

# ***SpeedyLoader: Efficient Pipelining of Data Preprocessing and Machine Learning Training***

**Rahma Nouaji, Stella Bitchebe, Oana Balmau**

*Dec 17, 2024, Women in CS Workshop Cameroon*



# What is Data preprocessing ?

- Transforming raw data into a clean and usable format to train machine learning models effectively.

# What is Data preprocessing ?

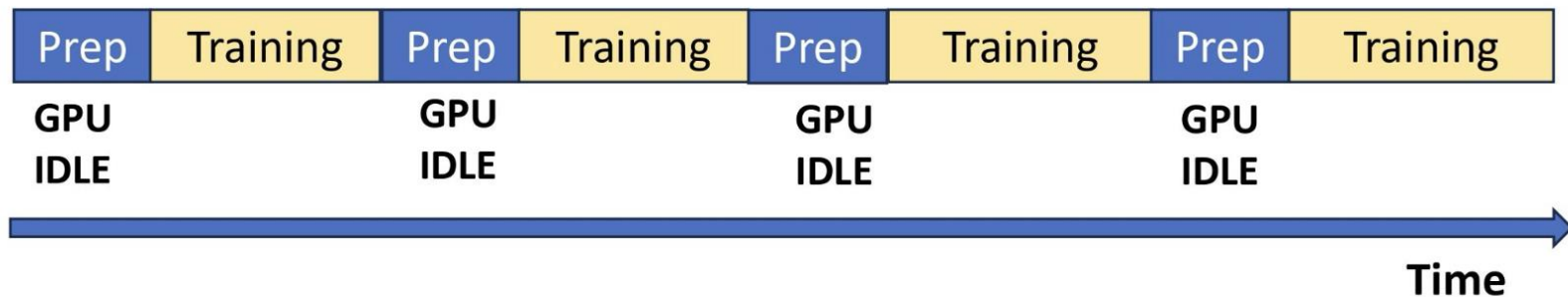
- Transforming raw data into a clean and usable format to train machine learning models effectively.
- 2 types:
  - Offline (batch) data preprocessing for static workloads (once before training).
  - **Online (real-time)** data preprocessing for dynamic workloads (with the training).

# What is Data preprocessing ?

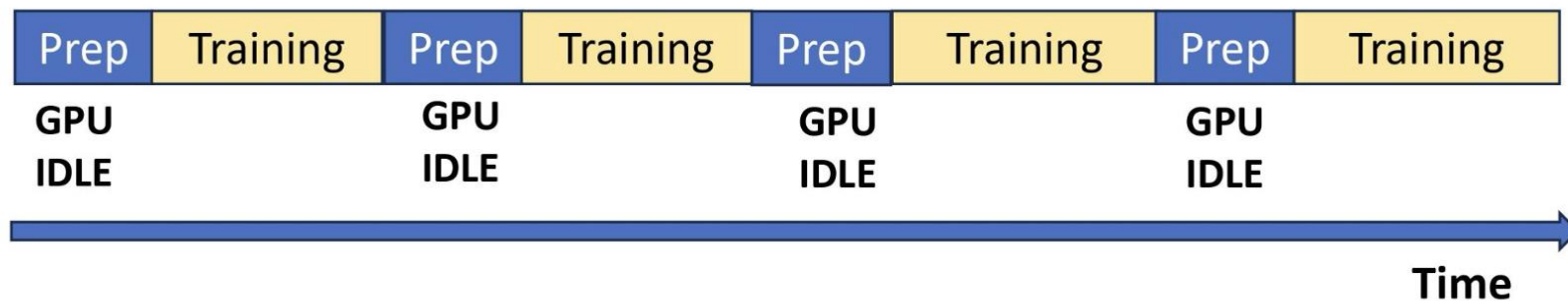
- Transforming raw data into a clean and usable format to train machine learning models effectively.
- 2 types:
  - Offline (batch) data preprocessing for static workloads (once before training).
  - Online (real-time) data preprocessing for dynamic workloads (with the training).



