# Refurbish Your Training Data: Reusing Partially Augmented Samples for Faster Deep Neural Network Training

Gyewon Lee[1,3], Irene Lee[2], Hyeonmin Ha[1],
Kyunggeun Lee[1], Hwarim Hyun[1], Ahnjae Shin[1,3], and Byung-Gon Chun[1,3]
*Seoul National University[1], Georgia Institute of Technology[2], FriendliAI[3]*

# DNN Training Pipeline

DNN Training =　　　　Data Preparation　　　　+　　　　Gradient Computation

# DNN Training Pipeline

DNN Training =        Data Preparation       +       Gradient Computation

- Data read and preprocessing
- On CPU

# DNN Training Pipeline

DNN Training =      Data Preparation      +      Gradient Computation

- Data read and preprocessing
- On CPU

- Forward and backward operations
- On DL accelerators (e.g., GPU, TPU)

# DNN Training Pipeline

DNN Training =   Data Preparation   +   Gradient Computation
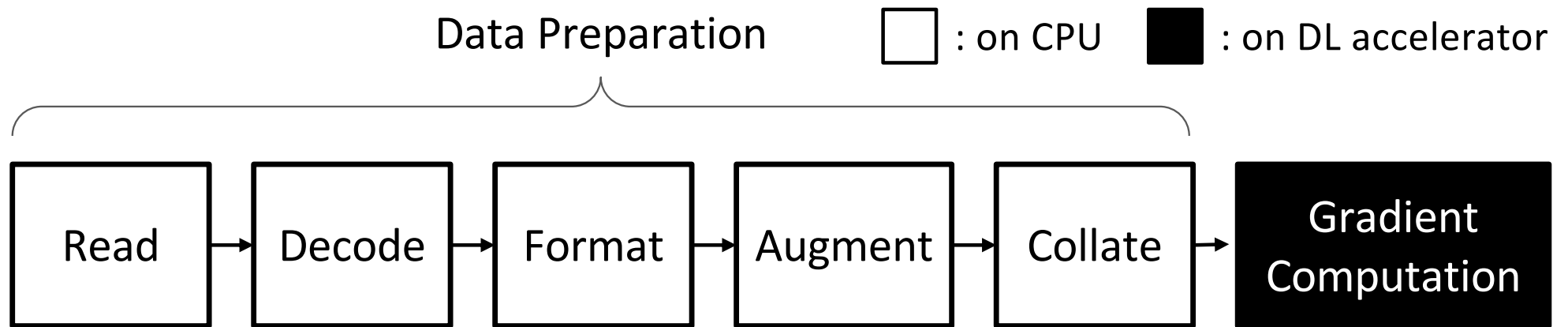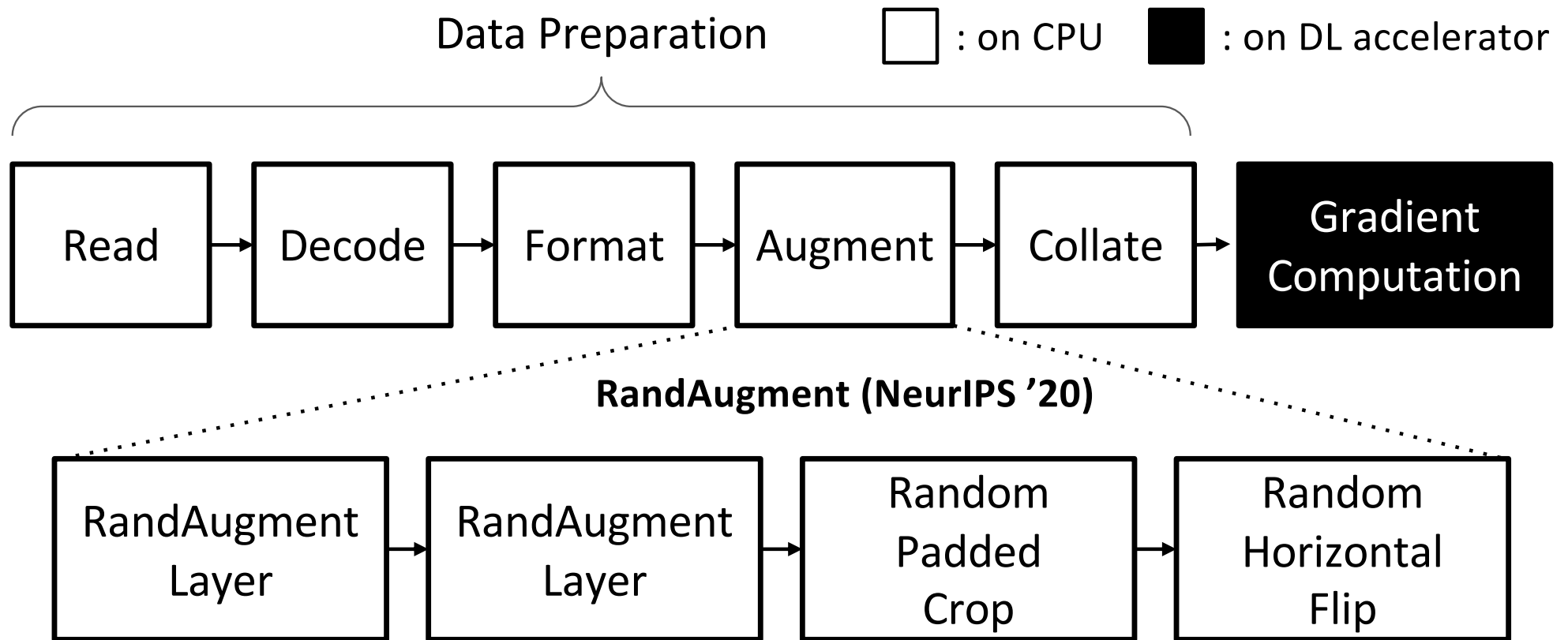
| Data Preparation box | Gradient Computation box |
|---|---|
| • Data read and preprocessing<br>• On CPU | • Forward and backward operations<br>• On DL accelerators (e.g., GPU, TPU) |

