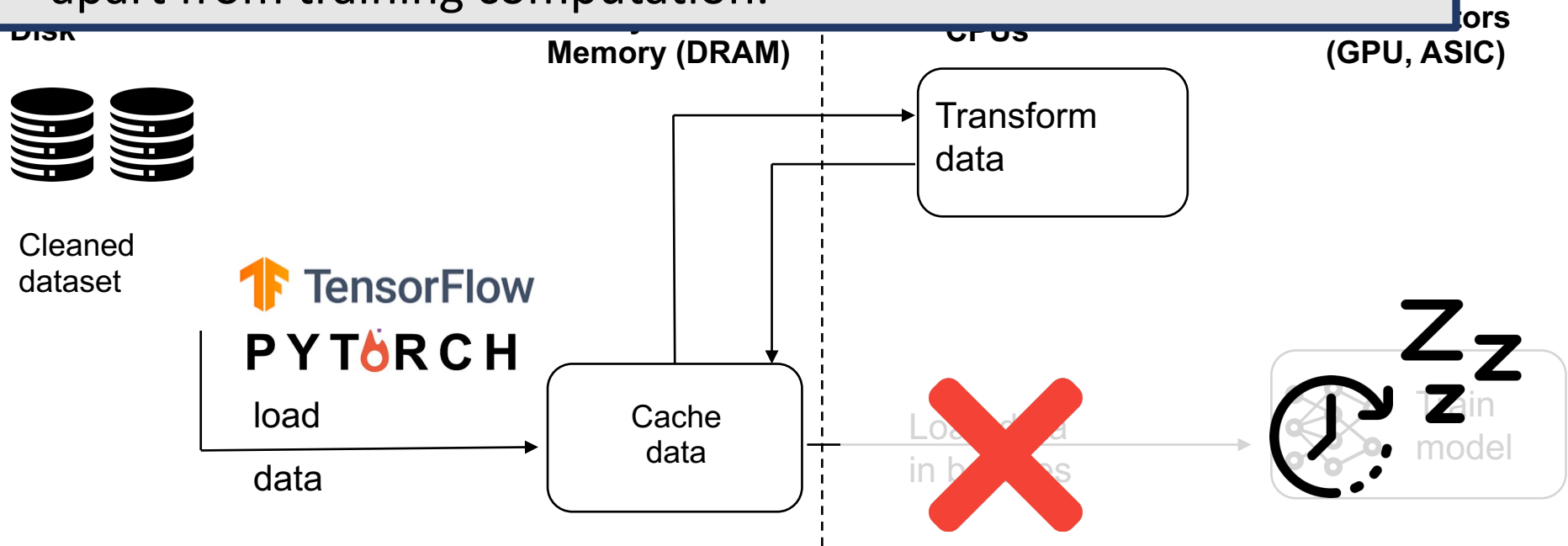data

Cache
data

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.*
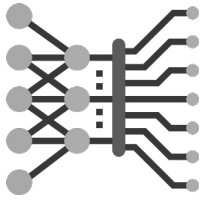
# Data pipeline in MLPerf Storage benchmark

