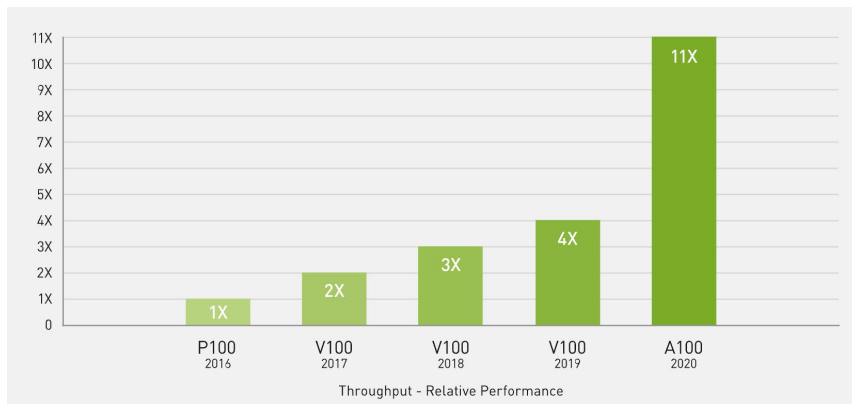


# FastFlow: Accelerating Deep Learning Model Training with Smart Offloading of Input Data Pipeline

*Taegeon Um, Byungsoo Oh, Byeongchan Seo,  
Minhyeok Kweun, Goeun Kim, Woo-Yeon Lee*

# GPUs are the Most Important Resources!



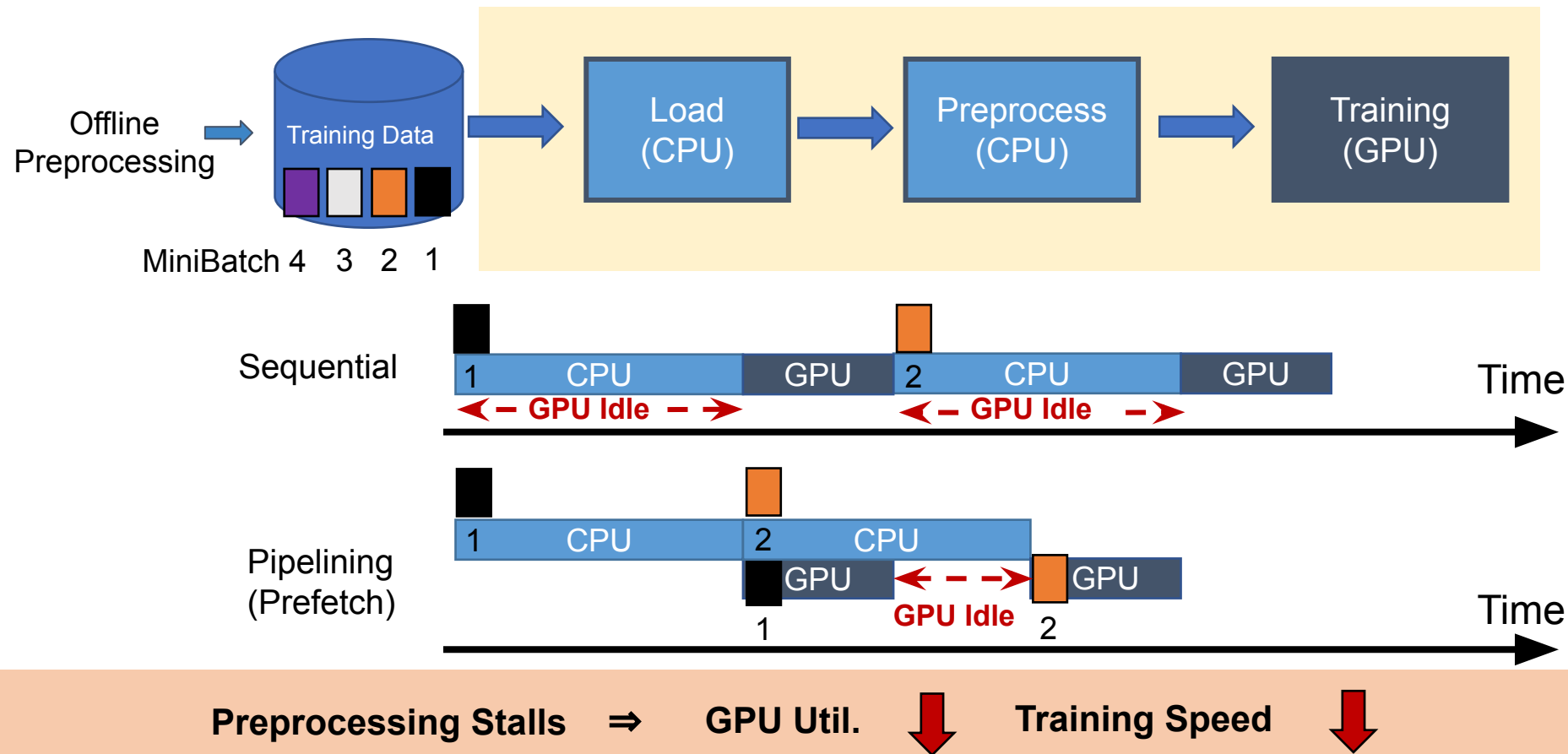
<https://www.nvidia.com/en-us/data-center/a100/>

## How to Reduce DL Training Time?

- **High-End, Expensive GPUs**
  - V100 << A100 << H100
- **More GPUs**
  - Training on 1 GPU << Training on 10 GPUs

