



# Seesaw: High-throughput LLM Inference via Model Re-sharding

Qidong Su<sup>1,2,3</sup>, Wei Zhao<sup>3,4</sup>, Xin Li<sup>3</sup>, Muralidhar Andoorveedu<sup>3</sup>,  
Chenhao Jiang<sup>1,2</sup>, Zhanda Zhu<sup>1,2,3</sup>, Kevin Song<sup>1,2</sup>,  
Christina Giannoula<sup>1,2,3</sup>, Gennady Pekhimenko<sup>1,2,3</sup>

1



2



3



4



# Overview

## Problem:

Throughput-oriented offline LLM text generation with multiple GPUs

## Observation:

Prefill and decode favor different *parallelization* strategies.

## Key Ideas:

Dynamically switch between different parallelization strategies

## Evaluation Results:

**1.36x** throughput increase (up to **1.78x**) compared with vLLM.

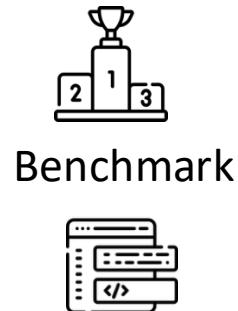
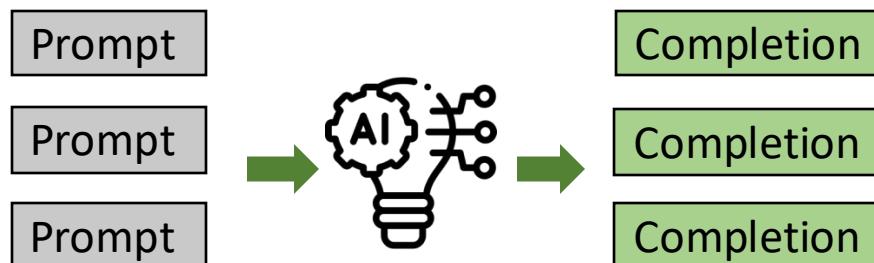
# LLM Offline Inference

Online:



**Metrics:**  
Latency,  
TTFT/TBT

Offline:

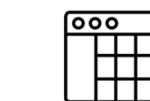


**Metrics:**  
Throughput  
E2E time



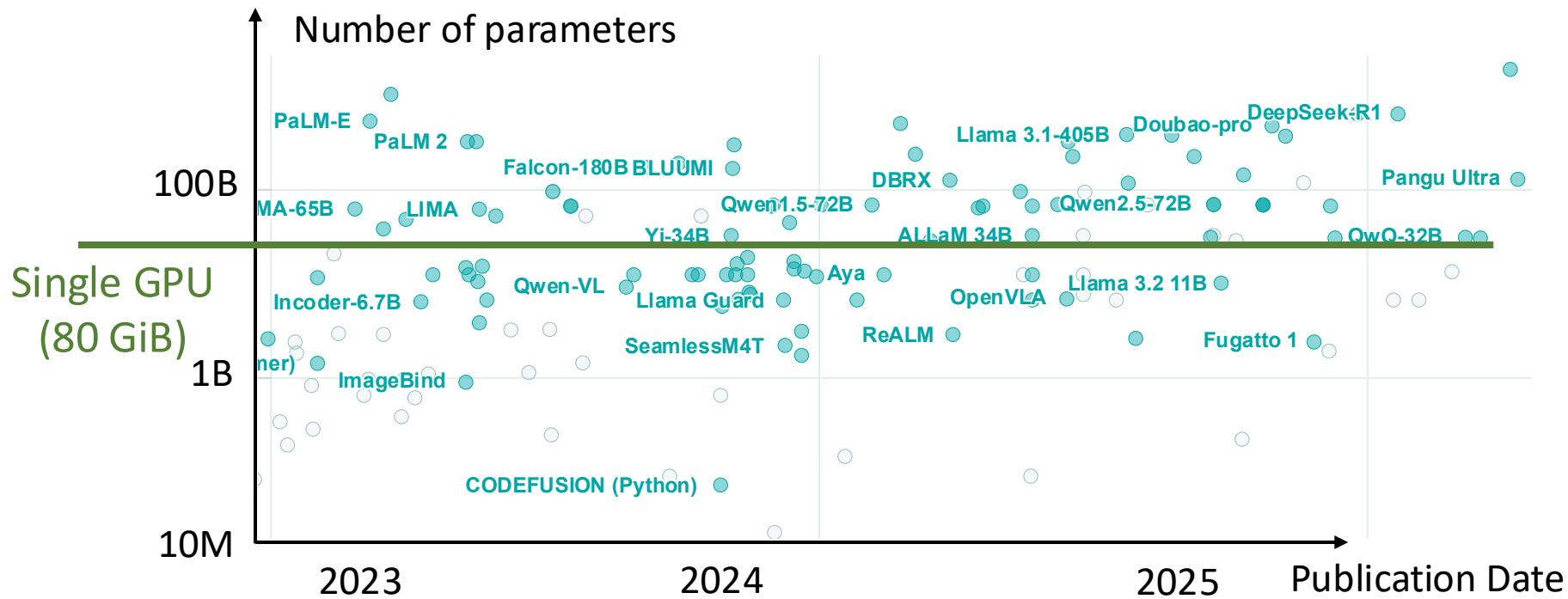
Batch API 50% cheaper

...