# Challenges with Online Data Preprocessing



# Challenges with Online Data Preprocessing

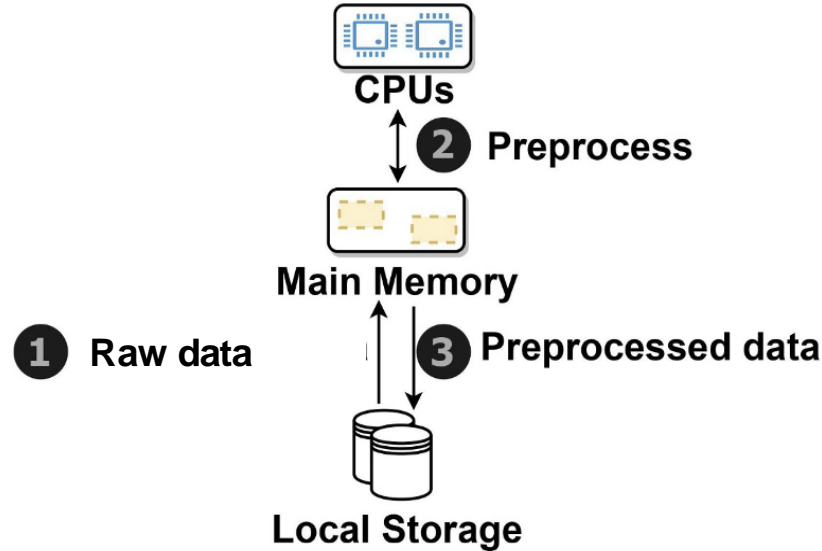


⊗ **75% GPU idleness!**

⊗ **Significant training delays!**

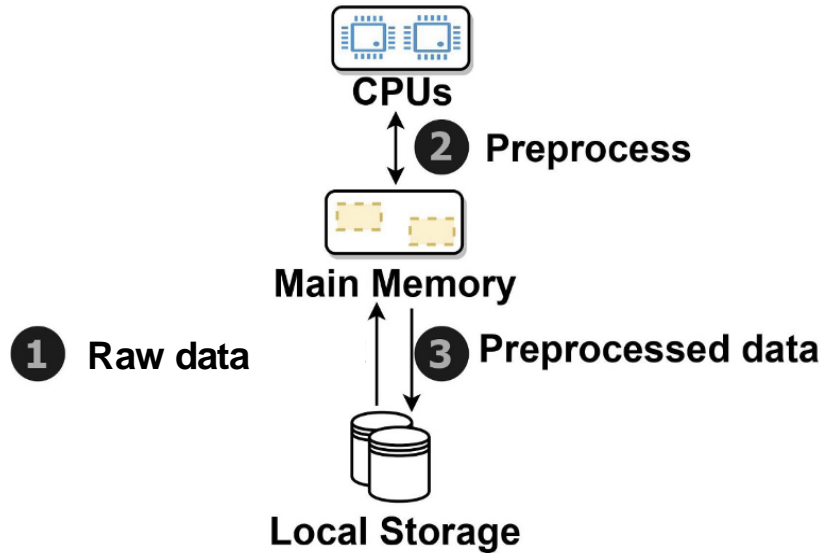
⊗ **GPU cost waste!**

# Overview of ML data preprocessing pipeline

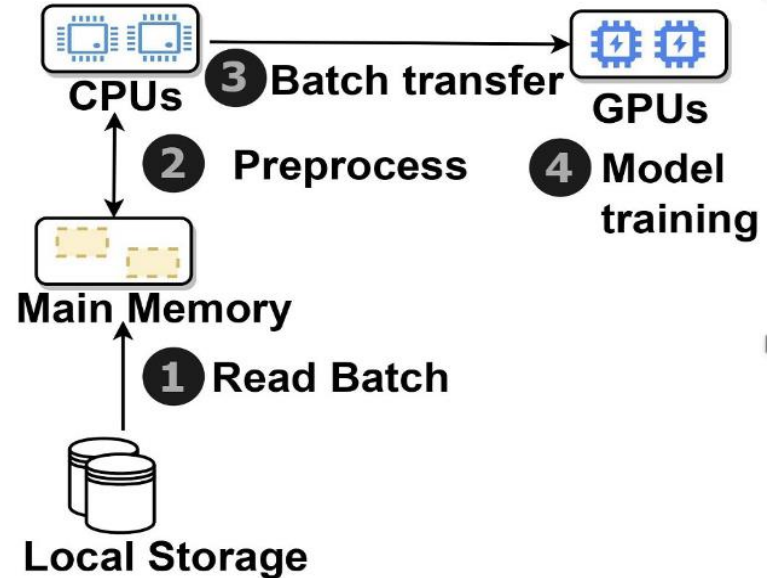


**A. Offline data preprocessing**  
Done once, before the training starts

# Overview of ML data preprocessing pipeline



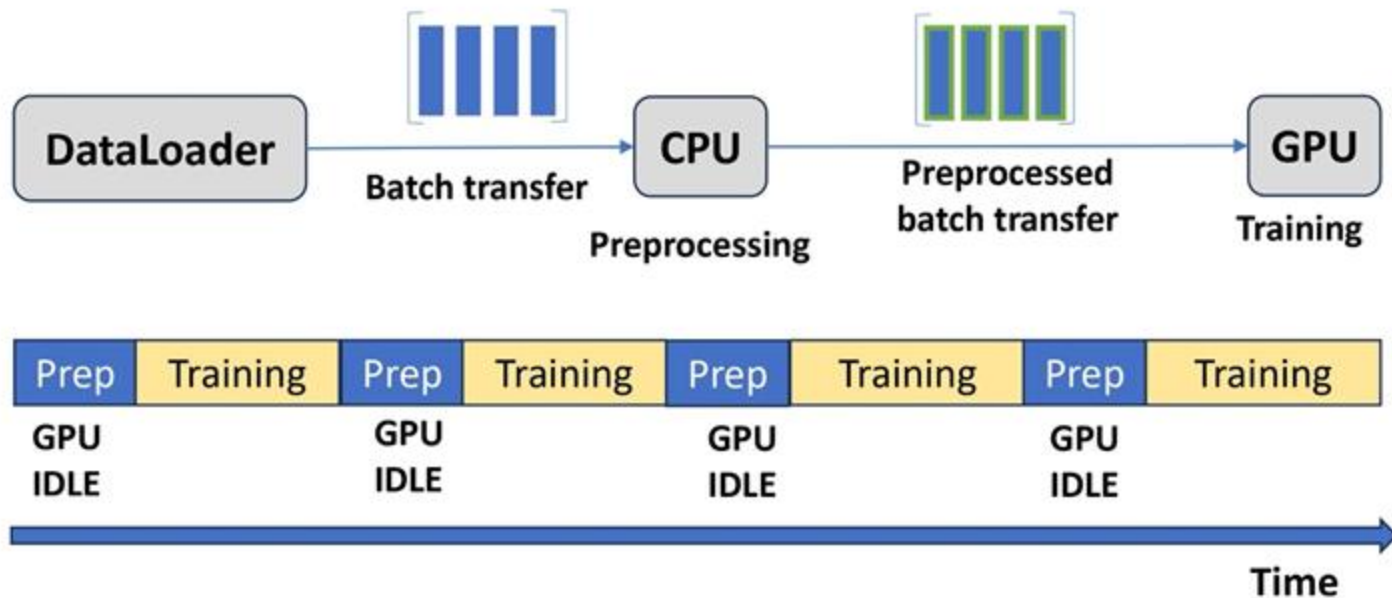
**A. Offline data preprocessing**  
Done once, before the training starts