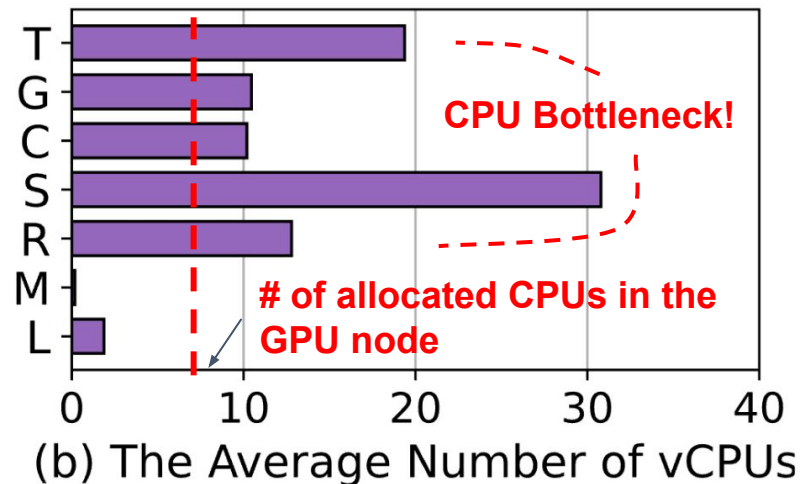
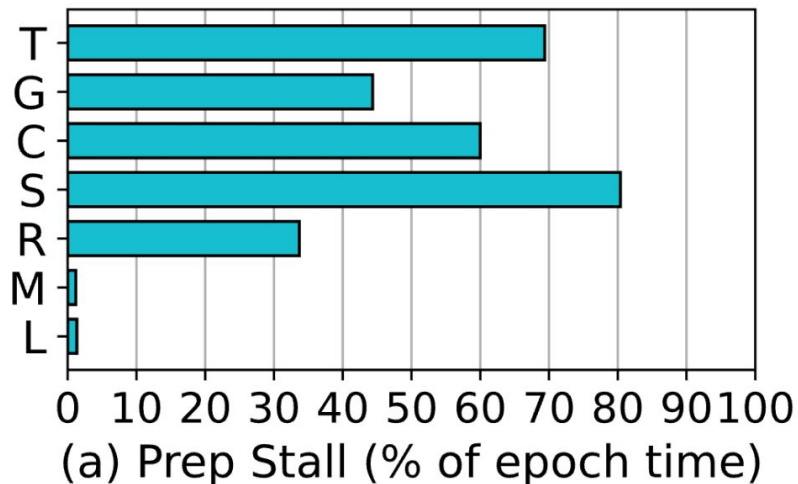
*Does this always guarantee training speed up?*

# Problem: CPU Bottlenecks and Prep Stalls



# Various Reasons of Prep Stalls

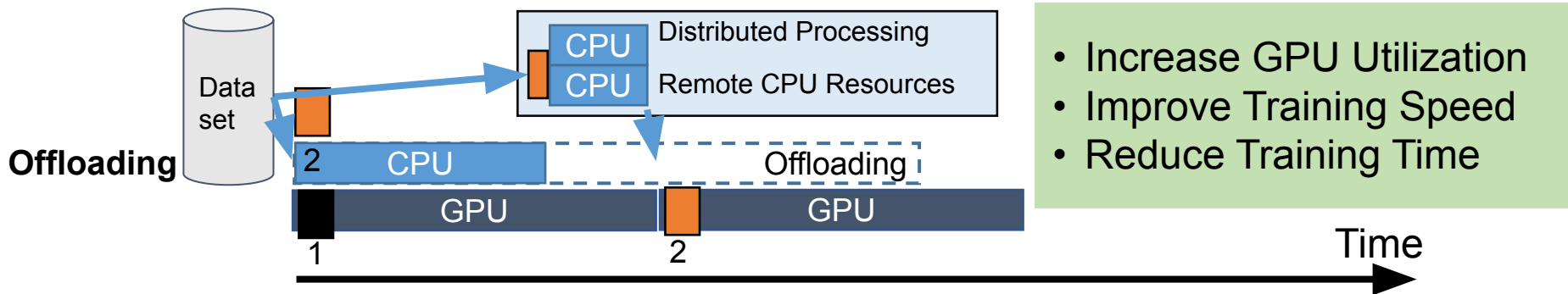
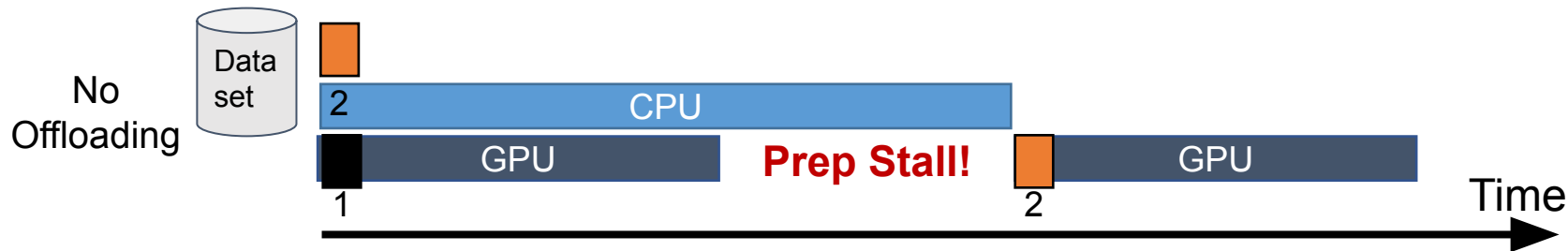
## Imbalance Computation Speed between GPU and CPU



- **Fixed CPU:GPU Ratio, but Various Workloads Require Diverse # of CPUs**
  - 8 vCPU cores : 1 V100 GPU (e.g., AWS P3.2xlarge)

- **Automatically Offloading Preprocessing to (Remote) CPU Resources**

- Automatic Decisions of **When** to Offload, **Which** Operator to Offload, **How** Much Data to Offload



- **Automatically Offloading Preprocessing to (Remote) CPU Resources**

- Automatic Decisions of **When** to Offload, **Which** Operator to Offload, **How** Much Data to Offload

