

# 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*

# Rise of LLMs

Technology

## ChatGPT sets record for fastest-growing user base - analyst note

By Krystal Hu

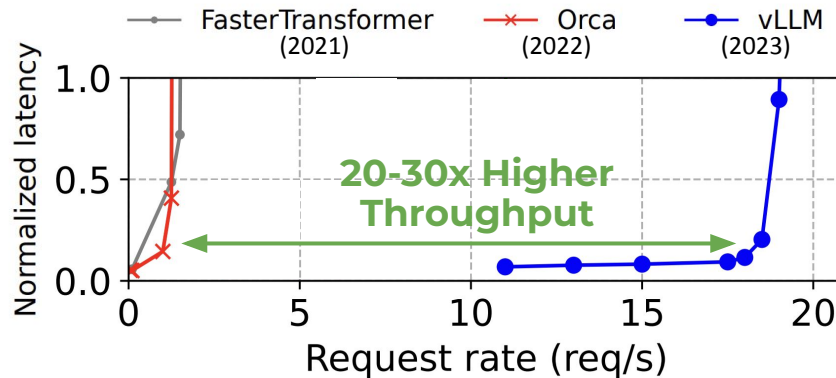
February 2, 2023 7:33 AM PST · Updated a year ago



CLIMATE

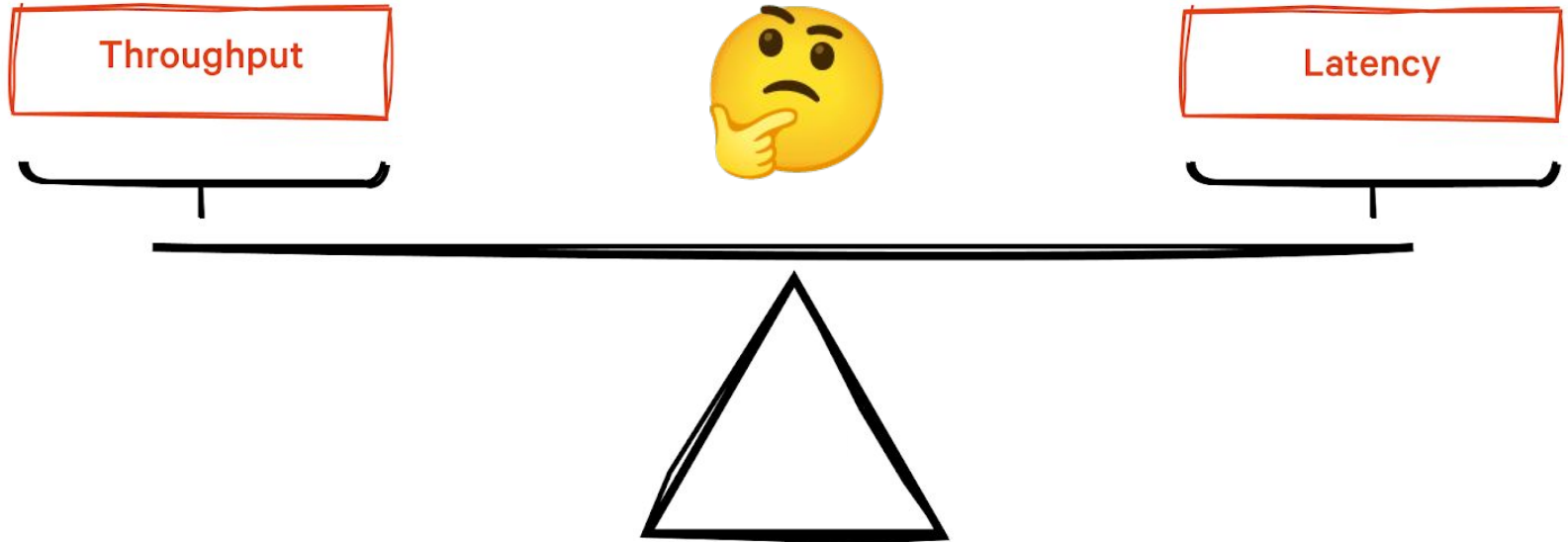
## Google's carbon emissions surge nearly 50% due to AI energy demand

PUBLISHED TUE, JUL 2 2024 3:41 PM EDT | UPDATED MON, JUL 8 2024 9:32 AM EDT



## & Inference Systems

Can we maintain low latency  
with high throughput?



