

# LLM Inference Serving Systems

Arvind Krishnamurthy

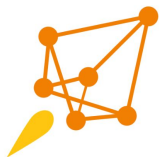
Material adapted from slides by Hao Zhang (UCSD) and Amey Agrawal (GTech)

# *Lecture Outline*

- Requirements of LLM Serving Systems
- Interleaved execution in Sarathi Serve
- Disaggregation in DistServe

# LLM Systems Today Optimize **Throughput**

vLLM



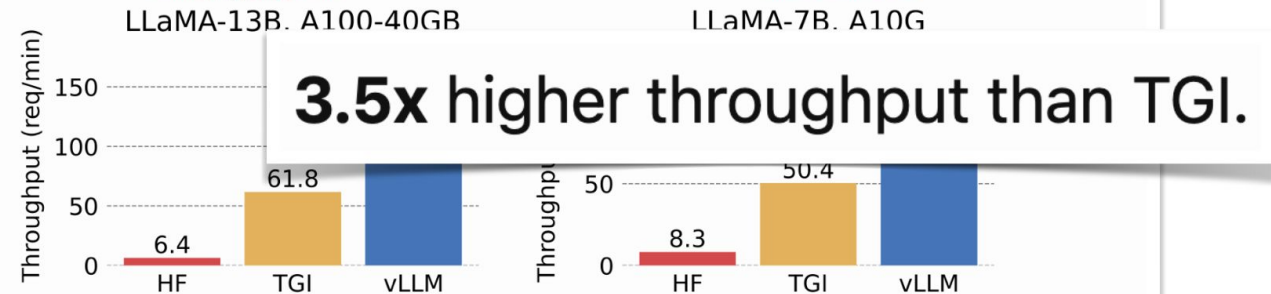
DeepSpeed MII



Beyond State-of-the-art Performance

**24x** higher throughput compared to HF

We sample the requests' input/output lengths from the ShareGPT dataset. In our experiments, vLLM achieves up to **24x** higher throughput compared to HF and up to **3.5x** higher throughput than TGI.



**3.5x** higher throughput than TGI.

Serving throughput when each request asks for *one output completion*. vLLM achieves 14x - 24x higher throughput than HF and 2.2x - 2.5x higher throughput than TGI.

# Motivation: Applications have Diverse SLO

## • TTFT

Time to first token

Initial response time



Chatbot



Fast initial response



Summarization



User can tolerate longer initial response

## • TPOT

Time per output token

Average time between two subsequent generated tokens

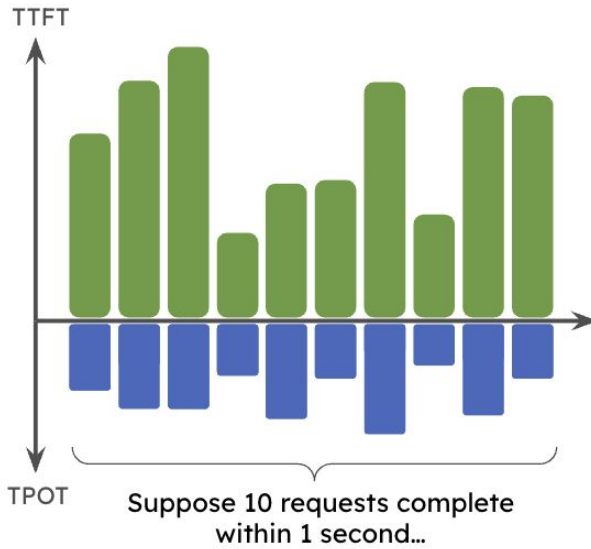


Human reading speed (P99 latency = 250ms)



Data output generation (P99 latency = 35ms)

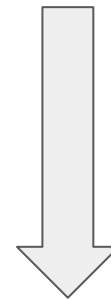
# High Throughput $\neq$ High Goodput



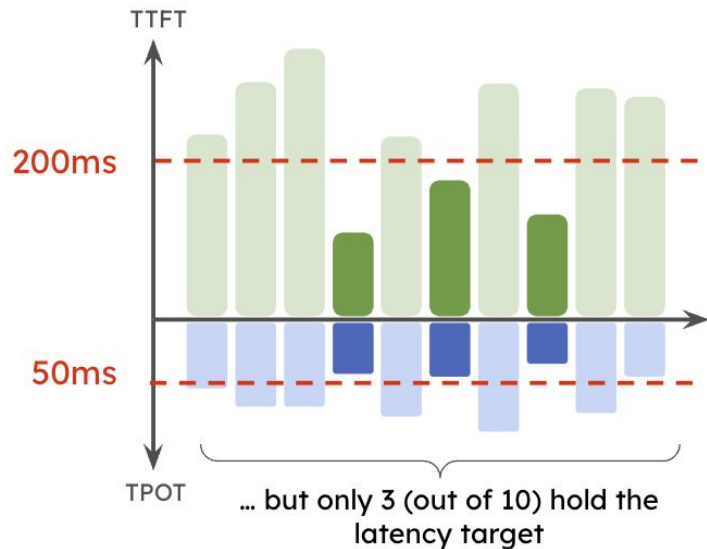
**Throughput = 10 rps**  
= completed request / time

**High  
Throughput  
System**

...



under SLO  
criteria



**Goodput = 3 rps** 🤯  
= completed request **within SLO** / time

can have  
**Low Goodput!**

# High Throughput $\neq$ High Goodput

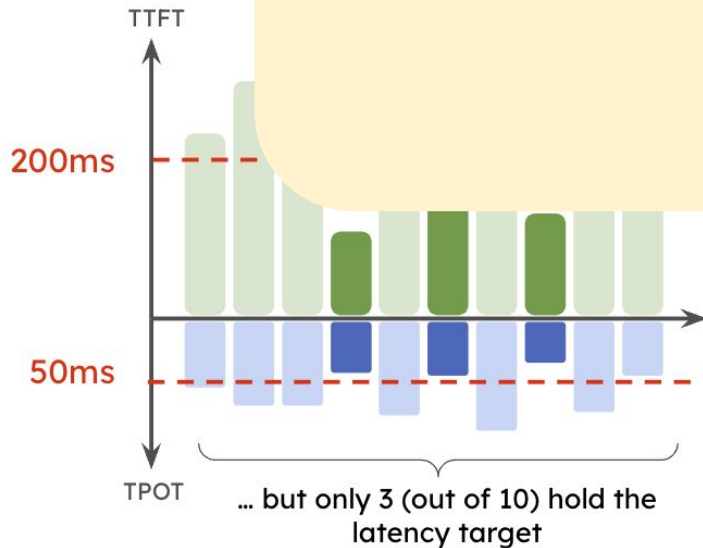
TTFT



High Throughput can still have **Low Goodput**

$\Rightarrow$  Poor UX 

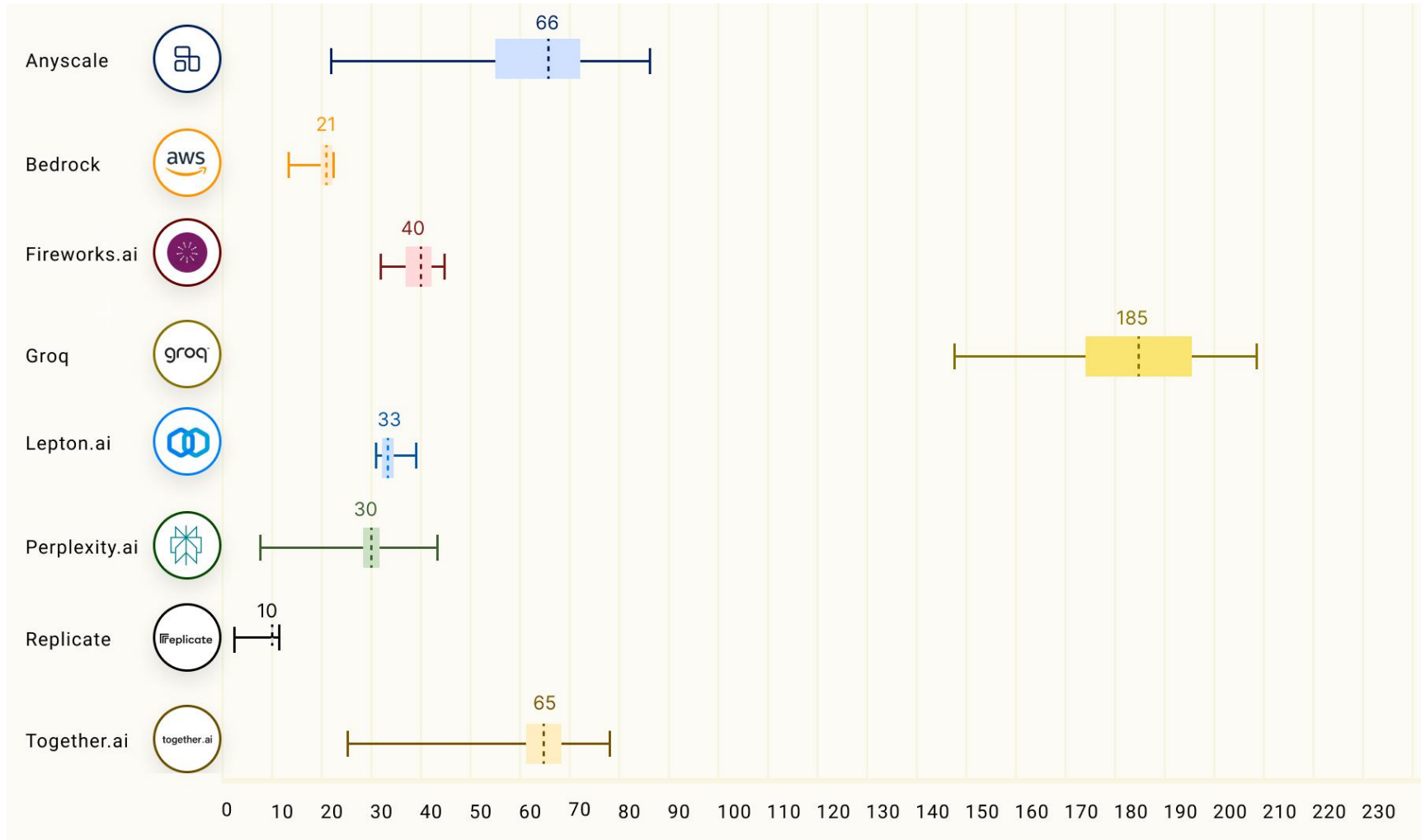
High  
out



**Goodput = 3 rps**  
= completed request **within SLO** / time

Low Goodput!