**B. Training for one epoch,**  
with online data preprocessing



# Inefficient pipelining using PyTorch DataLoader





# SpeedyLoader



## SpeedyLoader:

*A data loader that overlaps the preprocessing and training steps to enhance training time efficiency and GPU utilization .*



# SpeedyLoader

*A data loader that overlaps the preprocessing and training steps to enhance training time efficiency and GPU utilization .*

## How ?

1. Combines offline and online into one block.



# SpeedyLoader

*A data loader that overlaps the preprocessing and training steps to enhance training time efficiency and GPU utilization .*

## How ?

1. Combines offline and online into one block.
2. Introduces a load balancer.

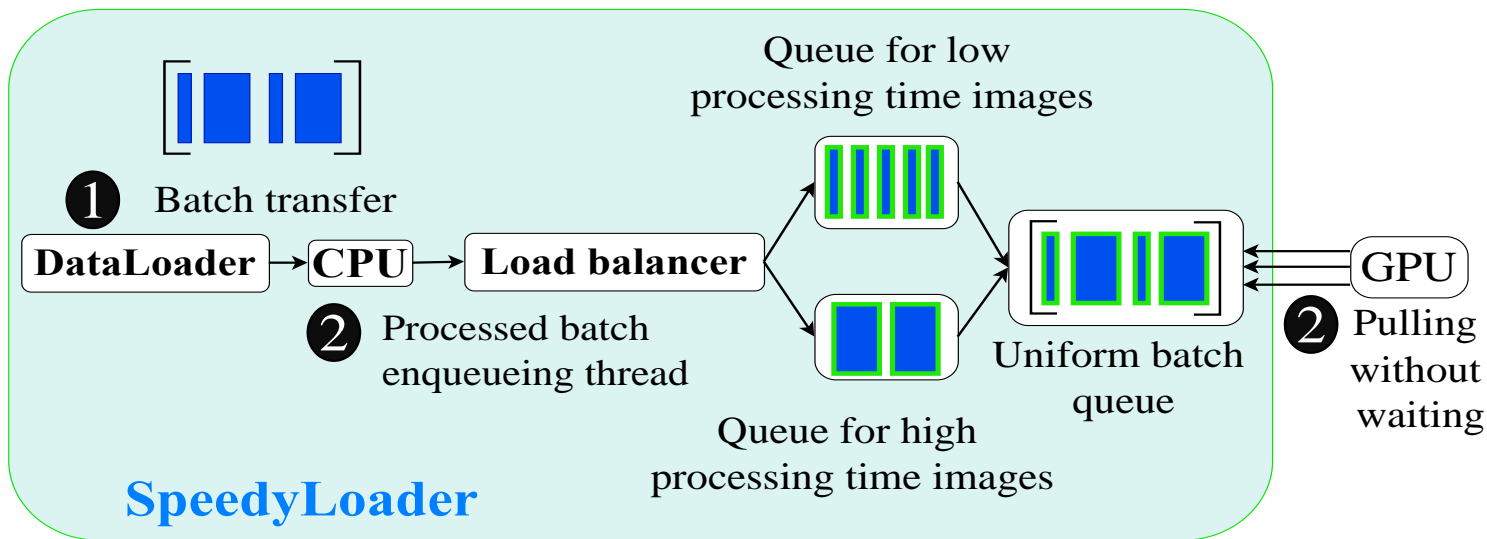


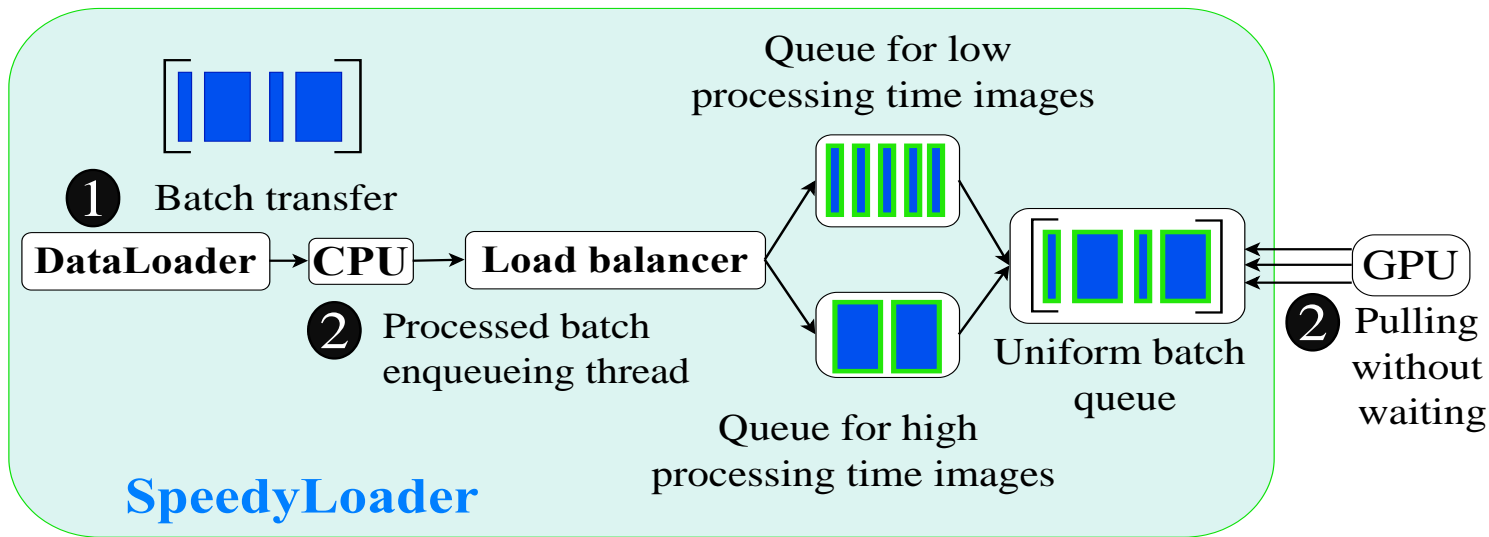
# SpeedyLoader

*A data loader that overlaps the preprocessing and training steps to enhance training time efficiency and GPU utilization .*

## How ?

1. Combines offline and online into one block.
2. Introduces a load balancer.
3. Enhances coordination between preprocessing and GPU threads by using a shared queue.





- ✓ **More than 90% GPU usage!**
- ✓ **More than 30% better training time!**
- ✓ **Optimized GPU usage efficiency!**



# Our workloads:

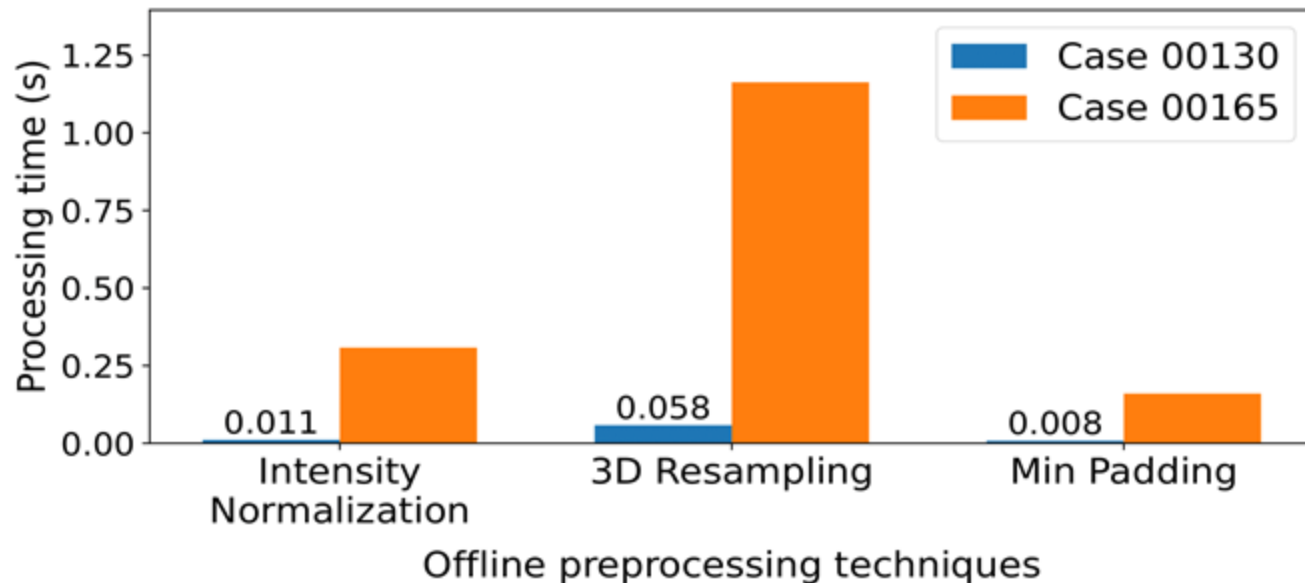
## ➤ Image Segmentation (3D-UNet)

- The KiTS19 challenge dataset with 210 images (**29GB**).
- 3D-UNet model.
- 8 online preprocessing techniques.

