清華大学
Tsinghua University

# Mobius: Fine Tuning Large-Scale Models on Commodity GPU Servers
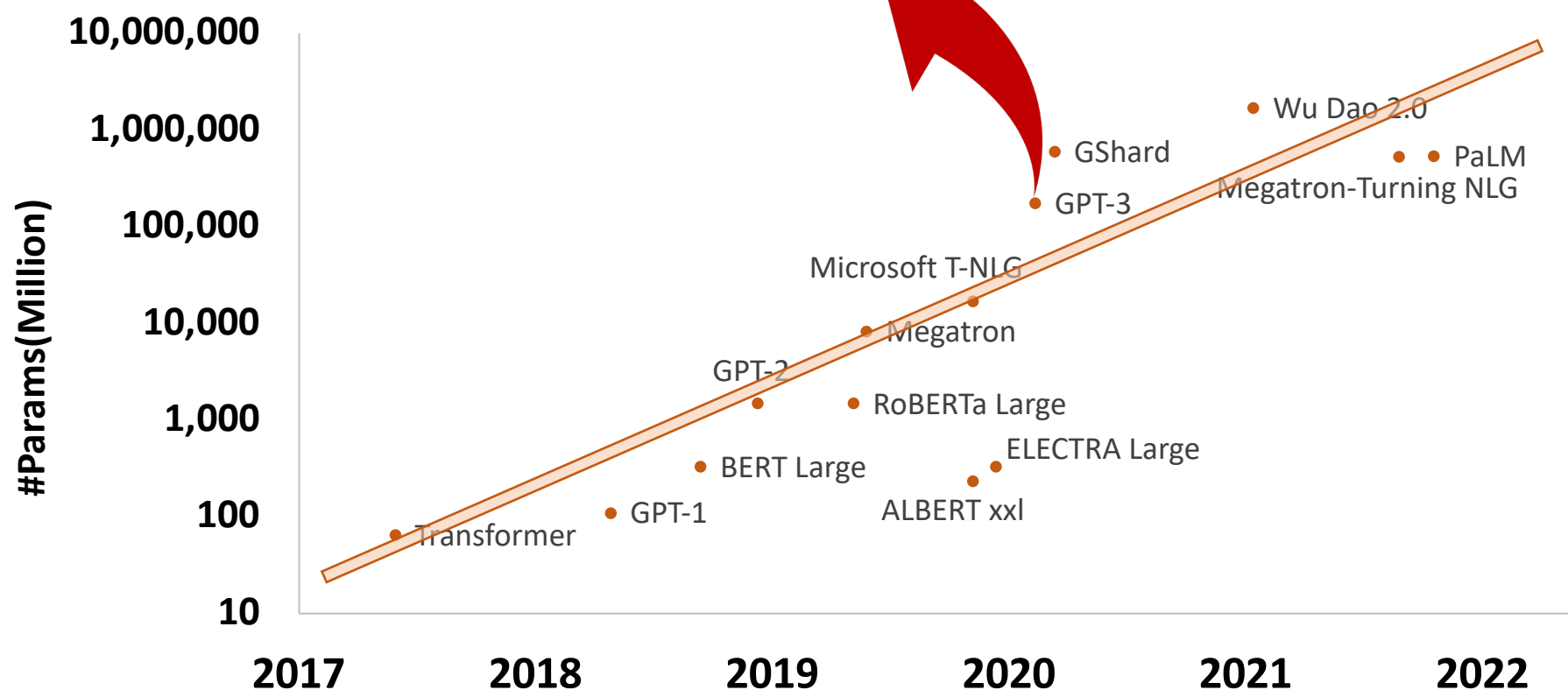
**Yangyang Feng**, Minhui Xie, Zijie Tian, Shuo Wang, Youyou Lu, and Jiwu Shu