**Bottleneck!**

**Getting faster: NVIDIA A100, Google TPU v3, …**
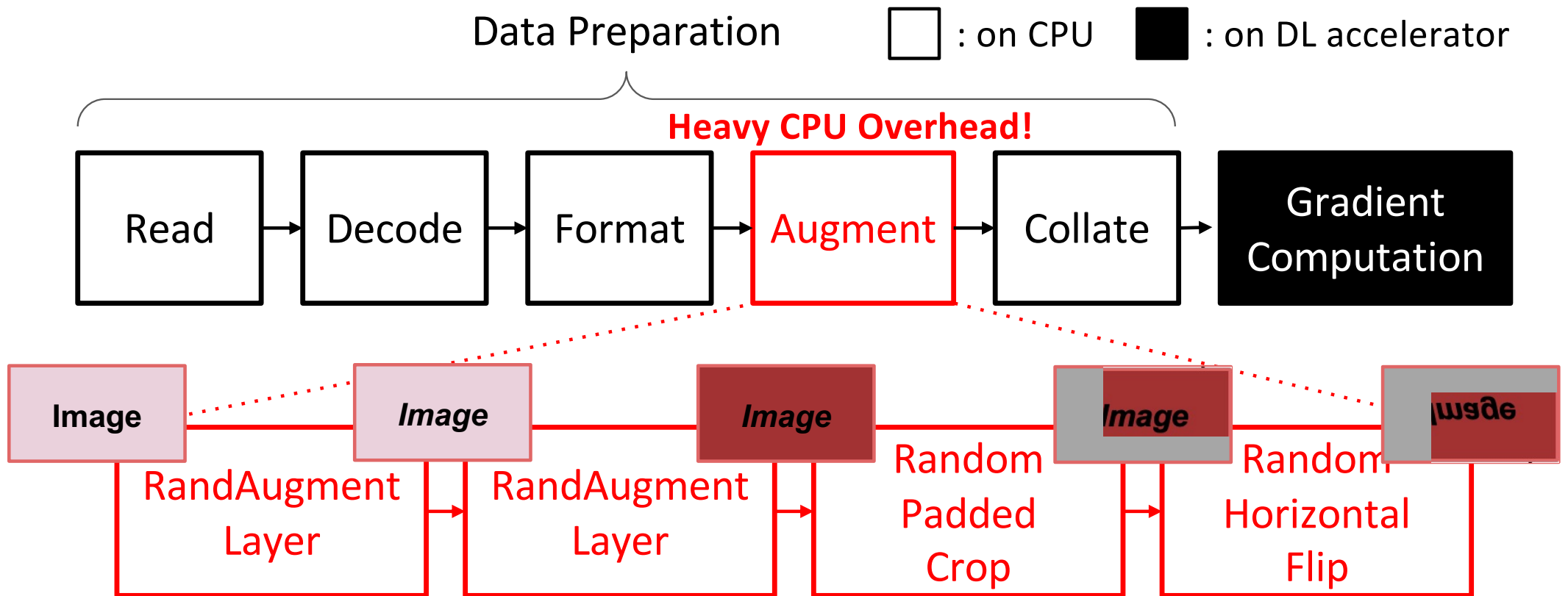
# DNN Training Pipeline

Data Preparation          ☐ : on CPU     ■ : on DL accelerator

```
Read → Decode → Format → Augment → Collate → Gradient Computation
```

# DNN Training Pipeline



Data Preparation

☐ : on CPU    ■ : on DL accelerator

Read → Decode → Format → Augment → Collate → **Gradient Computation**

**RandAugment (NeurIPS '20)**

RandAugment Layer → RandAugment Layer → Random Padded Crop → Random Horizontal Flip

# DNN Training Pipeline

Data Preparation

☐ : on CPU  ■ : on DL accelerator

Read → Decode → Format → Augment → Collate → **Gradient Computation**

Image — RandAugment Layer → Image — RandAugment Layer → Image — Random Padded Crop → Image — Random Horizontal Flip → Image
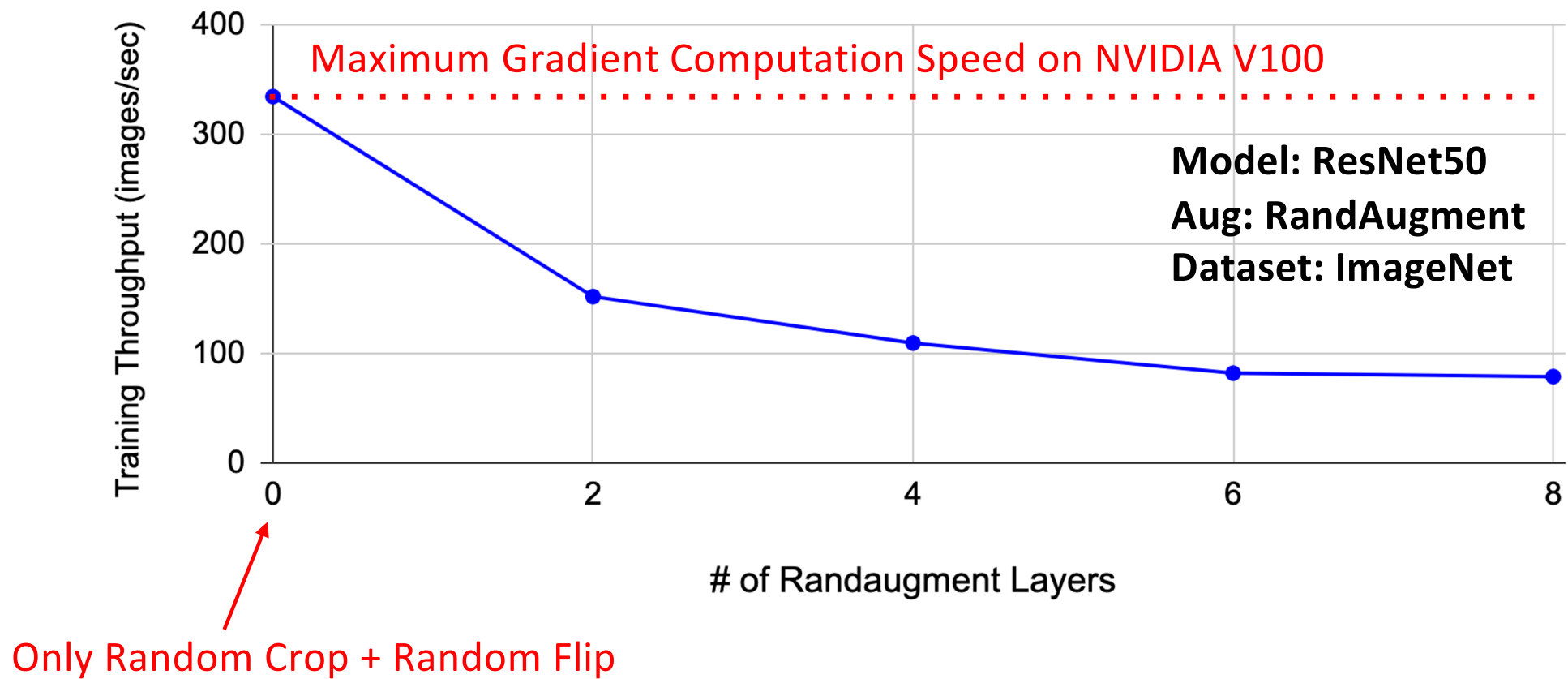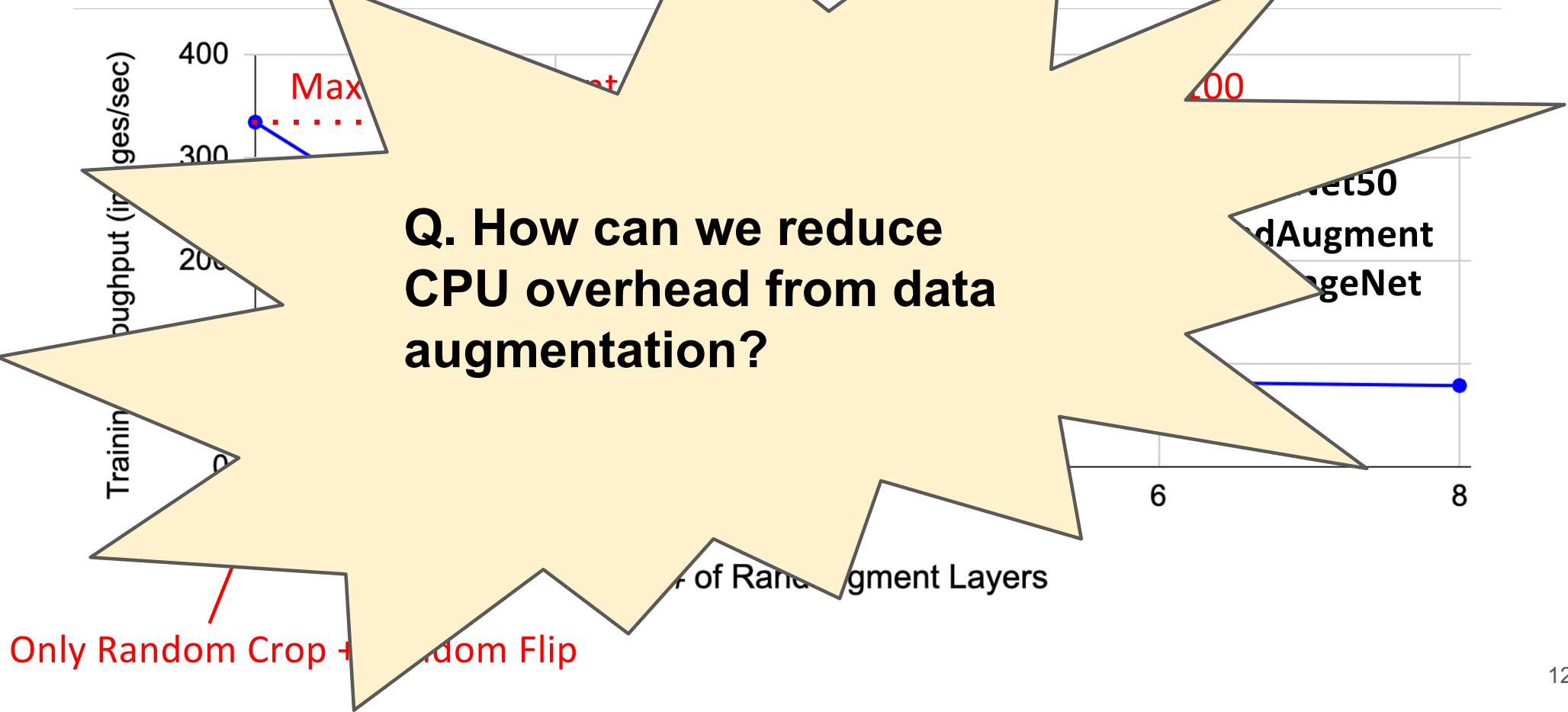
8

