

LLM SERVING AND ALIGNMENT

Team 6

Fengyu Gao, Shunqiang Feng, Wei Shen, Zihan Zhao



SGLang: Efficient Execution of Structured Language Model Programs

Zihan Zhao (rxy6cc)

```
# In this script, we demonstrate how to pass input to the chat method:
```

```
conversation = [  
    {  
        "role": "system",  
        "content": "You are a helpful assistant"  
    },  
    {  
        "role": "user",  
        "content": "Hello"  
    },  
    {  
        "role": "assistant",  
        "content": "Hello! How can I assist you today?"  
    },  
    {  
        "role": "user",  
        "content":  
            "Write an essay about the importance of higher education.",  
    },  
]  
  
outputs = llm.chat(conversation, sampling_params, use_tqdm=False)  
print_outputs(outputs)
```

Format the inputs **by hand**

Run the LLM with the inputs

Just like how you regularly code

SGLang

12.7k ★

```
import sglang as sgl
```

```
@sgl.function
```

```
def multi_turn_question(s, question_1, question_2):
```

```
    s += sgl.user(question_1)
```

```
    s += sgl.assistant(sgl.gen("answer_1", max_tokens=256))
```

```
    s += sgl.user(question_2)
```

```
    s += sgl.assistant(sgl.gen("answer_2", max_tokens=256))
```

```
def single():
```

```
    state = multi_turn_question.run(
```

```
        question_1="What is the capital of the United States?",
```

```
        question_2="List two local attractions.",
```

```
    )
```

```
    for m in state.messages():
```

```
        print(m["role"], ":", m["content"])
```

```
    print("\n-- answer_1 --\n", state["answer_1"])
```

Define the interaction flow

Run the LLM directly
from the flow!

A much better programming paradigm

vLLM vs. SGLang

- vLLM

- A framework focusing on system efficiency

- SGLang

- A framework focusing on system efficiency **and programming efficiency**

Motivations

- Lack of efficient systems
 - *vLLM was not built yet*
 - **Redundant** computations
 - **Redundant** memory usage
- Lack of programming efficiency
 - **Convolutd** code to start a server, expose an API, and run LLMs
 - LLM applications **!=** Any other applications (e.g. webapp)

Frontend Language

- Challenges
 - String manipulations
 - Prompt construction (e.g. roles, message, attachments, etc.)
 - Multimodality supports
 - Multimodal token placement (e.g. audio tokens, video tokens, image tokens, etc.)
 - Output parsing
 - Yes-or-no selection
 - Code extraction

Frontend Language

- Language primitives
 - Roles
 - `system` / `user` / `assistant`
 - Multimodal files
 - `image` / `video`
 - Control flow
 - `fork` / `select` / `gen`

```
dimensions = ["Clarity", "Originality", "Evidence"]
@function
def multi_dimensional_judge(s, path, essay):
    s += system("Evaluate an essay about an image.")
    s += user(image(path) + "Essay:" + essay)
    s += assistant("Sure!")

    # Return directly if it is not related
    s += user("Is the essay related to the image?")
    s += assistant(select("related", choices=["yes", "no"]))
    if s["related"] == "no": return

    # Judge multiple dimensions in parallel
    forks = s.fork(len(dimensions))
    for f, dim in zip(forks, dimensions):
        f += user("Evaluate based on the following dimension:" +
            dim + ". End your judgment with the word 'END'")
        f += assistant("Judgment:" + gen("judgment", stop="END"))

    # Merge the judgments
    judgment = "\n".join(f["judgment"] for f in forks)

    # Generate a summary and a grade. Return in the JSON format.
    s += user("Provide the judgment, summary, and a letter grade")
    s += assistant(judgment + "In summary," + gen("summary", stop=".")
        + "The grade of it is" + gen("grade"))

    schema = r'{"summary": "[\w\d\s]+\.", "grade": "[ABCD][+-]?"}'
    s += user("Return in the JSON format.")
    s += assistant(gen("output", regex=schema))

state = multi_dimensional_judge.run(...)
print(state["output"])
```

Handle chat template and multi-modal inputs

Select an option with the highest probability

Fetch result; Use Python control flow

Runtime optimization: KV Cache Reuse (Sec. 3)

Multiple generation calls run in parallel

Fetch generation results

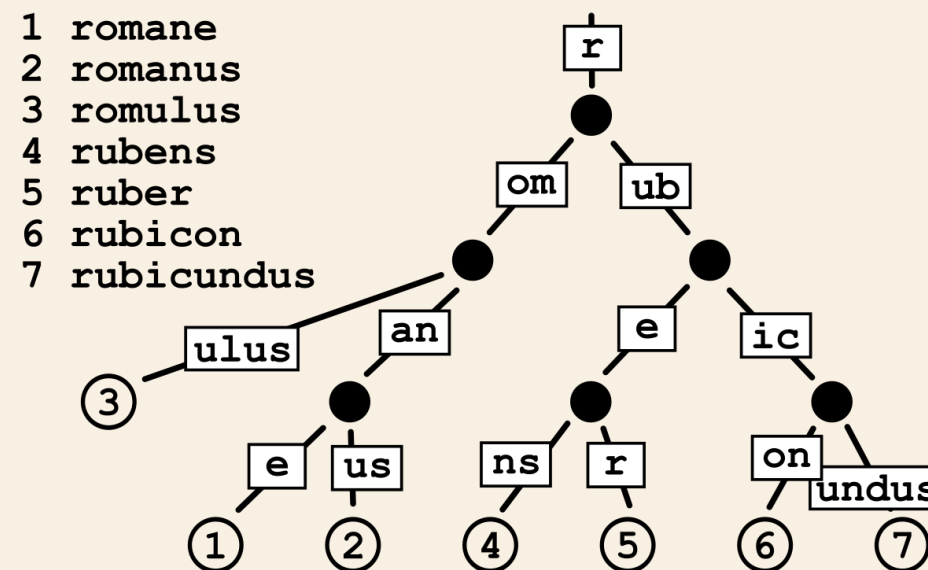
Runtime optimization: API speculative execution (Sec. 5)

Runtime optimization: fast constrained decoding (Sec. 4)

Run an SGLang program

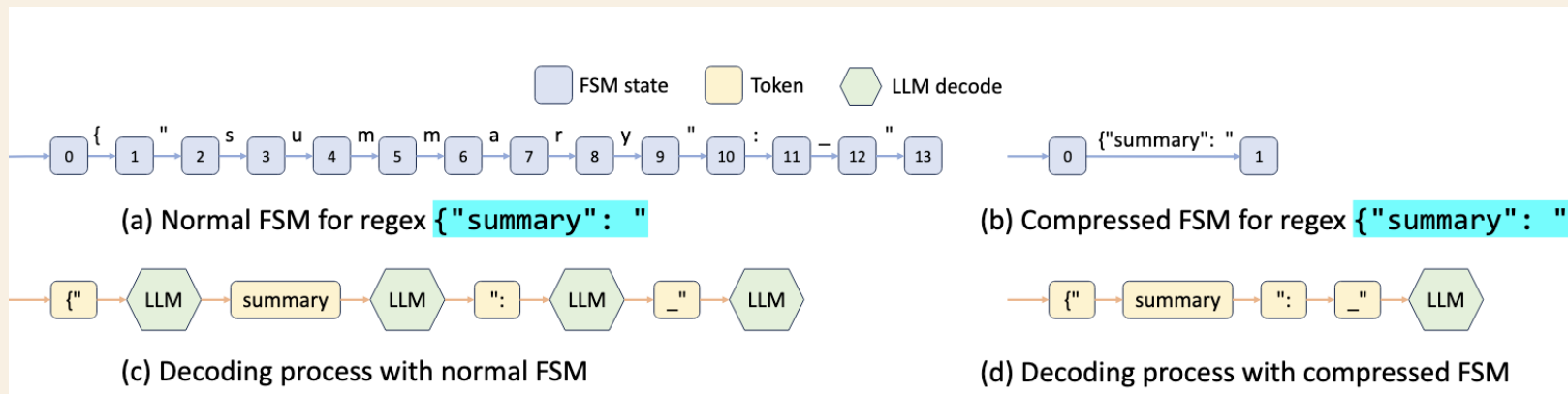
KV Cache Reuse

- RadixAttention
 - Prefix sharing
 - Radix tree + LRU cache
 - Prompt deduplication
 - Substring repetition discoveries



Accelerated Constrained Decoding

- Compressed finite state machines (CFSM)
 - FSM determines whether generated tokens met the constraints
 - *It's token-by-token*
 - CFSM supports multi-token processing



API Speculative Execution

- For black-box API
- Intuition: The first
- Imagine an LLM

(1d) ReAct (Reason + Act)

Thought 1: I need to search Apple Remote and find the program it was originally designed to interact with.