**Demo**



# In this talk...

 **Latency-throughput tradeoff:** Analyzing LLM batching policies

 **Finding a free lunch:** Arithmetic Intensity Slack in LLM Inference

 **Stall-free batching:** Leveraging chunked prefill to overcome the latency-throughput tradeoff

 **Evaluations:** Key results and analysis

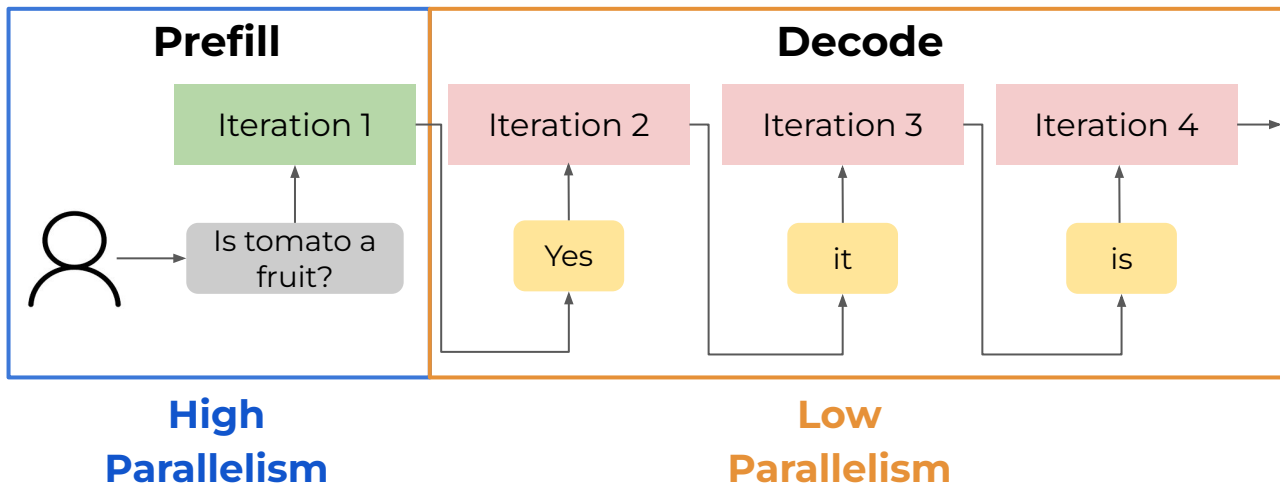
What causes the latency-throughput tradeoff in LLM inference systems? 🕵️