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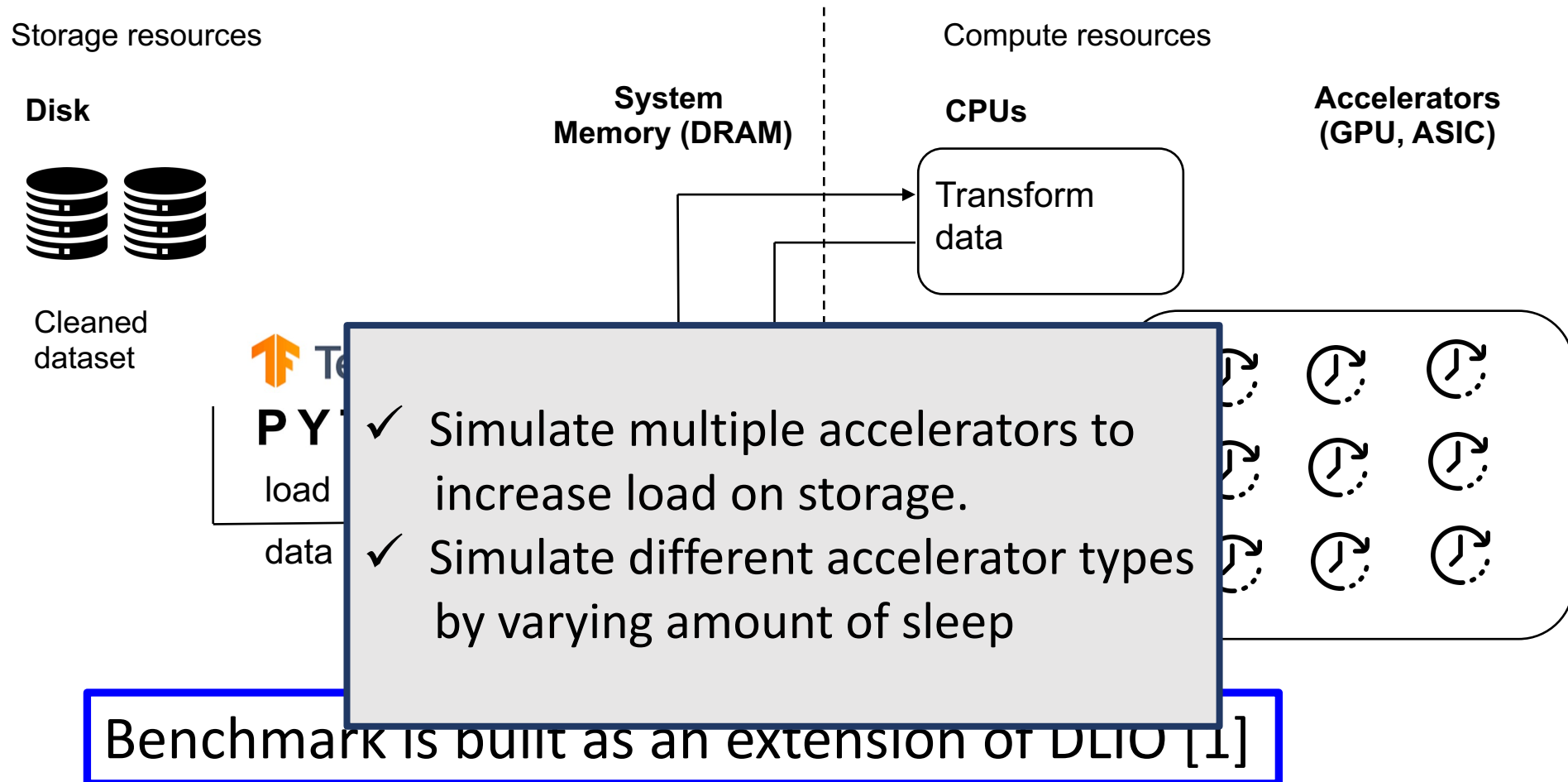


**Benchmark is built as an extension of DLIO [1]**

*[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

Storage resources

Compute resources

**Disk**

**System Memory (DRAM)**

**CPUs**

**Accelerators (GPU, ASIC)**

Cleaned dataset

Transform data

PY

load

data

✓ Simulate multiple accelerators to increase load on storage.
✓ Simulate different accelerator types by varying amount of sleep