# LLM Performance leaderboard (tokens/sec)



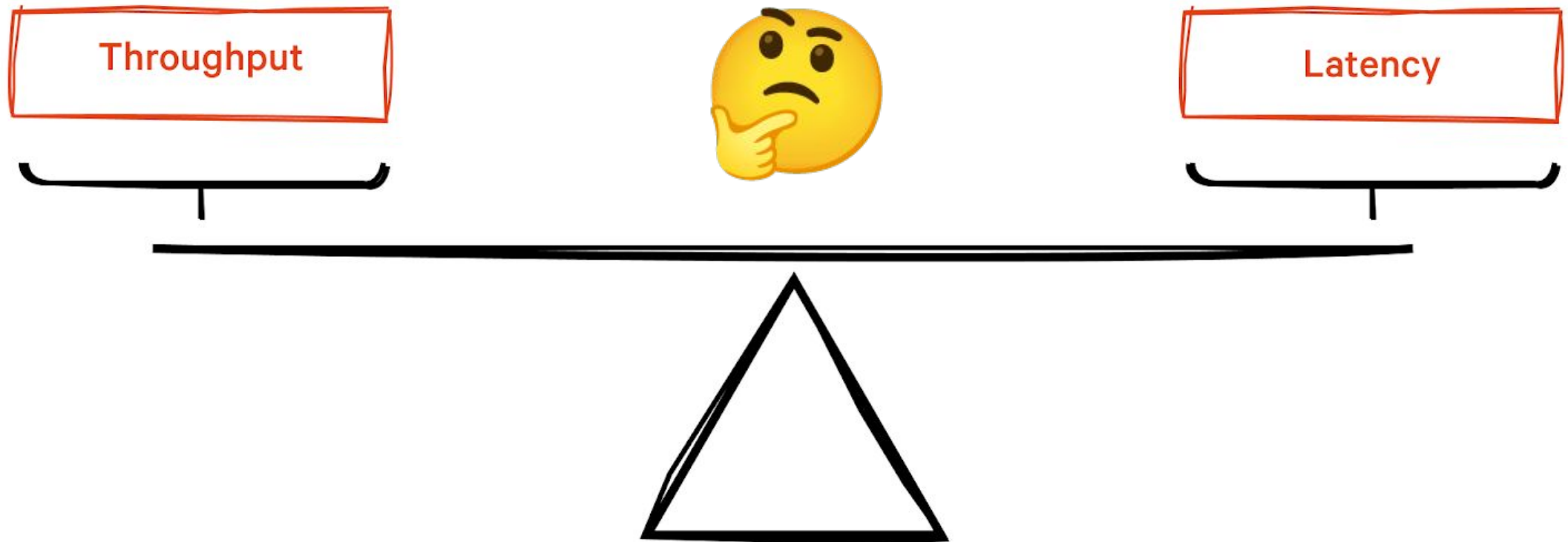
# Taming Throughput-Latency Tradeoff in LLM Inference with **Sarathi-Serve**

**Amey Agrawal<sup>1</sup>, Nitin Kedia<sup>2</sup>**, Ashish Panwar<sup>2</sup>, Jayashree Mohan<sup>2</sup>,  
Nipun Kwatra<sup>2</sup>, Bhargav Gulavani<sup>2</sup>, Alexey Tumanov<sup>1</sup>, Ramachandran  
Ramjee<sup>2</sup>

*<sup>1</sup>Georgia Institute of Technology, <sup>2</sup>Microsoft Research India*



Can we maintain low latency  
with high throughput?

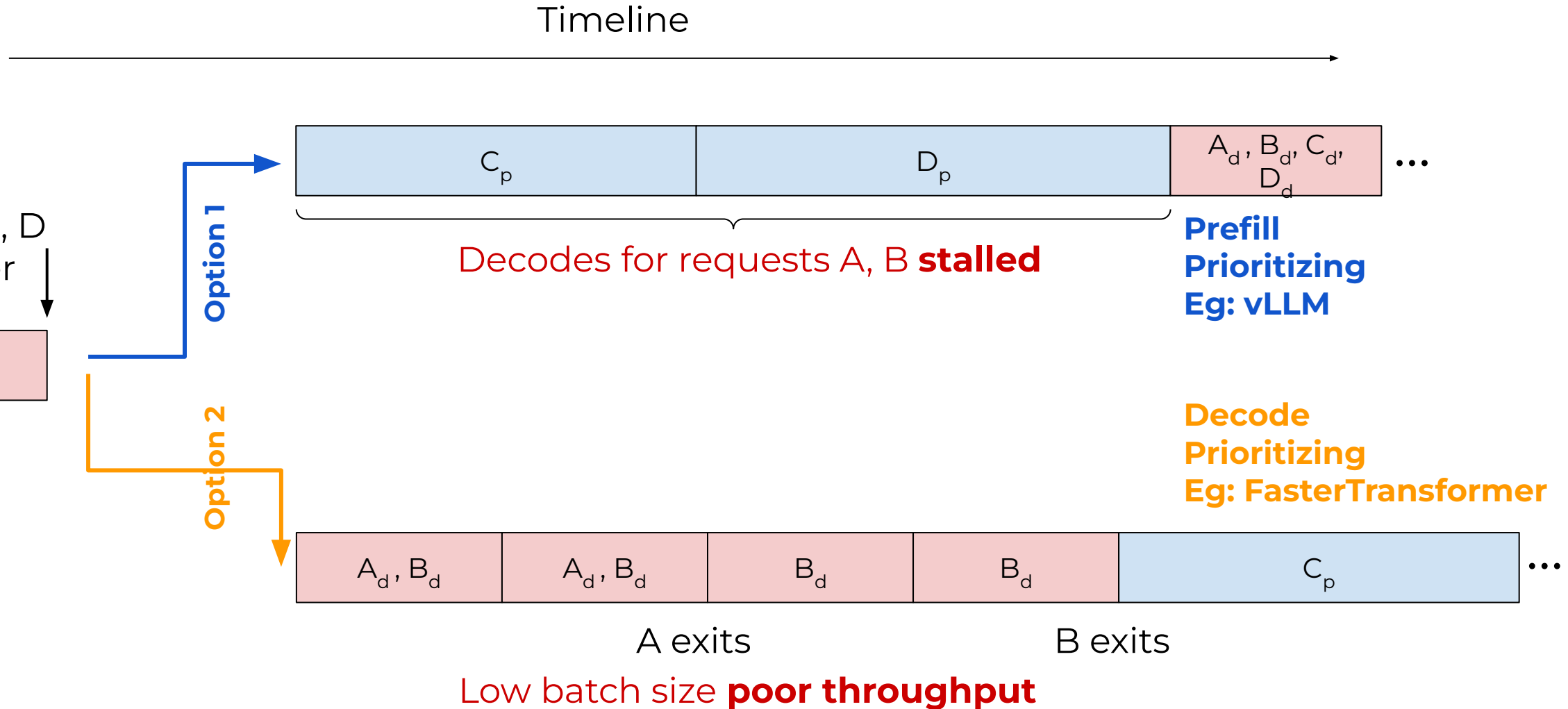


# Background: Continuous Batching

- Integrate new requests as old requests rotate off
- Need to perform prefill for new request
  - Results in a stall for existing requests
- Two queues in the system: prefill requests waiting to be integrated; ongoing decode requests

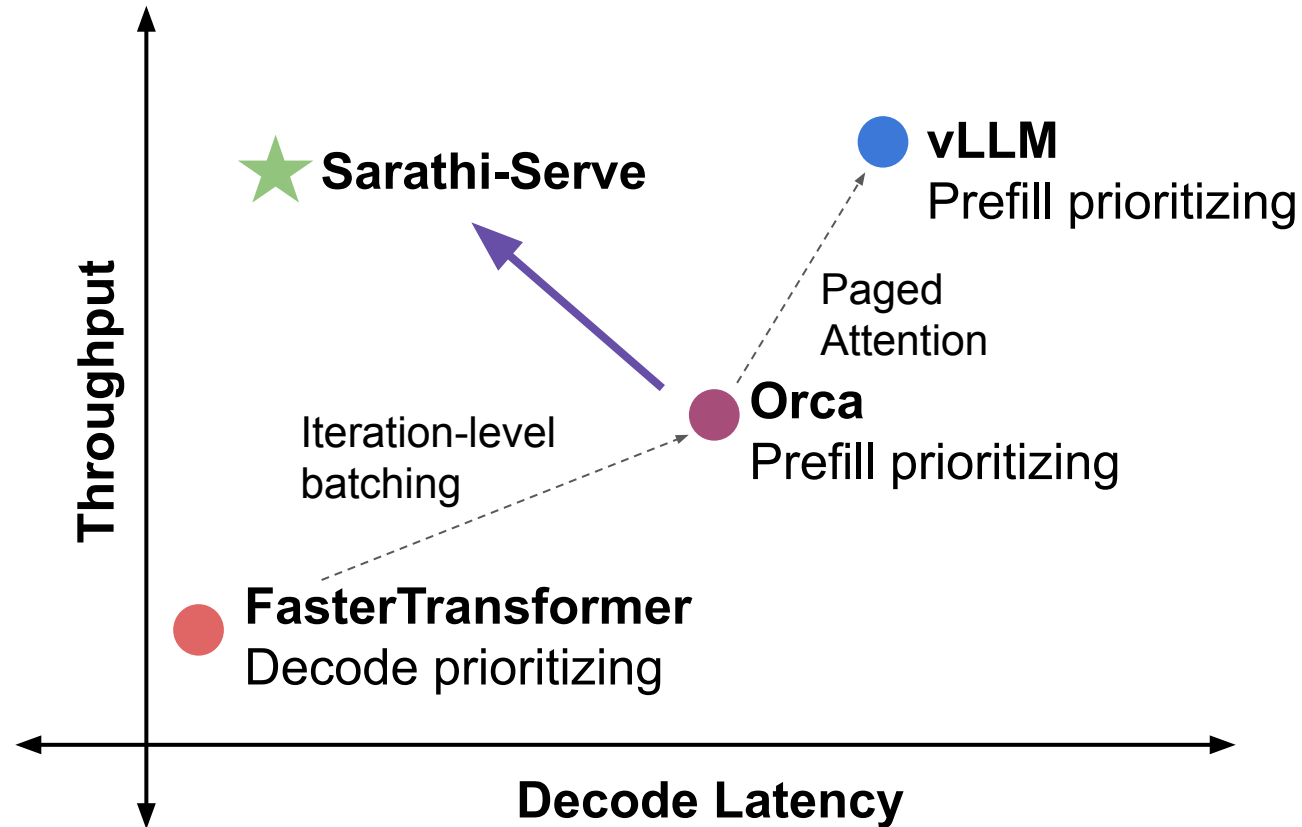
$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$	$S_1$	END	$S_6$	$S_6$
$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	END
$S_3$	$S_3$	$S_3$	$S_3$	END	$S_5$	$S_5$	$S_5$
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	END	$S_7$