Act 1: `Search[Apple Remote]`

Obs 1: The Apple Remote is a remote control introduced in October 2005 by Apple ... originally designed to control the Front Row media center program ...

Thought 2: Apple Remote was originally designed to control the Front Row media center program. I need to search Front Row next and find what other device can control it.

Act 2: `Search[Front Row]`

Obs 2: Could not find [Front Row]. Similar: ['Front Row Seat to Earth', 'Front Row Motorsports', 'Front Row (software)', ...]

Thought 3: Front Row is not found. I need to search Front Row (software) .

Act 3: `Search[Front Row (software)]`

Obs 3: Front Row is a discontinued media center software ...

Thought 4: Front Row (software) is controlled by an Apple Remote or the keyboard function keys. So the answer is keyboard function keys.

Act 4: `Finish[keyboard function keys]`



Evaluation

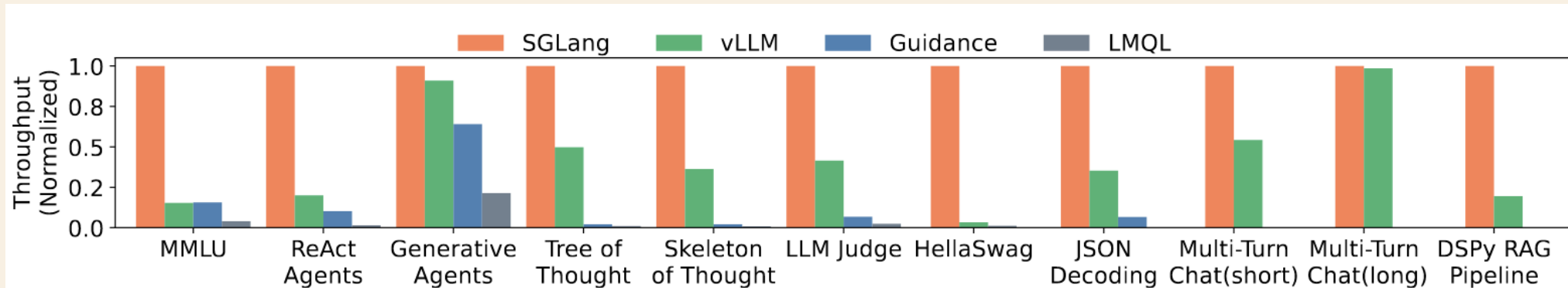


Figure 5: Normalized throughput on Llama-7B models. Higher is better.

Takeaways

- A developer-friendly frontend “language”
- An efficient backend runtime
 - **RadixAttention** to increase KV cache reusability
 - **Compressed FSM** to accelerate constrained decoding
 - **API speculative execution** to reduce E2E latency



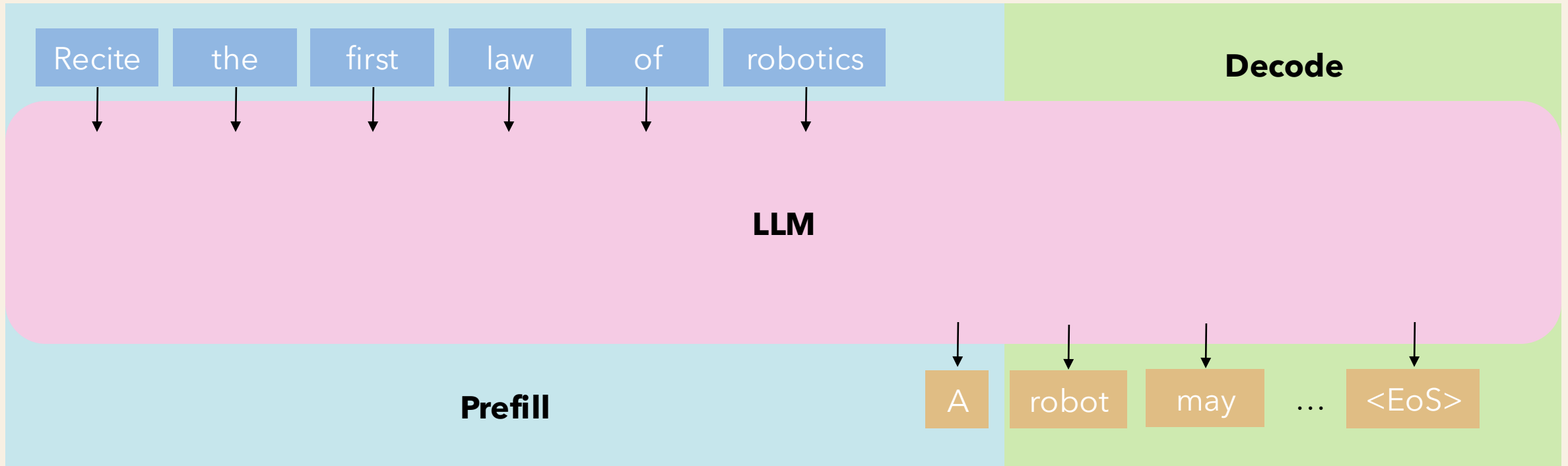
Taming Throughput-Latency Tradeoff in LLM Inference with Sarathi-Serve:

An efficient LLM inference scheduler that significantly improves throughput while maintaining low latency.

Shunqiang Feng (mpp7ez)

1. Background

(a) Prefill & Decode

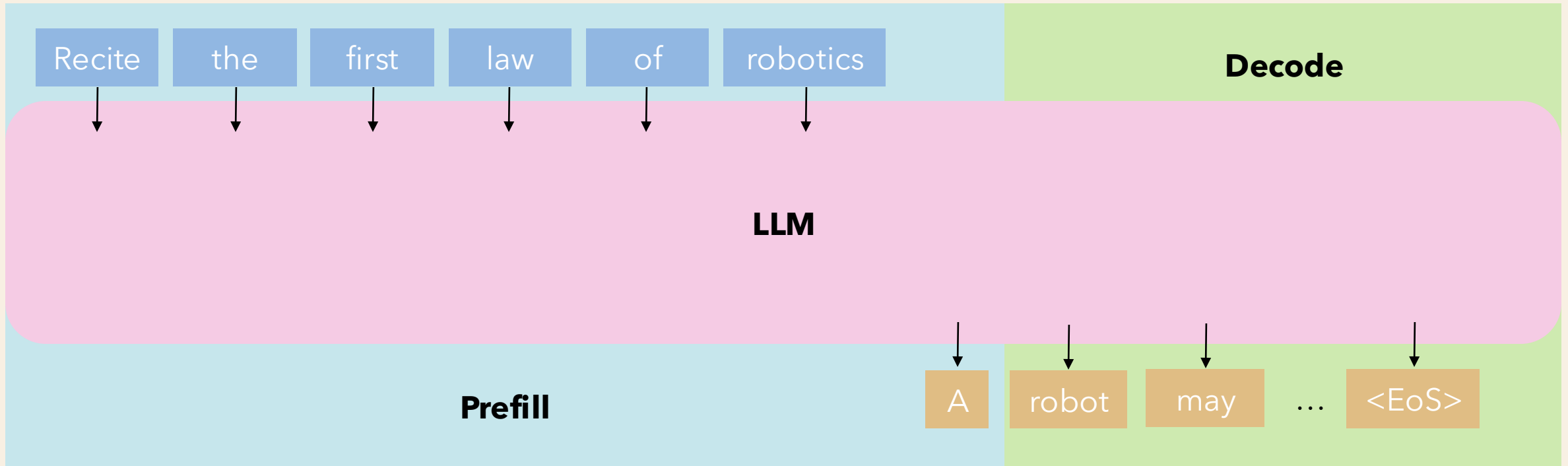


Each LLM serving request goes through two phases:

1. Prefill
To process the entire input prompt and produces the first output token
2. Decode
To generate the rest of output tokens, one-at-a-time.