## **Existing Work**

- tf.data.service: Manual offloading by users
- Cachew(ATC '22): Automatically decides the number of (CPU) workers and offloads all data and preprocessing operations to remote CPUs

## **Our Work**

- Automatically decides **when** to offload, **which** operators to offload, **how much data** to offload by considering diverse remote CPU environments (e.g., limited network bandwidth and CPU computations)

1

2

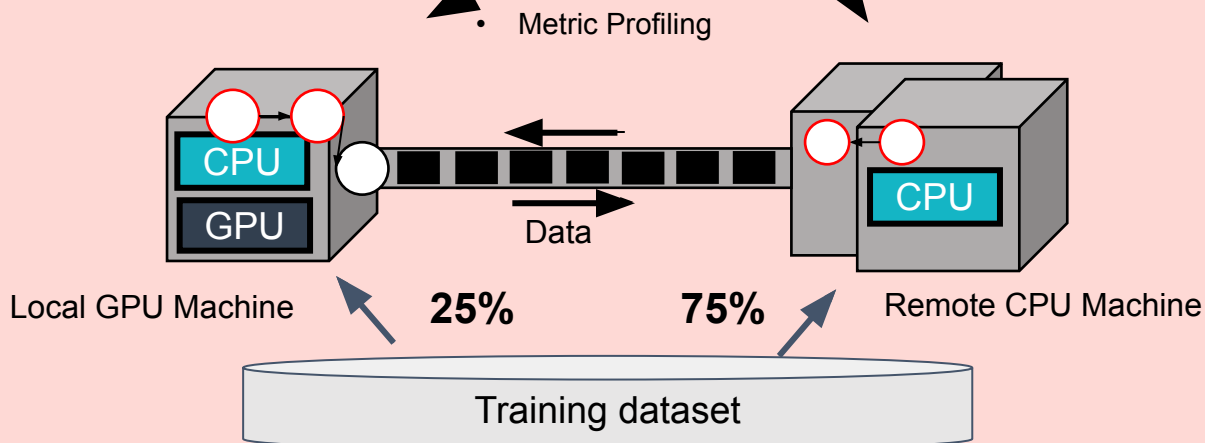
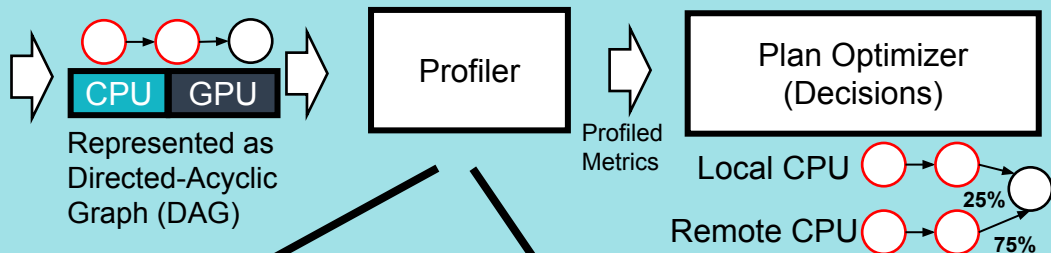
TIME

- Problem & Motivation
- Limitations of Existing Work and Our Approach
- **FastFlow Design**
- Demo
- Evaluation
- Conclusion

# FastFlow: A DLT System with Smart Offloading

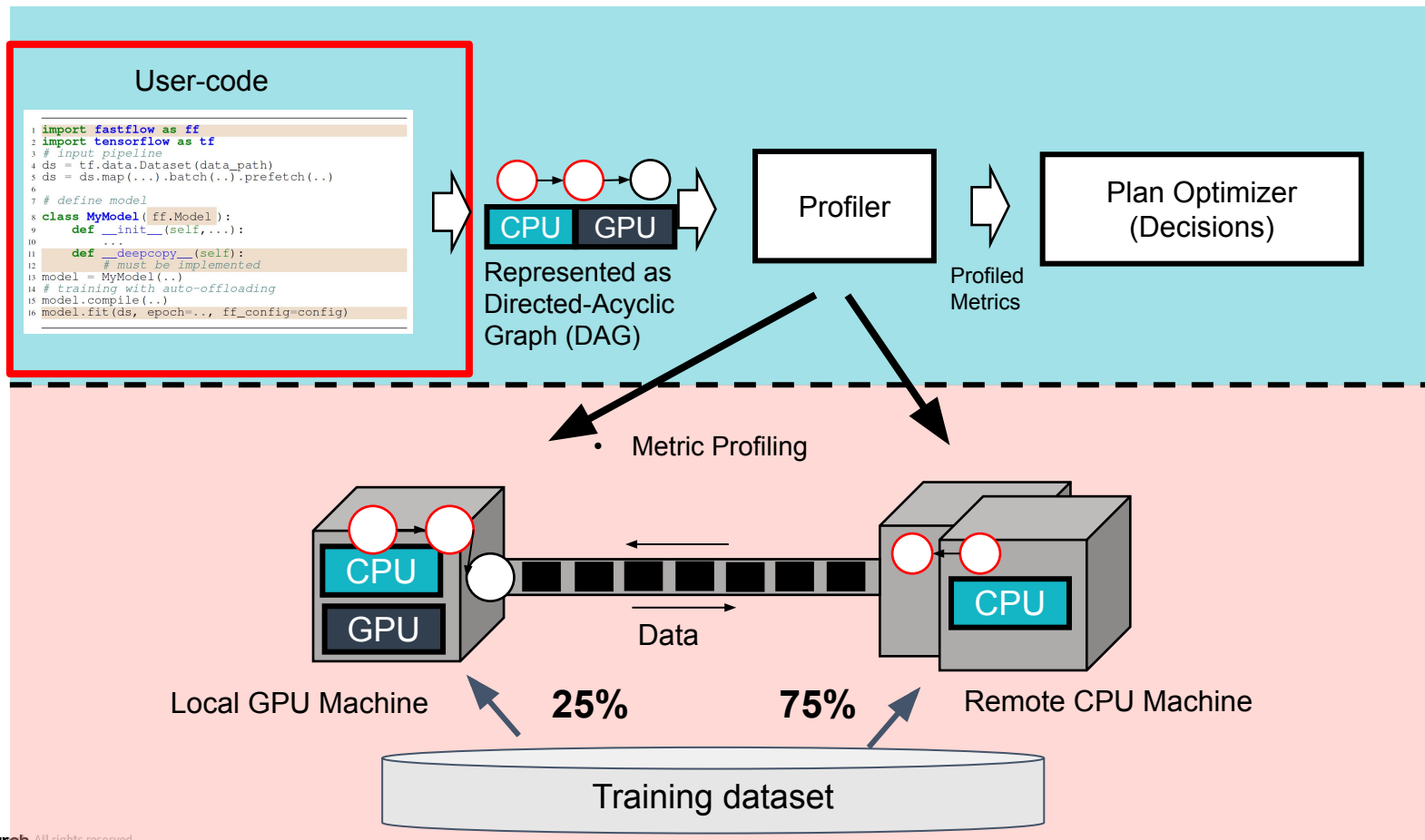
## User-code

```
1 import fastflow as ff
2 import tensorflow as tf
3 # input pipeline
4 ds = tf.data.Dataset(data_path)
5 ds = ds.map(...).batch(...).prefetch(...)
6
7 # define model
8 class MyModel(ff.Model):
9     def __init__(self,...):
10         ...
11     def __deepcopy__(self):
12         # must be implemented
13 model = MyModel(...)
14 # training with auto-offloading
15 model.compile(...)
16 model.fit(ds, epoch=..., ff_config=config)
```