Dataset Synthesis

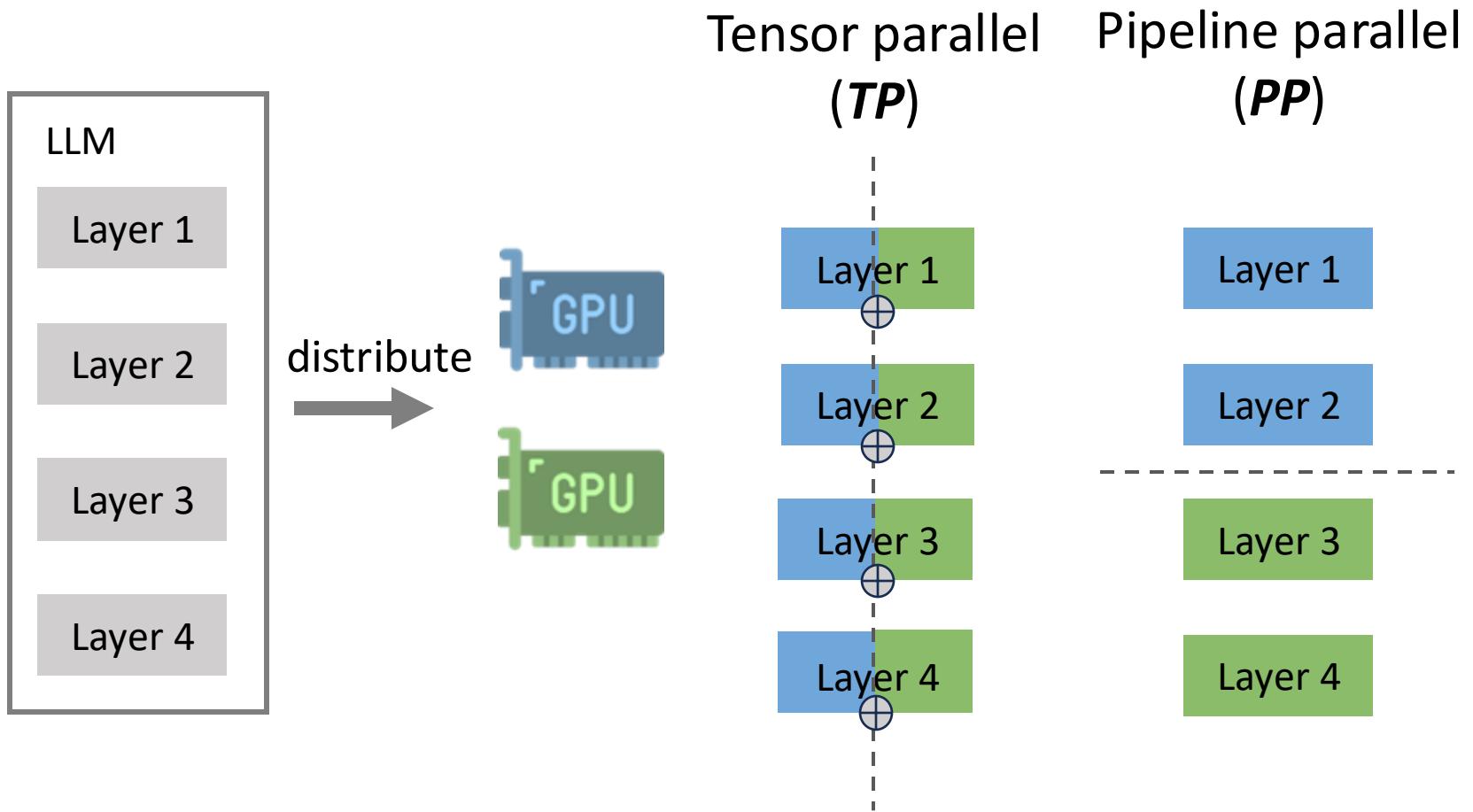
# Modern LLMs are Large



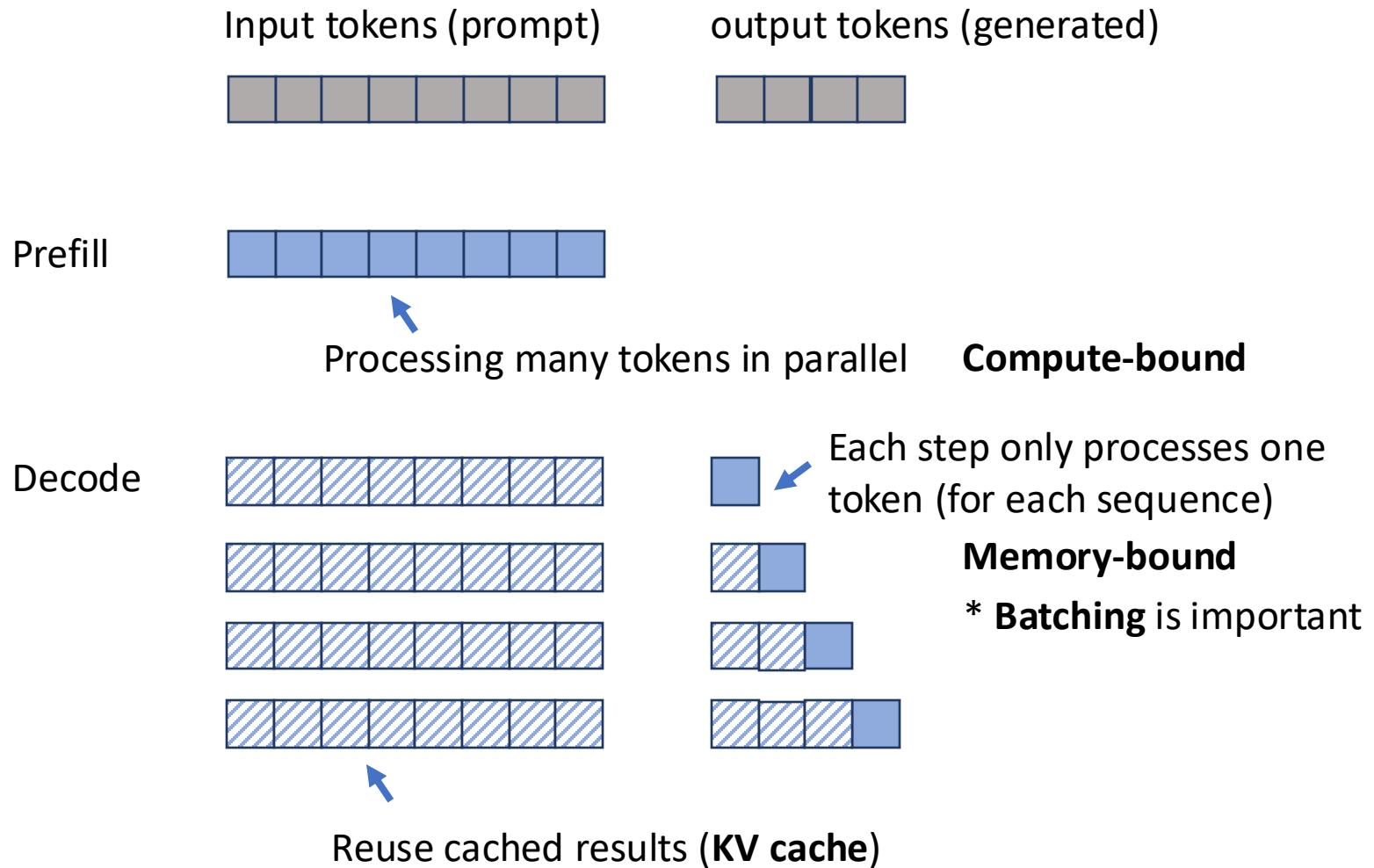
- We need multiple GPUs to deploy them.
- How to partition the model?