DNN Training Pipeline

# Overhead of Data Augmentation

- Investigate the impact of data augmentation overhead
- Workload: Training ResNet50 on ImageNet with RandAugment
  - Configuration: # of RandAugment Layers
- Environment: One NVIDIA V100 GPU with four physical CPU Cores
  - Same CPU-GPU ratio as cloud GPU VMs such as AWS P3 and GCP N1 instances
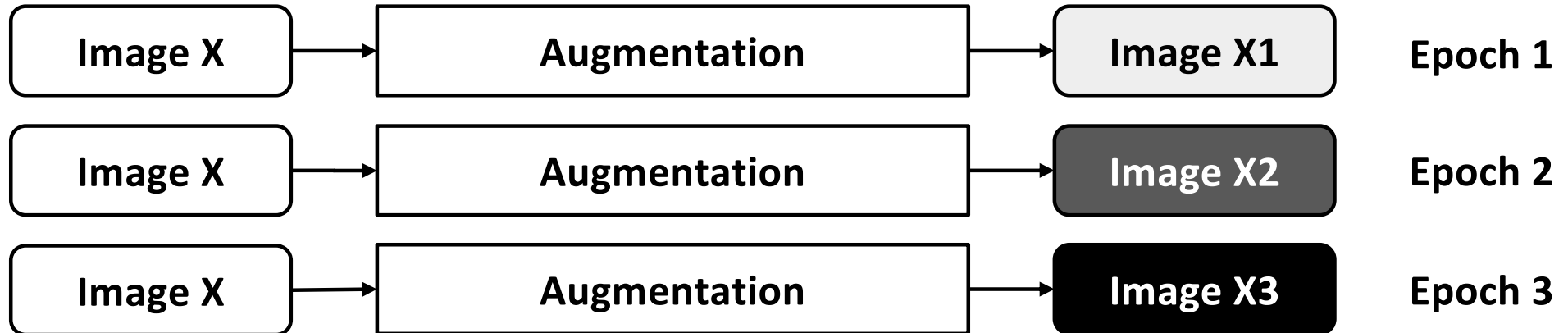
# Overhead of Data Augmentation

# Overhead of Data Augmentation



Max ...

Only Random Crop + Random Flip

Q. How can we reduce CPU overhead from data augmentation?

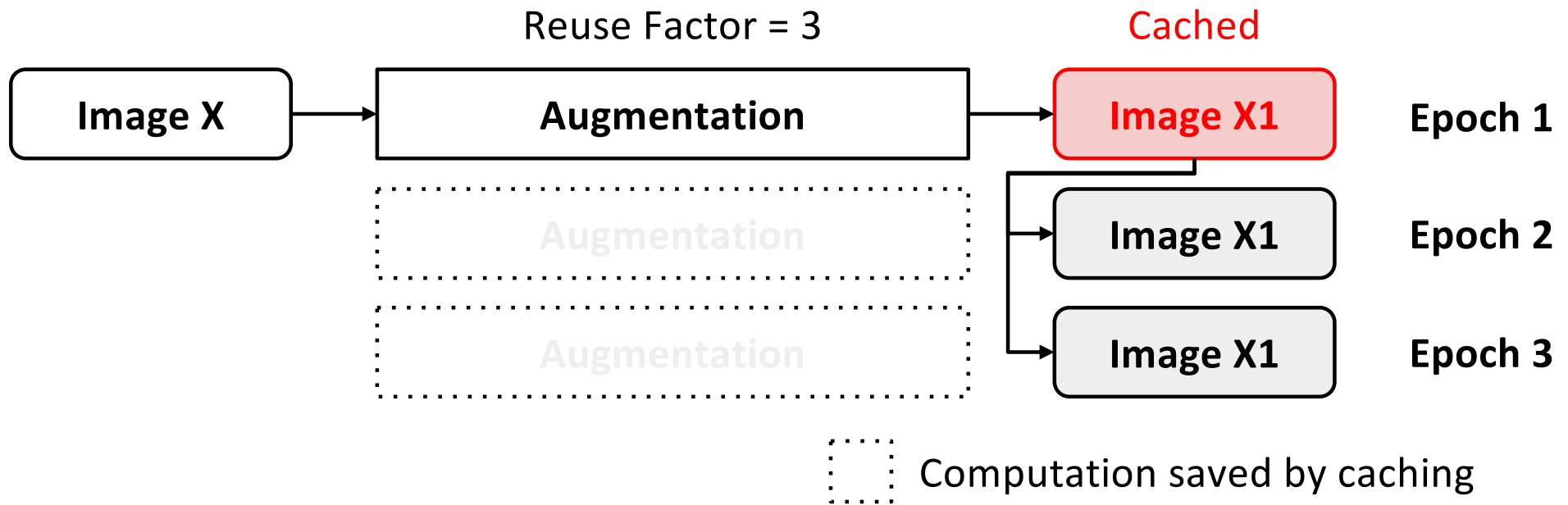# Existing Approach: Data Echoing

- Data echoing (arXiv '20, NeurIPS '20): Cache & reuse previously materialized samples
- Useful for training tasks with slow I/O
  - e.g., Training data on remote storage

# Standard Training

| | | | |
|---|---|---|---|
| Image X | → | Augmentation | → Image X1 |  Epoch 1 |
| Image X | → | Augmentation | → Image X2 |  Epoch 2 |
| Image X | → | Augmentation | → Image X3 |  Epoch 3 |

# Data Echoing

Problem: Sample diversity decreases to a great degree.
-> Low generalization of trained models