# FastFlow: A DLT System with Smart Offloading



- Goal: Minimize the Modification of Existing TensorFlow Code for Users to Easily Adopt FastFlow

## TensorFlow Code

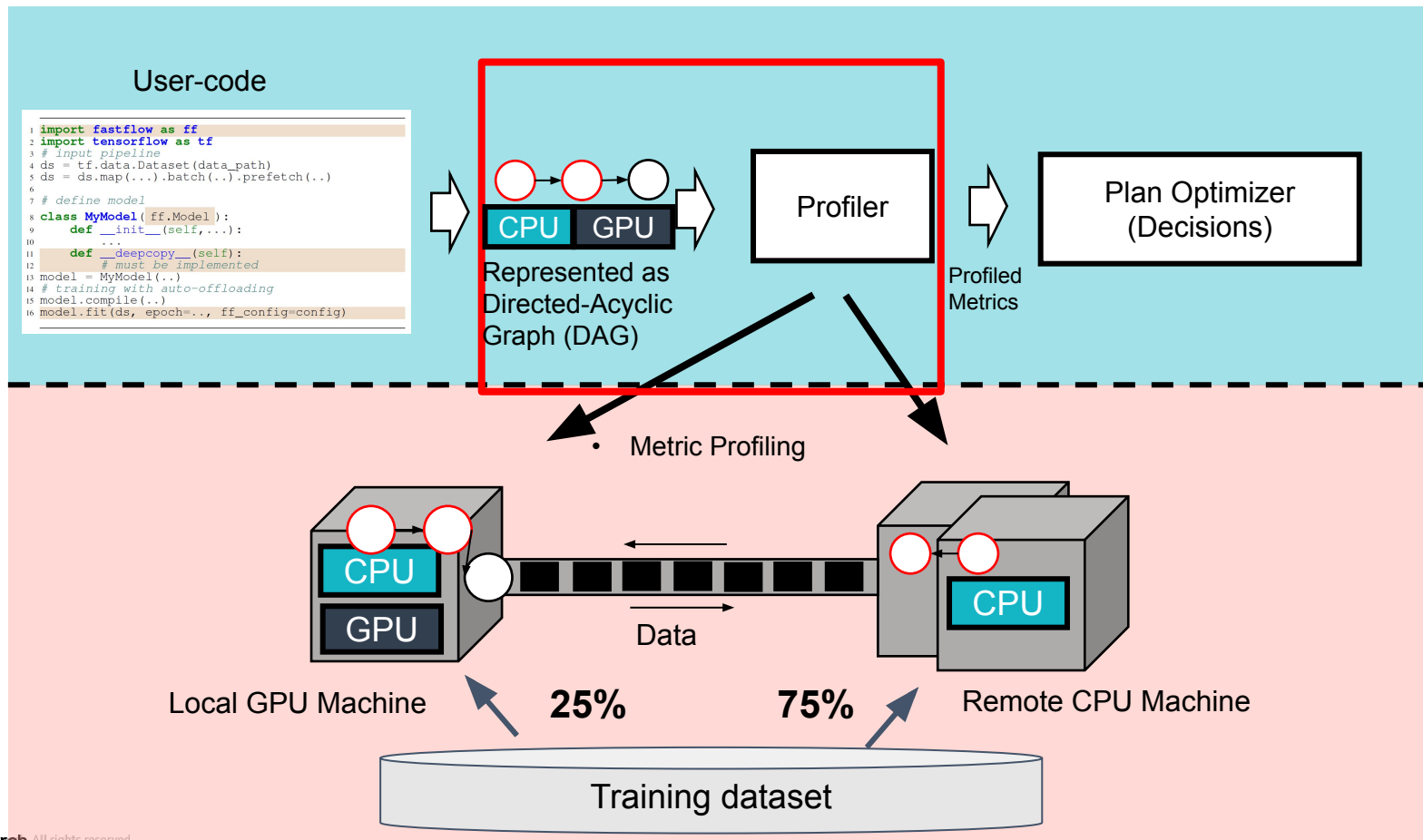
```
1 import tensorflow as tf
2 # Input pipeline
3 ds = tf.data.Dataset(data_path)
4 # Manual configuration for offloading
5 ds = ds.map(...)
6     .batch(...).prefetch(..)
7 # define model
8 class MyModel(tf.keras.Model):
9     def __init__(self,...):
10         ...
11 model = MyModel(..)
12 # training
13 model.compile(..)
14 model.fit(ds, epoch=..,)
```

