

## Agenda:

- Batching scheme we left off w/ last time
- Iteration-level scheduling ("continuous Batching")
- Different Inference Perf metrics: TTFT, TBT
- scheduling prioritization options (metric-dependent)
  - ↳ Decode vs. Prefill prioritizing
- chunked prefill

## \*RESOURCES:

- Any Scale Blog Post : <https://www.anyscale.com/blog/continuous-batching-llm-inference>
- ORCA OSDI Paper : <https://www.usenix.org/conference/osdi22/presentation/yu>
- Sarathi Serve Paper <https://www.usenix.org/system/files/osdi24-agrawal.pdf>
- vLLM paper: <https://arxiv.org/pdf/2309.06180>

## BATCHING SCHEME WE LEFT OFF WITH.

- Assumption: we are serving requests for a SINGLE model
- 7B model in fp16 takes up:  
$$7 \times 10^{12} \cdot 2 \approx 14 \text{ GB}$$
- Approx. remaining of GPU memory (40 GB) can be used by KV cache
- LLM serving is quite memory bound (wait to keep model weights loaded) - and serve

multiple (a batch of requests)

\* note - we still don't want to waste  
compute and recompute KV cache on each decode, so:

- when we serve a request, we want to  
make sure we can fit its entire KV

cache in GPU memory

↳ otherwise recompute would be  
too much

- maximum batch size we can fit is therefore  
determined by:

1. overhead per token in KV cache
2. max sequence length any request

(can have

$$\hookrightarrow \text{mem. for KV cache} = \frac{\text{mem. for KV cache}}{\text{mem. for KV cache}} = \frac{26 \text{ GB}}{26 \text{ GB}} = 1 \text{ GB}$$

# tokens that can fit =  $\frac{1 \text{ GB}}{0.0002 \text{ GB}} = 5000$

$$\text{max batchsize} = \frac{5000}{\text{max seq length}} = 128$$

"static" batching scheme we touched on:

1. wait until  $\text{batch size} > \text{requests admitted}$   
to the system

2. For all batches, run prompt-processing  
step \* to produce 1<sup>st</sup> token \*

3. Proceed one "iteration" at a time, where we run 1 decode step across all our examples in the batch to predict the next token in each sequence.

↳ if at any point, a sequence ends,

don't do anything w/ that example

until ~~all~~ batch size \* examples finish w/o

\*clarification: paged attn memory management scheme from last time, \*effective\* batch

size is lower (due to fragmentation, not all memory is available)

- say we have determined some batch size we know will fit, why is above scheme suboptimal?

\*not all requests will use max # of tokens!!

<https://www.anyscale.com/blog/continuous-batching-llm-inference>

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$				
$S_2$	$S_2$	$S_2$					
$S_3$	$S_3$	$S_3$	$S_3$				
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$			

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

white boxes =  
idle GPU cycles

\*how often is this a problem? likelihood of white squares?)

- ↳ depends on whether input sequences run together have \*similar/same size\*
- ↳ if all reqs produce 1 token (e.g., classification)
- ↳ if all prompts are same size

- in workloads like chatbots, there is a HIGH likelihood that sequences in some batch have diff. sizes

## ITERATION-LEVEL SCHEDULING / CONTINUOUS BATCHING

- introduced by orca (OSDI 2022) - vllm implements this now

→ higher GPU utilization

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$				
$S_2$	$S_2$	$S_2$					
$S_3$	$S_3$	$S_3$	$S_3$				
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$			

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

+ handles late-arriving / early-finish req more efficiently.

\* key idea: let new requests start immediately if there is space

- ↳ at each iteration - check if a seq. has ended and whether to admit a new request

# INFERENCE PERF METRICS

\* TTFT - time to 1<sup>st</sup> token → how long does prefill take

\* TBT - time between tokens → how long does it take to generate each subsequent token?

- there are 2 parts to serving (token prediction tasks):

→ can be run in parallel (in one step)