1. Background

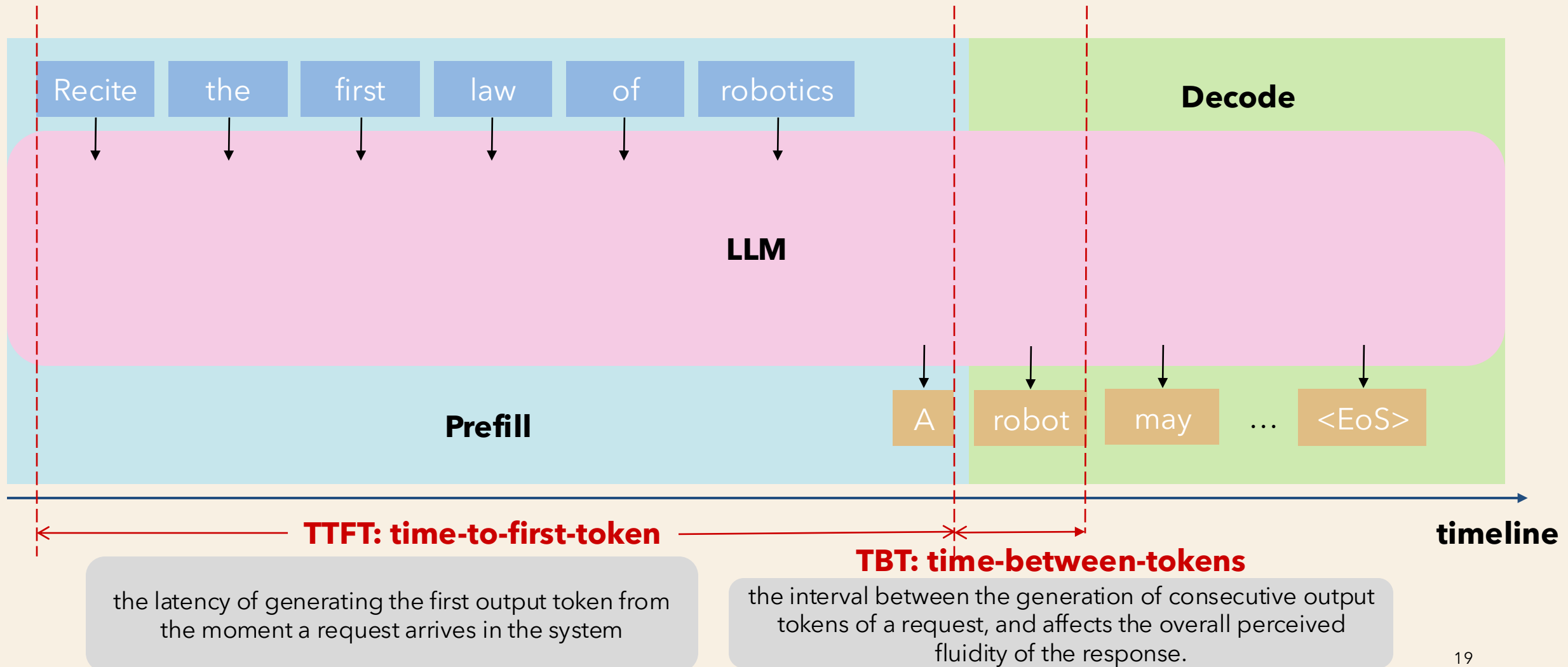
(a) Prefill & Decode



Stage	Resource Demand	Processing Style	benefit from batching?
Prefill	Computation-Intensive	Processes all input tokens in parallel	No
Decode	Memory-Intensive	Processes one token at a time	Yes

1. Background

(b) LLM Service Metrics: Latency



1. Background

(b) LLM Service Metrics: Throughput

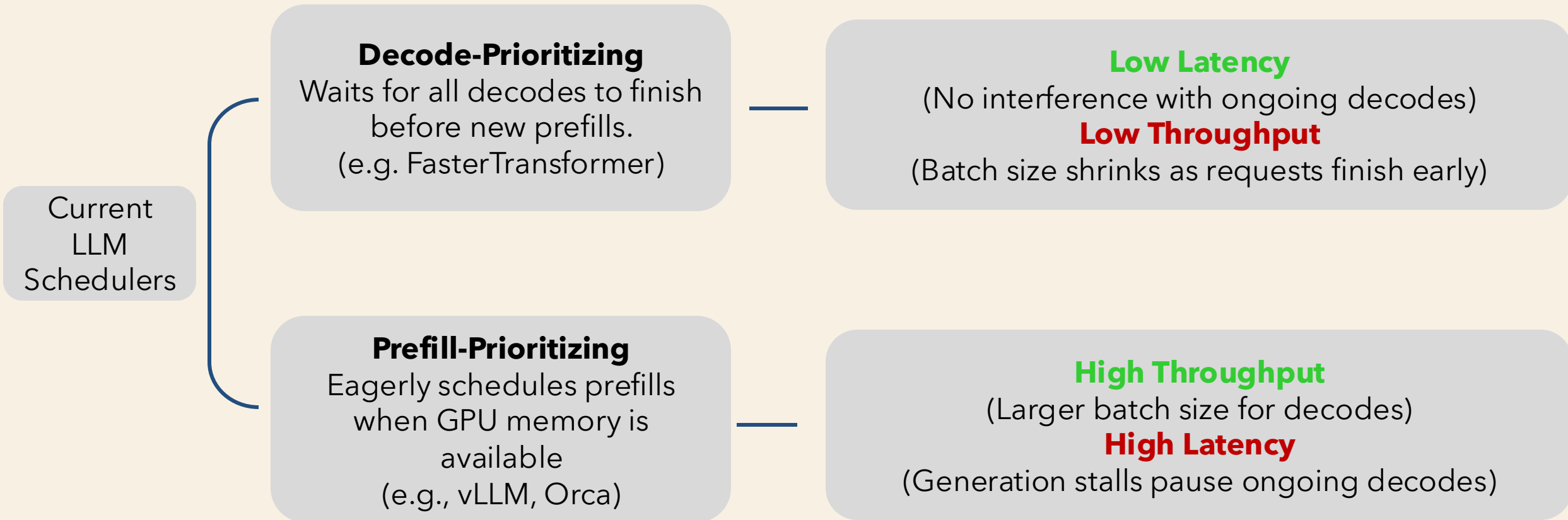
Capacity

- Definition:
the maximum request load (queries-per-second) a system can sustain while meeting certain latency targets.

Higher capacity is desirable because it reduces the cost of serving.