## FastFlow Code

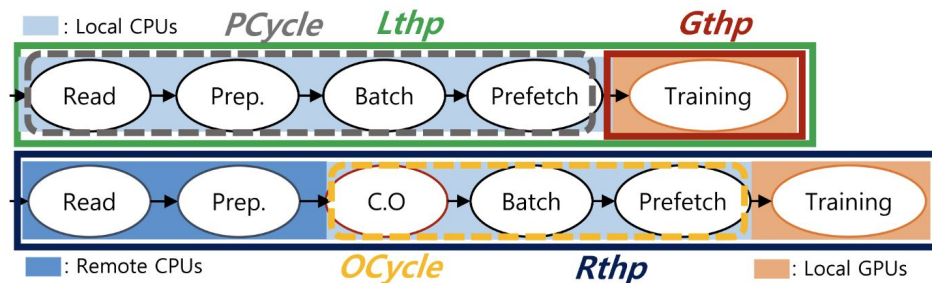
```
1 import fastflow as ff
2 import tensorflow as tf
3 # input pipeline
4 ds = tf.data.Dataset(data_path)
5 ds = ds.map(...).batch(...).prefetch(..)
6
7 # define model
8 class MyModel(ff.Model):
9     def __init__(self,...):
10         ...
11
12
13 model = MyModel(..)
14 # training with auto-offloading
15 model.compile(..)
16 model.fit(ds, epoch=.., ff_config=config)
```

- No Modification of Main TensorFlow Logic (Model Generation, Preprocessing)
- Only Need to Set Some Configuration Parameters for Smart Offloading

# FastFlow: A DLT System with Smart Offloading

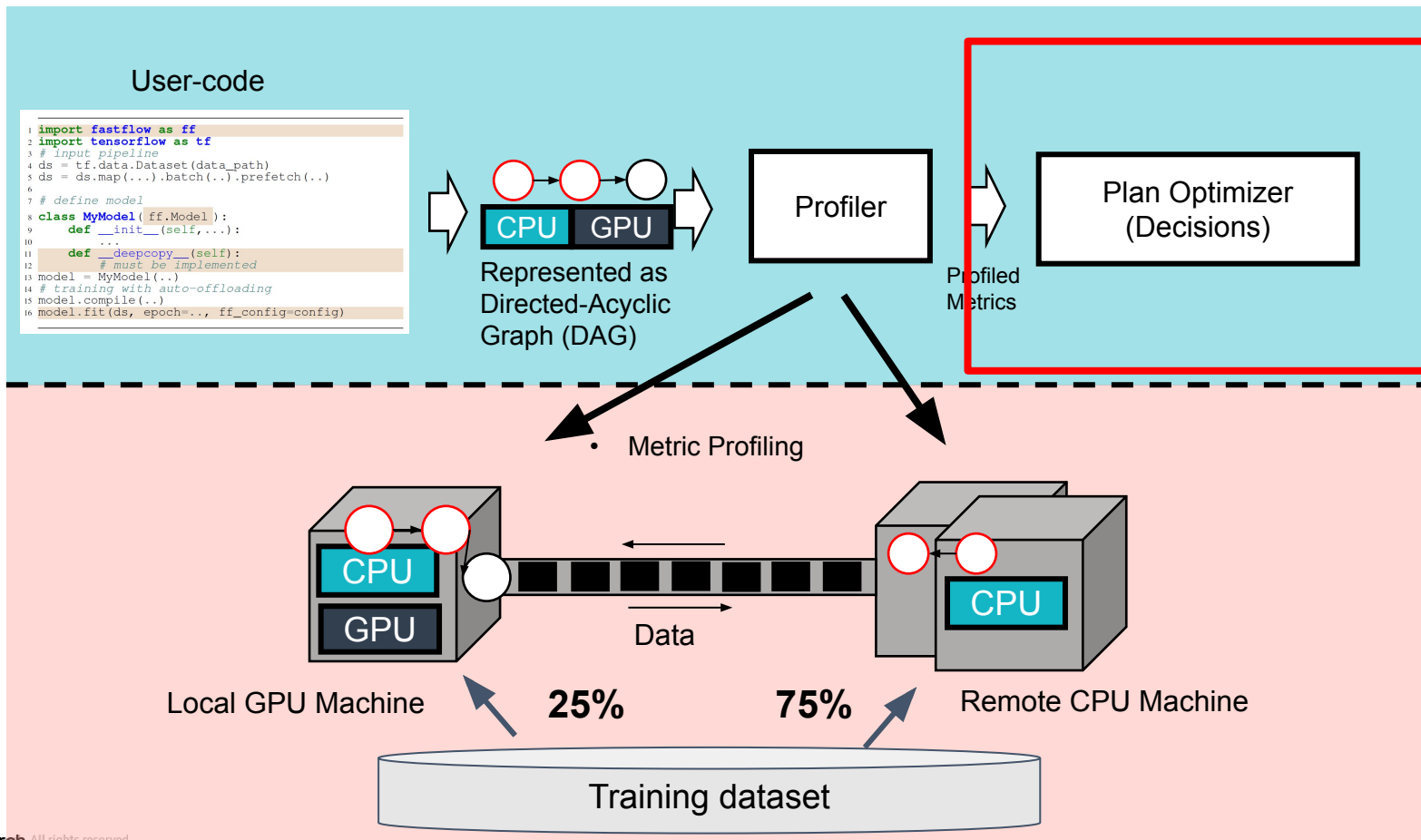


- Goal: Lightweight Profiling Overheads For Optimal Offloading Decisions
  - Characteristics of ML/DL Training; Iterative per Minibatch
  - Profiling A Few Steps (e.g., ~100)
    - Negligible Profiling Overheads (time), but High Accuracy of Profiled Metrics
  - Profile Metrics Related to Applications and Resource Environments
    - Ex) Training throughput with local or remote CPUs (Lthp/Rthp), GPU throughput (Gthp), Offloading overhead (PCycle, OCycle)