# Background: LLM Inference 101

👍 GPU Utilization

👎 GPU Utilization



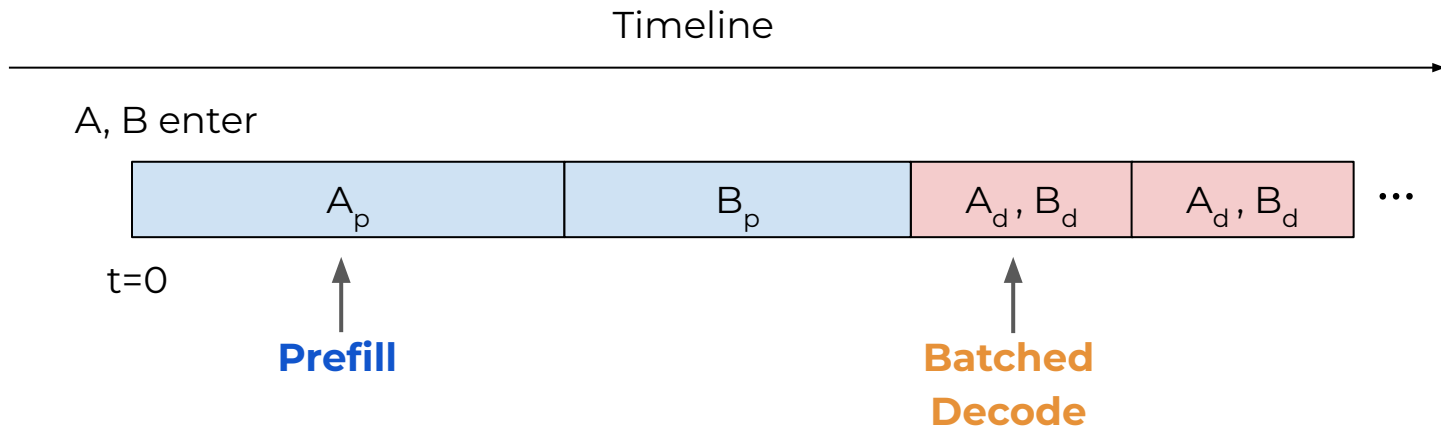
How to improve parallelism during  
decode phase? 🤔