# The Prefill-Decode Scheduling Conundrum





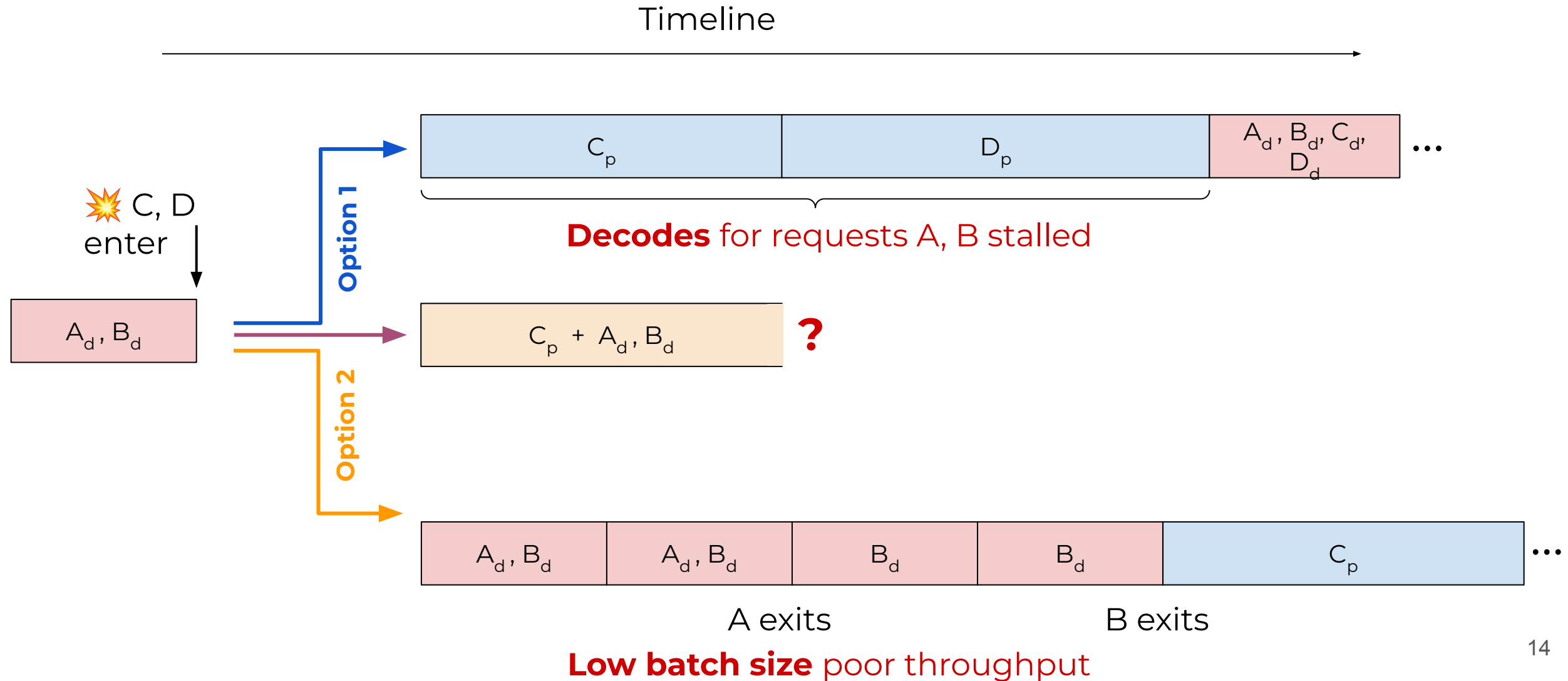
# The Latency-Throughput Tradeoff



Existing batching policies make a harsh latency-throughput tradeoff !

How can we achieve both high throughput and low-latency? 🤔

# The Prefill-Decode Scheduling Conundrum



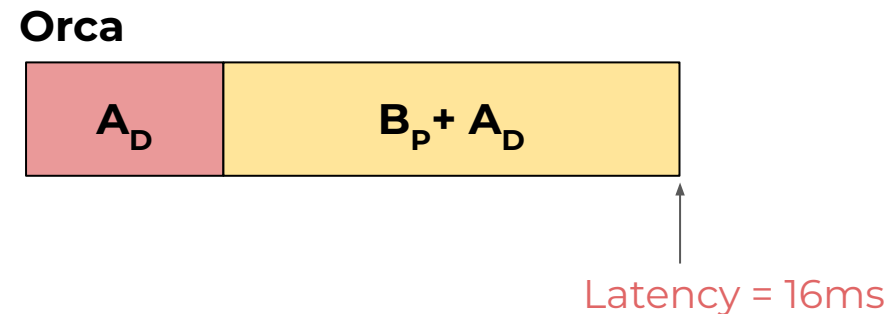
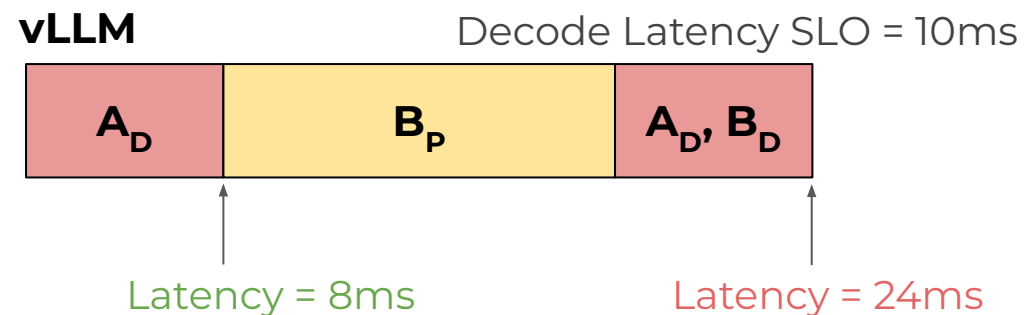
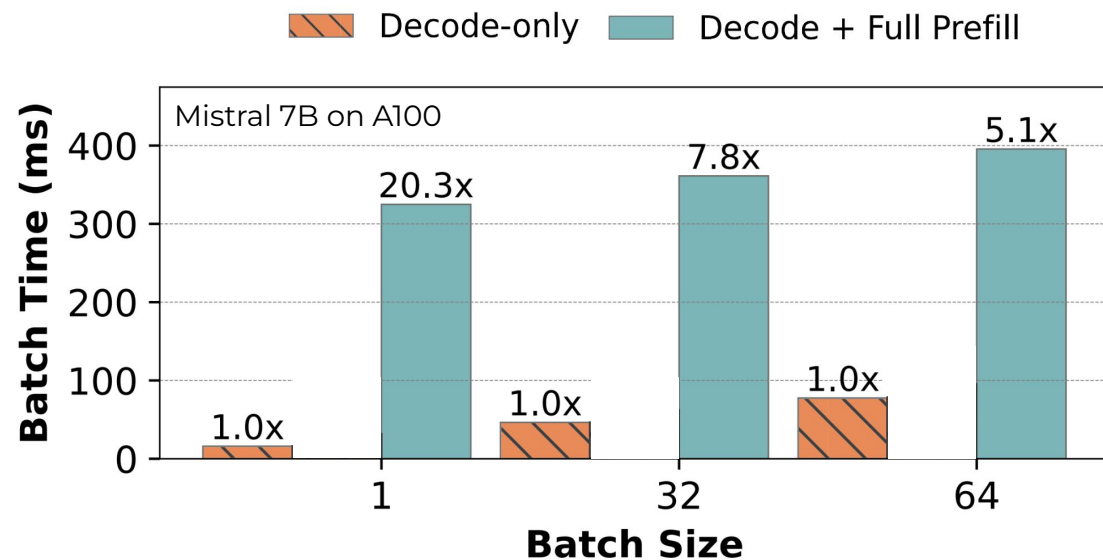
# 🔍 Mixed Batching