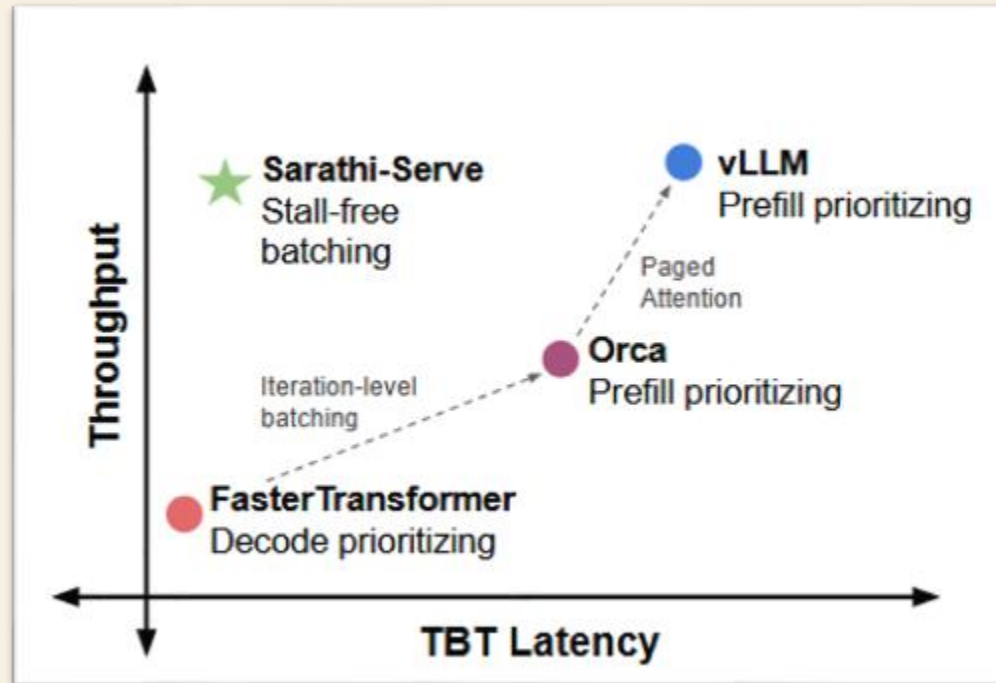
1. Background

(c) Current LLM Schedulers



1. Background

(c) Current LLM Schedulers



Tradeoff between throughput and latency in current schedulers

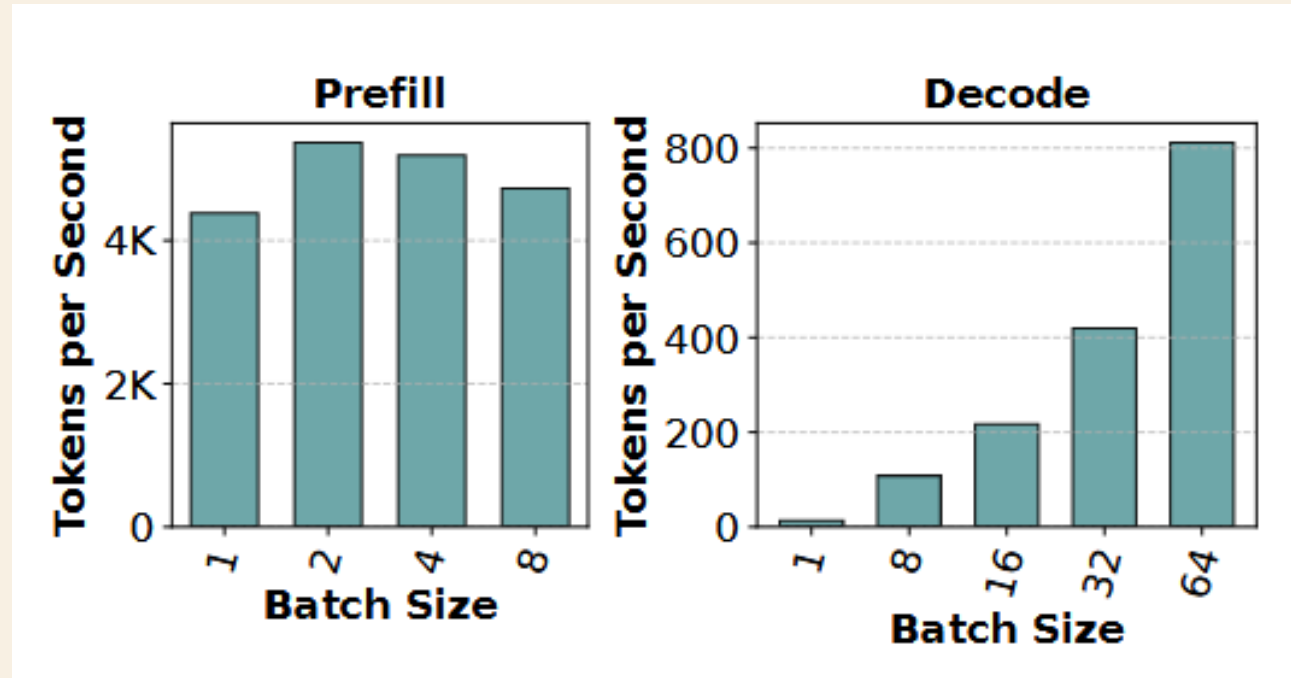
1. Background

(d) Another Challenge for Prefill-Prioritizing

- Background
 - Pipeline-parallelism (PP) used for cross-node LLM inference
- Issue: Pipeline bubbles waste GPU cycles
 - Caused by varying runtimes of prefill and decode micro-batches
- Impact: Degrades system throughput

2. Motivation

(a) Cost Analysis of Prefill and Decode

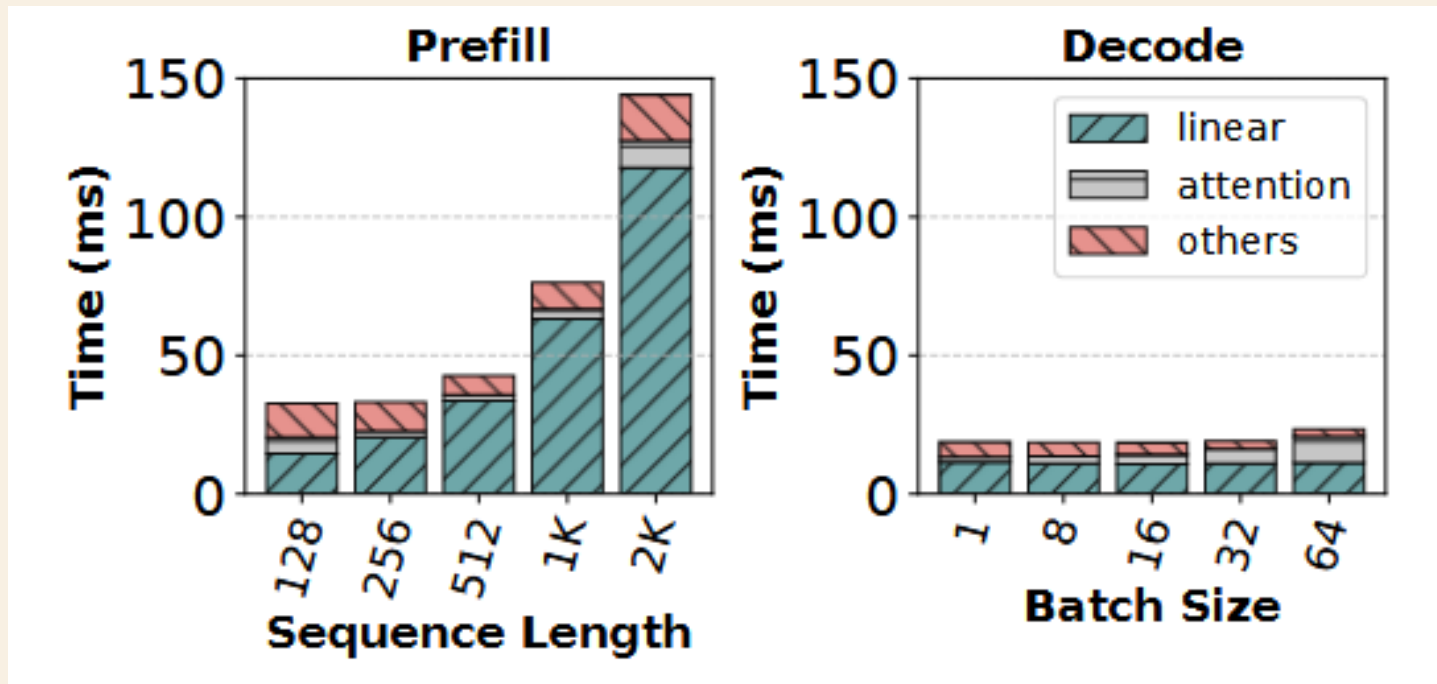


Takeaway-1

Prefill and Decode phase demonstrate contrasting behaviors wherein batching boosts decode phase throughput immensely but has little effect on prefill throughput.