**Batching** 🙄





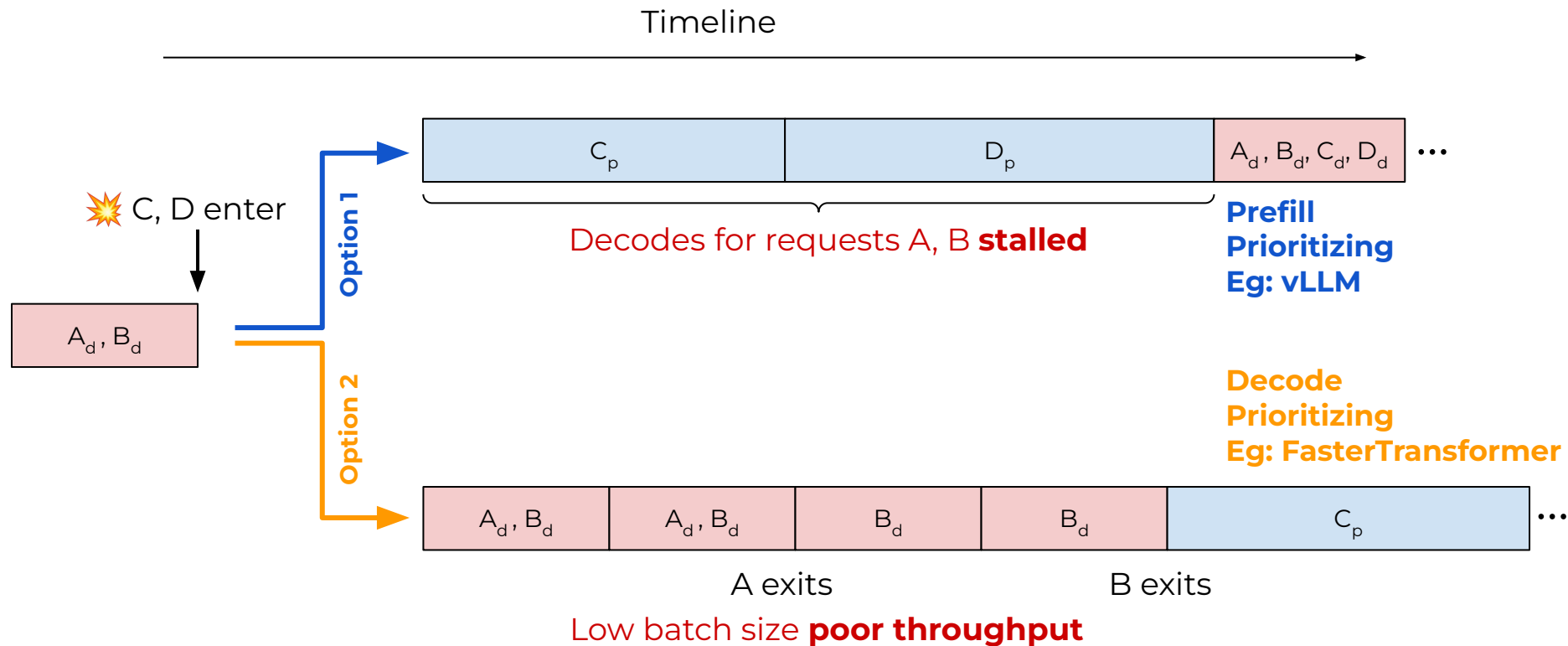
# Background: Batching LLM Inference



Decode efficiency increases linearly with batch size 🚀

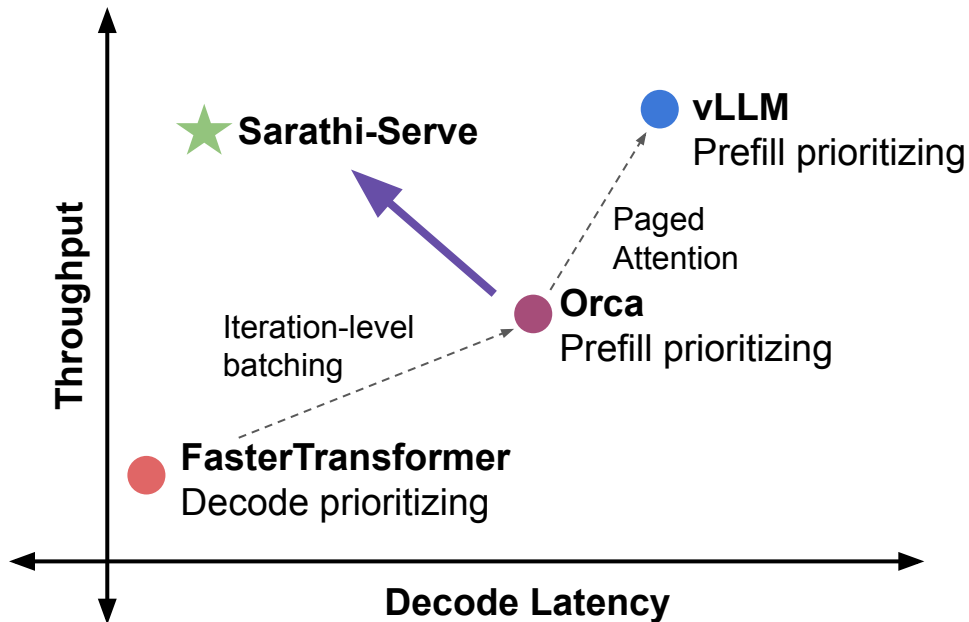
📦 Batch size  $\Rightarrow$  📦 Throughput