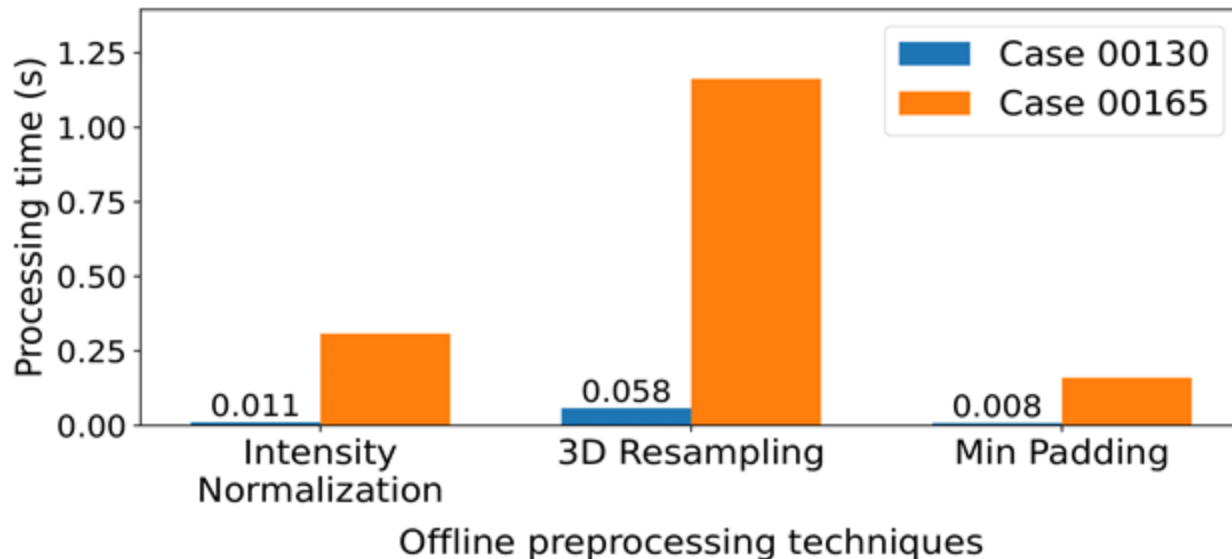
## ➤ Single Shot Detector (SSD) is an object detection network.

- Coco detection dataset (**352GB**).
- ResNeXt50\_32x4d Model.
- 5 online preprocessing techniques.

# Offline preprocessing (3D-UNet)



# Offline preprocessing (3D-UNet)



[1, 53, 512, 512]

[1, 734, 512, 512]

## Average time in ms:

Intensity norm 115ms

3D Resampling 380ms

Min Padding 55ms

# Online preprocessing (3D-UNet)

**Table 1.** Execution time (in ms) for Online Preprocessing Techniques on case\_00039 for two training runs.

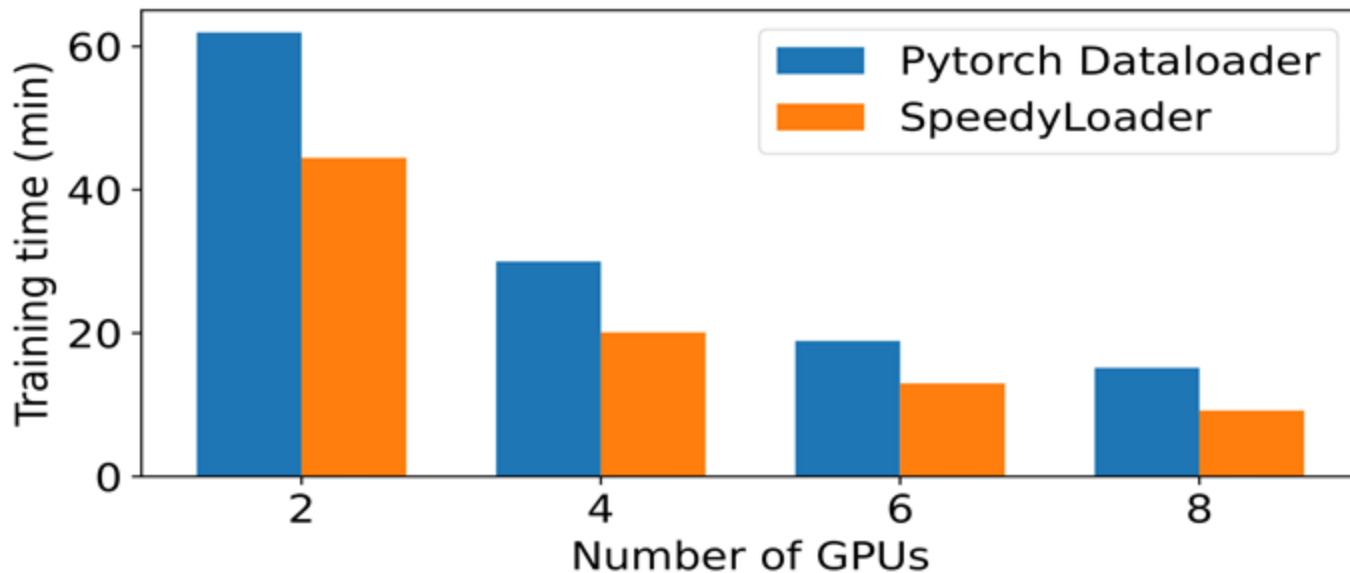
Technique	Time run 1 (ms)	Time run 2 (ms)
Random flip	$2.762 \times 10^{-3}$	32
Cast	3	3
Random Brightness Aug	6	$1.054 \times 10^{-3}$
Gaussian Noise	$2.005 \times 10^{-3}$	149
Random Balance Crop	965	0.172