Reuse Factor = 3

Cached

| | | |
|---|---|---|
| Image X | Augmentation | Image X1 | Epoch 1 |

Augmentation

Image X1 — Epoch 2

Augmentation
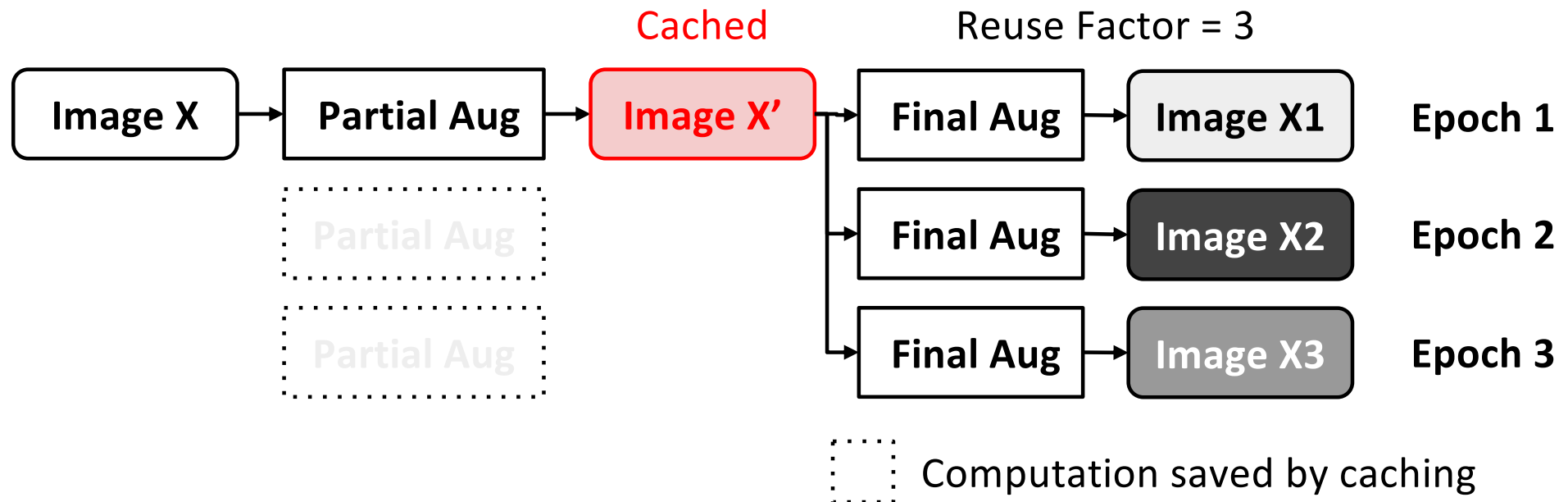
Image X1 — Epoch 3

Computation saved by caching

# Contents

- Background & Motivation
- **Data Refurbishing**
- Revamper
- Evaluation

# Our Approach: Data Refurbishing

**Solution:** Cache & reuse *partially augmented samples* by splitting augmentation pipelines
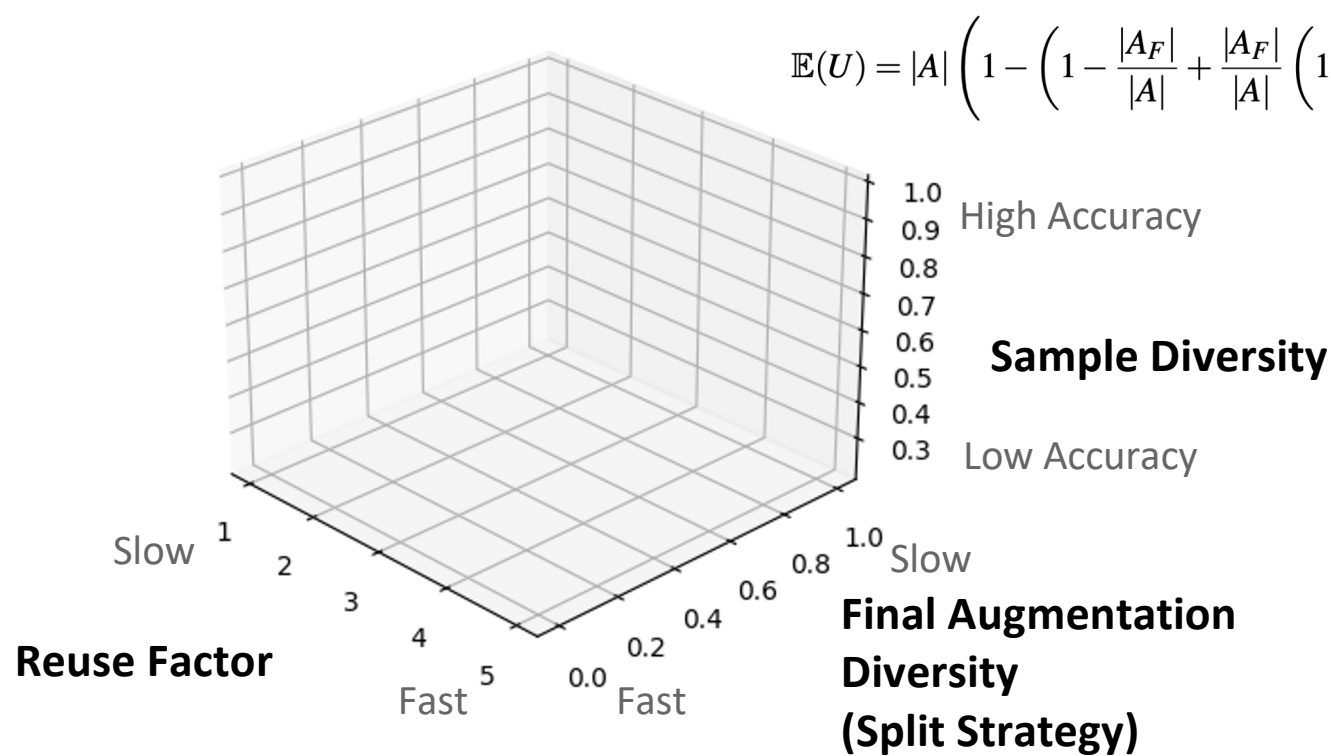