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.*
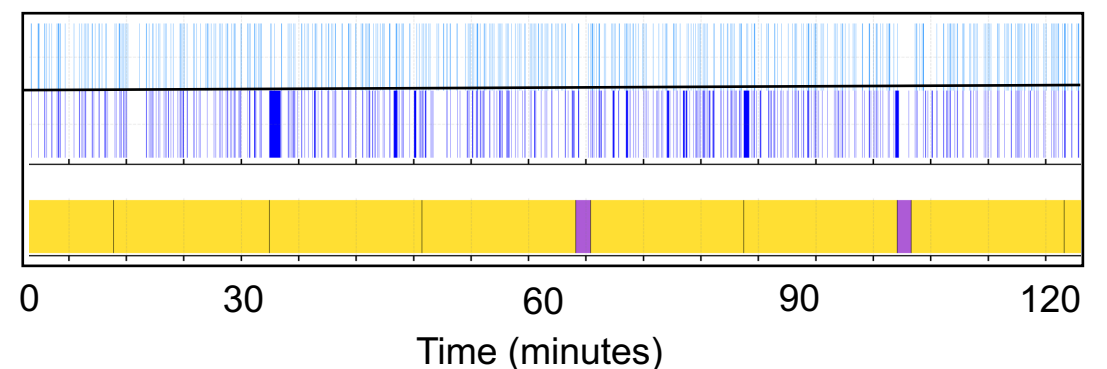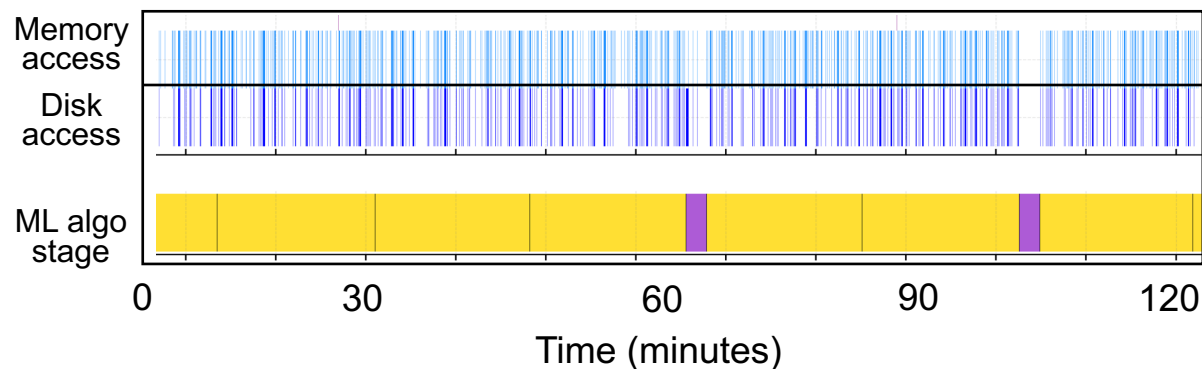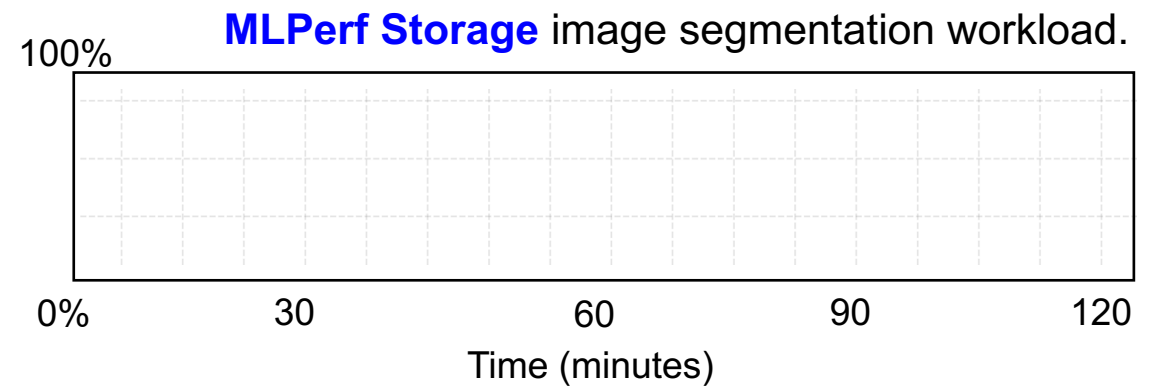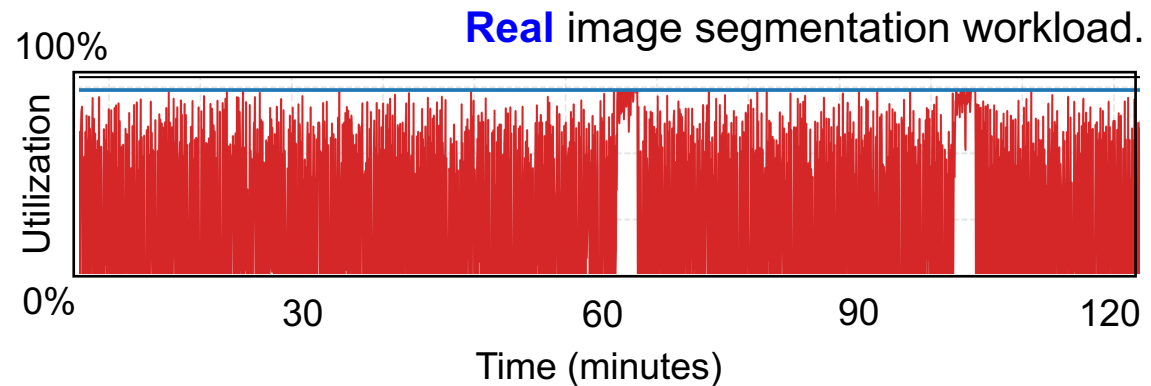
# Simulating training time does not impact I/O patterns