2. Motivation

(a) Cost Analysis of Prefill and Decode



From the figure, we see that linear operators contribute to the majority of the runtime cost. Therefore, optimizing linear operators is important for improving LLM inference.
(Let's focus on linear operators later)

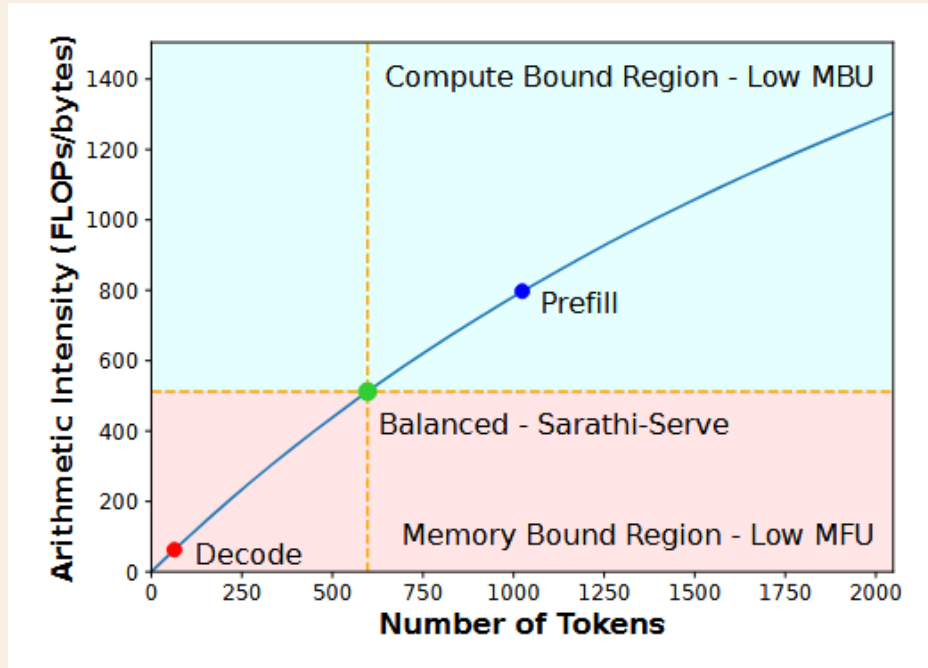
2. Motivation

(a) Cost Analysis of Prefill and Decode

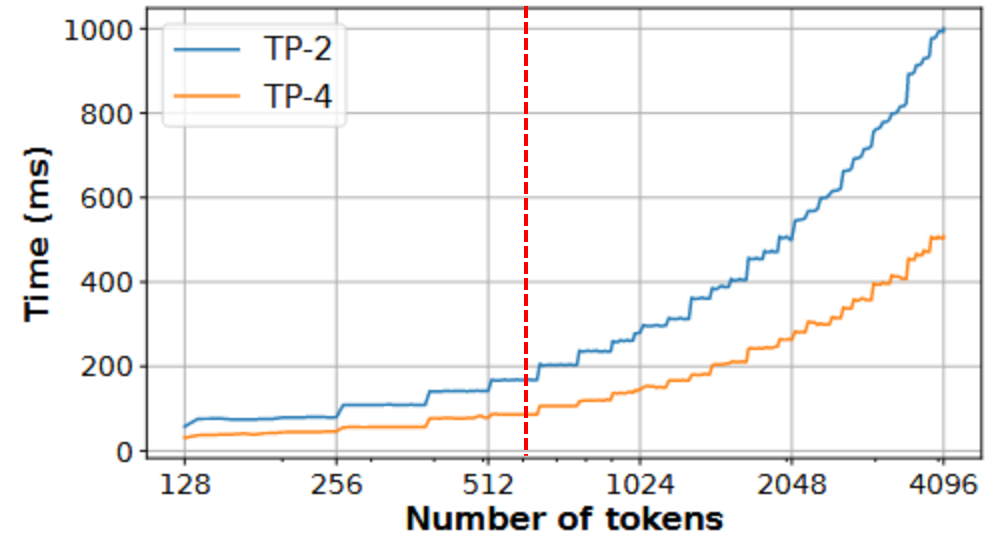
For an operation, Total Running Time = $\max(T_{\text{math}}, T_{\text{mem}})$
If $T_{\text{mem}} > T_{\text{math}}$: Operation is Memory-Bound
If $T_{\text{math}} > T_{\text{mem}}$: Operation is Compute-Bound
if $T_{\text{math}} = T_{\text{mem}}$, both compute and memory bandwidth utilization are maximized.

2. Motivation

(a) Cost Analysis of Prefill and Decode



Arithmetic intensity trend for LLaMA2-70B linear operations with different number of token running on four A100s.



Linear layer execution time as function of number of tokens in a batch for LLaMA2-70B on A100(s) with different tensor parallel degrees.

Takeaway-2

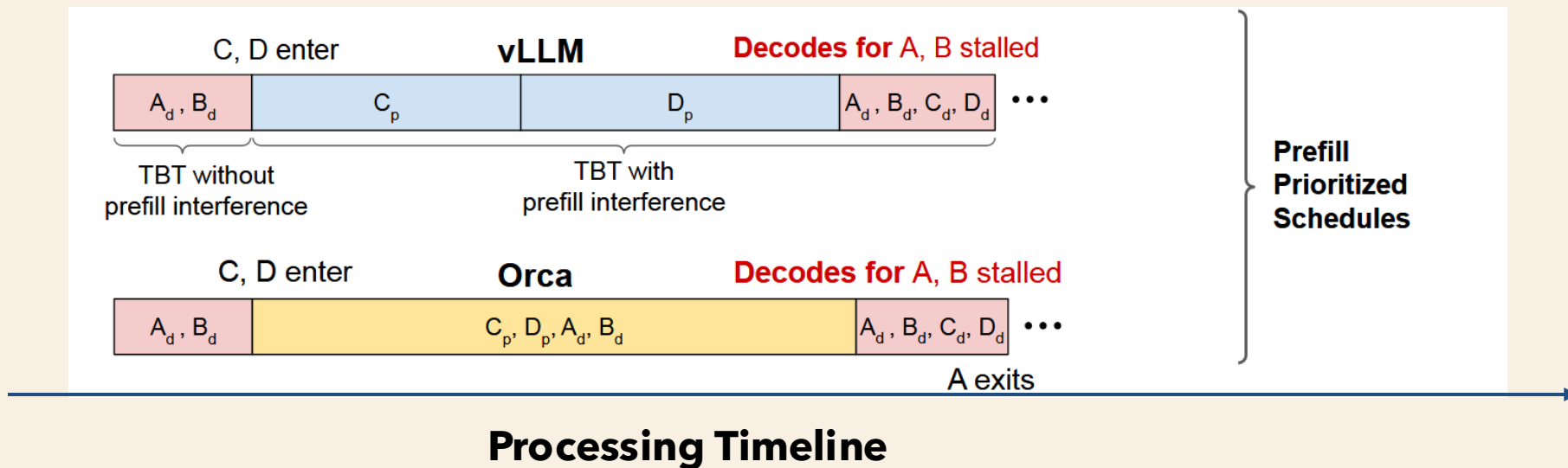
Decode batches operate in memory-bound regime leaving compute underutilized. This implies that more tokens can be processed along with a decode batch without significantly increasing its latency.

2. Motivation

(b) Throughput-Latency Trade-off

A, B: existing requests in decode

C, D: new requests



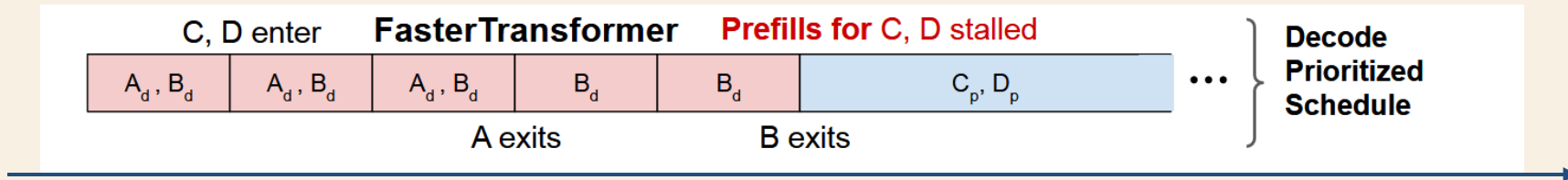
vLLM & Orca (Prefill-Prioritizing):

- Eagerly schedule prefills (C, D), pausing decodes (A, B)
- Result: Generation stalls (high TBT latency)