## Idea

🧑‍🔬 Fused computation of prefill and decodes

## Challenge

😭 Naively combining prefill and decode operations leads to increase in latency



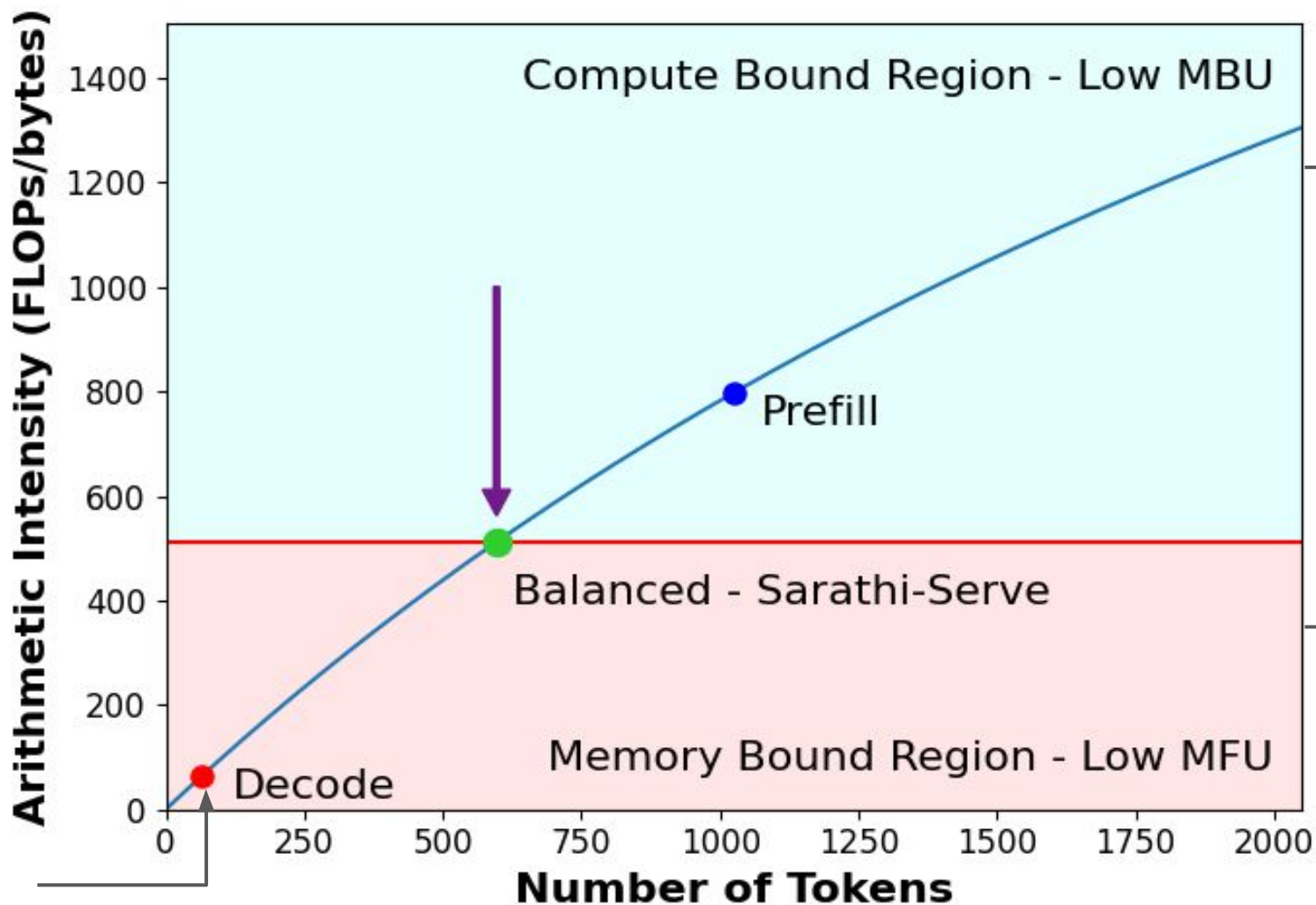
# Key Insight

Prefill computation can be done at a marginal cost with careful batching 





# Observation: Arithmetic Intensity Slack



Constrained due to memory overhead in decode phase

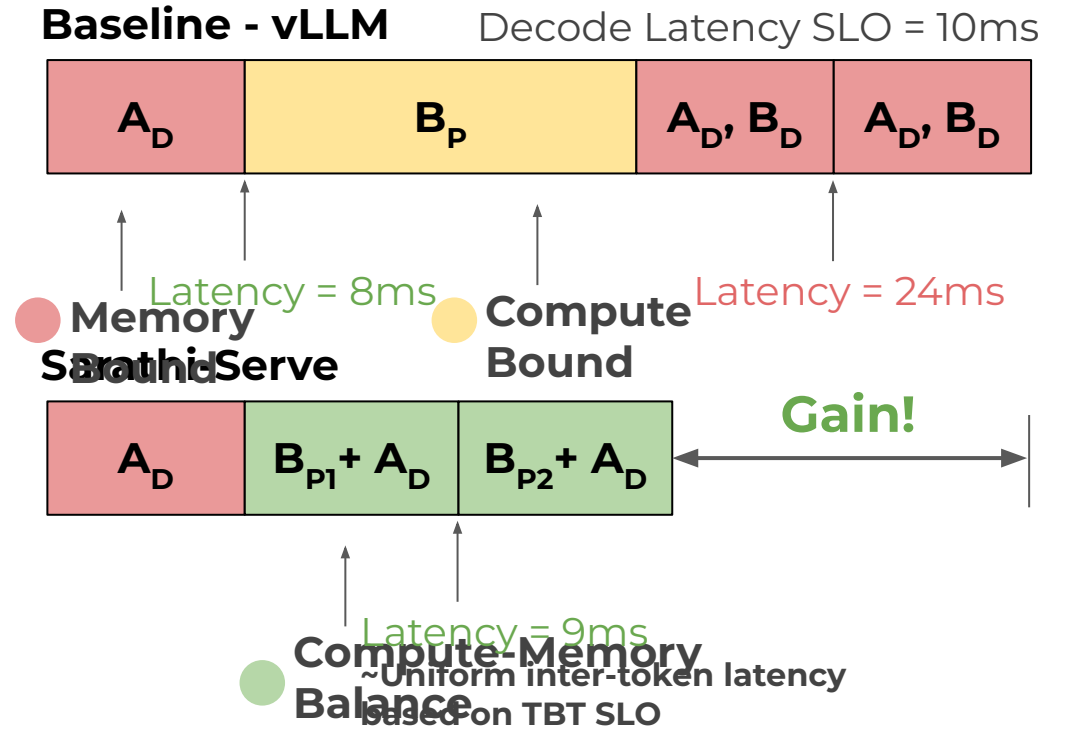
Latency dominated by compute  
Grows linearly with num tokens

Latency dominated by weight fetch time  
Independent of num tokens

# 🧠 Stall-free Batching

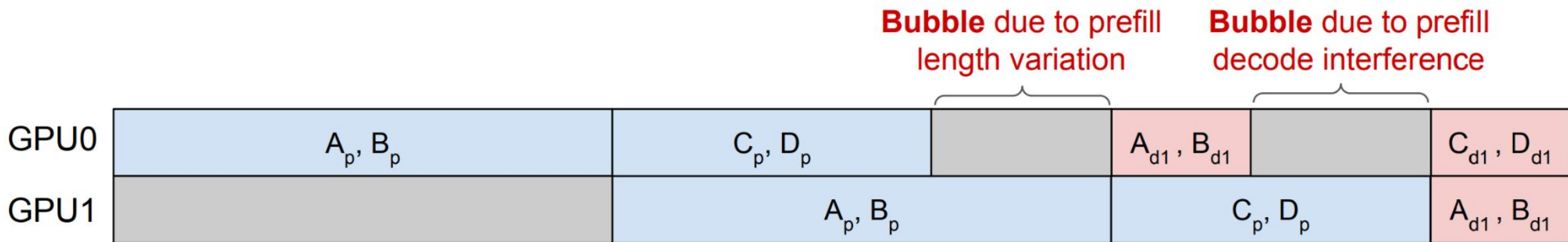
## Key Idea

🔪 Split large prefills into smaller chunks – just enough to consume the leftover compute budget in decode batches