# What are our key takeaways from the profiling study?

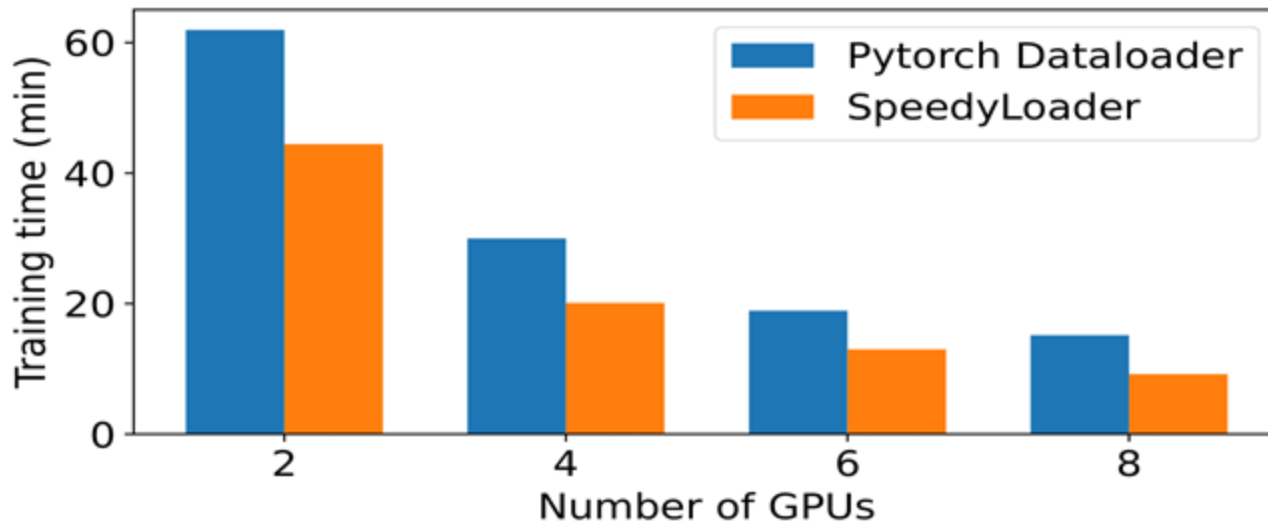
The image processing time is influenced by **two** factors:

1. **Image size** in the offline preprocessing.
2. The **randomness** in the online preprocessing transformations.

## Results: Total training time of 3D-UNet model.

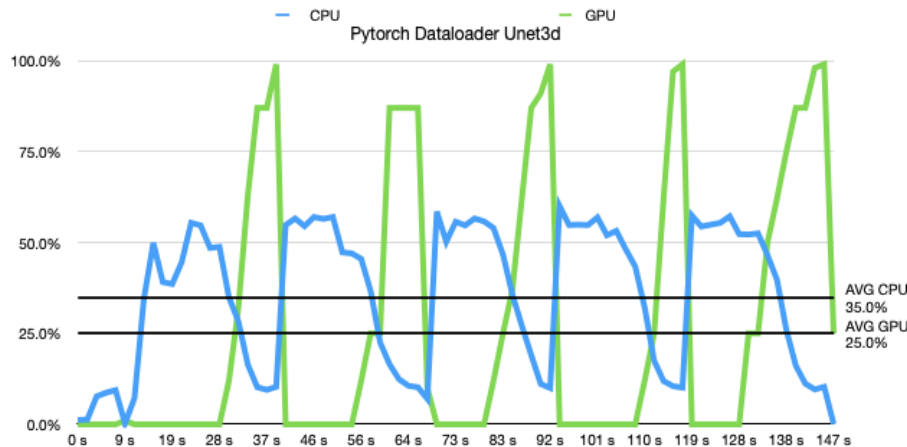


## Results: Total training time of 3D-UNet model.

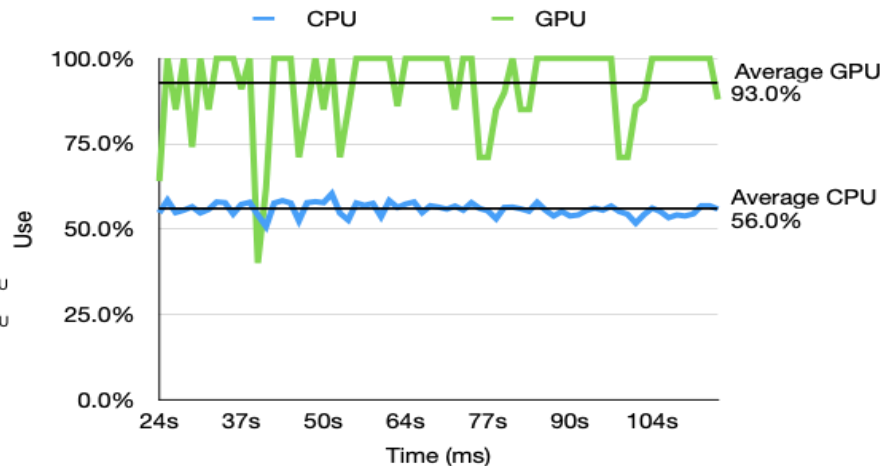


➡ ***SpeedyLoader*** provides up to **30%** better training time.