*Tsinghua University*

*http://storage.cs.tsinghua.edu.cn*

# Explosive Growth of Model Size

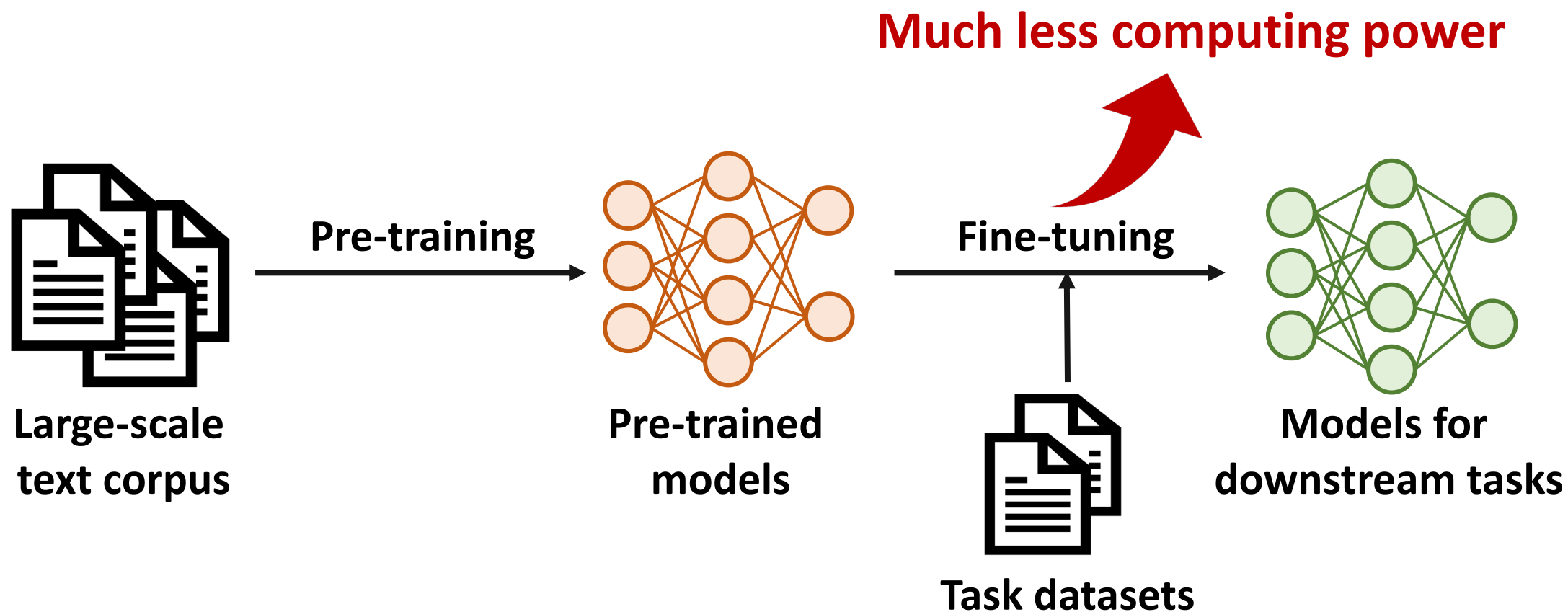**3640 petaflop/s-day ≈ A100 x 30 years**
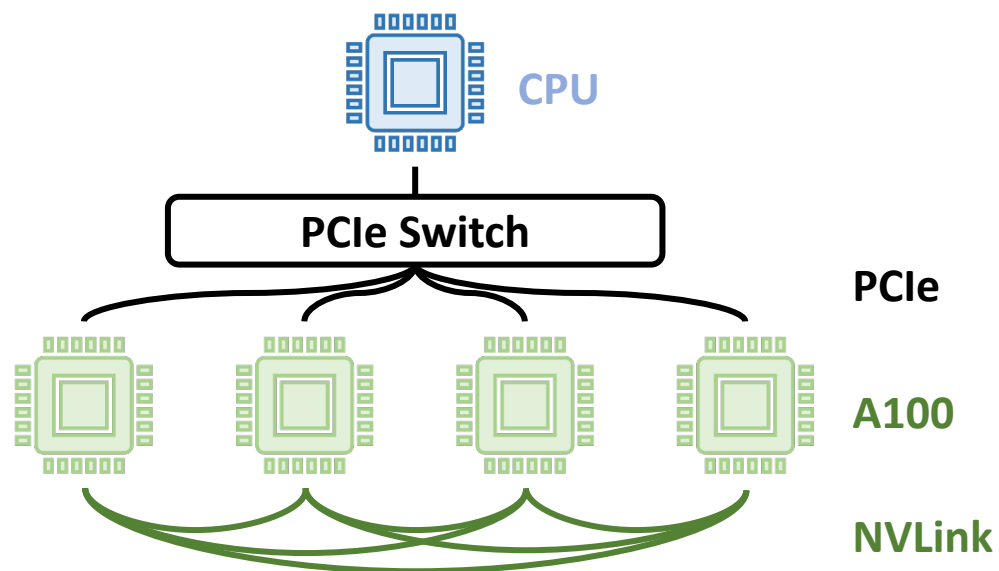
**≈ $4.5 million**



[1] Brown, Tom, et al. "Language models are few-shot learners." Advances in neural information processing systems 33 (2020): 1877-1901.