2. Motivation

(b) Throughput-Latency Trade-off

A, B: existing requests in decode
C, D: new requests



Processing Timeline

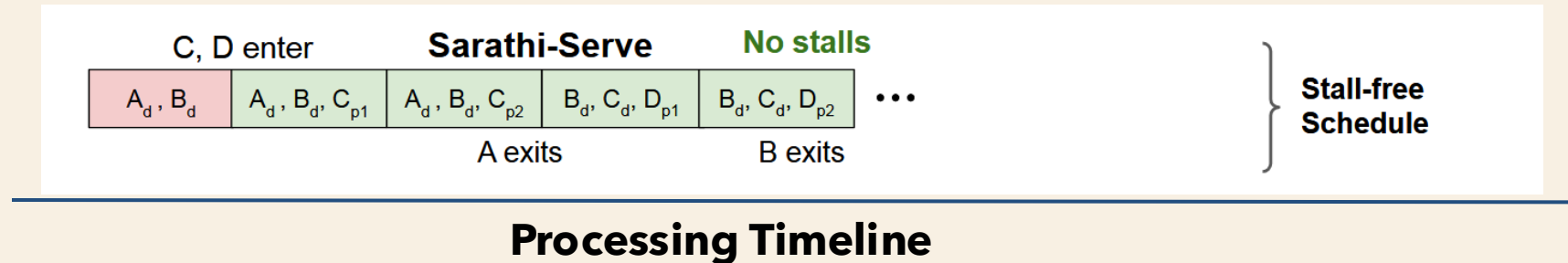
FasterTransformer (Decode-Prioritizing):

- Waits for decodes (A, B) to finish before prefills (C, D)
- Result: No stalls, but low throughput (small decode batch size)

2. Motivation

(b) Throughput-Latency Trade-off

A, B: existing requests in decode
C, D: new requests

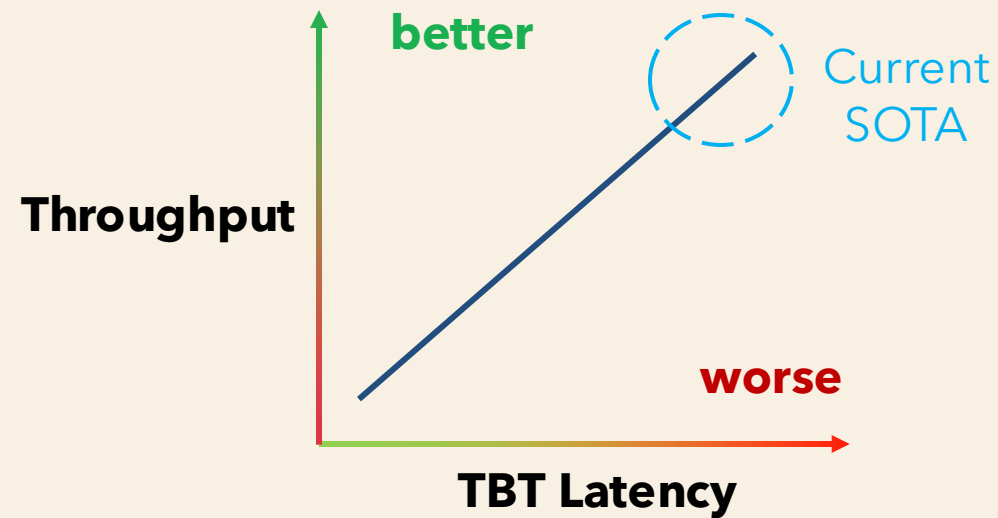


Proposed work Sarathi-Serve

- Stall-free scheduling: No pauses, high throughput

2. Motivation

(b) Throughput-Latency Trade-off



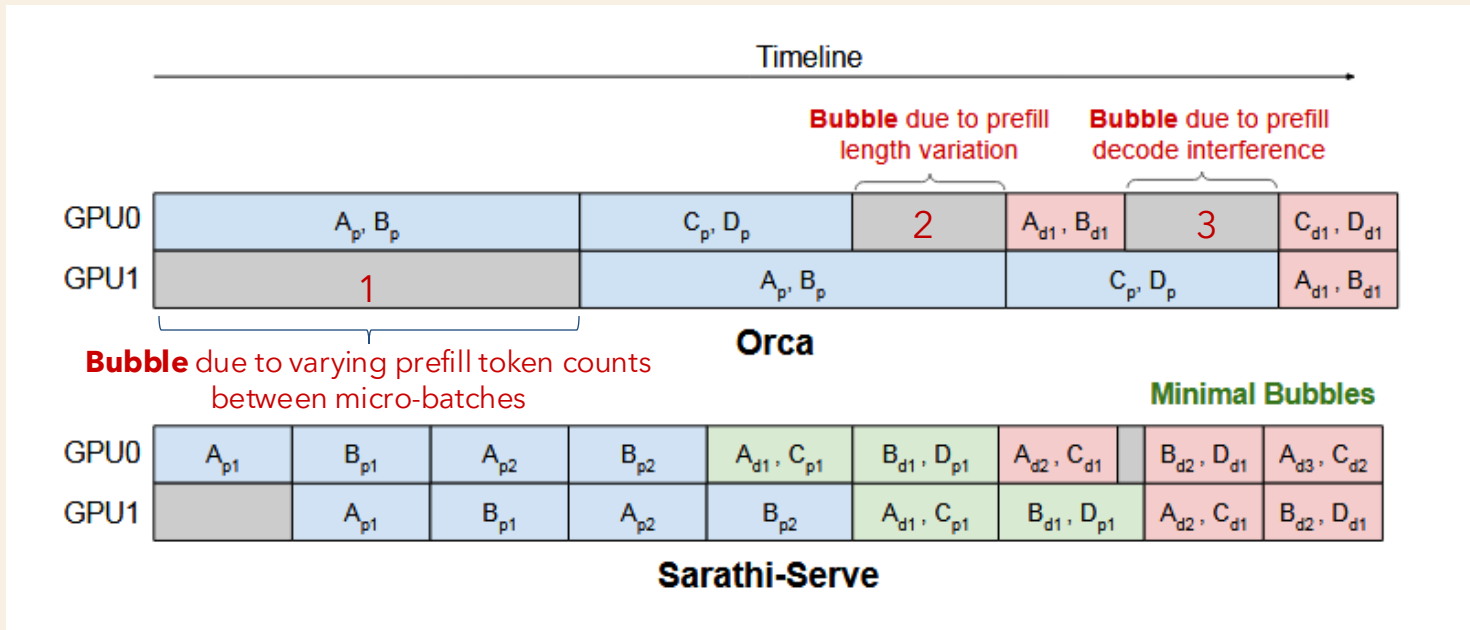
Takeaway-3

The interleaving of prefills and decodes involves **a trade-off between throughput and latency** for current LLM inference schedulers.

State-of-the-art systems today use prefill-prioritizing schedules that **trade TBT latency for high throughput**.

2. Motivation

(c) Pipeline Bubbles waste GPU Cycles



- A 2-way pipeline parallel iteration-level schedule in Orca across 4 requests (A,B,C,D) shows **the existence of pipeline bubbles due to non-uniform batch execution times.**

Takeaway-4

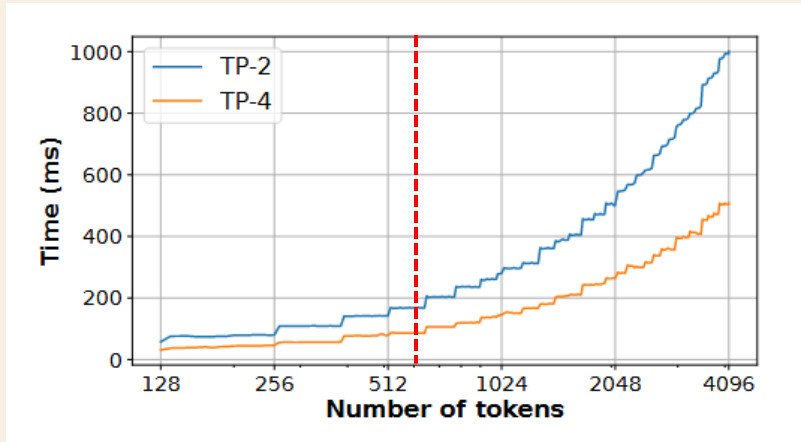
There can be a large variance in compute time of LLM iterations depending on composition of prefill- and decode-tokens in the batch. This can lead to significant bubbles when using pipeline-parallelism.