\*source: <https://epoch.ai/data/notable-ai-models>

# How to Partition a Model?

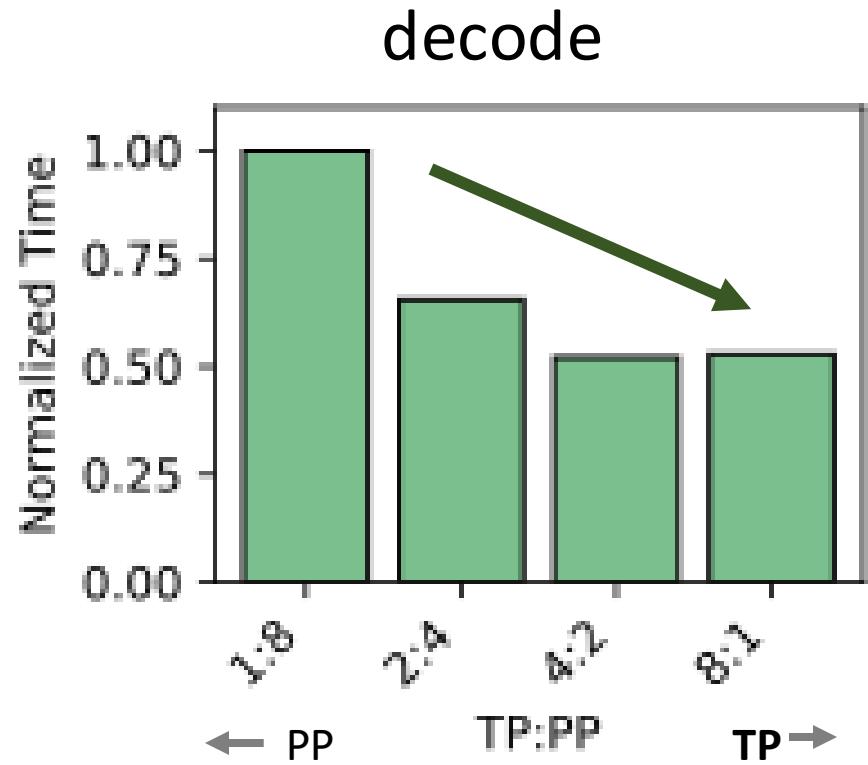
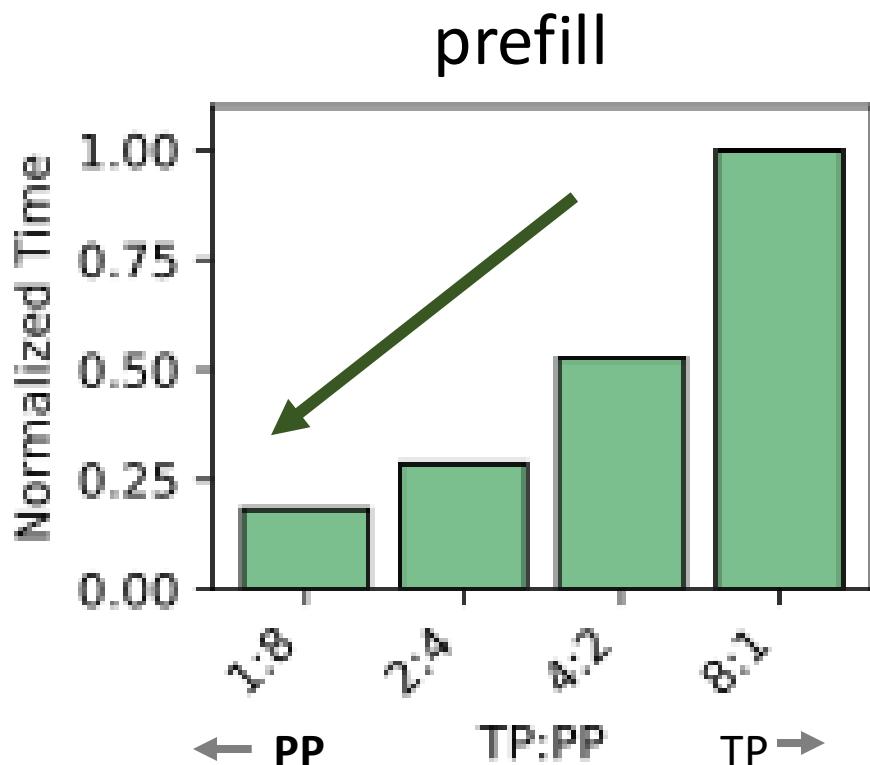


# Prefill and Decode



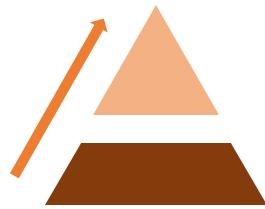
# PP for prefilling, TP for decoding

(lower is better)

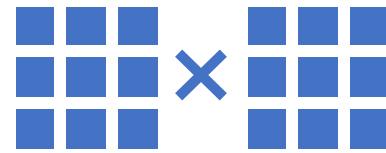


(Measured on 8xL4)

# How are TP and PP different?



Memory Access



Computation



Network  
Communication

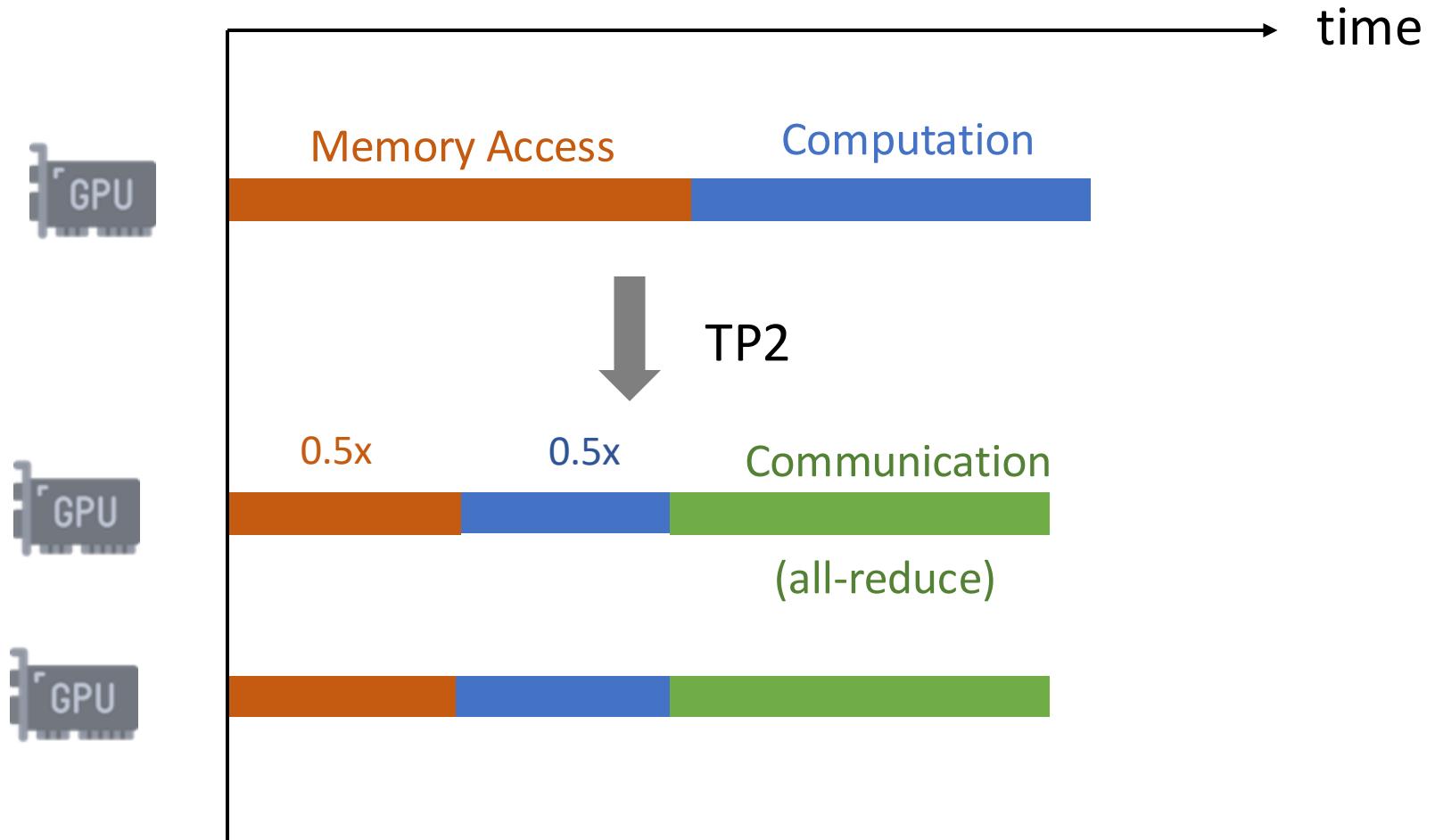
Tensor  
Parallel (TP)