# Pre-training and Then Fine Tuning

**Much less computing power**



Large-scale text corpus → **Pre-training** → Pre-trained models → **Fine-tuning** → Models for downstream tasks

Task datasets
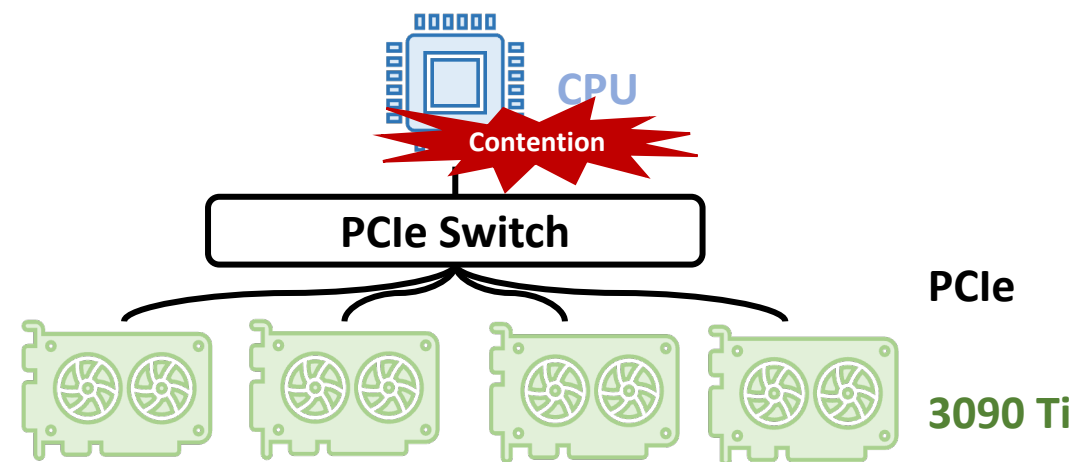
# Commodity GPU Servers



**Data Center GPU Server**

$100, 000

432 tensor cores

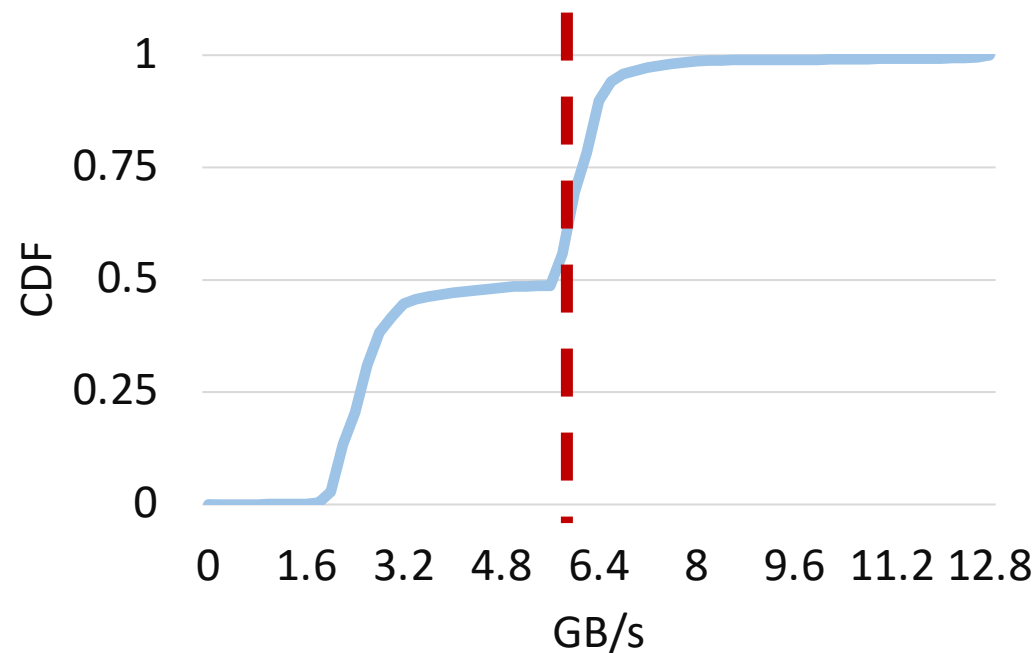900 GB/s inter-GPUs comm.