**Legend:** 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
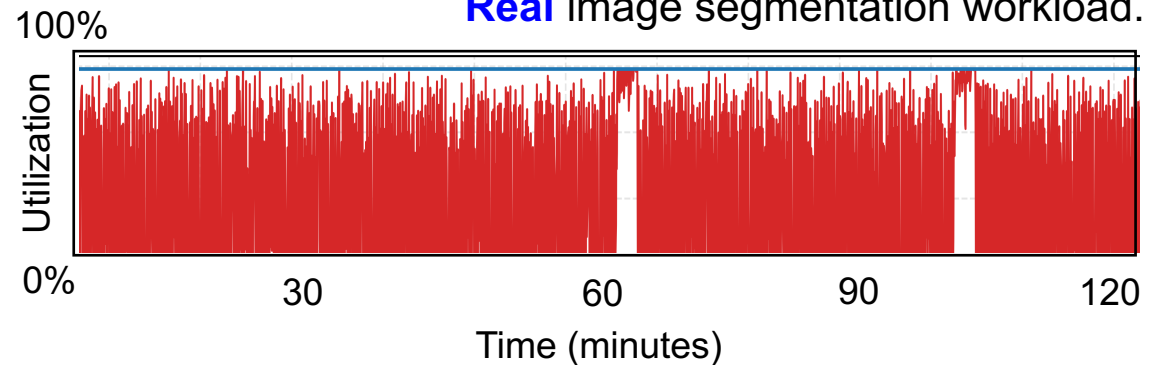


**Experiment setup:** DGX-1 with 8xV100 GPUs, 512GB DRAM. Dataset : KiTS19, Dataset size:Memory size ratio 2:1