Pipeline  
Parallel (PP)

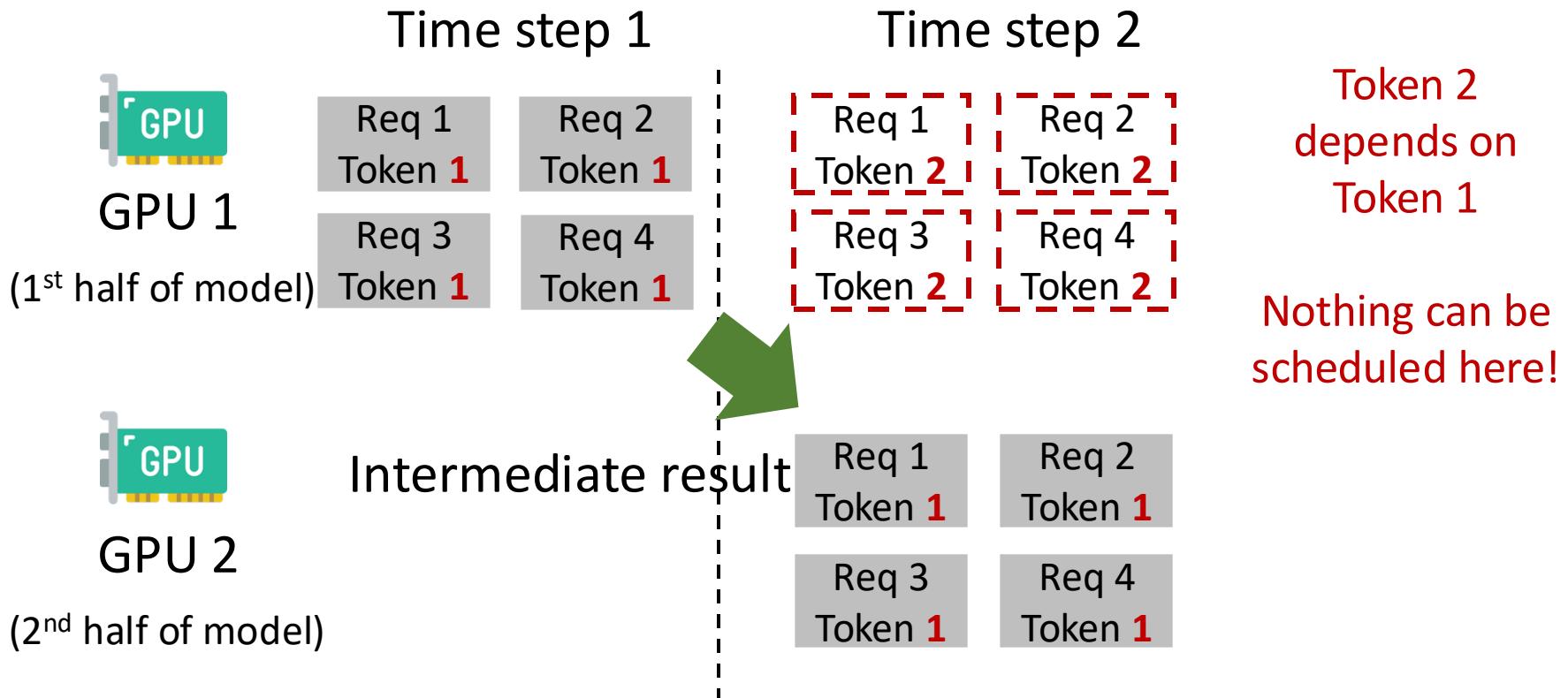


# TP: communication overhead



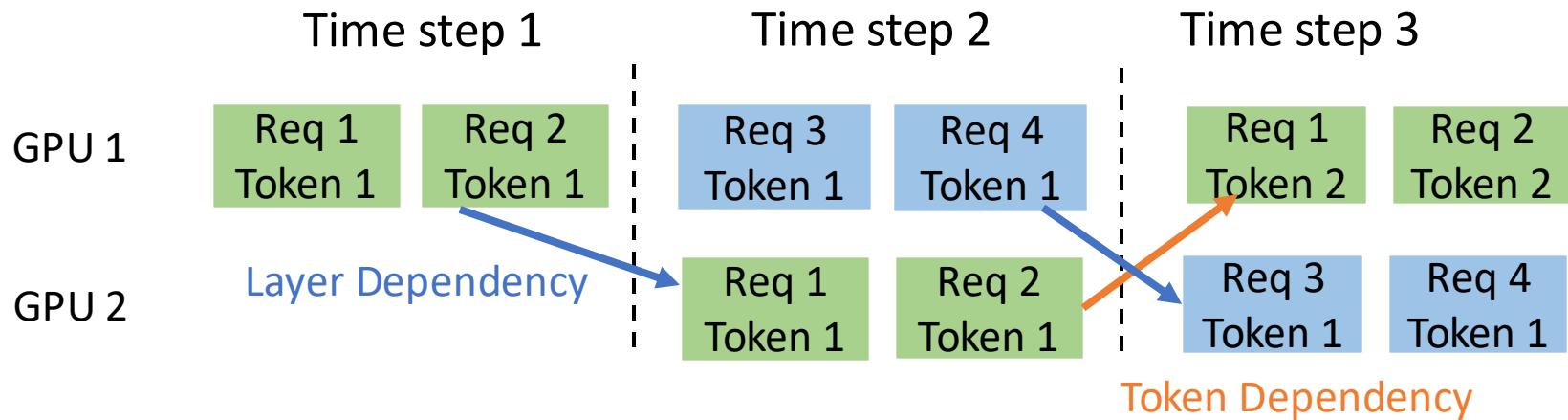
# Pipeline Parallel: Micro-batching

Example: 2 GPUs, 4 requests