GPU Direct P2P

**Commodity GPU Server**
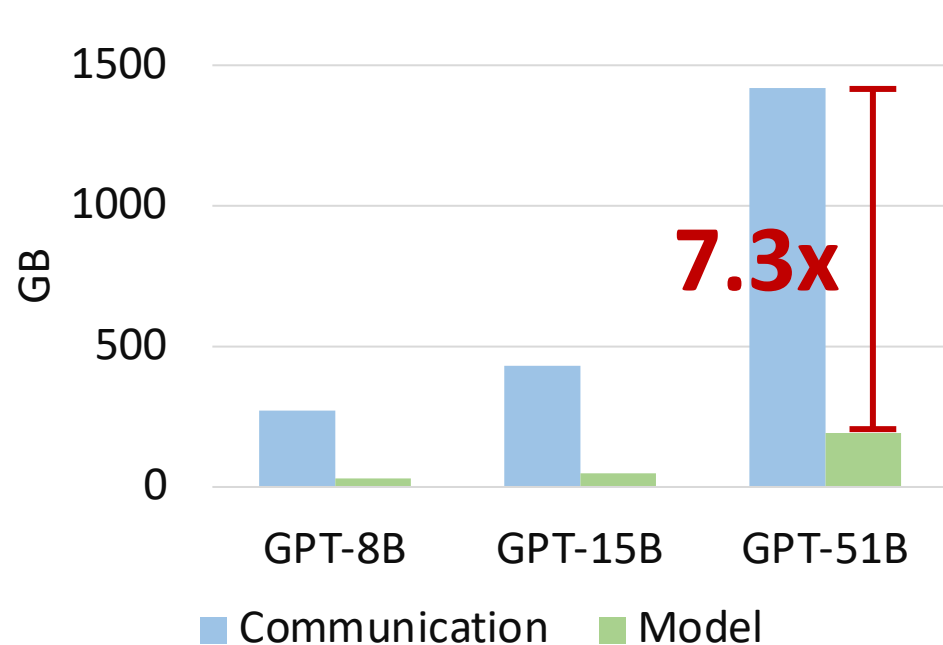
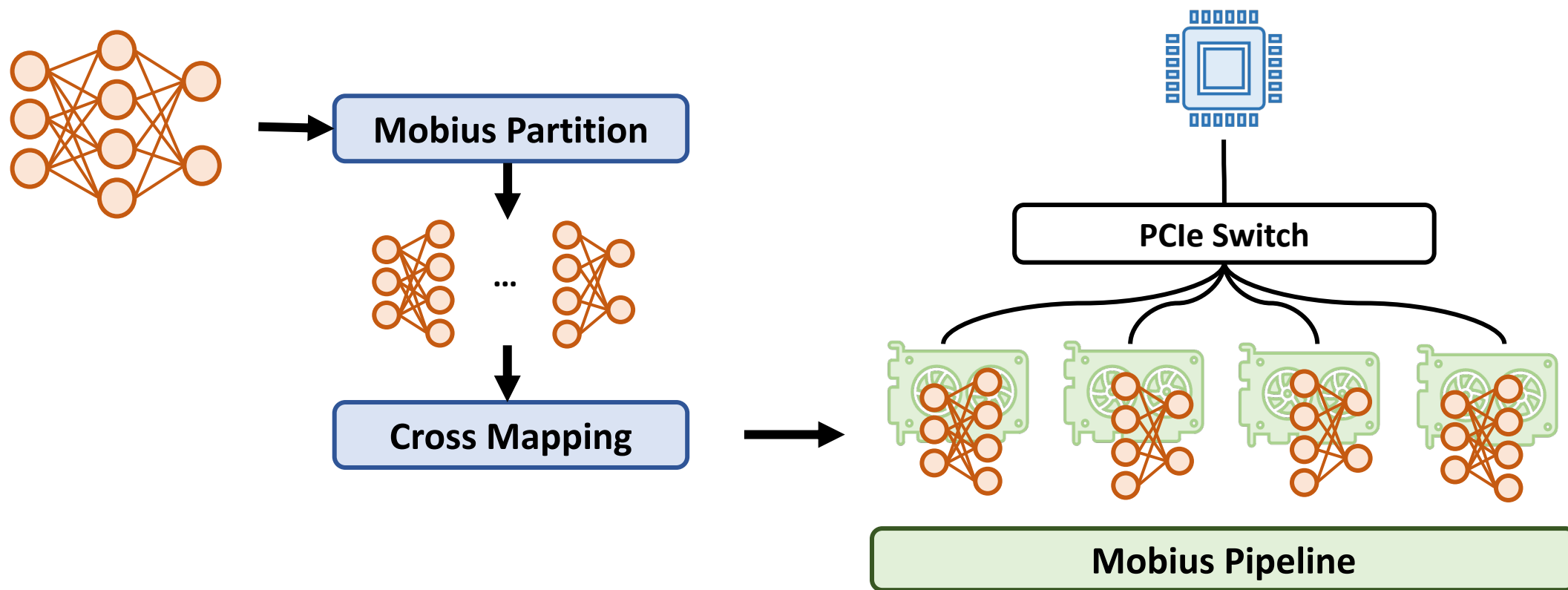$10, 000

336 tensor cores

16 GB/s inter-GPUs comm.