# Results: CPU and GPU usage 3D-UNet training for 5 epochs, 8 V100 GPUs



Pytorch Dataloader

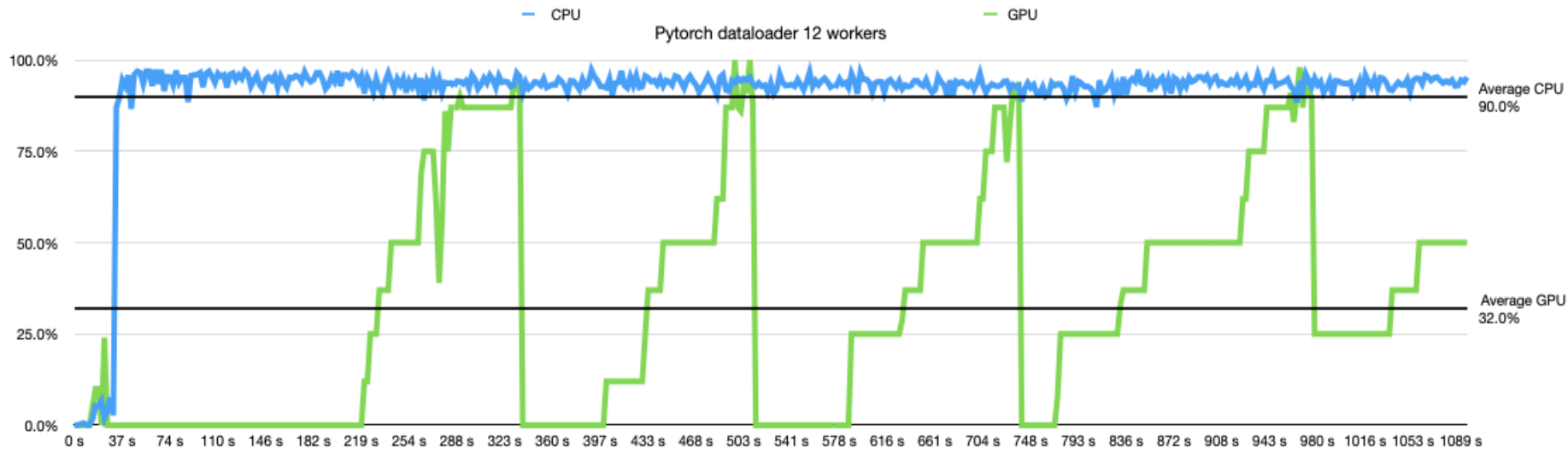


SpeedyLoader

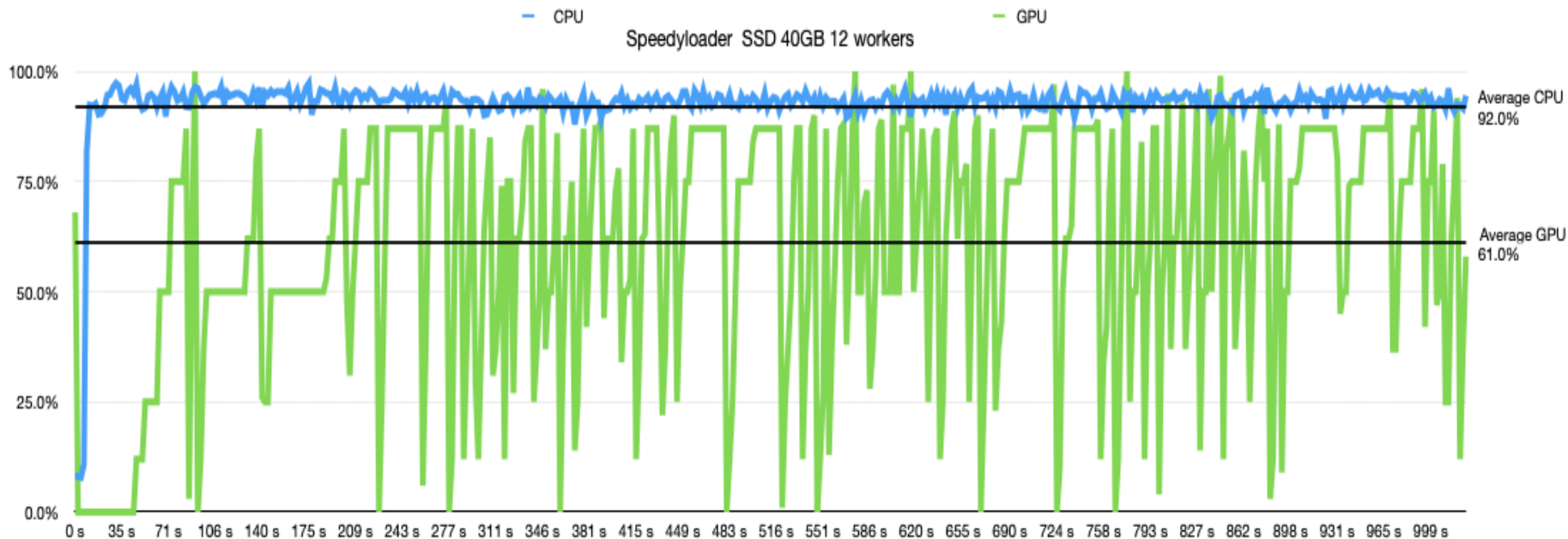
➡ *SpeedyLoader provides up to **3.8x** better GPU usage.*



## Results: SSD with Pytorch DataLoader training for 1 epoch, 8 V100 GPUs



## Results: SSD with SpeedyLoader training for 1 epoch, 8 V100 GPUs



➡ **SpeedyLoader** provides up to **2x** better GPU usage.