# The Prefill-Decode Scheduling Conundrum





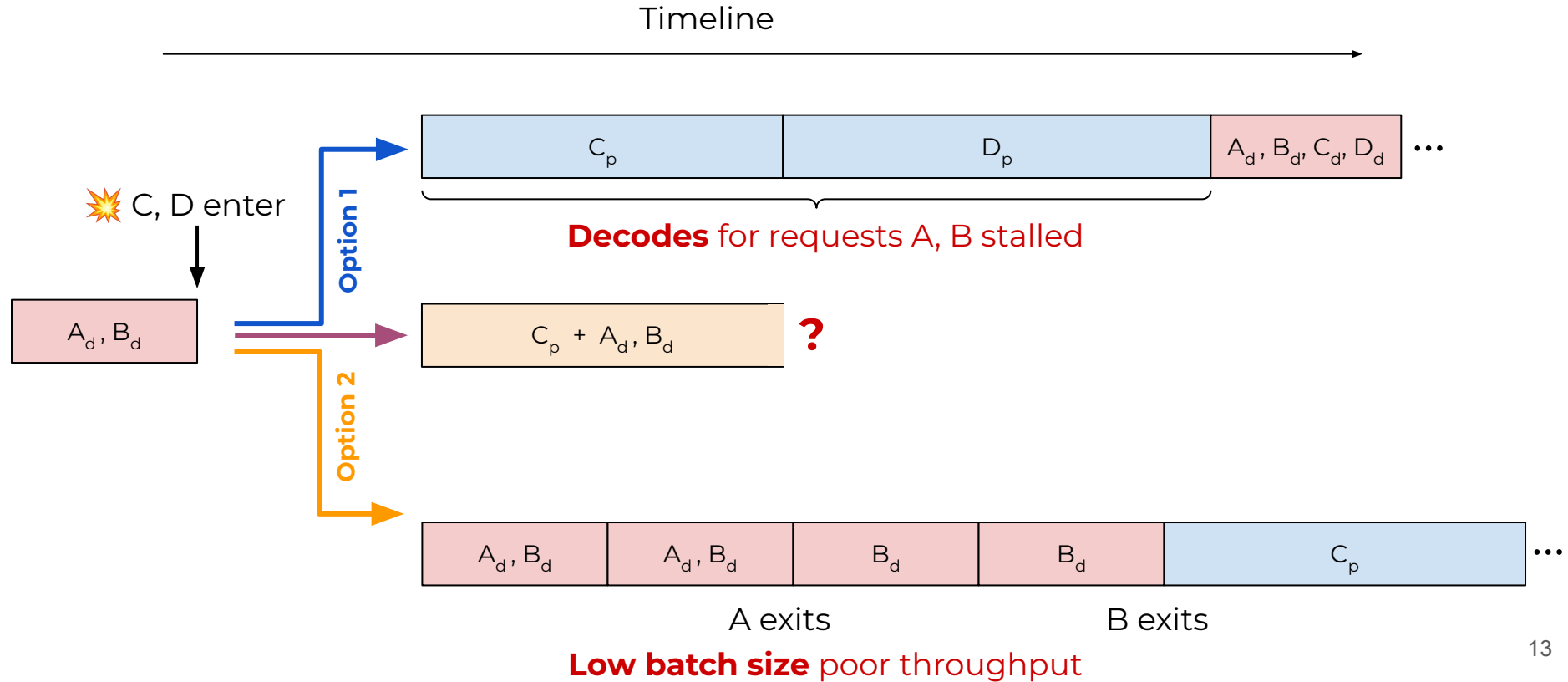
# The Latency-Throughput Tradeoff



Existing batching policies make a harsh latency-throughput tradeoff !