# Pipeline Parallel: Micro-batching

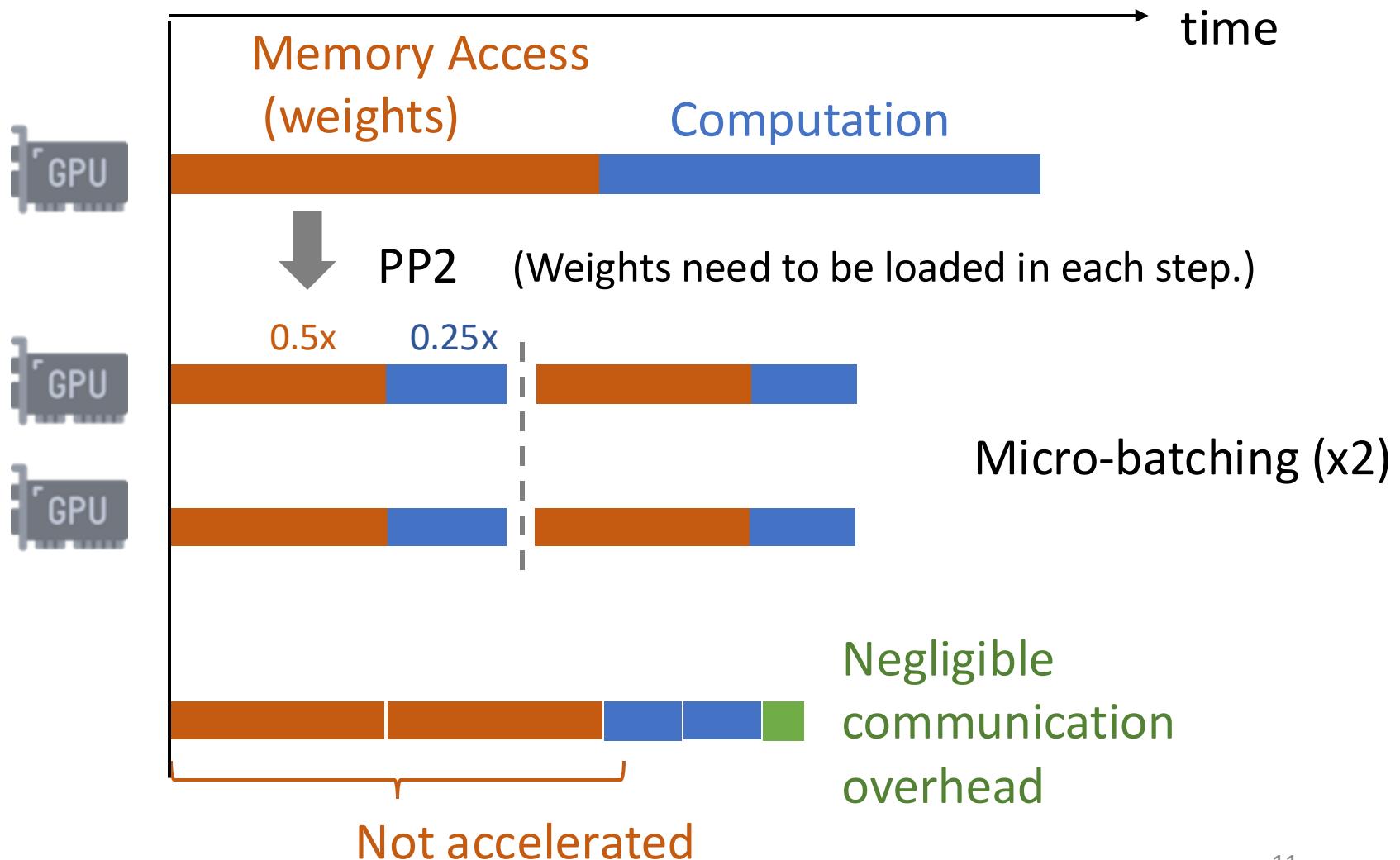
Example: 2 GPUs, 4 requests



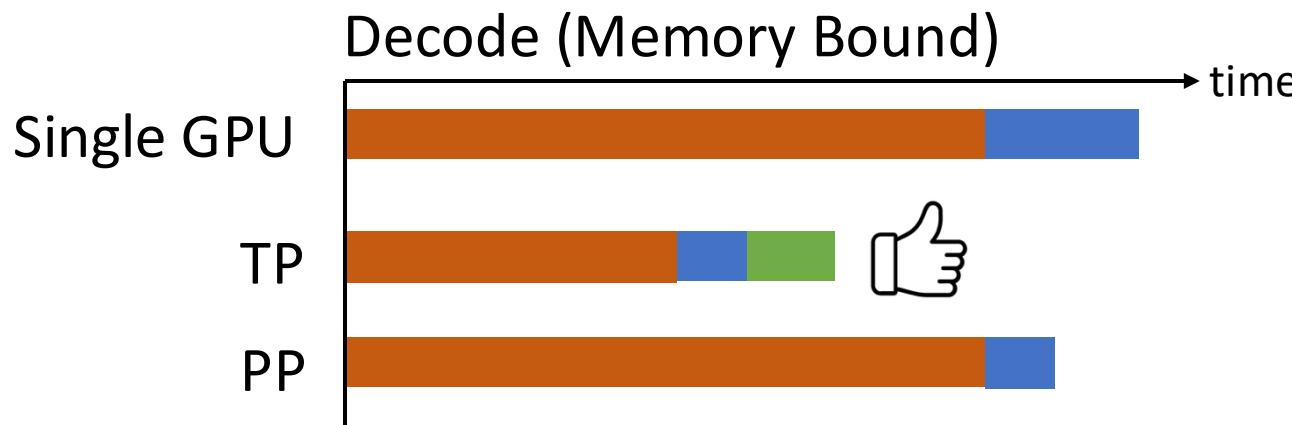
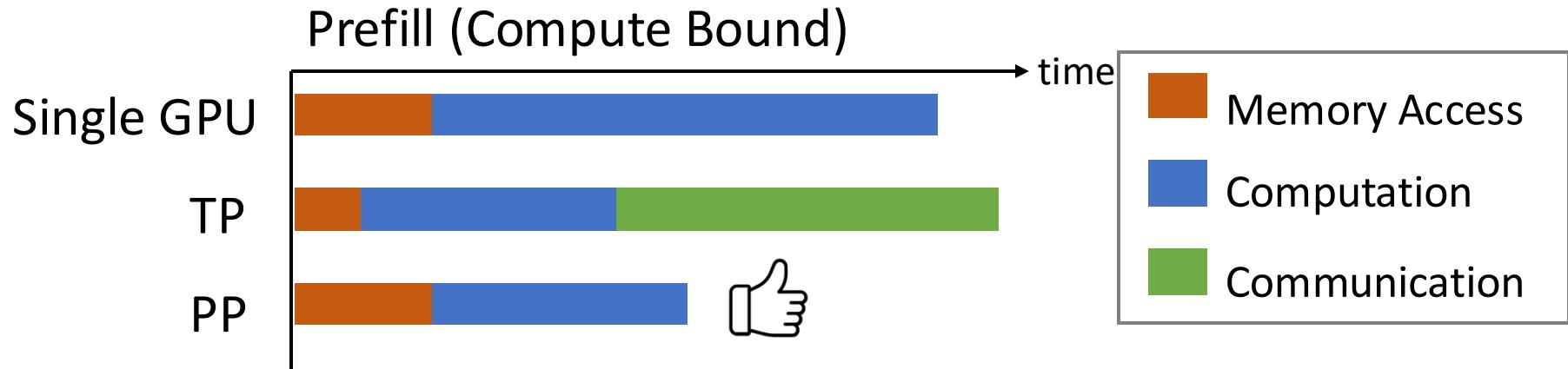
A step can only process  $\frac{1}{PP}$  of all requests!

We need  $PPx$  more steps!