Cached

Reuse Factor = 3

Image X → Partial Aug → Image X' → Final Aug → Image X1    Epoch 1

Final Aug → Image X2    Epoch 2

Final Aug → Image X3    Epoch 3

Computation saved by caching

# Analysis on Sample Diversity

- Notations
  - Given a sample,
    - U (Sample Diversity): # of unique augmented samples during training
    - |A| (Augmentation Diversity): # of possible unique augmented samples by an augmentation pipeline A
    - $|A_F|$: The augmentation diversity of the final augmentation
  - r (Reuse Factor): # of reuses for each cached sample
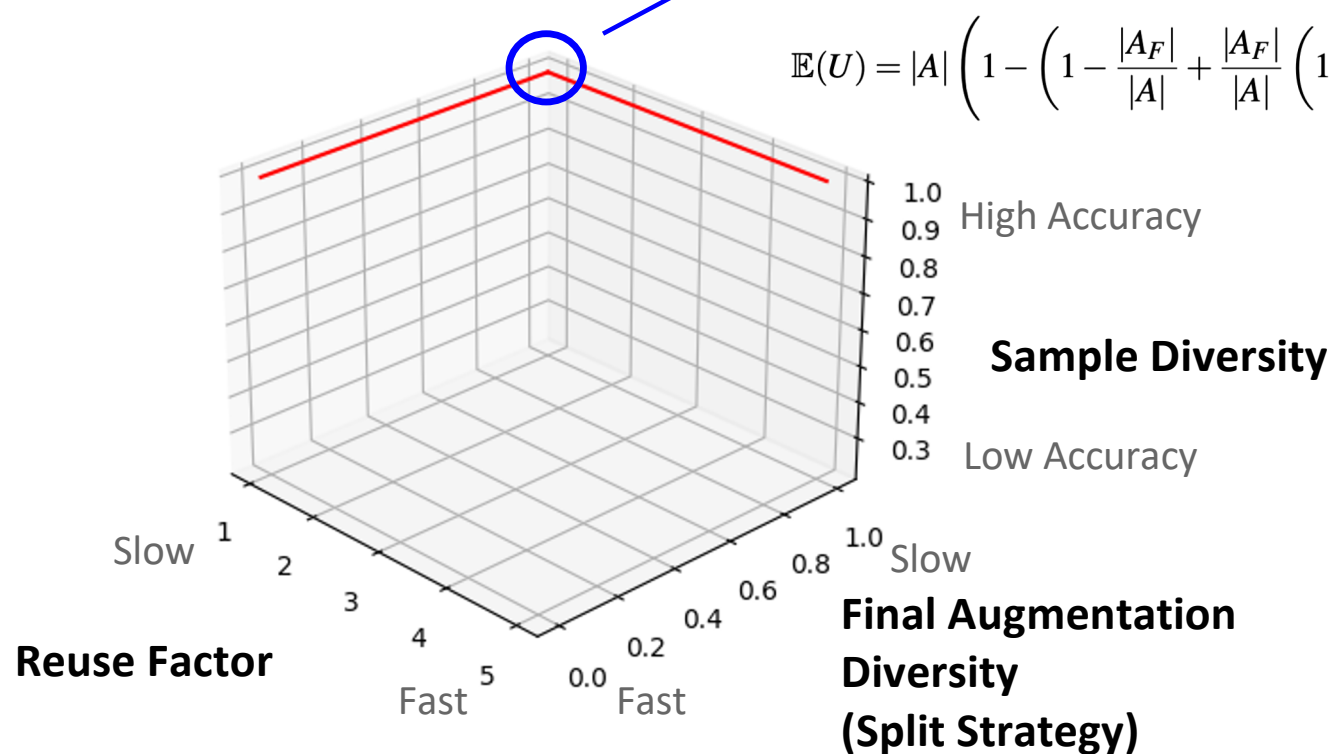  - k: The total number of training epochs

$$\mathbb{E}(U) = |A| \left( 1 - \left( 1 - \frac{|A_F|}{|A|} + \frac{|A_F|}{|A|} \left( 1 - \frac{1}{|A_F|} \right)^r \right)^{\frac{k}{r}} \right)$$

# Analysis on Sample Diversity

**Aug: RandAugment**
**k (# of epochs) = 300**

$$\mathbb{E}(U) = |A| \left( 1 - \left( 1 - \frac{|A_F|}{|A|} + \frac{|A_F|}{|A|} \left( 1 - \frac{1}{|A_F|} \right)^r \right)^{\frac{k}{r}} \right)$$

# Case #1: Standard Training

**$|A_F| = |A|$ or $r = 1$**

**Aug: RandAugment**
**k (# of epochs) = 300**
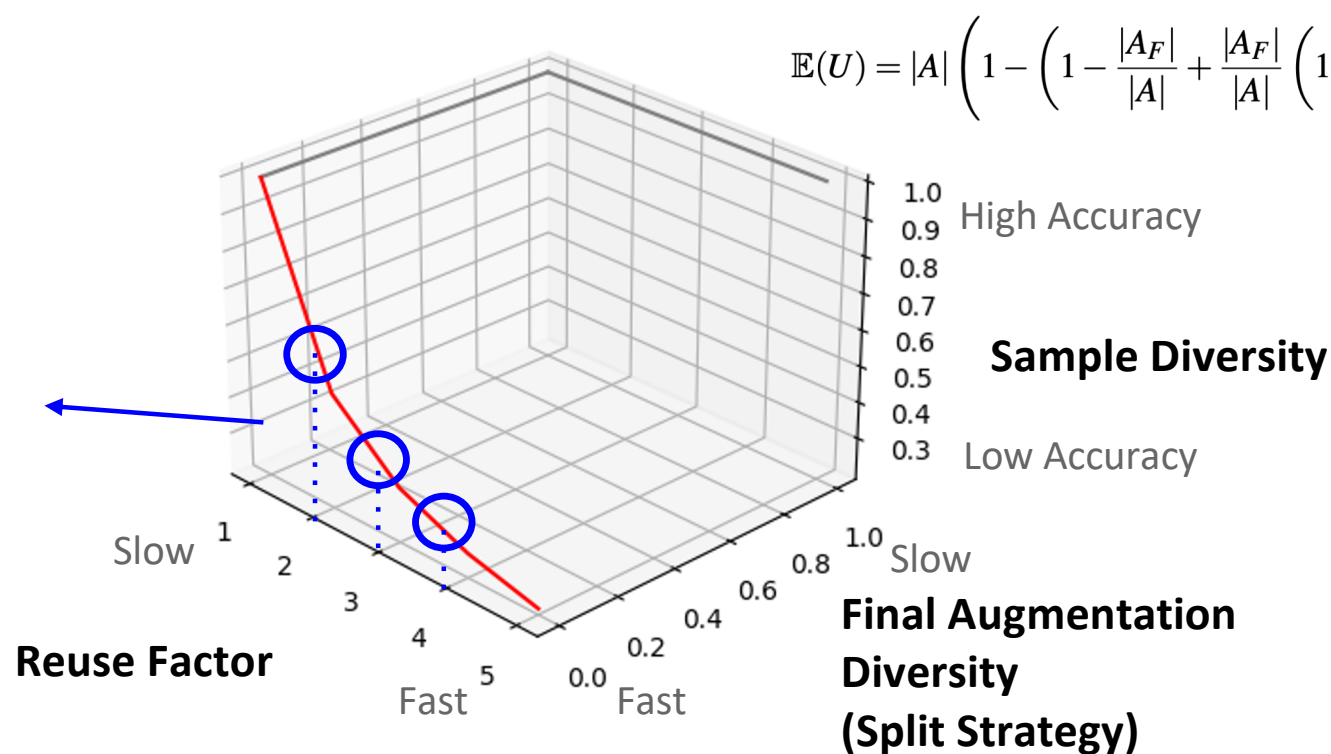
**High sample diversity
but low throughput**

$$\mathbb{E}(U) = |A| \left( 1 - \left( 1 - \frac{|A_F|}{|A|} + \frac{|A_F|}{|A|} \left( 1 - \frac{1}{|A_F|} \right)^r \right)^{\frac{k}{r}} \right)$$



High Accuracy

**Sample Diversity**

Low Accuracy

**Reuse Factor**

**Final Augmentation
Diversity
(Split Strategy)**

# Case #2: Data Echoing

**$|A_F| = 1$ and $r > 1$**

**Aug: RandAugment**
**k (# of epochs) = 300**

$$\mathbb{E}(U) = |A| \left( 1 - \left( 1 - \frac{|A_F|}{|A|} + \frac{|A_F|}{|A|} \left( 1 - \frac{1}{|A_F|} \right)^r \right)^{\frac{k}{r}} \right)$$
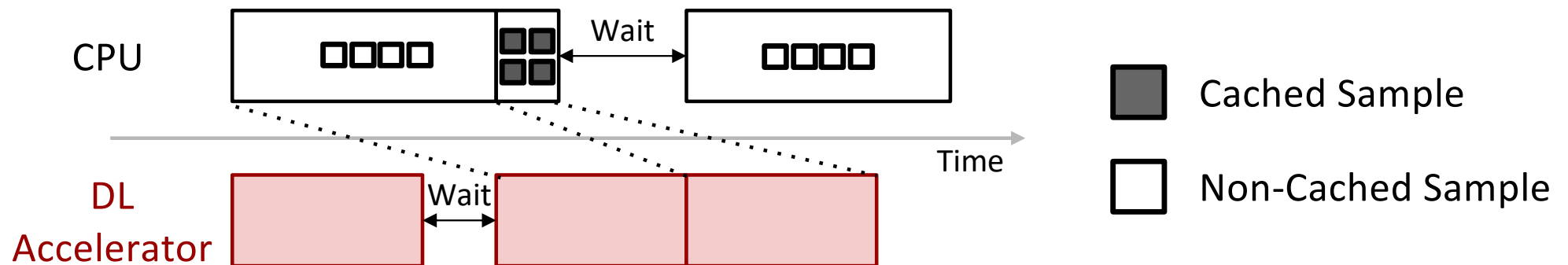
**High throughput but low sample diversity**