**Figure 4: An illustration of measuring metrics for offloading decisions.**

# FastFlow: A DLT System with Smart Offloading



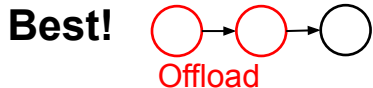
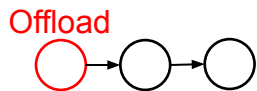
- Goal: Optimal Automatic Decisions for Offloading

1. When to Offload? (Do We Have to Offload?)

- If Preprocessing Speed on CPUs < Computation Speed on GPUs  $\Rightarrow$  Try to Offload! (CPU is the Bottleneck)  
Otherwise, Do Not Offload

2. Which Operations to Offload?

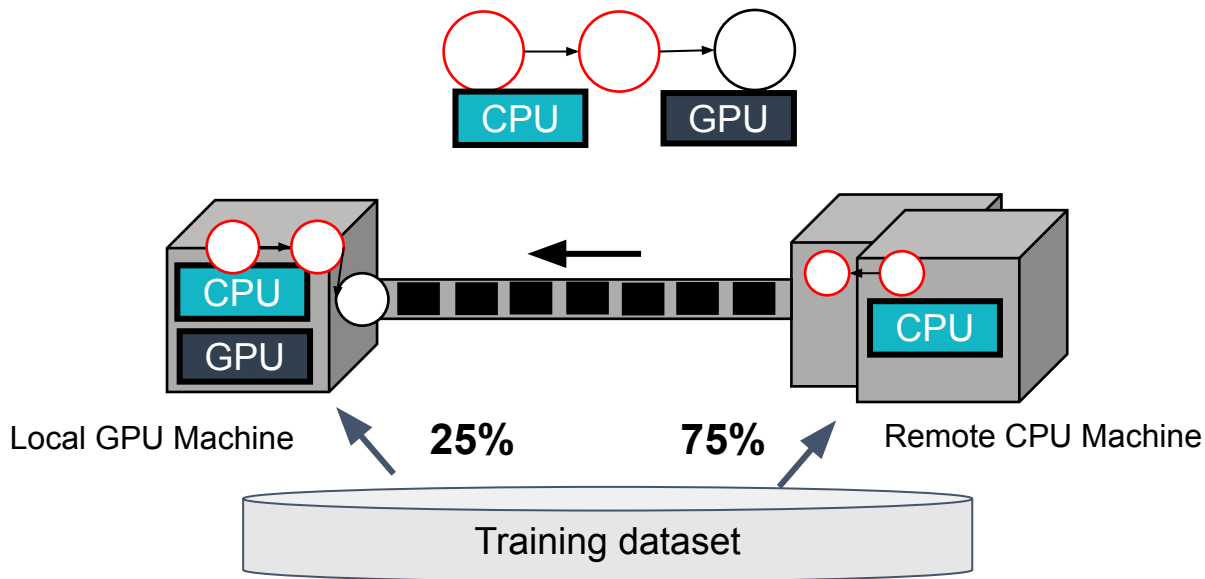
- Compare several candidate pipelines for offloading and choose the one that leads to maximum training throughput



- Goal: Optimal Automatic Decisions for Offloading

## 3. How Much Data to Offload?

- Considers Computing Capacity and Network Bandwidth of Remote Machines
- Ex) If Remote CPU Throughput = 3x Local CPU Throughput, then 75% in Remote 25% in Local



- Problem & Motivation
- Limitations of Existing Work and Our Approach
- FastFlow Design
- **Demo**
- Evaluation
- Conclusion



(tg-tf) byungsoo.oh@test:~/data-engineering-framework/fastflow/examples\$

test (debian buster/sid 64bit / Linux 3.10.0-957.27.2.el7.x86\_64) Uptime: 176 days, 4:54:59

- 2.80/2.80GHz CPU - 18.7% GPU Tesla V100-SX MEM - 7.4% SWAP - 0.0% LOAD 64-core  
CPU [ 18.7%] user 15.3% proc: 0% total 753G  
MEM [ 7.4%] system 3.0% mem: 0% used 55.8G  
SWAP [ 0.0%] iowait 0.0% temperature: 32C free 698G

NETWORK Rx/s Tx/s TASKS 24 (61 thr), 1 run, 23 slp, 0 oth sorted automatically by CPU consumption  
eth0 2Kb 3Kb  
lo 0b 0b

	CPU%	MEM%	VIRT	RES	PID	USER	TIME+	THR	NI	S	R/s	W/s
TCP CONNECTIONS	>2.4	0.0	351M	92.1M	4968	byungsoo.	1:11	1	0	R	0	0
Listen	0.0	0.0	199M	51.0M	4424	byungsoo.	0:00	1	0	S	0	0
Initiated	0	0.0	2.45G	29.2M	4389	byungsoo.	0:00	33	0	S	0	0
Established	0	0.0	101M	4.24M	4460	root	0:00	1	0	S	?	?
Terminated	3	0.0	375M	3.91M	4443	root	0:00	6	0	S	?	?
Tracked	0	0.0	21.9M	3.67M	4674	byungsoo.	0:03	1	0	S	0	0
	0.0	0.0	46.2M	3.33M	9135	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	46.2M	3.32M	9140	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	19.2M	3.25M	4675	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	19.1M	3.16M	4472	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	19.1M	3.14M	4695	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	19.1M	3.14M	4917	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	19.0M	2.86M	4957	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	56.3M	1.87M	4374	root	0:00	1	0	S	0	0
	0.0	0.0	56.3M	1.87M	4445	root	0:00	1	0	S	0	0
	0.0	0.0	101M	1.85M	4471	byungsoo.	0:02	1	0	S	0	0