# PP doesn't accelerate loading weights



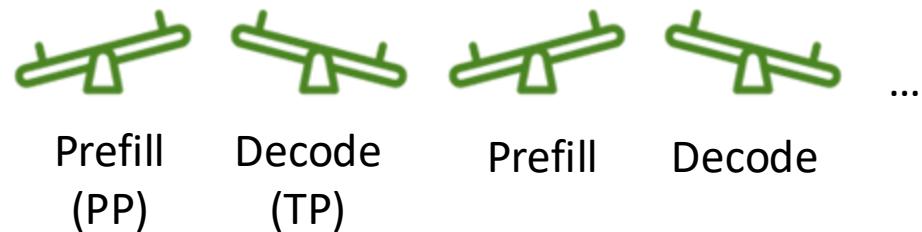
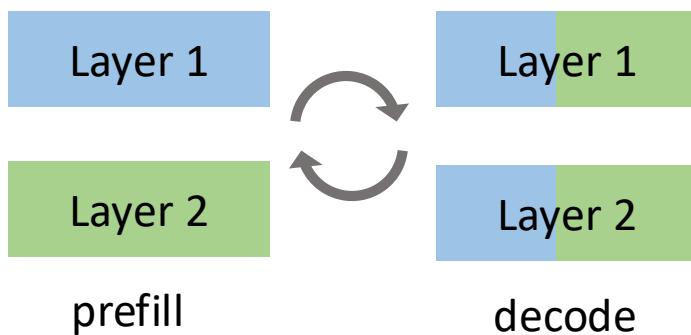
TP=efficient memory access  
PP=reduce communication





# Seesaw: Key Ideas

- Pipeline Parallelism for Prefilling
- Tensor Parallelism for Decoding
- Re-Shard the model during the stage transition.



# Challenge: Reshard Overhead

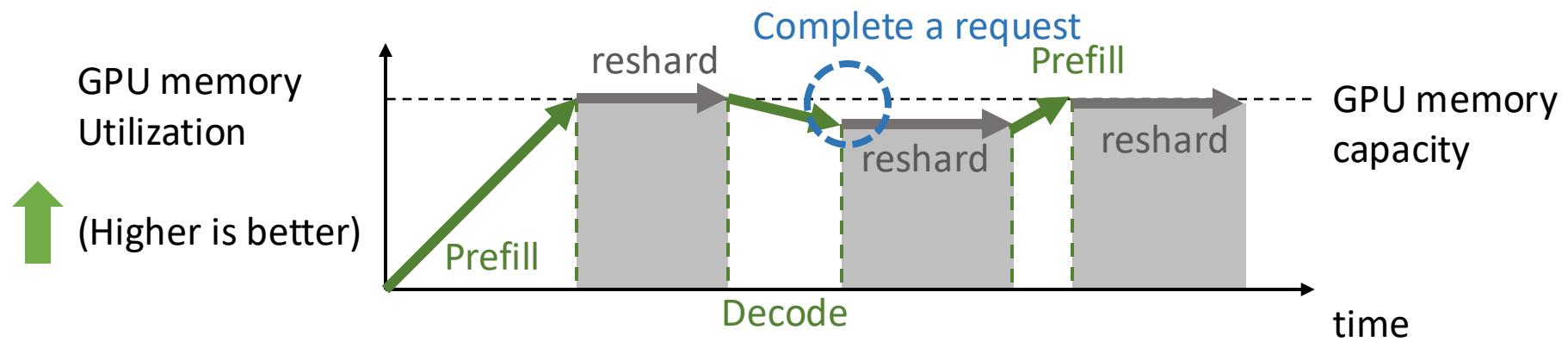
- Re-sharding moves tens of gigabytes of tensors
- We need to reduce the reshading frequency

However, **Continuous Batching** causes high reshard frequency.

# Challenge: Reshard Overhead

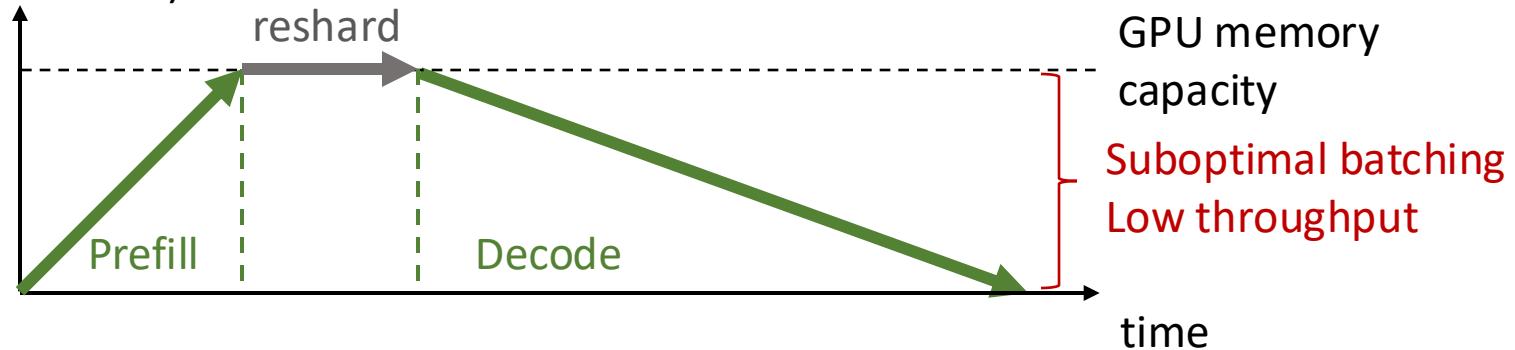
## Recall:

- In decoding, large batch sizes = higher throughput
- Higher memory utilization gives larger batch sizes

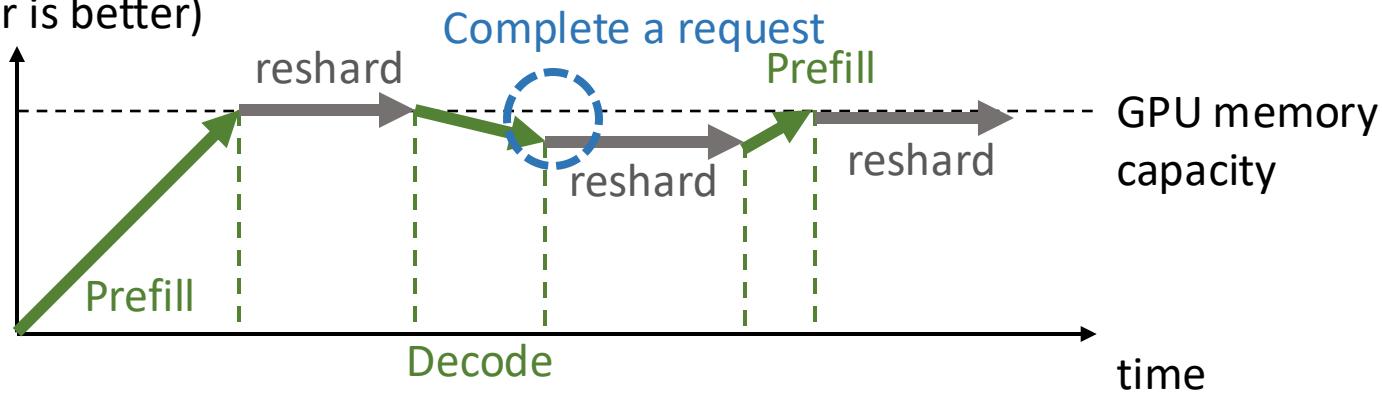


# Alternative: Decode-Prioritizing

GPU memory utilization  
(Higher is better)



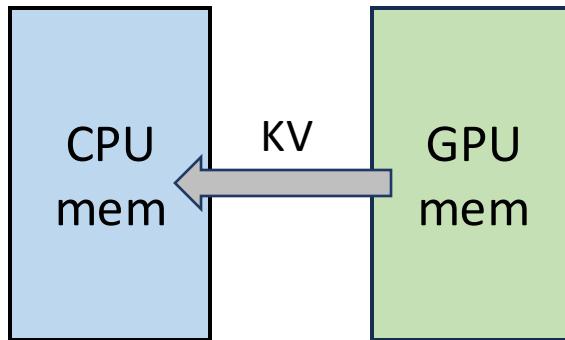
GPU memory utilization  
(Higher is better)