# Case #3: Data Refurbishing

**1 < $|A_F|$ < $|A|$ and r > 1**

**Aug: RandAugment**
**k (# of epochs) = 300**

$$\mathbb{E}(U) = |A| \left( 1 - \left( 1 - \frac{|A_F|}{|A|} + \frac{|A_F|}{|A|} \left( 1 - \frac{1}{|A_F|} \right)^r \right)^{\frac{k}{r}} \right)$$

**Exploit "sweet spot"**
**=> High throughput &**
**high sample diversity**



High Accuracy

1.0
0.9
0.8
0.7
0.6
0.5
0.4
0.3

**Sample Diversity**

Low Accuracy

Slow  1
      2
      3
      4
Fast  5

Slow
1.0
0.8
0.6
0.4
0.2
0.0
Fast

**Reuse Factor**

**Final Augmentation**
**Diversity**
**(Split Strategy)**

# Case #3: Data Refurbishing

**1 < |A$_F$| < |A| and r > 1**

**Aug: RandAugment**
**k (# of epochs) = 300**

**Exploit "sweet sp**
**=> High throughp**
**high sample dive**

$$\mathbb{E}(U) = |A| \left( 1 - \left( 1 - \frac{|A_F|}{|A|} + \frac{|A_F|}{|A|} \left( 1 - \frac{1}{|A_F|} \right)^r \right)^{\frac{k}{r}} \right)$$

Accuracy

mple Diversity

Accuracy

**Good Split Strategy**

1. Final augmentation has "enough" diversity
2. Final augmentation has low computation overhead

**Reuse Factor**

3
4
5
Fast

0.4
0.2
0.0
Fast

**Final Augmentation Diversity (Split Strategy)**

# Contents

- Background & Motivation
- Data Refurbishing
- **Revamper**
- Evaluation

# Challenge: Inconsistent Batch Time

- ## Within a mini-batch,
  - CPU processing time fluctuates according to the # of cache misses
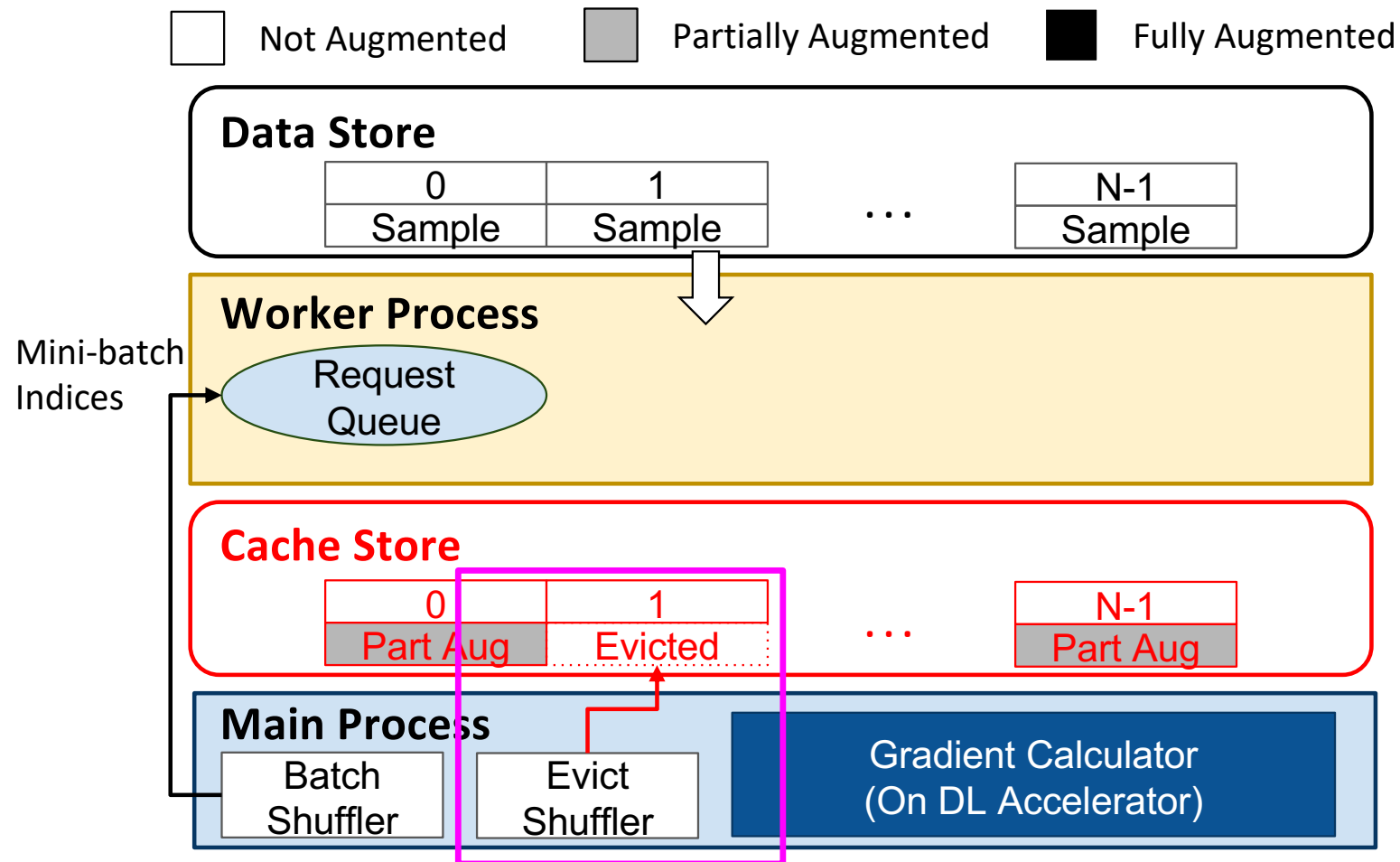  - Gradient computation time on DL accelerator remains the same

=> Poor computation overlap

# Challenge: Inconsistent Batch Time

- Within a mini-batch,
    - CPU processing time fluctuates according to the # of cache misses
    - Gradient computation time on DL accelerator remains the same

=> Poor computat

CPU

DL
Accelerator

Cached Sample

Non-Cached Sample

**Solution: Revamper**

1. Balanced Eviction: Balance # of cache misses across epochs
2. Cache-Aware Shuffle: Balance # of cache misses within an epoch