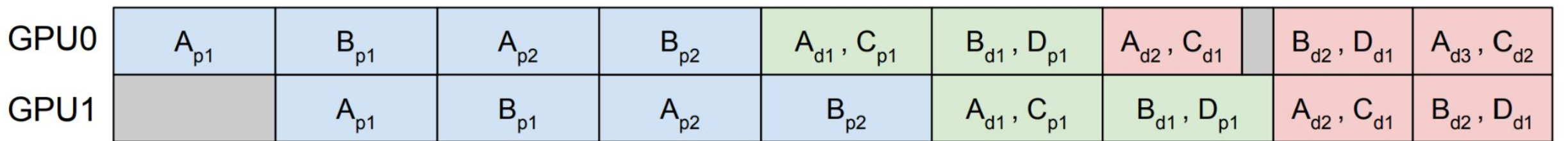
# Avoid Pipeline Bubbles

Timeline



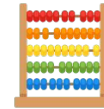
**Orca**

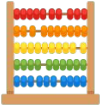
**Minimal Bubbles**



**Sarathi-Serve**

# Evaluations



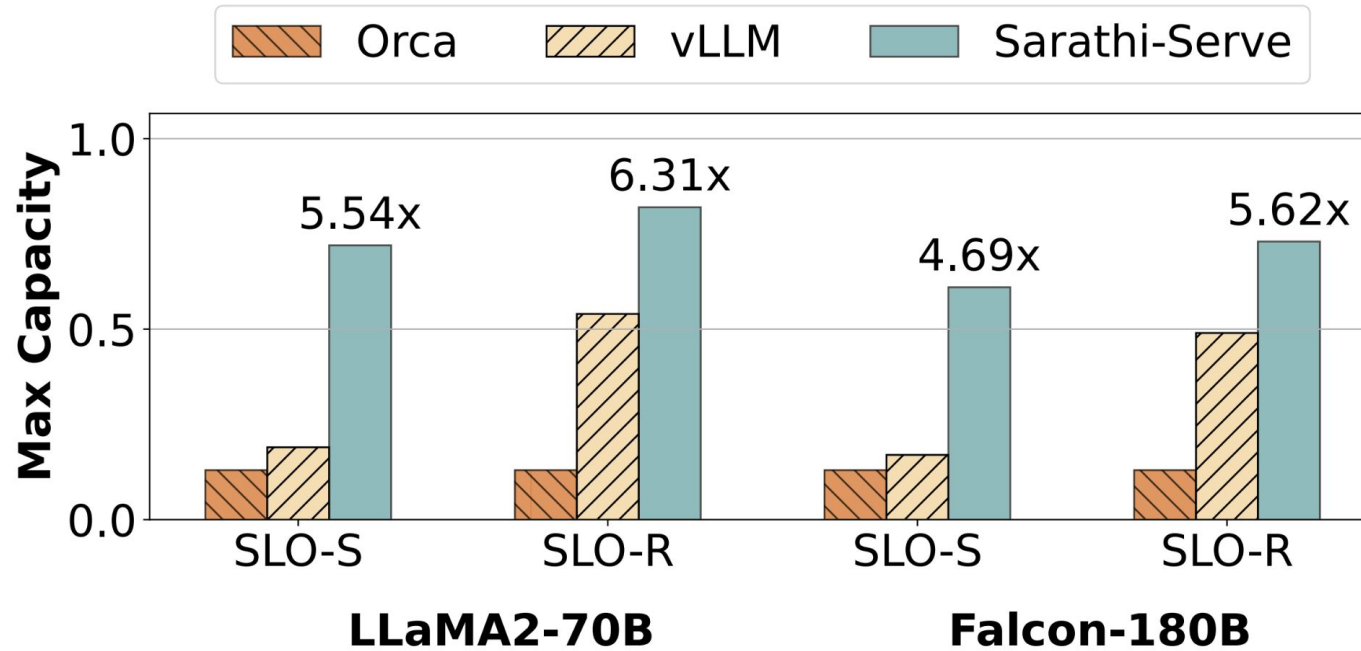


# Serving Capacity under SLOs

## Setup

ShareGPT4 trace on on A100 GPUs with strict (S) and relaxed (R) latency SLOs

adapt using different chunk sizes



**5-6x higher capacity 🦑**



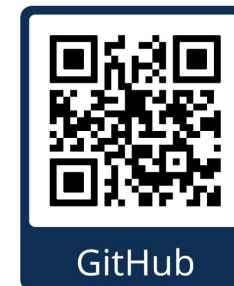
# Summary

 **Problem:** State-of-the-art systems sacrifice decode latency to achieve higher throughput

 **Key Insight** - Low arithmetic intensity of decodes allows for adding compute intensive prefills with negligible decode latency cost

 **Key Results** - We achieve optimality in both latency and throughput simultaneously leading up to 6x higher capacity under SLO constraints

 **Industry Adoption** - Available in all major serving frameworks and more.



# Course Logistics

- Assignment 2 out soon
  - Performance modeling of LLM prefill and decode
  - Use AMD clusters for the assignment
  - Read document on AMD cluster usage
- Information on project logistics will be released next week

# Throughput is Not All You Need

Maximizing **Goodput** in LLM Serving  
using **Prefill-Decode Disaggregation**

Hao Zhang @ Hao AI Lab & vLLM Team



# Prefill and Decode have Distinct Characteristics

- **Prefill**

Compute-bound

One prefill saturates compute. Batching across context tokens.



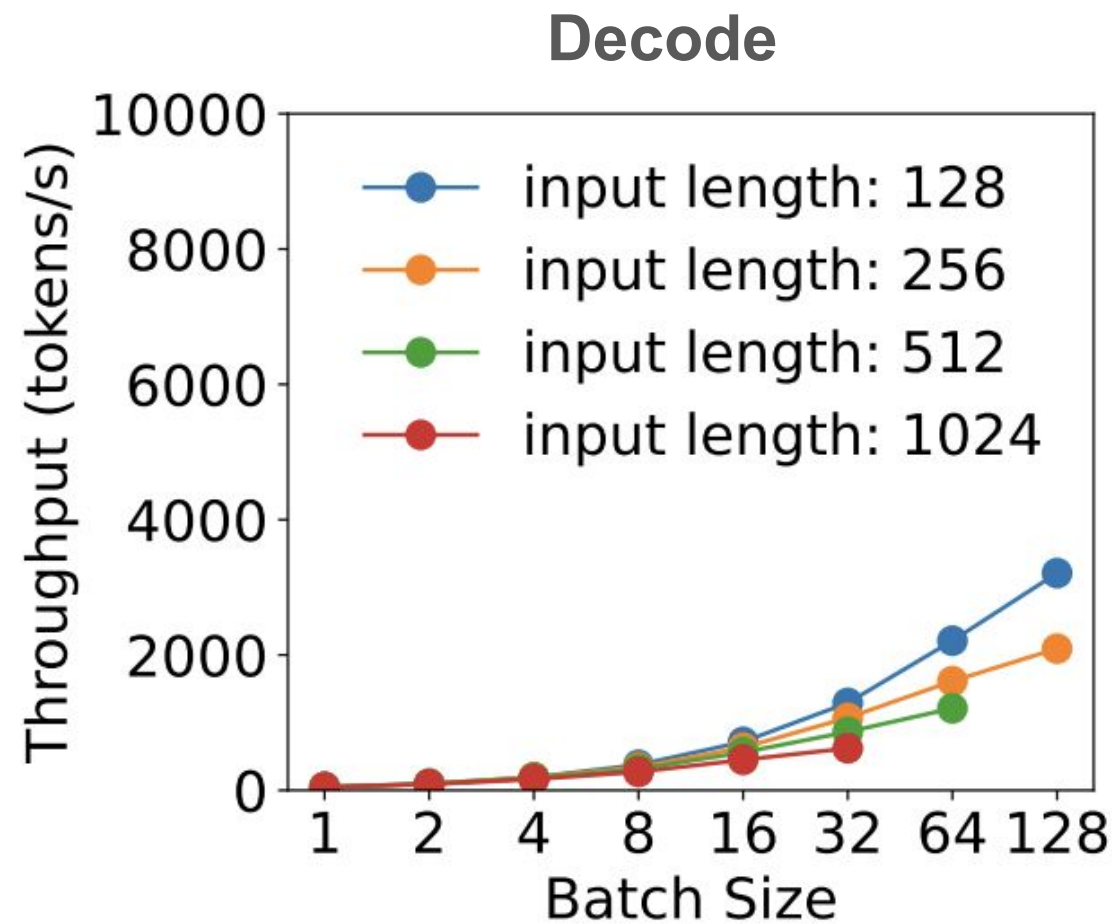
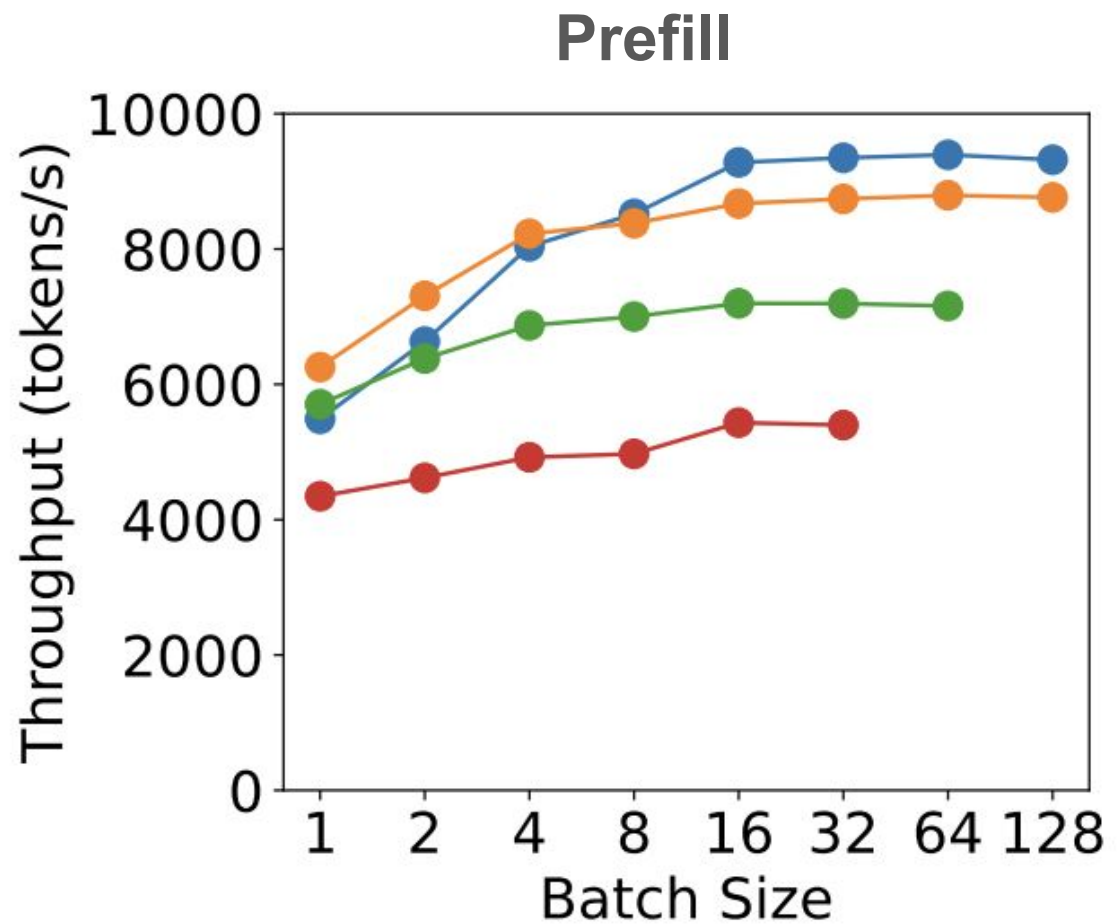
- **Decode**