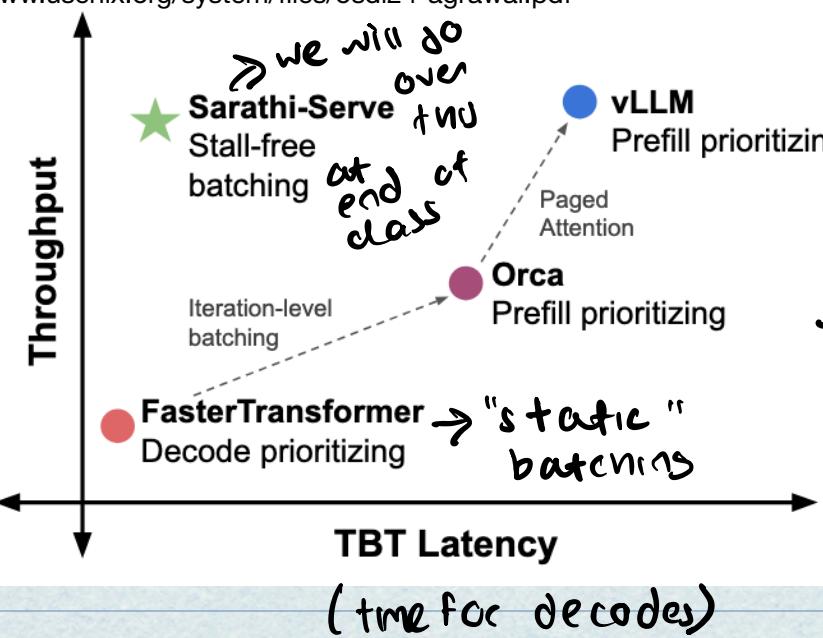
1. Prefill: prompt processing to get 1<sup>st</sup> token
2. Decode: each subsequent iteration after

\* prioritizing prefill will get you lower TTFT, <sup>high throughput</sup> (faster before)

\* prioritizing ongoing decodes will get you lower TBT, also will compromise throughput

\* what does it mean to prioritize? Do these things come at a cost / is there a fundamental tradeoff?

<https://www.usenix.org/system/files/osdi24-agrawal.pdf>



w/ iteration-level scheduling - (continuous batches), we are:  
- prioritizing prefill lets you fill in the empty spaces in batch, increasing throughput

\* why does prefill prioritizing sacrifice

## TBT / decode latency?

\* let's look more closely at "static" batching

and "continuousbatching" scheduling algorithm:

request-level

• STATIC BATCHING: optimizes TBT, but

wastes GPU utilization  
(Faster Transformer)

<https://www.usenix.org/system/files/osdi24-agrawal.pdf>

```
1: Initialize current batch  $B \leftarrow \emptyset$ 
2: while True do
3:   if  $B = \emptyset$  then → only when batch is empty, do we admit a new request
4:      $R_{new} \leftarrow \text{get\_next\_request()}$ 
5:     while can_allocate_request( $R_{new}$ ) do
6:        $B \leftarrow B + R_{new}$  ↪ if "fits" in batch
7:        $R_{new} \leftarrow \text{get\_next\_request()}$ 
8:     prefill( $B$ )
9:   else → if  $B$  has ANYTHING at all, decode( $B$ ) finish ongoing decodes
10:    decode( $B$ )
11:     $B \leftarrow \text{filter\_finished\_requests}(B) \rightarrow \text{check if } B \neq \emptyset$ 
```

→ DECODE

PRIORITYING =

only admits new req if current decoder can do all done

As we

saw before, lead to low throughput!

## What abt continuous batching?

```
1: Initialize current batch  $B \leftarrow \emptyset$ 
2: while True do → at each iter
3:    $B_{new} \leftarrow \emptyset$  → initialize set of new reqs to add
4:    $R_{new} \leftarrow \text{get\_next\_request()}$ 
5:   while can_allocate_request( $R_{new}$ ) do
6:      $B_{new} \leftarrow B_{new} + R_{new}$  ↪ checks there is enough space for prefill of  $R_{new}$ 
7:      $R_{new} \leftarrow \text{get\_next\_request()}$ 
8:   if  $B_{new} \neq \emptyset$  then
9:     prefill( $B_{new}$ )
10:     $B \leftarrow B + B_{new}$  ↪ run all NEW prefills
11:   else
12:     decode( $B$ ) → only run decodes once we've finished new prefills
13:      $B \leftarrow \text{filter\_finished\_requests}(B)$ 
```

PREFILL PRIORITYZ

INT:  
at each iteration,  
run all  
prefills we  
can

↳ this makes TBT longer because we can have "generative stalls":

- before executing line 12 (decode),

we see if we have new prefill to run and run those first.

→ so there is a fundamental tradeoff b/w TBT latency and Throughput in existing schedulers.

- at each iteration, we have the choice to:

① run existing requests decode cycles -

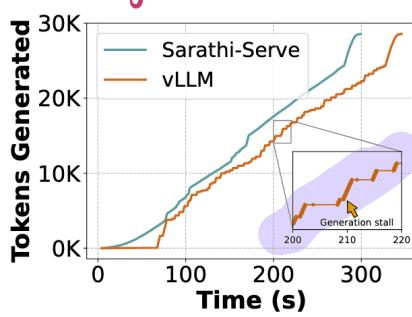
prioritizes TBT at expense of TBT

② run prefills - which will "stall ongoing decodes"

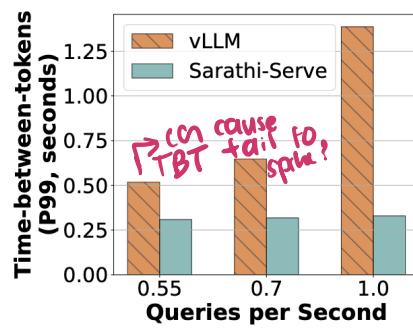
- what does the stalling look like?