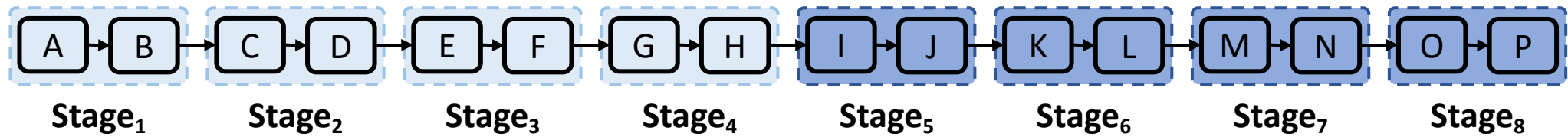
NO GPU Direct P2P

# Training on Commodity GPU Servers



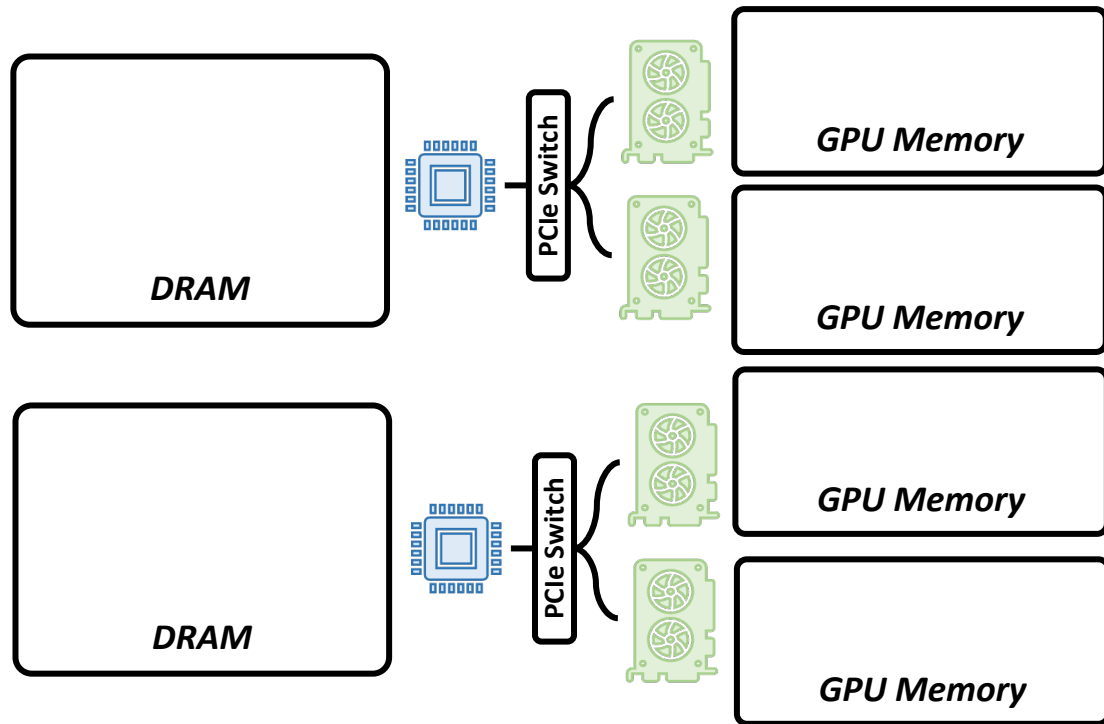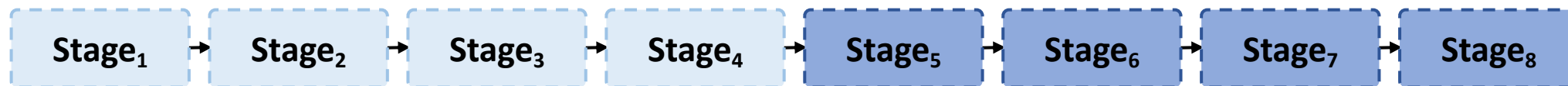