Memory-bound

Must batch many requests together to saturate compute

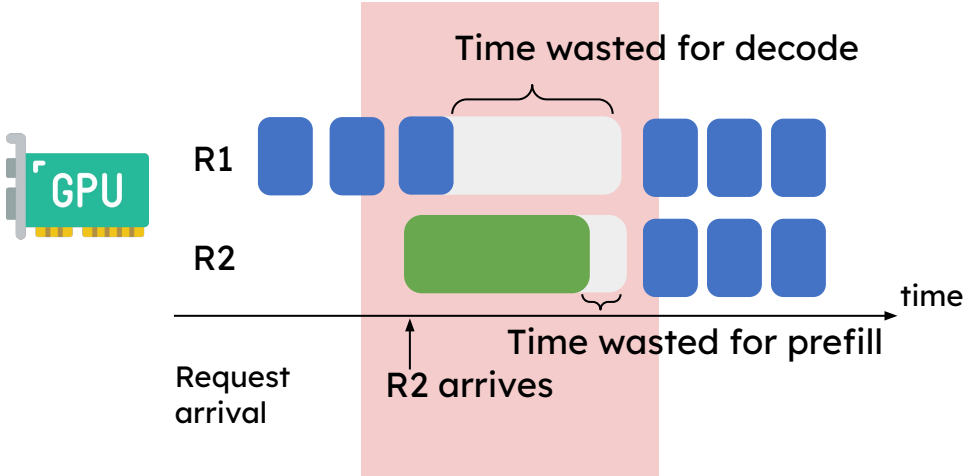


# Prefill and Decode have Distinct Characteristics



# Continuous Batching Causes Interference

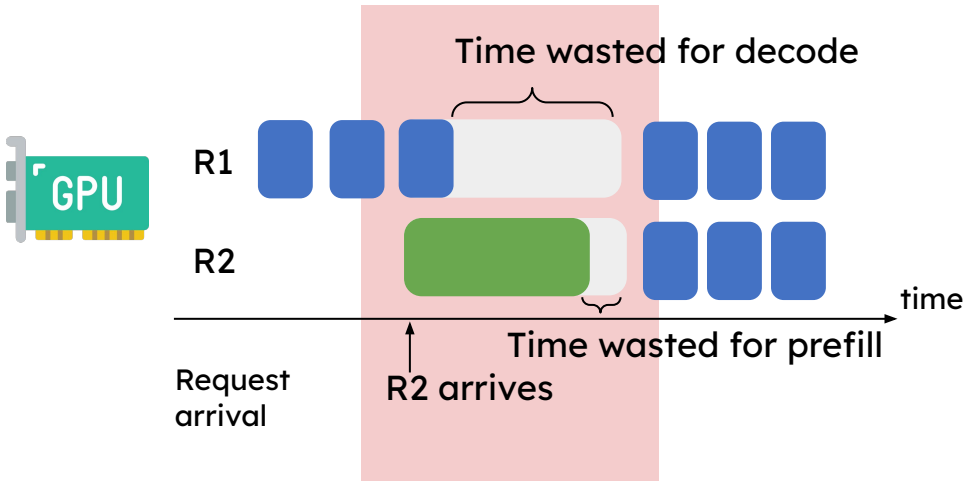
Continuous Batching  
Batch R1 and R2 together in 1 GPU



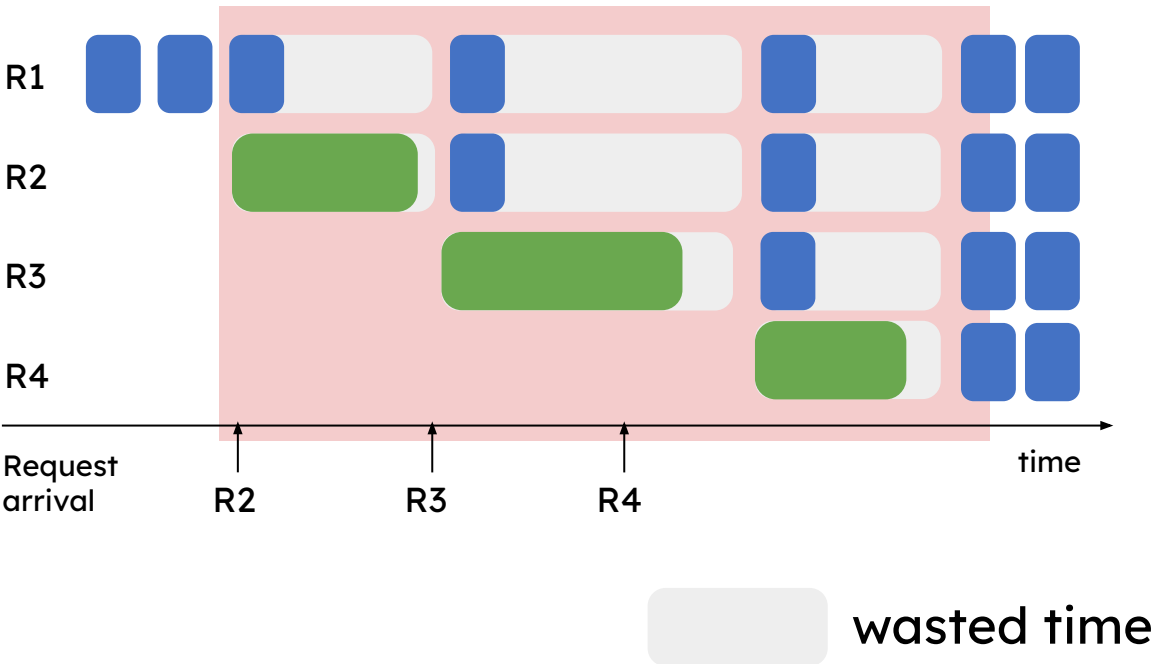
 wasted time

# Continuous Batching Causes Interference

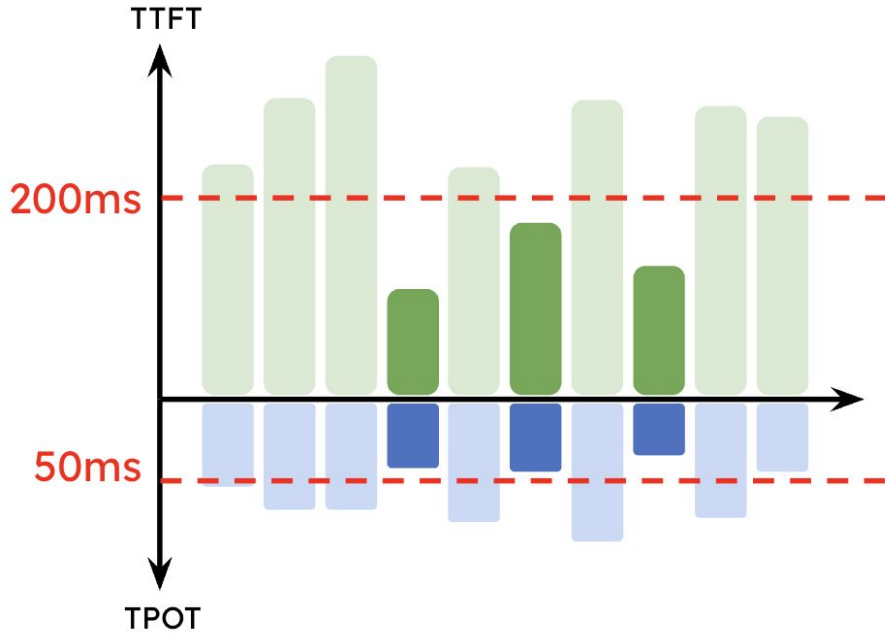
## Continuous Batching Batch R1 and R2 together in 1 GPU



## Continuous Batching Batch R1~R4 together in 1 GPU

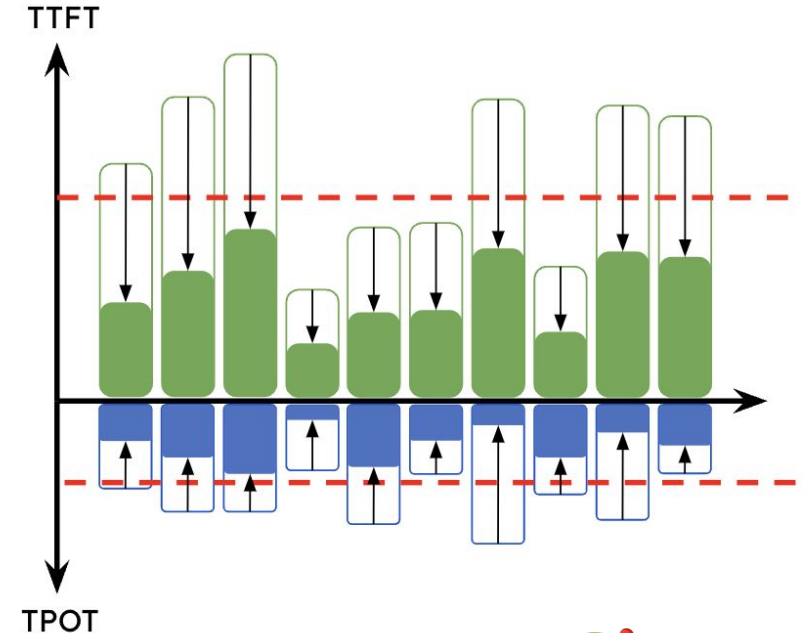
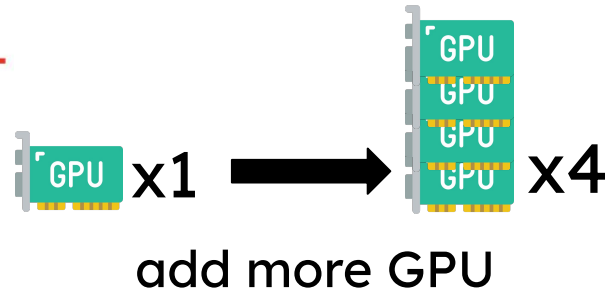