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

# The Prefill-Decode Scheduling Conundrum



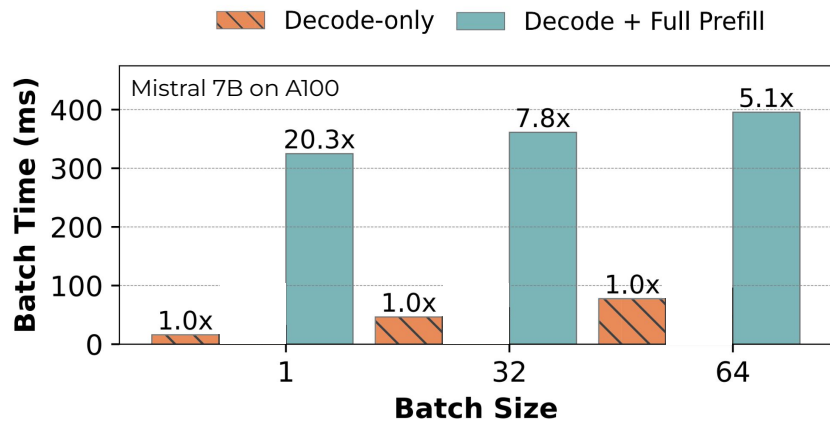
# Mixed Batching

## Idea

👤 Fused computation of prefill and decodes

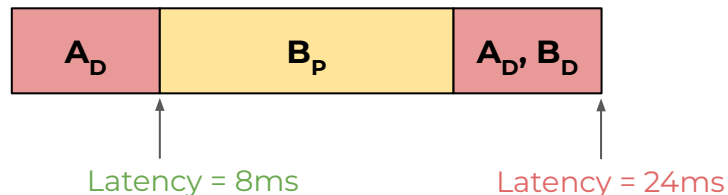
## Challenge

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

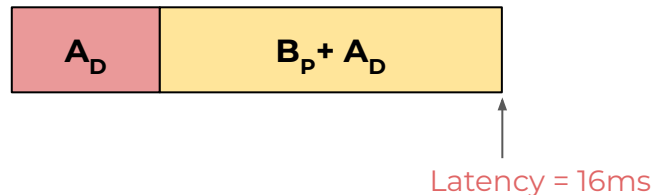


### vLLM

Decode Latency SLO = 10ms



### Orca

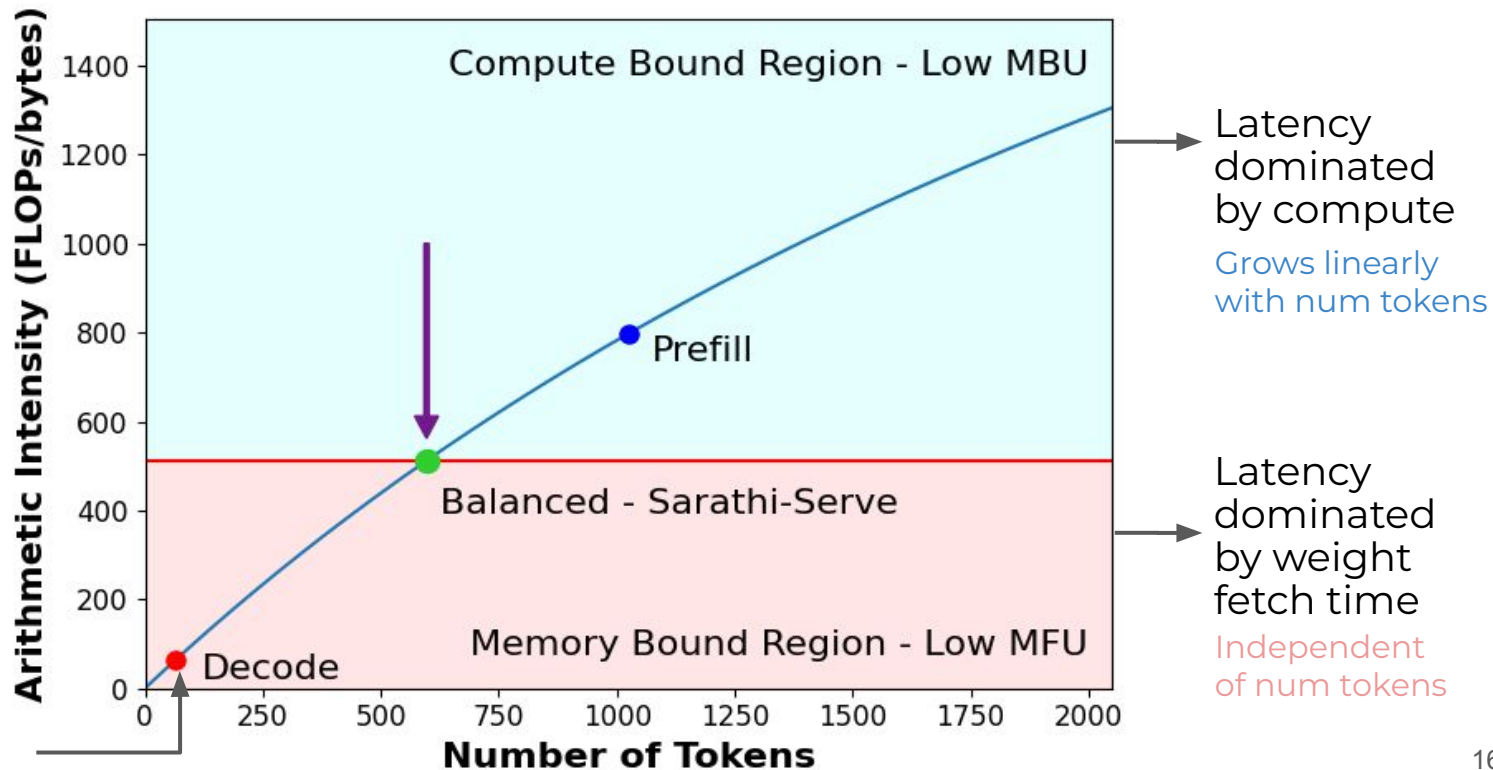


## Key Insight

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



# Observation: Arithmetic Intensity Slack



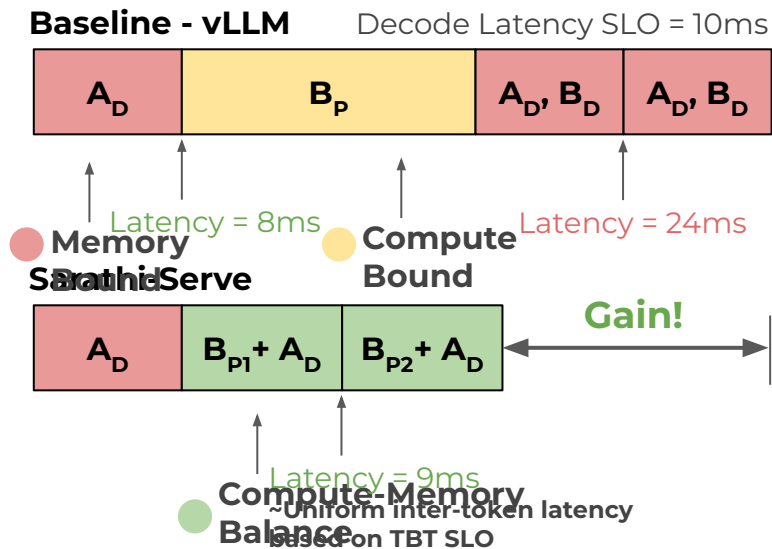