**70% of training time is spent on communication**
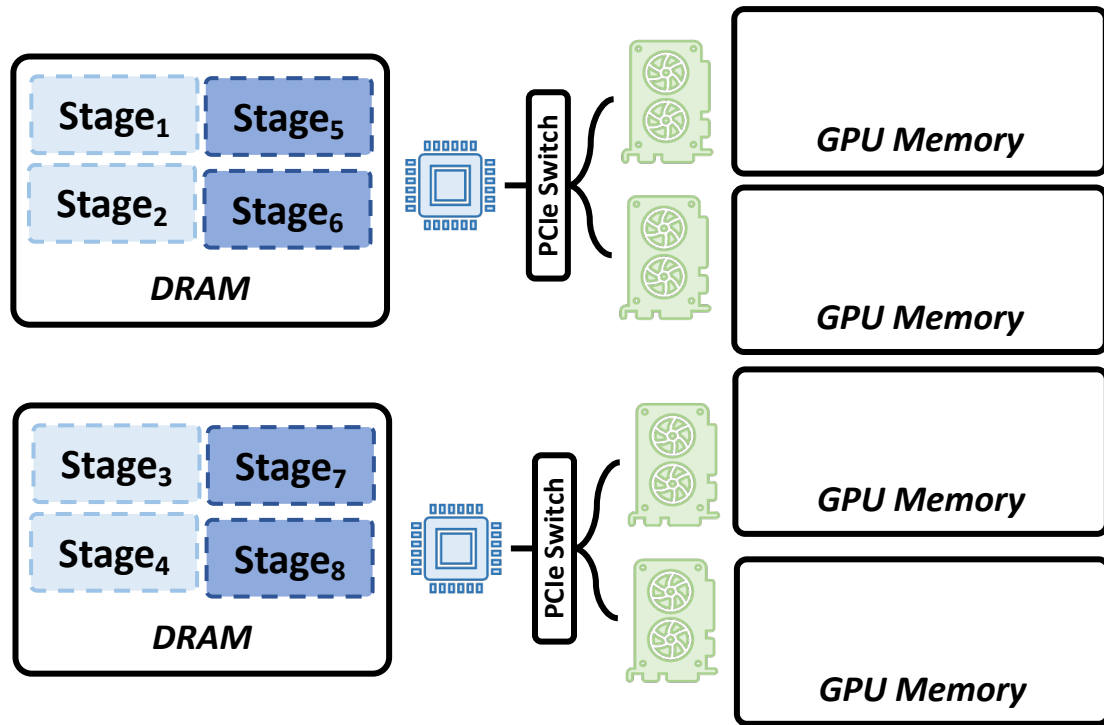
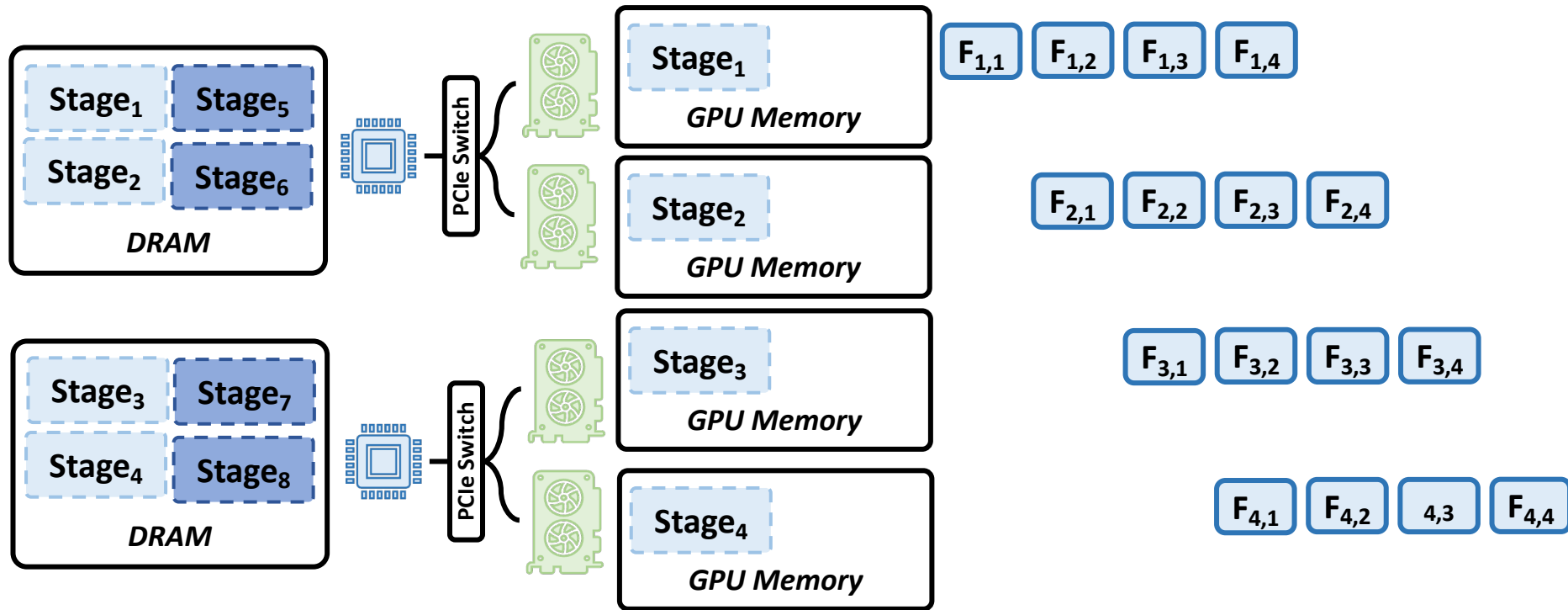# Mobius Overall