# Colocation → Overprovision Resource to meet SLO



Poor UX 🥰

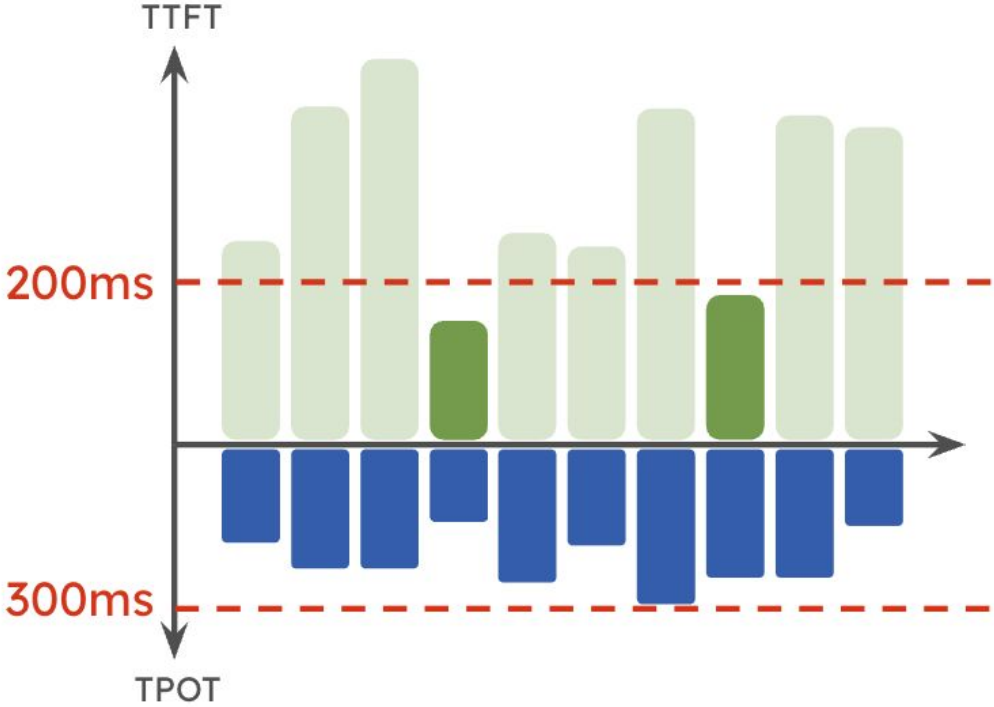
lower cost 💰



Good UX 🥰

Higher cost 💰💰💰💰

# Colocation → Coupled Parallelism

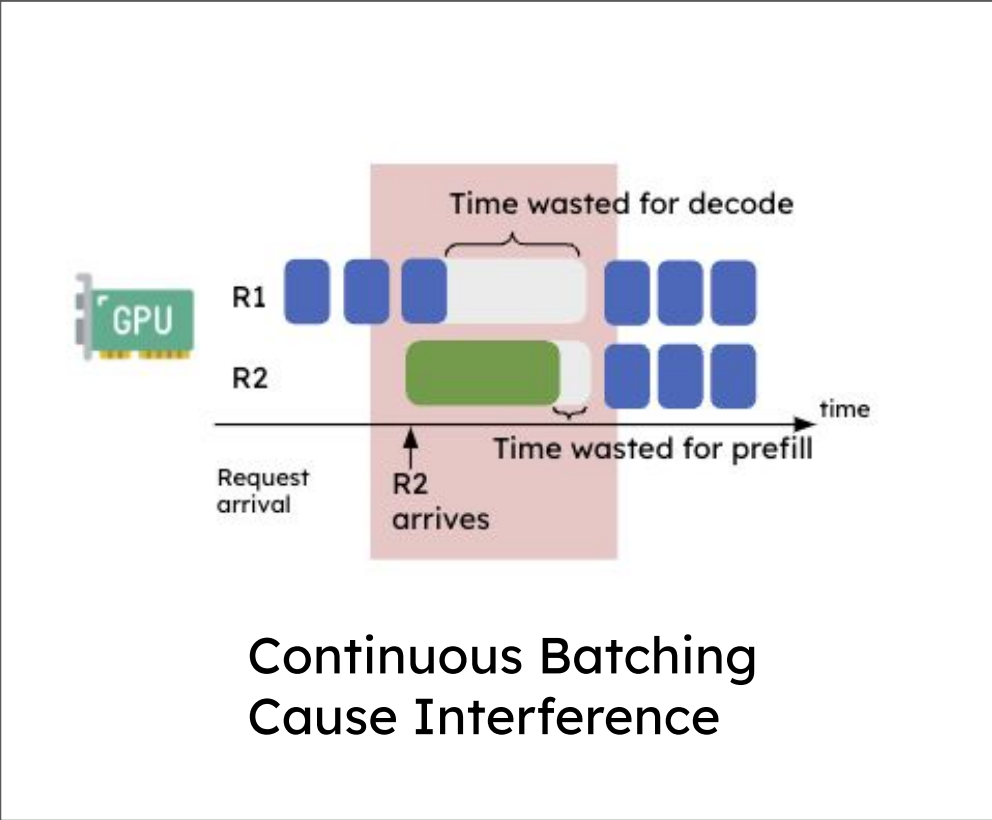


TTFT tight, TPOT loose

	PP	TP
● Prefill	😐	❤️
● Decode	❤️	😐

Prefill and Decode have different preferences

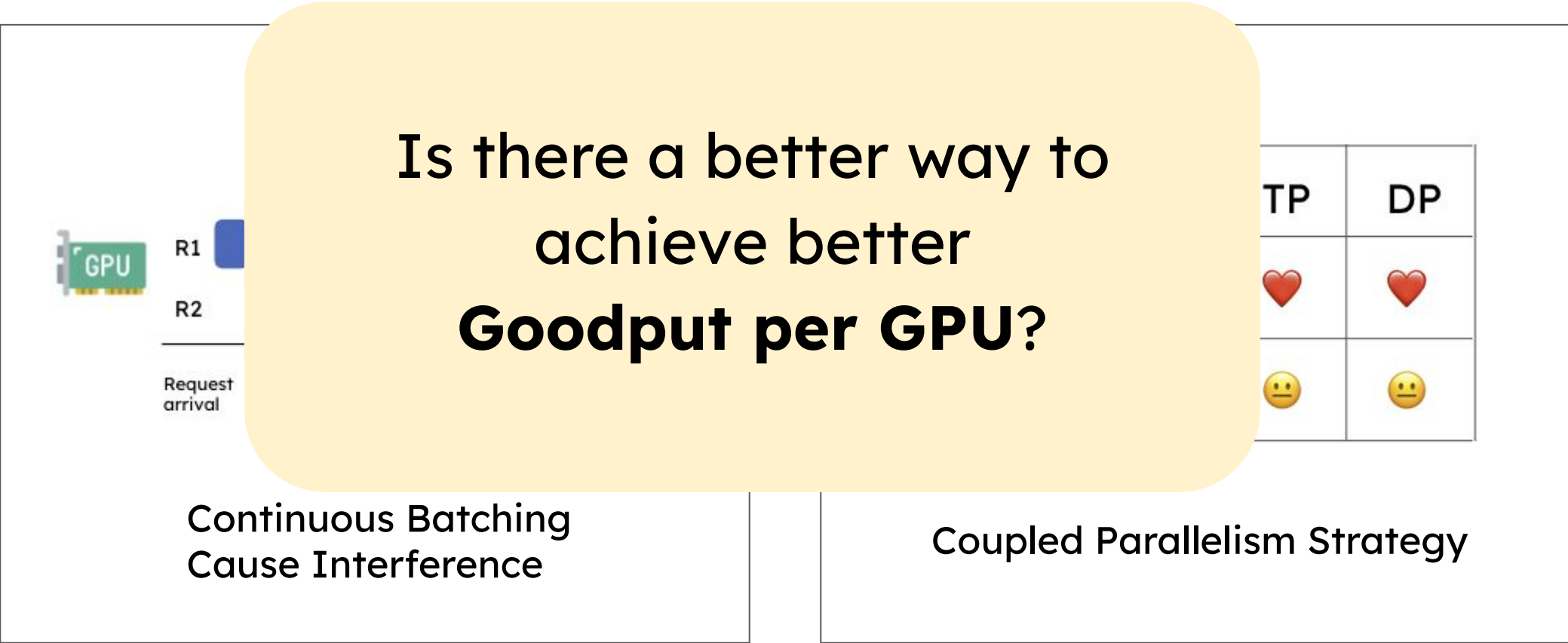
# Summary: Problems caused by Colocation



	PP	TP	DP
● Prefill	😐	❤️	❤️
● Decode	❤️	😐	😐

Coupled Parallelism Strategy

# Summary: Problems caused by Colocation

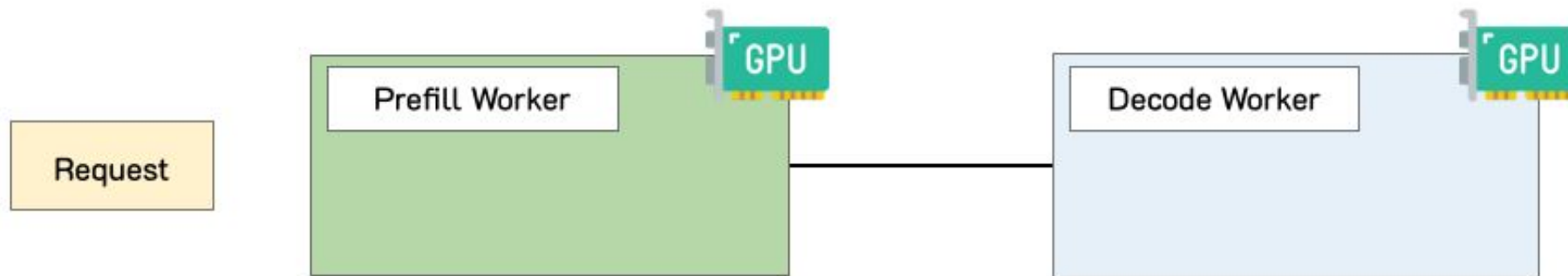




# DistServe: Disaggregating Prefill and Decode

Disaggregation is a technique that

## Request Arrived

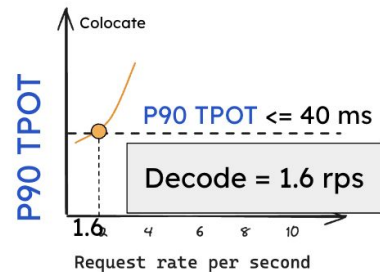
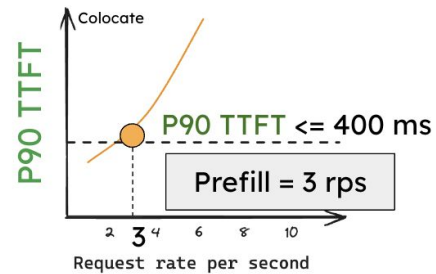