# PyTorch Dataloader

# Revamper

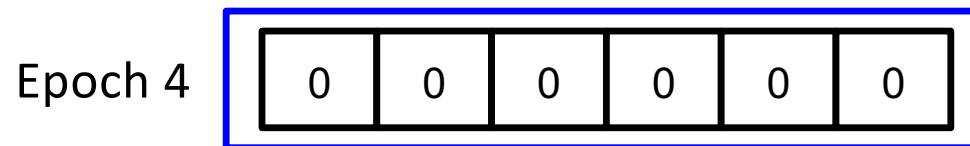# Revamper

# Revamper

# Balanced Eviction

Reuse Factor = 3

▮ Cached Sample       ☐ Non-Cached Sample

| Epoch 1 | 0 | 0 | 0 | 0 | 0 | 0 |

| Epoch 2 | 2 | 2 | 2 | 2 | 2 | 2 |

| Epoch 3 | 1 | 1 | 1 | 1 | 1 | 1 |

Fast: Possibly bottlenecked by DL accelerators

| Epoch 4 | 0 | 0 | 0 | 0 | 0 | 0 |

Slow: Possibly bottlenecked by CPU

**Naive (Reference Count)**

# Balanced Eviction

# Cache-Aware Shuffle
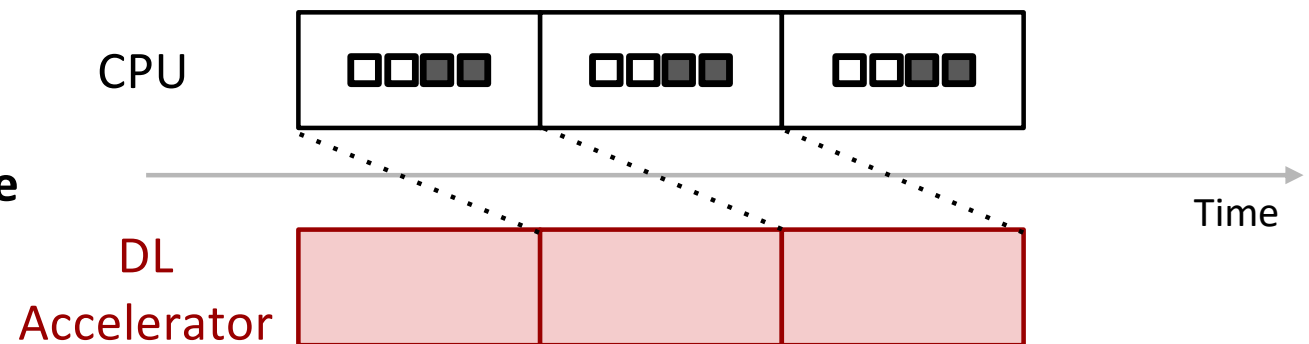
# Cache-Aware Shuffle

# Contents

- Background & Motivation
- Data Refurbishing
- Revamper
- **Evaluation**

# Implementation

- Implemented in 2000+ lines of Python code based on PyTorch 1.6
- Identical interface to the PyTorch dataloader except for some additional parameters
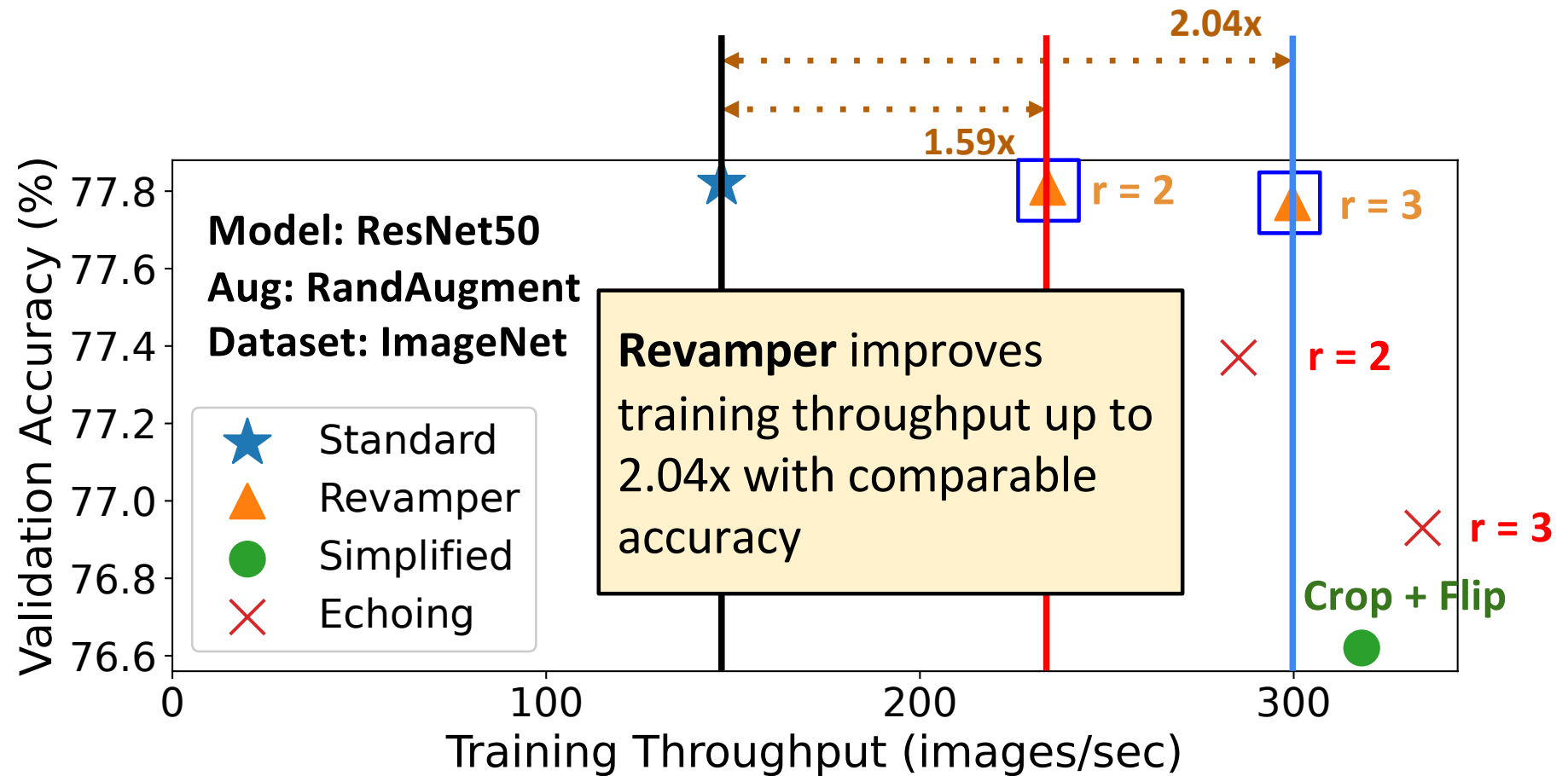  - e.g., reuse factor and split strategy