- consider vLLM in the graph below

↳ when decodes are interrupted, there will be a "stall" in new tokens generated



(a) Generation stall.



(b) High tail latency.

→ SARATHI SERVE / CHUNKED PREFILLS -

→ paper that analyzes existing inference

system and shows mismatch in input / latency

→ proposes 2 techniques to solve:

1) chunked prefills

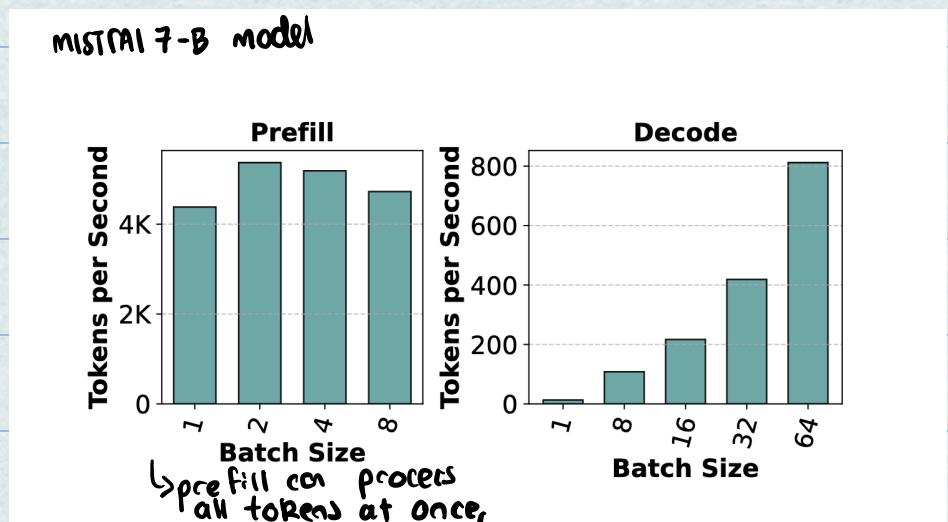
2) stall-free batching

- we will try to build up these techniques

• 3 (related) observations:

→ OBSERVATION 1:

- batching boosts throughput for decode, but not so much for prefill



To get more tput in decode, we need to run more decodes.

Observation 2:

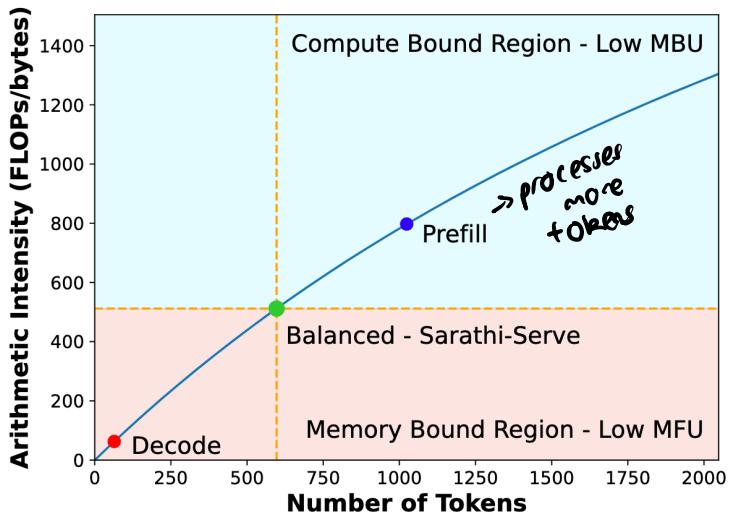
- Decode has low AI (memory bound)  
low compute utilization

- Prefill has high AI (compute bound,  
low memory bw utilization)

- Neither at ridge point

llam2-708, running on 4A100s

these schemes use Pipeline parallelism - more thoughts in paper



- this graph JUST focuses on linear layers so cost of fetching layer weights vs. doing computation on linear layers  
(another graph in paper shows linear ops as largest timespent)

- when we are in memory bound regions (e.g., decode) - we can do more compute for same data fetched for free

↳ so \*more tokens can be processed along w/ decodes w/o ruining TBT for ongoing decodes