Constrained due to memory overhead in decode phase



# 🧠 Stall-free Batching

## Key Idea

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



**Demo**



# Evaluations





## Background: Performance Metrics

 **Time to first token (TTFT):** Time required for the first token to show up from the time user submits a request

 **Time between tokens (TBT):** Latency between each output token

 **Capacity:** Maximum QPS that can be served while satisfying latency SLOs

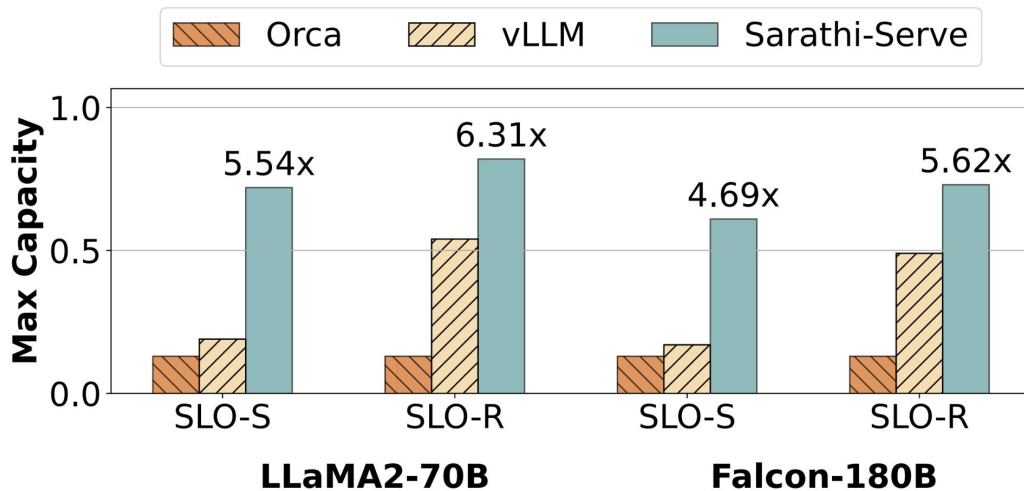


# Serving Capacity under SLOs

## 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.