Sarathi-Serve is able to minimize these stalls by creating **uniform-compute** batches.

3. Method

(a) Chunked-prefills



Insight

A prefill request with **modest sequence length** can effectively saturate GPU compute

Dataset	Prompt Tokens			Output Tokens		
	Median	P90	Std.	Median	P90	Std.
<i>openchat_sharegpt4</i>	1730	5696	2088	415	834	101
<i>arxiv_summarization</i>	7059	12985	3638	208	371	265

Fact

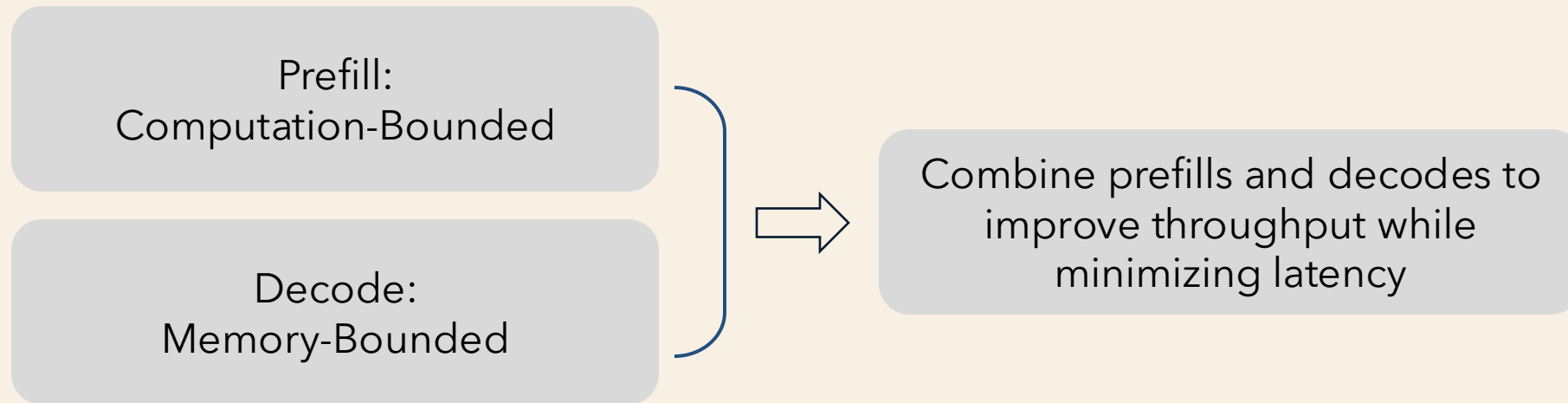
Prompts are often very long, and long prefills with decode iterations increase TBT latency.

Solution: split large prefills to small chunks

Break large prefill requests into smaller units of compute which are still large enough to saturate GPU compute.

3. Method

(b) Stall-free batching



Naïve Method in Orca & vLLM

- stall existing decodes to execute prefills

Sarathi-Serve: stall-free batching

- leverages the arithmetic intensity slack in decode iterations to execute prefills without delaying the execution of decode requests in the system.

3. Method

(b) Stall-free batching - Algorithm

```
1: Input:  $T_{max}$ , Application TBT SLO.  
2: Initialize  $token\_budget$ ,  $\tau \leftarrow compute\_token\_budget(T_{max})$   
3: Initialize  $batch\_num\_tokens$ ,  $n_t \leftarrow 0$   
4: Initialize current batch  $B \leftarrow \emptyset$   
5: while True do  
6:   for  $R$  in  $B$  do  
7:     if  $is\_prefill\_complete(R)$  then  
8:        $n_t \leftarrow n_t + 1$   
9:   for  $R$  in  $B$  do  
10:    if not  $is\_prefill\_complete(R)$  then  
11:       $c \leftarrow get\_next\_chunk\_size(R, \tau, n_t)$   
12:       $n_t \leftarrow n_t + c$   
13:    $R_{new} \leftarrow get\_next\_request()$   
14:   while  $can\_allocate\_request(R_{new}) \wedge n_t < \tau$  do  
15:      $c \leftarrow get\_next\_chunk\_size(R_{new}, \tau, n_t)$   
16:     if  $c > 0$  then  
17:        $n_t \leftarrow n_t + c$   
18:        $B \leftarrow R_{new}$   
19:     else  
20:       break  
21:  
22:    $process\_hybrid\_batch(B)$   
23:    $B \leftarrow filter\_finished\_requests(B)$   
24:    $n_t \leftarrow 0$ 
```

1. Calculates token budget based on user-specified SLO [Service Level Objective] (we will introduce later)

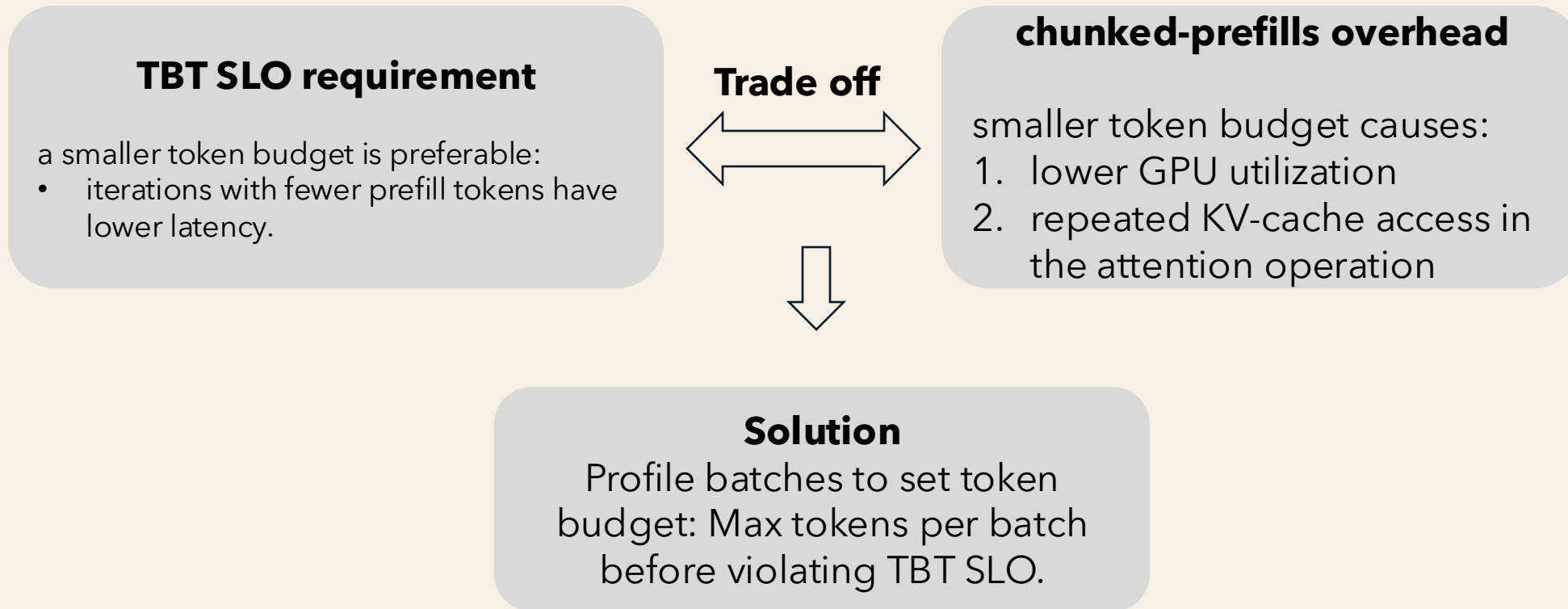
2. the batch is filled with ongoing decode tokens