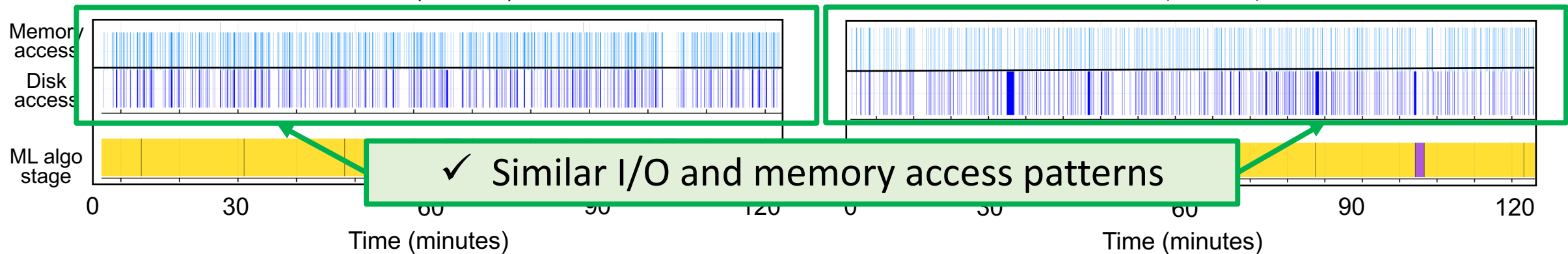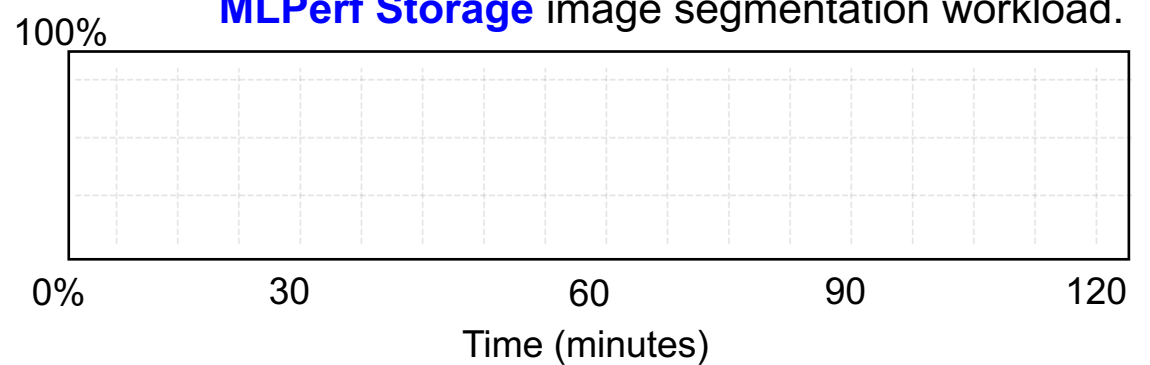
# Simulating training time does not impact I/O patterns

ML Training ■ ML Evaluation ■ Disk I/O Read ■ In-memory Read ■ GPU
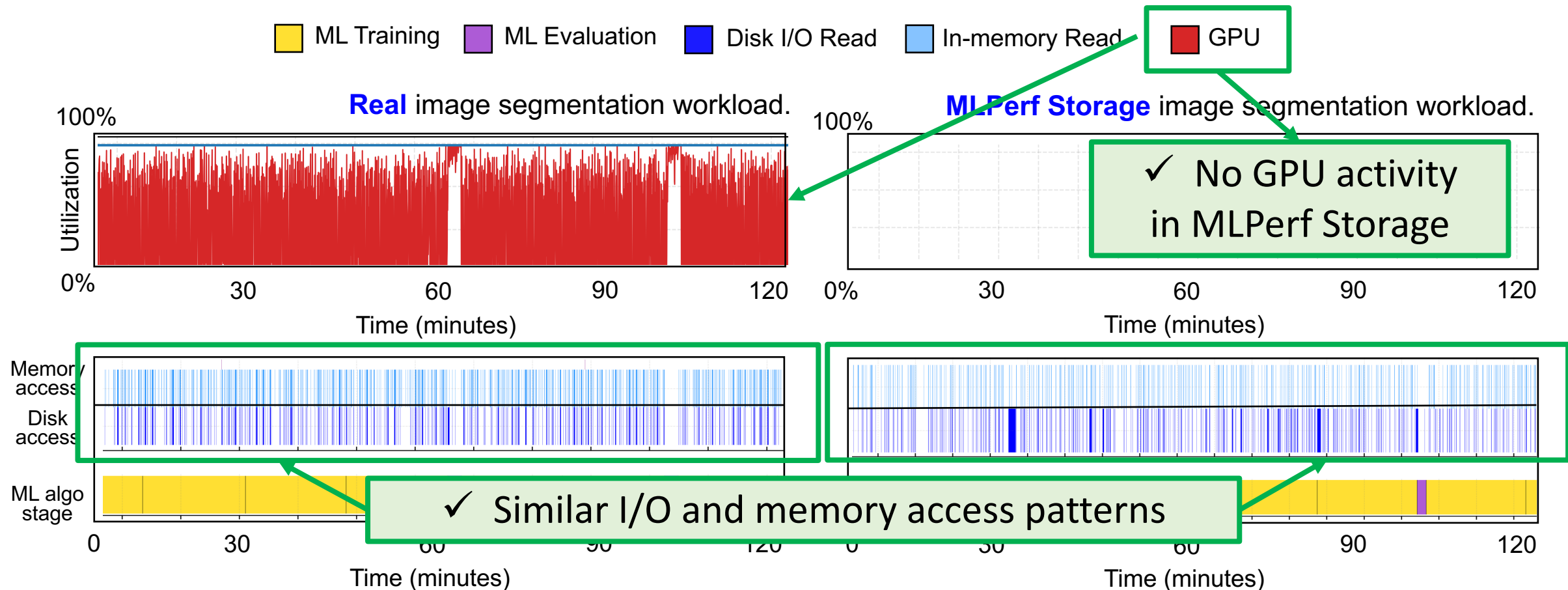
**Real** image segmentation workload.