Timeline

# Disaggregation achieves better goodput

## Colocate

1 GPU for both Prefill and Decode



Max System goodput

= Min(Prefill, Decode)

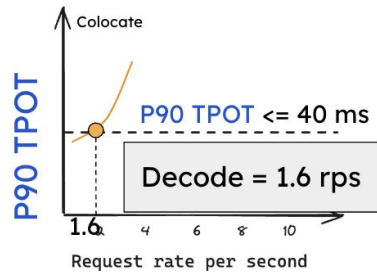
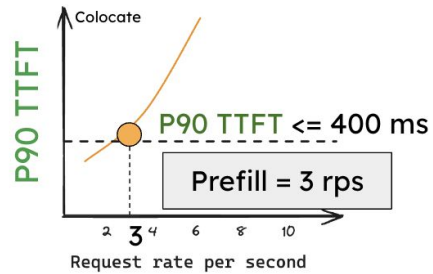
= 1.6 rps / GPU



# Disaggregation achieves better goodput

## Colocate

1 GPU for both Prefill and Decode



Max System goodput

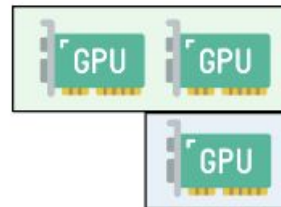
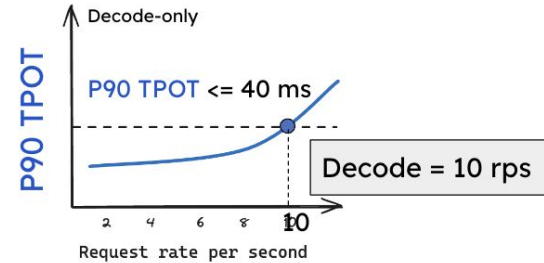
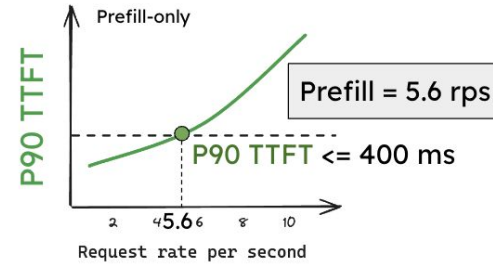
=  $\text{Min}(\text{Prefill}, \text{Decode})$

= 1.6 rps / GPU



## Disaggregate (2P1D)

2 GPU for Prefill + 1 GPU for Decode



Disaggregate (2P1D) goodput

=  $\text{Min}(5.6 \times 2, 10)$  rps / 3 GPU

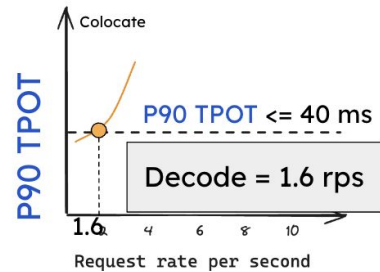
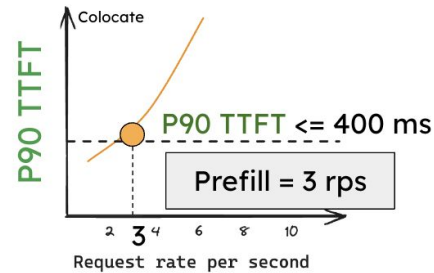
= 3.3 rps / GPU



# Disaggregation achieves better goodput

## Colocate

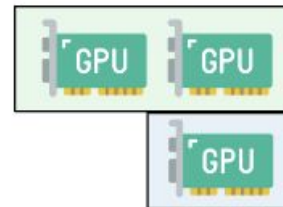
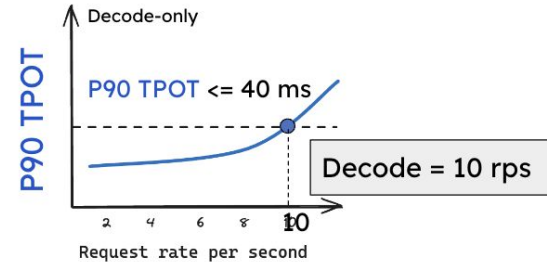
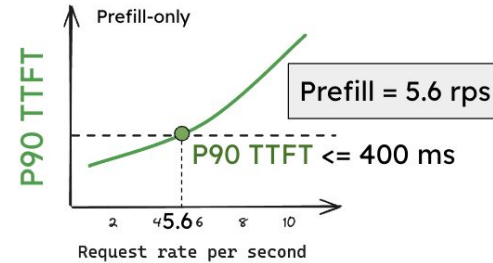
1 GPU for both Prefill and Decode



Max System goodput  
= Min(Prefill, Decode)  
= 1.6 rps / GPU 🙄

## Disaggregate (2P1D)

2 GPU for Prefill + 1 GPU for Decode

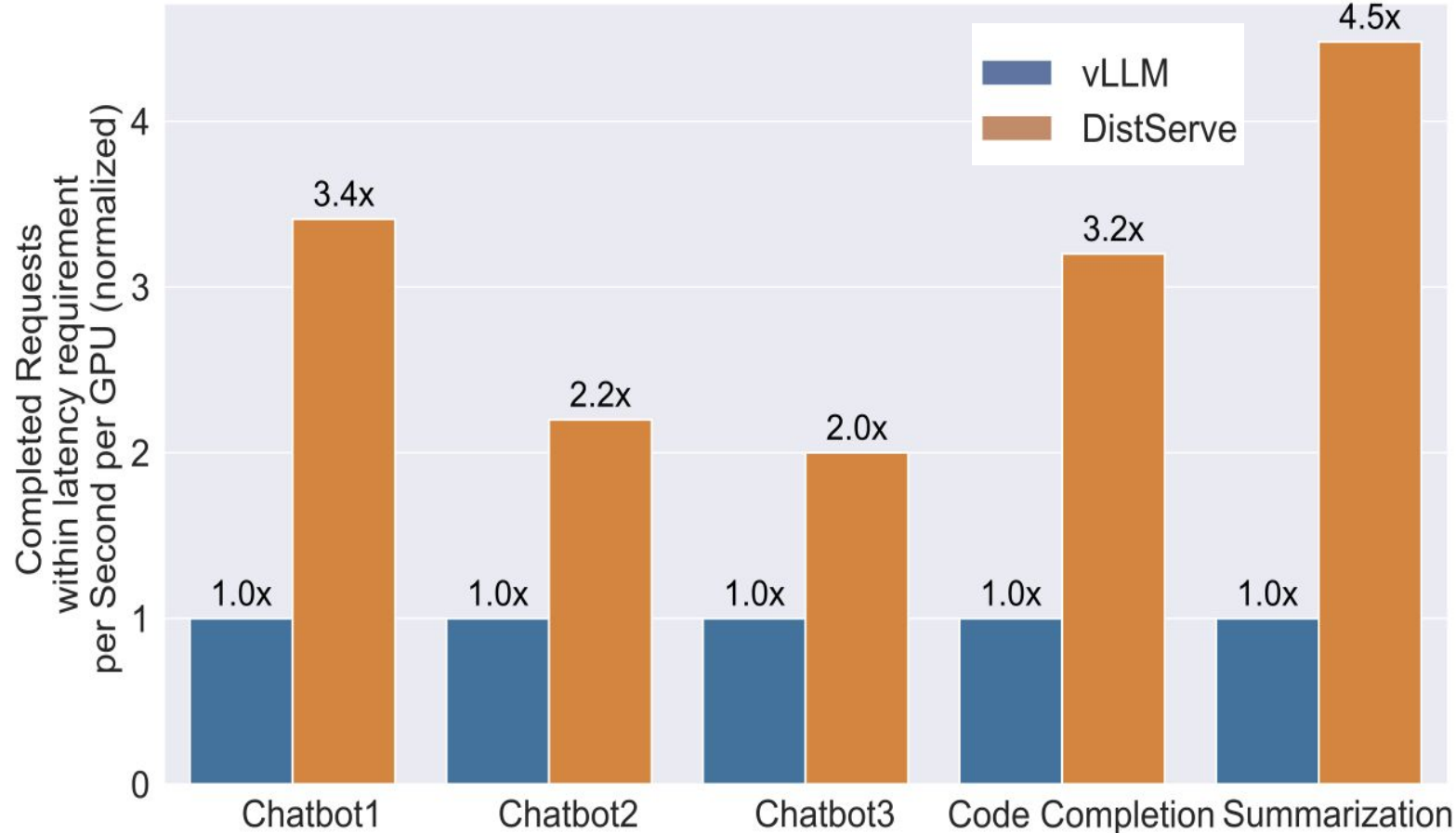


Disaggregate (2P1D) goodput  
= Min (5.6 x 2, 10) rps / 3 GPU  
= 3.3 rps / GPU 😎

Simple Disaggregation  
achieves **2x** goodput  
(per GPU)

What are issues or potential downsides with disaggregation?


# Evaluation



**Achieves 2.0x - 4.48x  
compared to vanilla vLLM**

- **Chatbot: 2.0 - 3.4x**
- **Code Completion: 3.2x**
- **Summarization: 4.5x**

# Summary

- **Goodput** instead of Throughput
- **Disaggregation** is effective to optimize **goodput!**
- **DistServe** achieves 2.0x - 4.48x compared to vLLM
- Integrating into  vLLM
- Already adopted by companies like **anyscale**