# Mobius Pipeline

# Mobius Pipeline

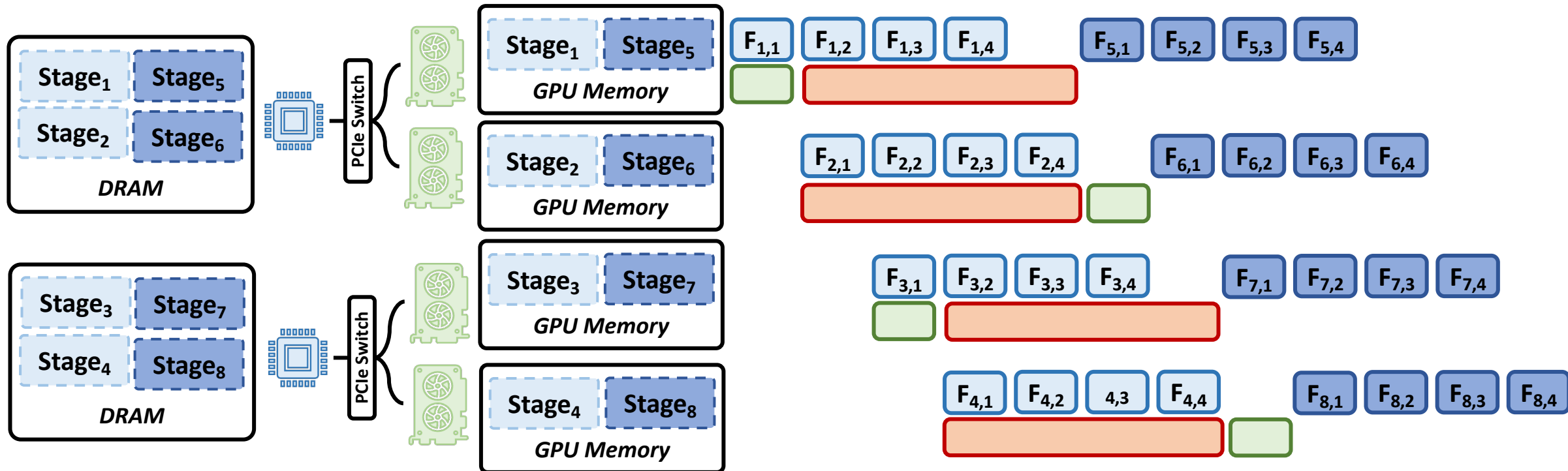# Mobius Pipeline

# Mobius Pipeline

# Mobius Pipeline
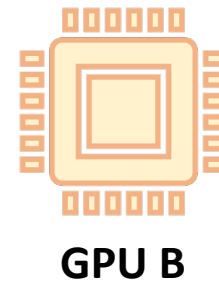


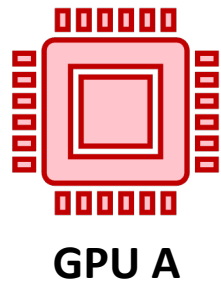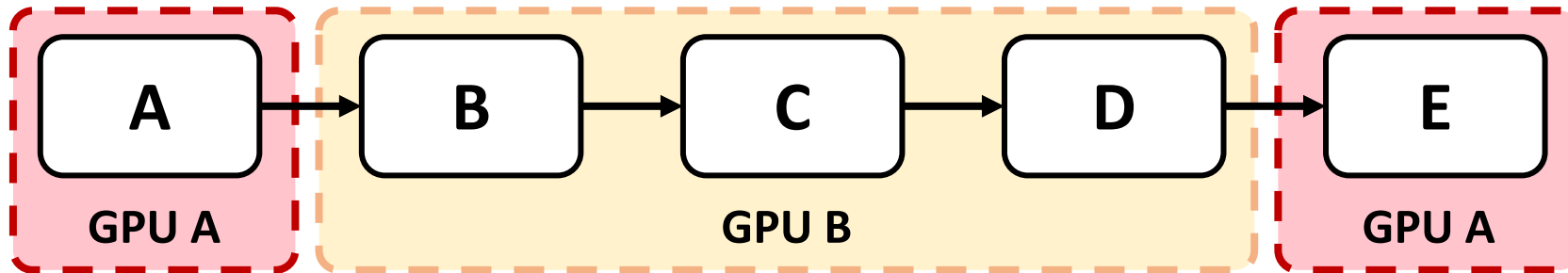$F_{i,j}$ — Stage$_i$'s execution on $j$th microbatch

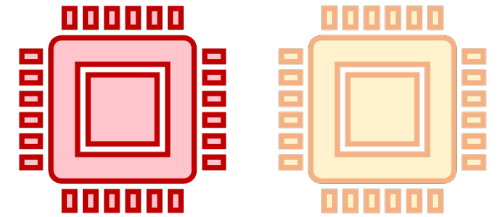Communication without contention

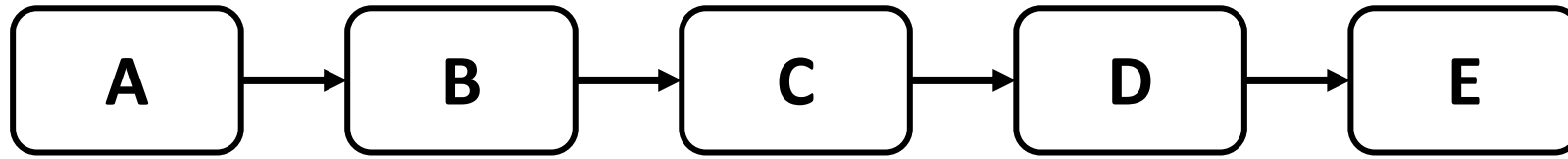Communication with contention

# Two Partition Questions

- How many stages are in each GPU?
- How many layers are in each stage?

# Mobius Partition

- Profile each layer's **memory footprint and computation overhead**
- Profile **hardware performance**, i.e. bandwidth
- Use mixed integer program (MIP) to fine the **optimal partition scheme**



# MIP

# Mobius Partition

**minimize**    *Training time of one step*

**subject to**    *Memory constraints*
- *Memory required by computation*
- *Memory required by prefetching*

*Pipeline order constraints*
- *Stage execution order*
- *Microbatch execution order*

# Communication Contention

# Cross Mapping



**Map adjacent stages to GPUs not under the same CPU root complex**

# Cross Mapping

**Number of GPUs under the same CPU root complex**

$$contention\left(stage_i, stage_j\right) = \frac{shared(i, j)}{|i - j|}$$

**Time difference to upload the two stages' data**

# Experimental Setup

TOPO 2+2

TOPO 3+1



CPU

PCIe

3090 Ti