# Evaluation: Environments

- Training server specification
  - CPU: Intel Xeon E5-2695v4 (18 cores, 2.10GHz, 45MB Cache)
  - RAM: 256GB DRAM
  - GPU: NVIDIA V100
  - Disk: Samsung 970 Pro 1TB NVMe SSDs
- We adjust CPU-GPU ratios using a Docker container (Default = 4:1)
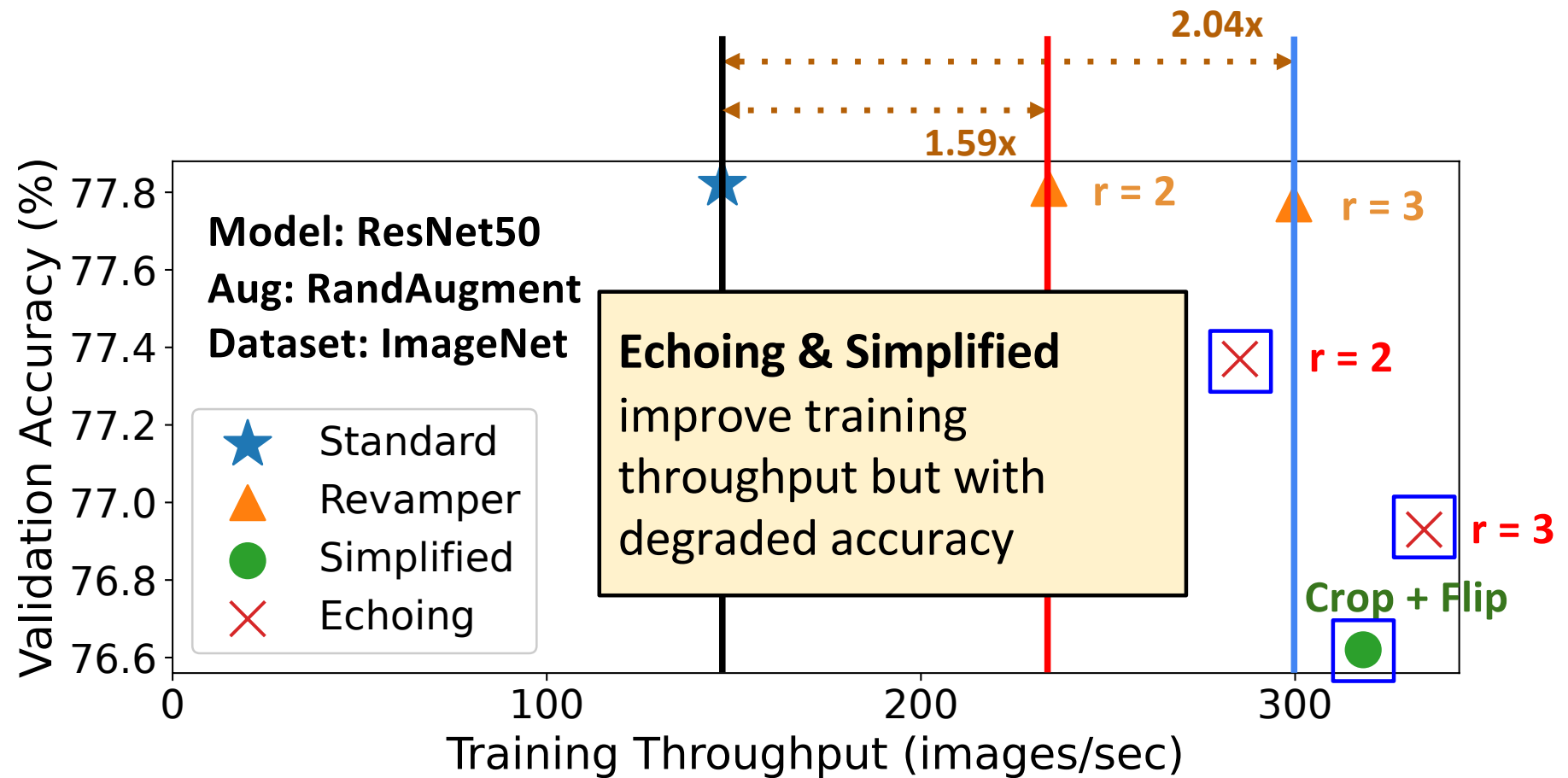- Workload: Image Classification

# Evaluation: Baselines

- Standard: DNN training without adopting data reusing mechanism
- Data Echoing: Cache & reuse fully augmented samples
- Simplified: Simply removing one or more transformation layers
- Same hyperparameters for all the training settings
  - Revamper does not require additional hyperparameter tuning
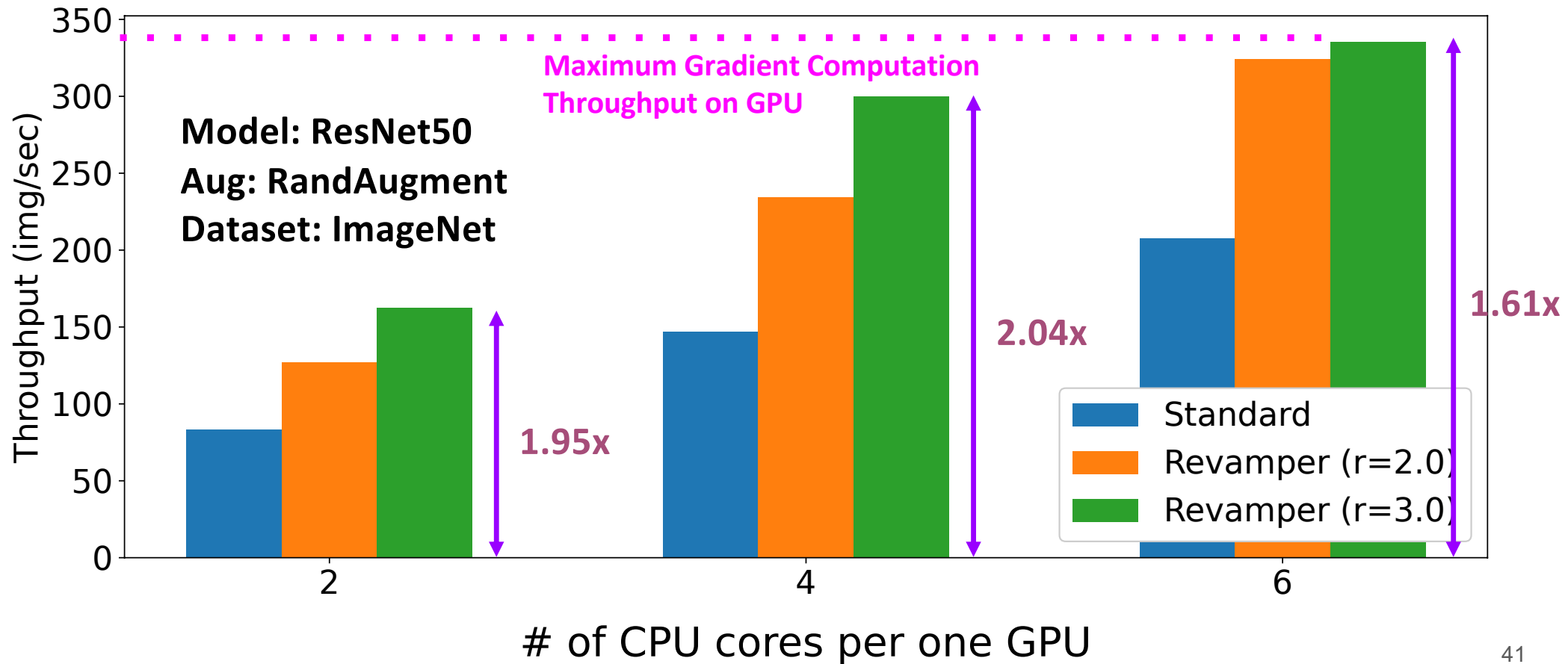
# Evaluation: Accuracy & Throughput

# Evaluation: Accuracy & Throughput

# Evaluation: CPU-GPU Ratio



Fewer CPUs -> Bigger Thp Gain

Maximum Gradient Computation Throughput on GPU

Model: ResNet50
Aug: RandAugment
Dataset: ImageNet

Throughput (img/sec)

1.95x

2.04x

1.61x

Standard
Revamper (r=2.0)
Revamper (r=3.0)

# of CPU cores per one GPU

41

# Conclusion

- **Data refurbishing** is a new intermediate data caching technique for DNN training that accelerates data augmentation while preserving diversity of augmented samples.

- **Revamper** realizes data refurbishing by maximizing computation overlap between CPU and DL accelerators with carefully-designed cache eviction and shuffle strategies.

- Revamper improves training throughput of DNN models by **1.03x-2.04x** while maintaining comparable accuracy.