3. include any partially completed prefill

4. Admit new requests within leftover token budget

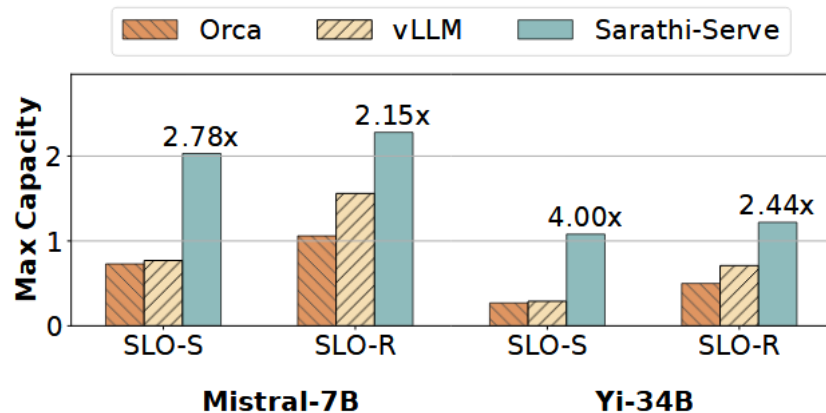
3. Method

(c) Determining Token Budget

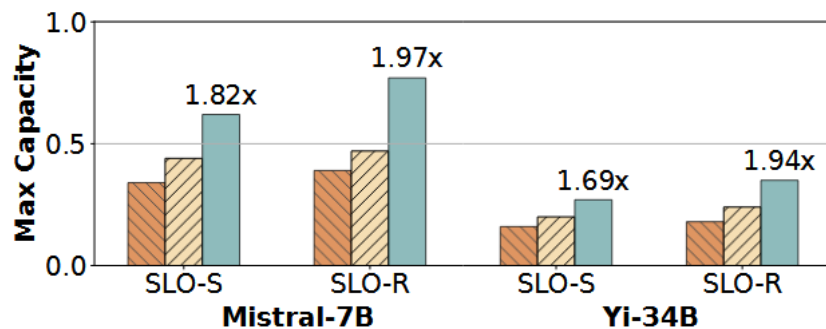


4. Evaluation

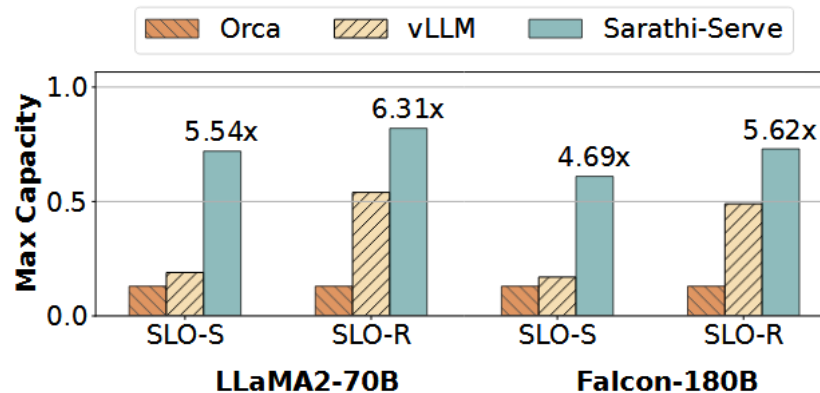
(a) Capacity



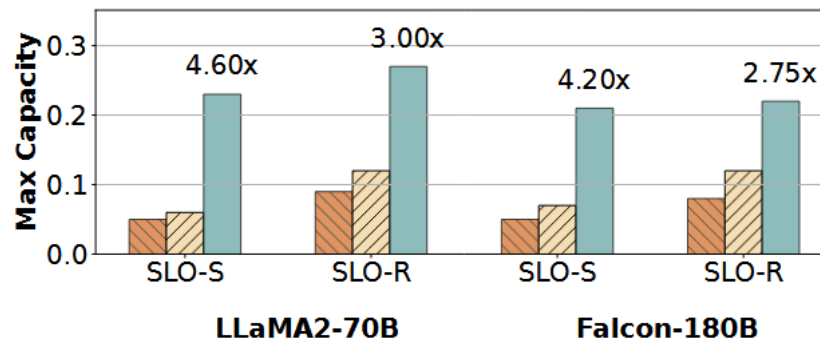
(a) Dataset: *openchat_sharegpt4*.



(b) Dataset: *arxiv_summarization*.



(a) Dataset: *openchat_sharegpt4*.



(b) Dataset: *arxiv_summarization*.

Note:

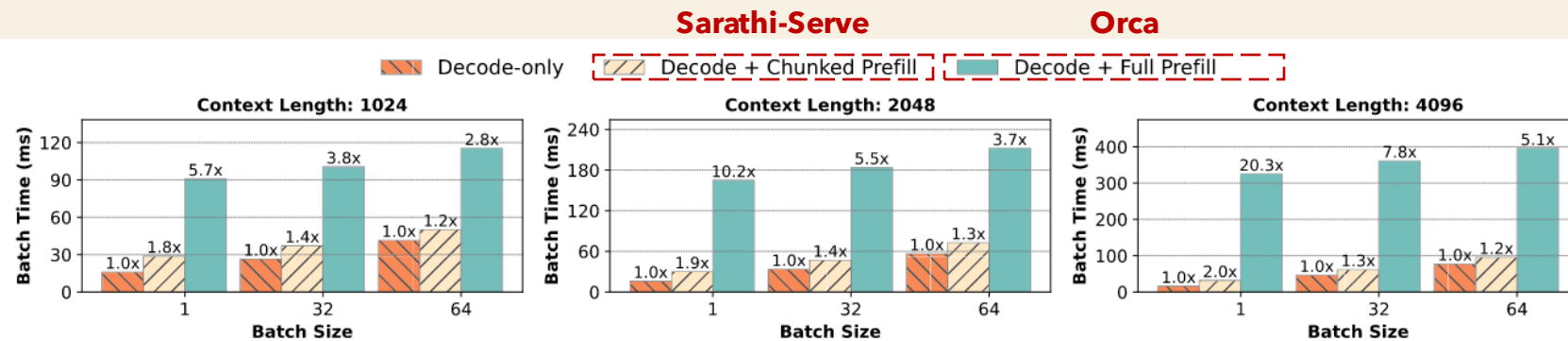
- SLO: Service Level Objective
- SLO-S: Strict SLO
- SLO-R: Relaxed SLO

Capacity
(in queries per second)
improves **up to 6x**

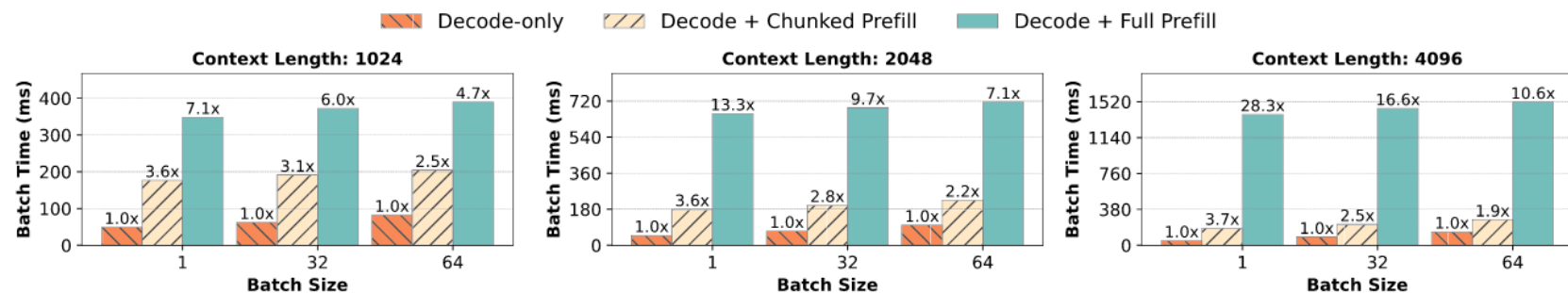
Capacity evaluation across different models and datasets

4. Evaluation

(b) Latency



(a) Mistral-7B on one A100s with token budget of 256.



(b) LLaMA2-70B on four A100s with token budget of 512.

Note:

- Decode + Full Prefill the hybrid batching of Orca
- Decode + Chunked Prefill Proposed work

Sarathi-Serve processes prefill tokens with much lower impact on the latency of decodes.

4. Evaluation

(c) Making Pipeline Parallel Viable

Sarathi-Serve optimizes pipeline parallelism by creating hybrid batches with uniform computational demands.

Key Benefits:

- Reduces pipeline bubbles, improving GPU utilization.
- Enables efficient operation in multi-node deployments over standard Ethernet connections.



Thank you!