**GPU Topologies**

# Overall Results

- Mobius and DeepSpeed with heterogeneous memory mode are able to **train larger models**
- Mobius **decreases per-step training time**
- Mobius brings more **significant performance improvement** when the GPU topology has **more severe communication contention**

# Communication Analysis

- DeepSpeed with heterogeneous memory mode requires frequent GPU all-to-all collective communications, while Mobius pipeline only transfers **small activations and activation gradients**
- **More than half of the data** is transferred at a bandwidth of more than 12 GB/s in Mobius

# Conclusion

- Commodity GPU server is an **affordable** option for fine-tuning large-scale models However, communication resources on commodity GPU servers are scarce
- We propose **Mobius** to reduce communication traffic and mitigate communication contention problem
  - **Mobius pipeline**: heterogeneous memory-based pipeline training scheme
  - **Mobius partition**: find the optimal partition scheme
  - **Cross mapping**: mitigate communication contention
- Mobius significantly reduces the training time by **3.8-5.1 times** compared with the prior art

# Thanks

**Mobius: Fine Tuning Large-Scale Models on Commodity GPU Servers**

**Yangyang Feng**, Minhui Xie, Zijie Tian, Shuo Wang, Youyou Lu, and Jiwu Shu

*Tsinghua University*

*http://storage.cs.tsinghua.edu.cn*
Email: fyy21@mails.tsinghua.edu.cn