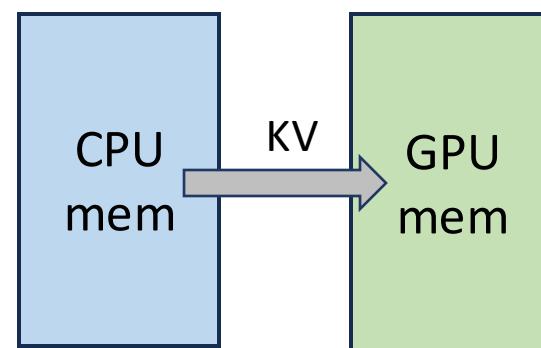
# Seesaw: Tiered KV Cache Buffering

## Prefill



Offload KV cache to CPU  
(Overlap with computation)

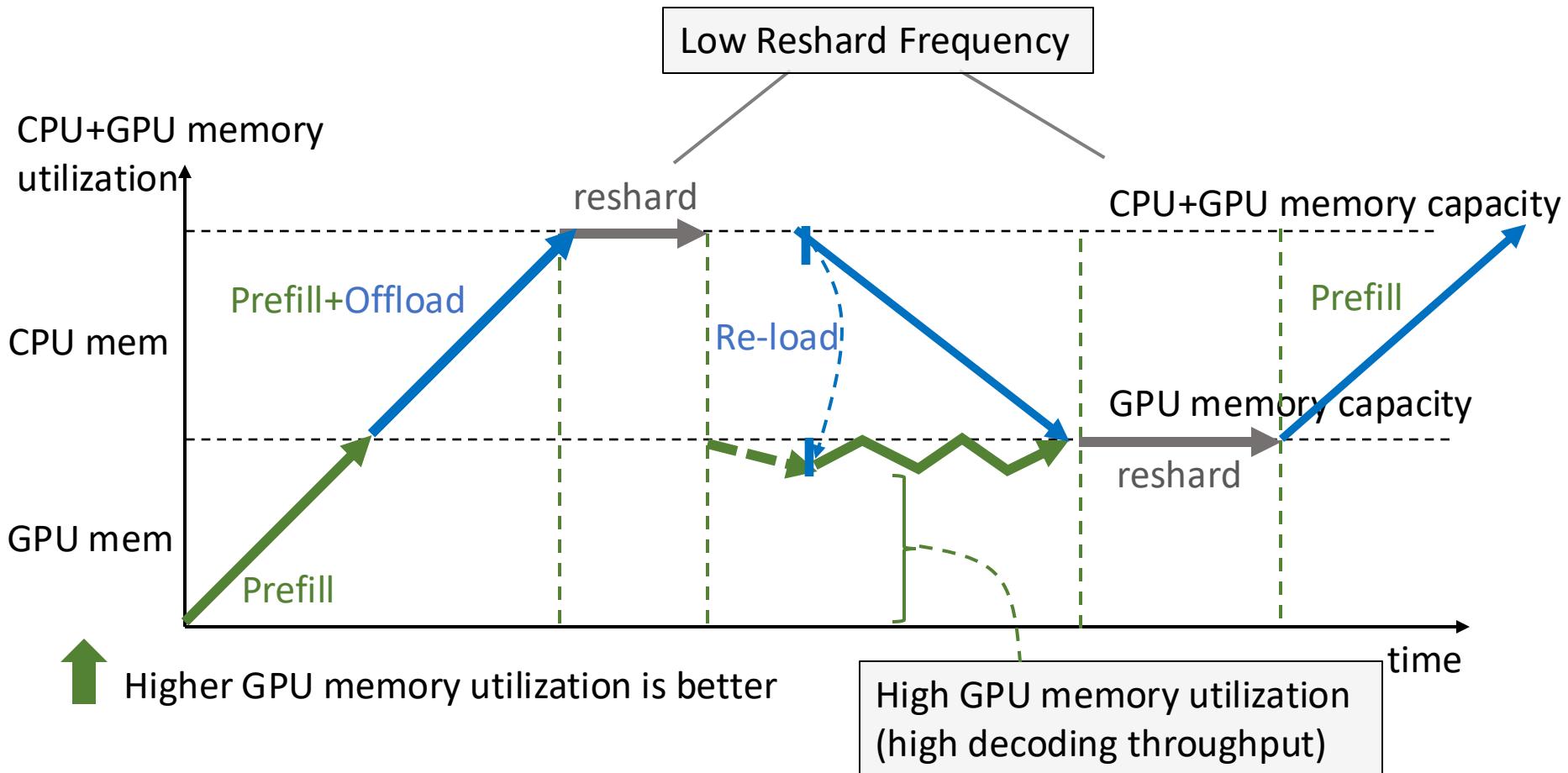
## Decode



Reload KV cache from CPU  
memory (asynchronously).



# Seesaw: Tiered KV Cache Buffering



# Evaluation Methodology

## Hardware:

- A10 (24G, PCIe)
- L4 (24G, PCIe)
- A100 (40G, PCIe/NVLink)



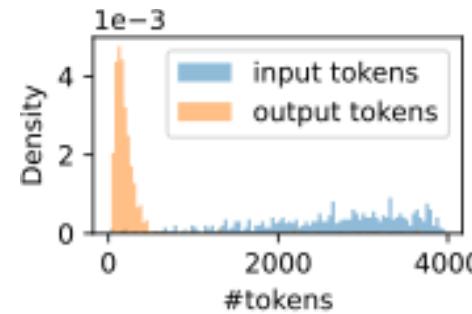
## Model:

- LLaMA3 15B
- CodeLLaMA 34B
- LLaMA2 70B

## Baseline:

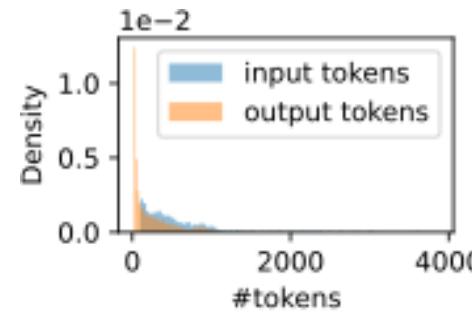
- vLLM
- + chunked prefill
- + tuned chunk size