DISK I/O R/s W/s  
dm-0 0 0  
nvme0n1 0 0  
nvme0n1p1 0 0  
sda 0 4K  
sda1 0 4K  
sda2 0 0

FILE SYS Used Total  
2022-11-30 11:01:09 UTC9G

[D] 0:00:00

Local  
GPU: 1 (V100)  
CPU: 7 cores

(tg-tf) byungsoo.oh@test:~\$

test (debian buster/sid 64bit / Linux 3.10.0-957.27.2.el7.x86\_64) Uptime: 176 days, 4:54:59  
Cloud flavor-baremetal instance i-0000102b (nova)

- 0.00/2.80GHz CPU - 19.1% 3 GPU Tesla V100- MEM - 7.4% SWAP - 0.0% LOAD 64-core  
CPU [ 19.1%] user 15.3% 0: 0% mem: 0% total 753G  
MEM [ 7.4%] system 3.0% 1: 0% mem: 0% used 55.8G  
SWAP [ 0.0%] iowait 0.0% 2: 0% mem: 0% free 698G

NETWORK Rx/s Tx/s TASKS 18 (55 thr), 1 run, 17 slp, 0 oth sorted automatically by CPU consumption  
eth0 256b 2Kb  
lo 0b 0b

	CPU%	MEM%	VIRT	RES	PID	USER	TIME+	THR	NI	S	R/s	W/s
TCP CONNECTIONS	>2.4	0.0	314M	55.1M	10636	byungsoo.	0:28	1	0	R	0	0
Listen	0.0	0.0	199M	51.0M	6428	byungsoo.	0:00	1	0	S	0	0
Initiated	0	0.0	2.45G	29.0M	6393	byungsoo.	0:00	33	0	S	0	0
Established	0	0.0	101M	4.25M	10596	root	0:00	1	0	S	?	?
Terminated	2	0.0	101M	4.25M	10616	root	0:00	1	0	S	?	?
Tracked	0	0.0	375M	3.66M	6447	root	0:00	6	0	S	?	?
	0.0	0.0	19.1M	3.21M	10609	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	18.9M	2.88M	10628	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	56.3M	1.87M	6378	root	0:00	1	0	S	0	0
	0.0	0.0	56.3M	1.87M	6448	root	0:00	1	0	S	0	0
	0.0	0.0	101M	1.86M	10607	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	101M	1.85M	10627	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	9.68M	1.43M	6453	byungsoo.	0:00	1	0	S	0	0
	0.0	0.0	70.6M	1.36M	6445	root	0:00	1	0	S	0	0
	0.0	0.0	4.53M	740K	1	root	0:00	1	0	S	0	0

DISK I/O R/s W/s  
dm-0 0 0  
nvme0n1 0 0  
nvme0n1p1 0 0  
sda 0 4K  
sda1 0 4K  
sda2 0 0

FILE SYS Used Total  
2022-11-30 11:01:09 UTC1G

[D] 0:00:00

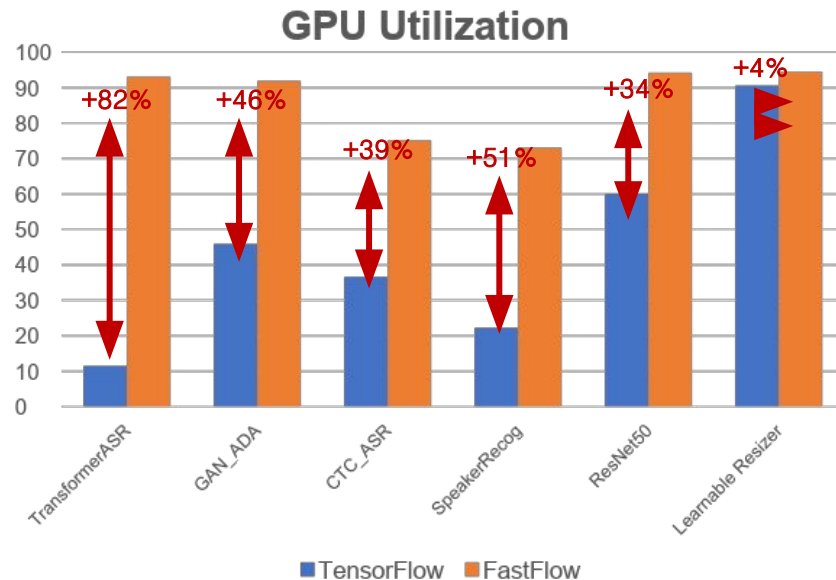
Remote  
CPU: 21 cores

- Problem & Motivation
- Limitations of Existing Work and Our Approach
- FastFlow Design
- Demo
- **Evaluation**
- Conclusion

- Q) How Much Does FastFlow Improve Performance And GPU Util. Compared to TensorFlow w/o Offloading? (21 remote CPUs)

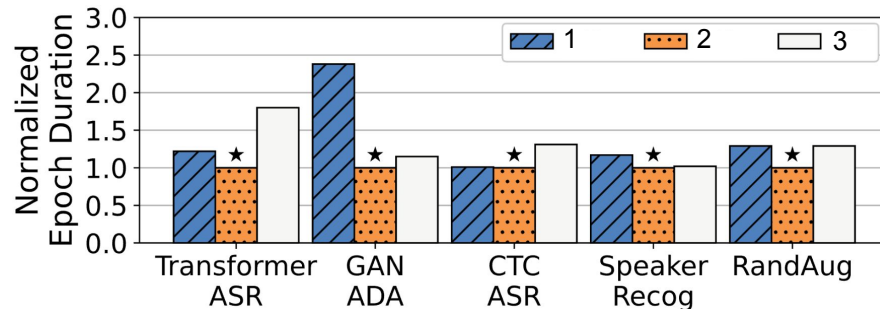
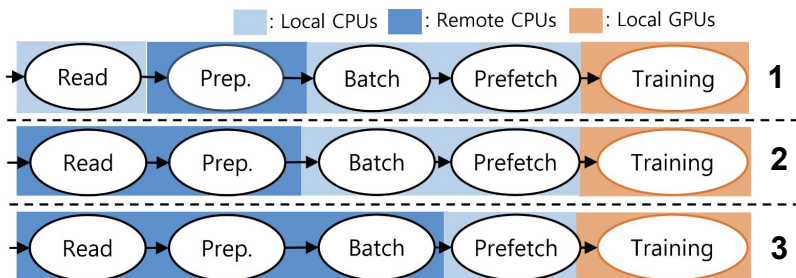