# Summary

- A study of data preprocessing techniques.

# Summary

- A study of data preprocessing techniques.
- **Bottleneck**: Inefficient pipelining of preprocessing and training.

# Summary

- A study of data preprocessing techniques.
- **Bottleneck**: Inefficient pipelining of preprocessing and training.
- **Solution**: SpeedyLoader
  - Relies on shared queue between loading threads and GPU threads.
  - Implements a Load Balancer to mitigate head-of-line blocking.

# Summary

- A study of data preprocessing techniques.
- **Bottleneck**: Inefficient pipelining of preprocessing and training.
- **Solution**: SpeedyLoader
  - Relies on shared queue between loading threads and GPU threads.
  - Implements a Load Balancer to mitigate head-of-line blocking.
- **30%** decrease in training time and up to **3.8x** increase in GPU usage with 91% accuracy.



# SpeedyLoader: Efficient Pipelining of Data Preprocessing and Machine Learning Training

Rahma Nouaji  
rahma.nouaji@mail.mcgill.ca  
McGill University  
Montreal, Quebec, Canada

Stella Bitchebe  
stella.bitchebe@mcgill.ca  
McGill University  
Montreal, Quebec, Canada

Oana Balmou  
oana.balmou@cs.mcgill.ca  
McGill University  
Montreal, Quebec, Canada

## Abstract

Data preprocessing consisting of tasks like sample reusing, cropping, and filtering, is a crucial step in machine learning (ML) workflows. Even though the preprocessing step is largely ignored by work that focuses on optimizing training algorithms, in practice for many workflows preprocessing and training are pipelined. Popular ML frameworks like PyTorch use data loaders to feed data into model training. If the pipeline between preprocessing and training is not done carefully, it can cause significant waiting times on the GPU side. To address this limitation, we introduce **SPEEDY-LOADER**, a system that overlaps preprocessing and training by leveraging asynchronous data preprocessing and avoiding data loading blocking. **SPEEDY-LOADER** incorporates dedicated data loading threads, which organize preprocessed samples into queues based on their predicted processing times. Contrary to the default PyTorch **DataLoader**, **SPEEDY-LOADER** reduces training time by up to 30% and increases GPU usage by 4.3x, all while maintaining a consistent evaluation accuracy of 91%.

**CCS Concepts** • Computing methodologies → Machine learning; Parallel algorithms; • Computer systems organization → Data flow architectures.

**Keywords** Machine learning, DataLoader, GPU-CPU overlap, Data preprocessing, Training, Pipelining

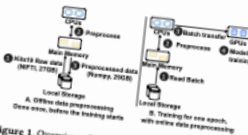
## ACM Reference Format:

Rahma Nouaji, Stella Bitchebe, and Oana Balmou. 2024. SpeedyLoader: Efficient Pipelining of Data Preprocessing and Machine

Learning Training. In *10th Workshop on Machine Learning and Systems (EuroMLSys '24)*, April 22, 2024, Athens, Greece. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3642970.3653824>

## 1 Introduction

The efficacy of Machine Learning (ML) deployments relies on high-quality data—obtained through data preprocessing and high-quality algorithms. The latter has attracted significant attention, leading to numerous techniques [16, 17, 23], (e.g., GPU, TPU, DPU, and other hardware accelerators [4, 5, 16]), and software frameworks [4, 5, 16]. Though data preprocessing has not received much attention relative to the work and processing efficiency (e.g., via operations like cropping, resampling, filtering, etc.) are crucial to the training process. Recent work shows that preprocessing has a significant impact on learning speed, prediction accuracy, energy efficiency, and scalability [14, 18, 27].



**Figure 1.** Overview of the ML data preprocessing pipeline. Step A involves one-time offline preprocessing. Step B shows the training phase with online preprocessing executed at the beginning of each epoch.

Figure 1 shows a typical workflow for data preprocessing in a computer vision application selected from the MLPert Training Benchmark suite [18]. Data preprocessing is done in two stages: offline preprocessing (Figure 1A) and online preprocessing (Figure 1B). Both online and offline preprocessing load data into system main memory and then perform transformations in the CPU. Offline preprocessing occurs before the training begins, whereas online preprocessing occurs on each batch of images during the training process. Depending on the dataset size, offline preprocessing can span several hours to several days worth of CPU time [6].

Find out more in our [paper!!](#)

Contact me:

[rahma.nouaji@mail.mcgill.ca](mailto:rahma.nouaji@mail.mcgill.ca)

Thank you for your attention!