### - OBSERVATION 3:

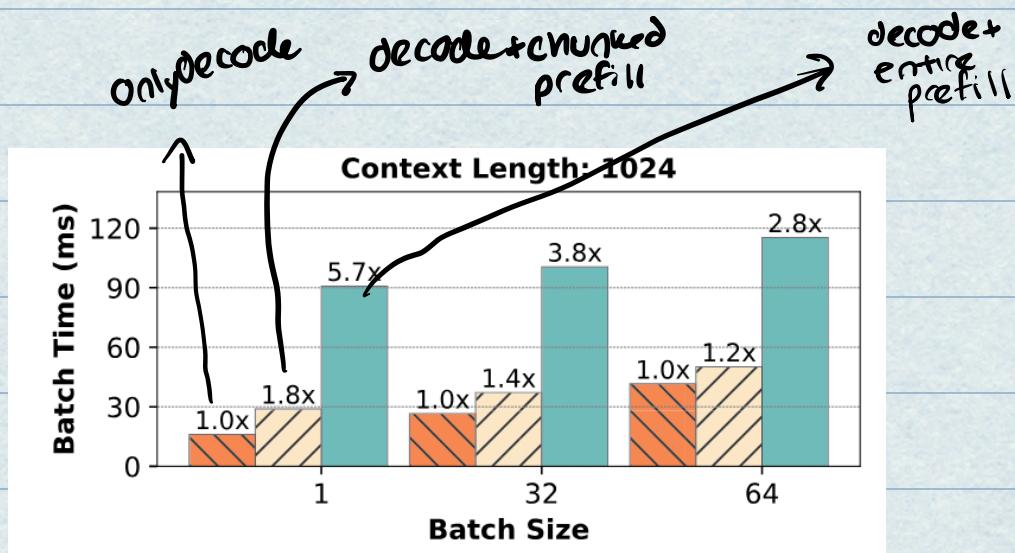
- as we discussed before, existing systems choose at each iteration to run new prefills (to incr. throughput) OR finish existing requests

- running prefills can stall decoders

## IDEA 1: chunked prefill

- when we have extra capacity from decodes, we'd like to do more work (decode has low AI, and w memory bound)
- if we try to add in a ENTIRE prefill - and we schedule a single iteration w/ existing decodes and a full prefill, latency will go up
- KEY IDEA: lets CHUNK the prefill.

↳ and execute small blocks at a time



\* we can execute a small portion of prefill along w/ decodes, to increase throughput w/o affecting latency that much

## IDEA2: STALL -Free Batches

↳ w/ chunked prefills, how do we schedule each iteration?

```
1: Input:  $T_{max}$ , Application TBT SLO.
2: Initialize token_budget,  $\tau \leftarrow \text{compute\_token\_buget}(T_{max})$ 
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:        $n_t \leftarrow n_t + 1 \rightarrow 1^{\text{st}} \text{ decode}$ 
9:   for  $R$  in  $B$  do
10:    if not is_prefill_complete( $R$ ) then
11:       $c \leftarrow \text{get\_next\_chunk\_size}(R, \tau, n_t)$ 
12:       $n_t \leftarrow n_t + c \rightarrow \text{add next prefills for existing requests}$ 
13:      |  $R_{new} \leftarrow \text{get\_next\_request()}$ 
14:      while can_allocate_request( $R_{new}) \wedge n_t < \tau$  do
15:         $c \leftarrow \text{get\_next\_chunk\_size}(R_{new}, \tau, n_t)$ 
16:        if  $c > 0$  then
17:           $n_t \leftarrow n_t + c \rightarrow \text{only add new prefills chunks after}$ 
18:           $B \leftarrow R_{new}$ 
19:        else
20:          break
21:
22:      process_hybrid_batch( $B$ )
23:       $B \leftarrow \text{filter\_finished\_requests}(B)$ 
24:       $n_t \leftarrow 0$ 
```

paper has more details on calculations than but we can think of as additional capacity until we become compute bound

→ prioritizes ongoing decodes: Always schedule these 1<sup>st</sup> to run

→ THEN schedule next chunk of existing batches  
→ THEN only admit new requests

\*process hybrid batch.

- can execute batch of decodes w/  
chunked prefills at some time
- previous systems were only executing  
a batch of decodes OR batch of  
prefills
- paper has ablations on hybrid batch  
technique

#### \* some related themes of work:

- because prefill and decode behave so differently,  
why not put them on dif. GPUs? schedules  
will be easier
  - ↳ can send KV cache in between
- lets cache KV cache in CPU memory  
(when it needs to be reused) - e.g. document  
and transfer it
  - ↳ but compress it in a smart  
way so its not so large to send