- FastFlow Accelerates Training Up to 3.8x Compared to TensorFlow



- FastFlow Increases GPU Utilization from 4%~82% Points Compared to TensorFlow

- Q) Does FastFlow Make Optimal Decisions for Offloading?
  - Decision of Which Operations to Offload**

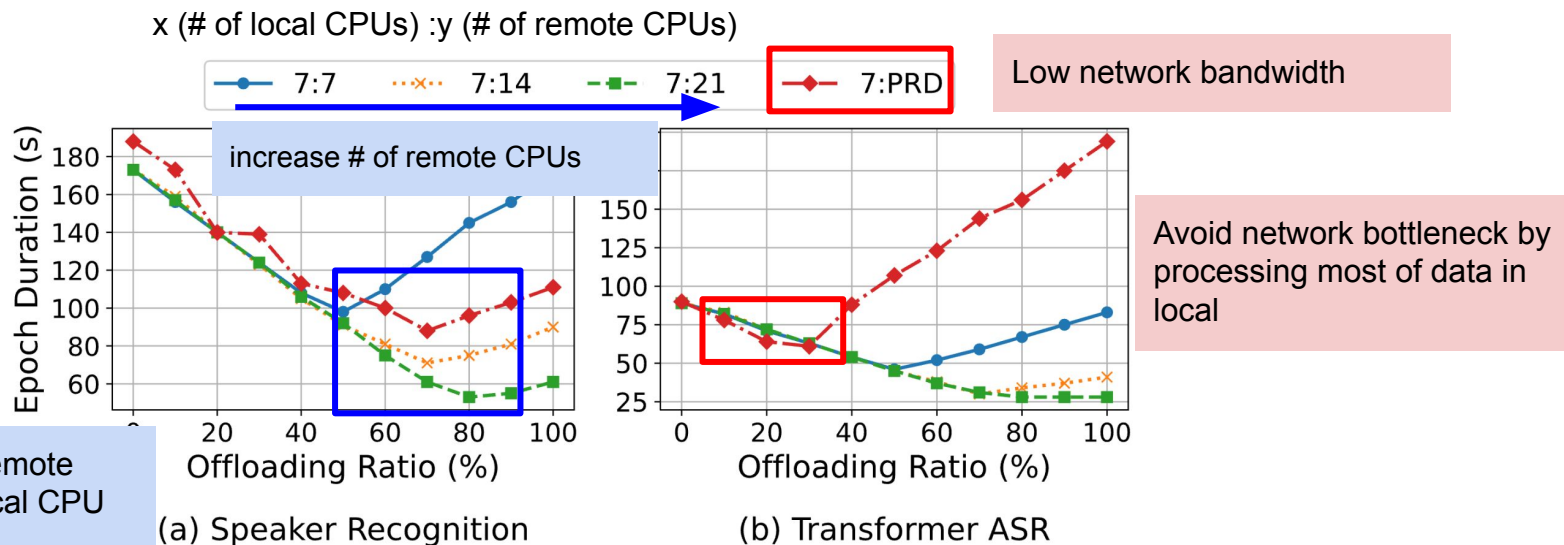


<Evaluation on remote CPUs with high disk and network bandwidth>

**Figure 8: Impact of the offloading operator selection of FastFlow.** \* mark shows the selection of FastFlow.

- FastFlow chooses **2** when remote nodes have high disk and network bandwidth
- If remote CPU node has low disk I/O, FastFlow chooses **1** to prevent disk bottleneck.

- Q) Does FastFlow Make Optimal Decisions for Offloading?
  - Decision of How Much Data to Offload



**Figure 9: Epoch duration in various offloading ratios on different remote workers. (a) is compute-intensive, and (b) is data-intensive preprocessing workload.**

## ▪ Q) How Long Does the Profiling Take Time ?

Job	Profile time (min)	Train time (min)	Fraction (%)
T	1.91 / 0.08*	68.37	2.80 / 0.11*
G	0.61 / 0.08*	162.22	0.37 / 0.05*
C	9.27 / 0.20*	299.32	3.10 / 0.07*
S	1.28 / 0.02*	150.58	0.85 / 0.01*
R	7.01 / 0.24*	11396.50	0.06 / 0.002*
M	0.78 / 0.53*	933.35	0.08 / 0.06*
L	0.65 / 0.26*	49.23	1.31 / 0.53*

## ▪ Profiling Overhead is Negligible (Only ~3% of Total Training Time)

 : Profiling Without Metric Store/Reload

 : Profiling With Metric Store/Reload.

FastFlow Further Reduces Profiling Time By Reusing Metrics Measured Before

- Problem & Motivation
- Limitations of Existing Work and Our Approach
- FastFlow Design
- Demo
- Evaluation
- **Conclusion**

- Although GPUs are the most important and expensive resource for DL training, it is under-utilized because of CPU bottlenecks
- To address the CPU bottlenecks, we designed and implemented FastFlow. We evaluated that FastFlow can significantly improve training speed with automatic decisions of offloading
- FastFlow can be applied to diverse applications on various resource environments without the modification of main training logic



# Thank you