## Workload:



arXiv-  
Summarization

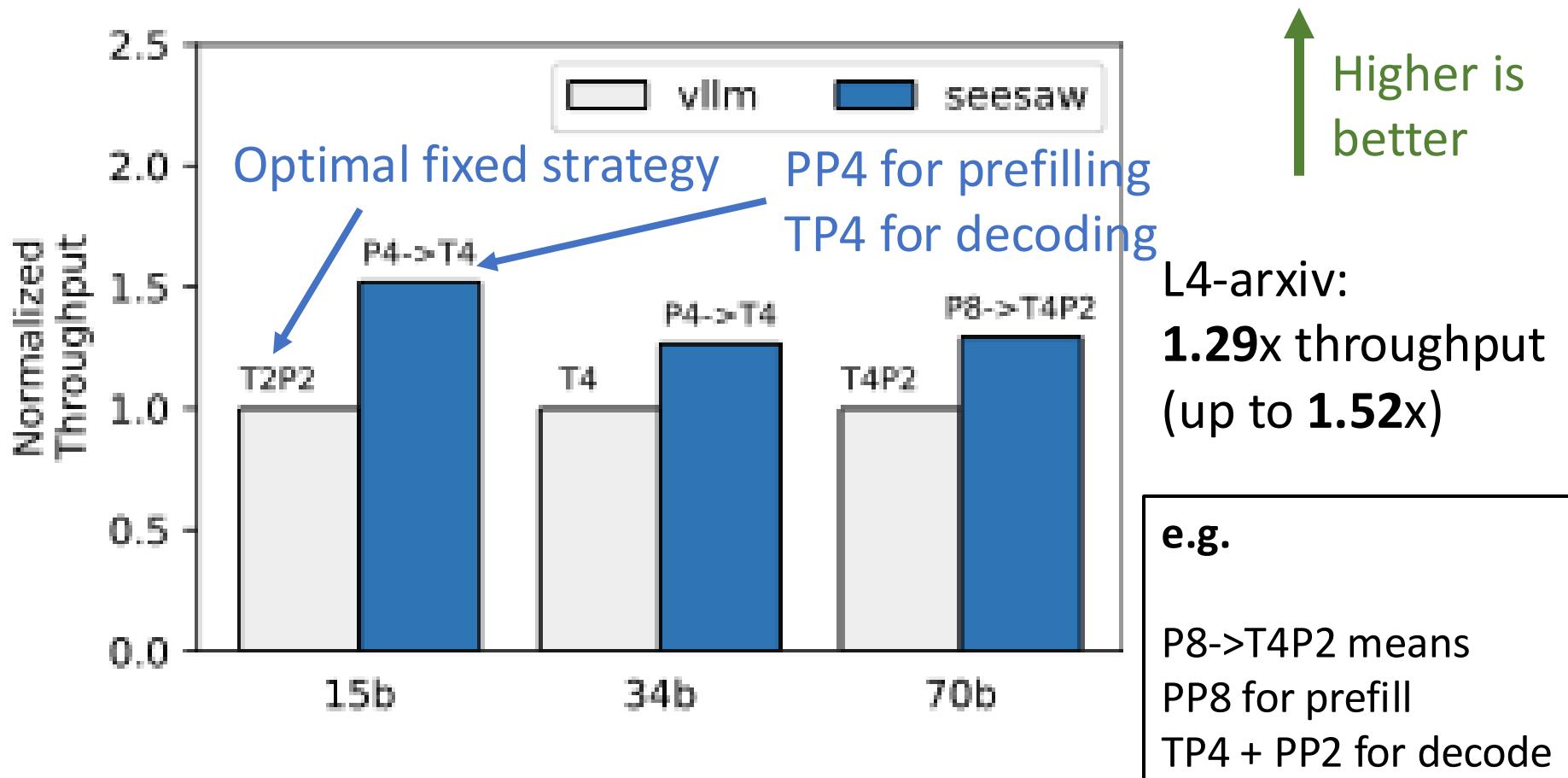
input > output



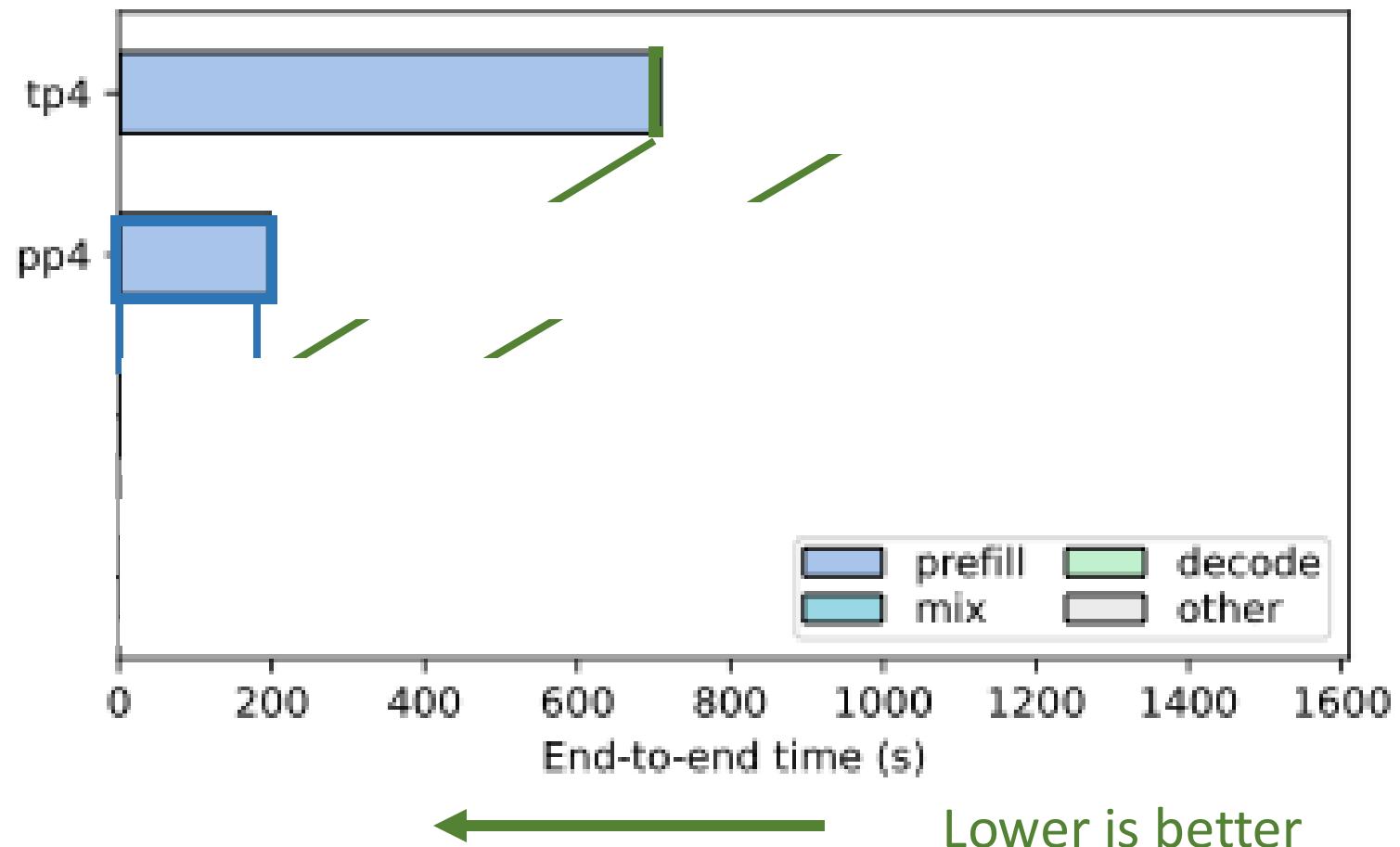
ShareGPT

Input ≈ output

# End-to-end Throughput Compare



# Speedup breakdown



# Summary

## Problem:

Throughput-oriented offline LLM text generation

## Observation:

Prefill and decode favor different parallelization strategies.

## Key Ideas:

Dynamically switch between different parallelization strategies.

Tiered KV cache buffering to reduce transition overhead

## Evaluation Results:

**1.36x** throughput increase (up to **1.78x**) compared with vLLM.