**MLPerf Storage** image segmentation workload.

✓ No GPU activity in MLPerf Storage

Utilization

100%

0%

30    60    90    120

Time (minutes)

Memory access

Disk access

ML algo stage

✓ Similar I/O and memory access patterns

0    30    60    90    120

Time (minutes)

**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/tree/v0.5-branch*

Parallelism

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

# Key Takeaways – MLPerf Storage

**MLPerf Storage is a new benchmark**

Realistic **storage** settings

**No accelerators required** to run

Follow MLPerf Storage repository for updates:

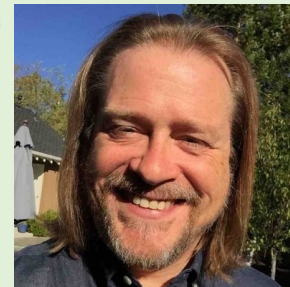https://github.com/mlcommons/storage

Get involved
mlcommons.org/en/get-involved/

**We appreciate your feedback**

Share your thoughts
Email oana.balmau@cs.mcgill.ca

Thanks to all working group co-chairs!



**Curtis Anderson**
**Panasas**

**Huihuo Zheng**
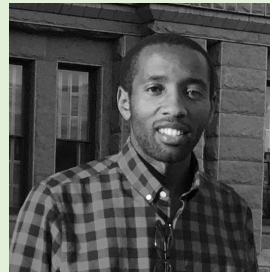**Argonne National Labs**

**Johnu George,**
**Nutanix**

# McGill DISCS Lab

**DISCS**

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

Postdoctoral
Researcher

*Dr. Stella Bitchebe*

PhD
Candidates:

*Nelson Bore*

Masters
Students

*Sebastian Rolon*   *Loïc Ho-Von*   *Aayush Kapur*   *Aidan Goldfarb*   *Rahma Nouaji*

Undergraduate
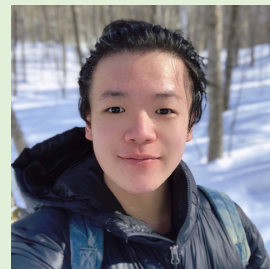Students

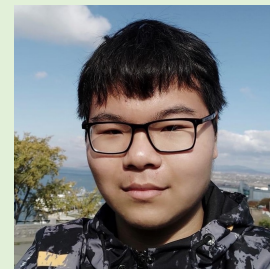*Zachary Doucet*   *Christian Zhao*   *Zhongjie Wu*   *Jiaxuan Chen*   *Changjun Zhou*   *Olivier Michaud*