

# SGLang FLPM

Yichuan Wang <[yichuan\\_wang@berkeley.edu](mailto:yichuan_wang@berkeley.edu)> & LMSYS  
Arxiv : <https://arxiv.org/pdf/2501.14312>

# Locality-aware Fair Scheduling in LLM Serving

Shiyi Cao<sup>\*1 2</sup>, Yichuan Wang<sup>\*1 2</sup>, Ziming Mao<sup>1</sup>, Pin-Lun Hsu<sup>2</sup>, Liangsheng Yin<sup>1 2</sup>, Tian Xia<sup>1</sup>, Dacheng Li<sup>1 2</sup>, Shu Liu<sup>1</sup>, Yineng Zhang<sup>2</sup>,  
Yang Zhou<sup>1</sup>, Ying Sheng<sup>2</sup>, Joseph E. Gonzalez<sup>1</sup>, Ion Stoica<sup>1</sup>

<sup>1</sup>UC Berkeley <sup>2</sup>LMSYS \*indicates equal contribution.

# Outlines

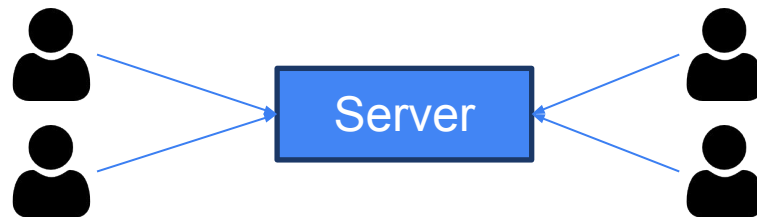
1. Background and motivation
2. DLPM and its distribute version D2LPM
3. Evaluation
4. Future Work

## **→ 1. Background and motivation**

# Background:

## 1. What is fairness in LLM serving?

- One server serves multiple users concurrently.
- The server's capacity is limited.

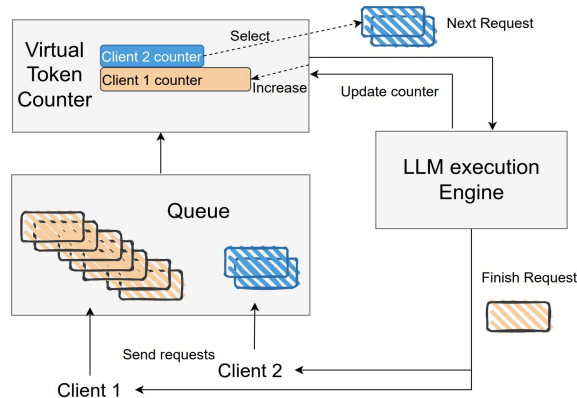


# VTC scheduler(Done by Ying)



Simple algorithm:

- ❑ Virtual counter for each user
- ❑ Update and schedule at token granularity
- ❑ Update according to the cost per token



Virtual counters:

100

80

120

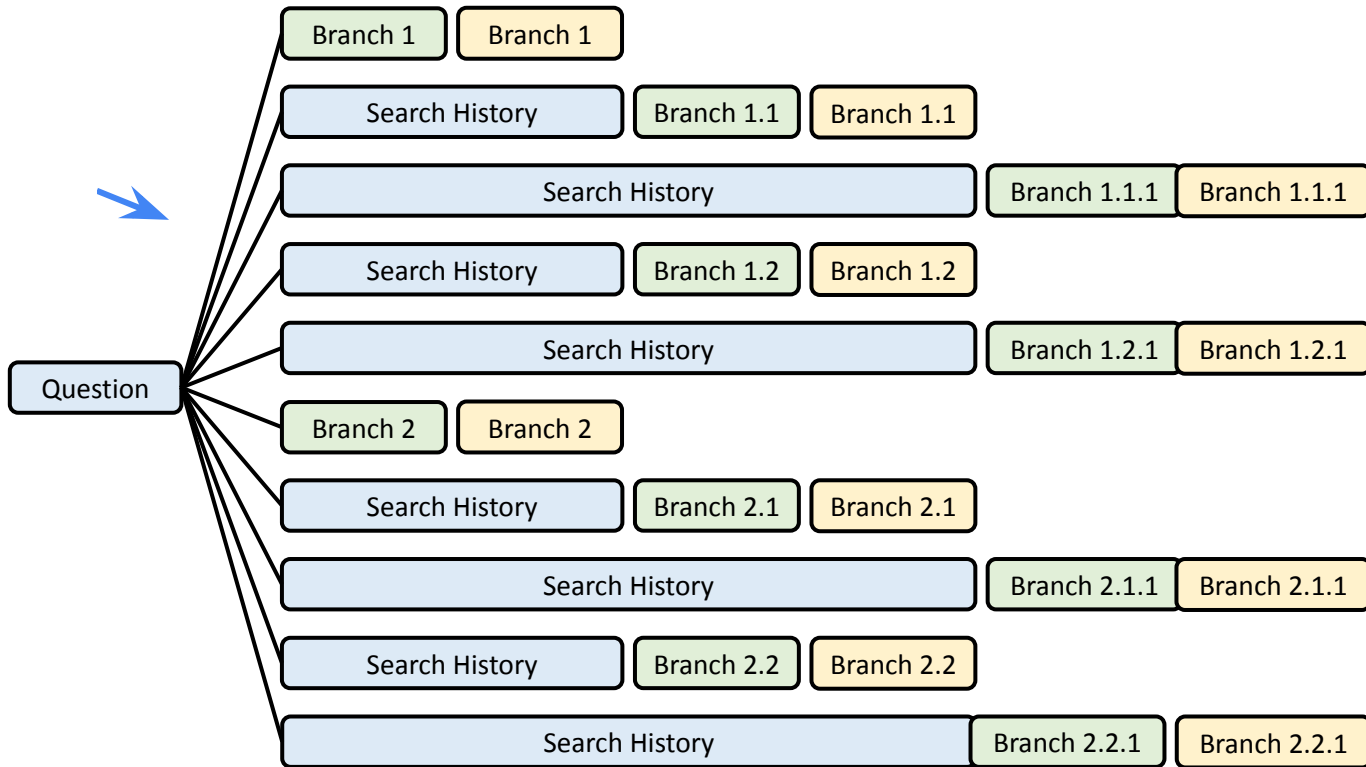
← New request

# Background:

## 2. Structure in LLM Calls and Prefix Cache

Prefix Cache

Tree of Thought



# Background:

Cache-aware schedule: LPM (Longest Prefix Match)



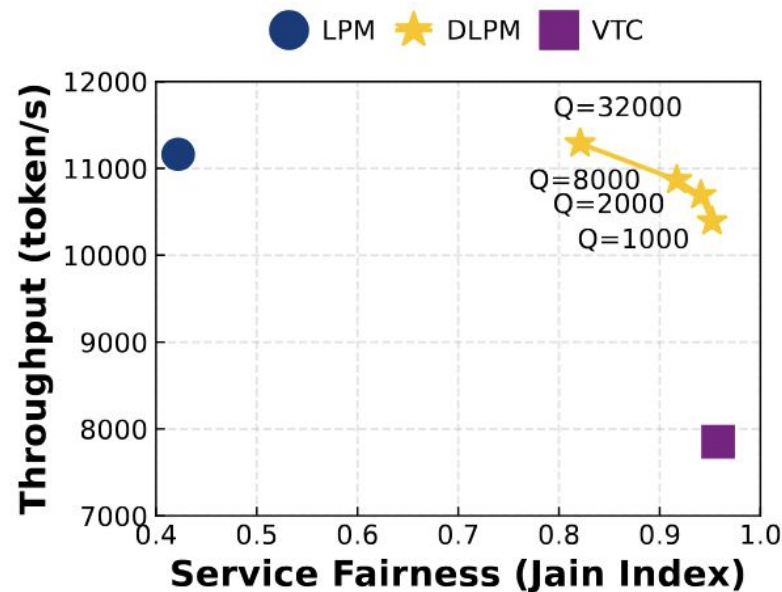
Even more severe fairness issue



# Motivation:

**LPM:** Prioritize locality but ignore fairness

**VTC:** Ensure fairness but ignore efficiency

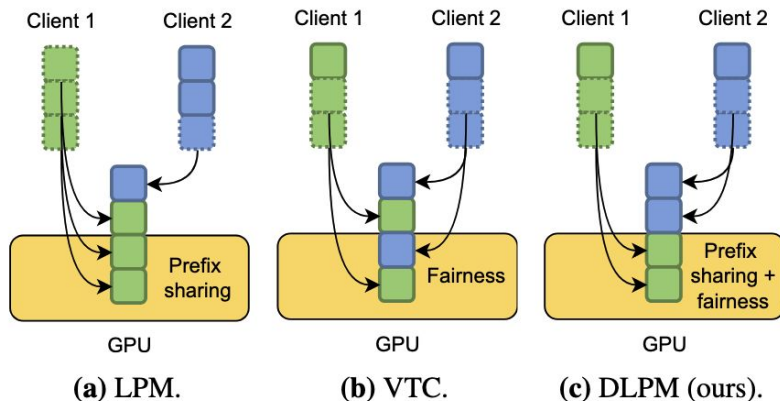


**Figure 1:** DLPM achieves a *new Pareto frontier* considering locality and fairness in LLM serving.  $Q$  is a hyper-parameter in DLPM, indicating how much we relax the fairness bound of DLPM. Results are obtained in a single A10 GPU.

\* Jain fairness index:  $J(x_1, x_2, \dots, x_n) = \frac{(\sum_{i=1}^n x_i)^2}{n \sum_{i=1}^n x_i^2}$

## **→ 2. DLPM and its distribute version D2LPM**

# DLPM



## Algorithm 1 Deficit Longest Prefix Match (DLPM)

---

```

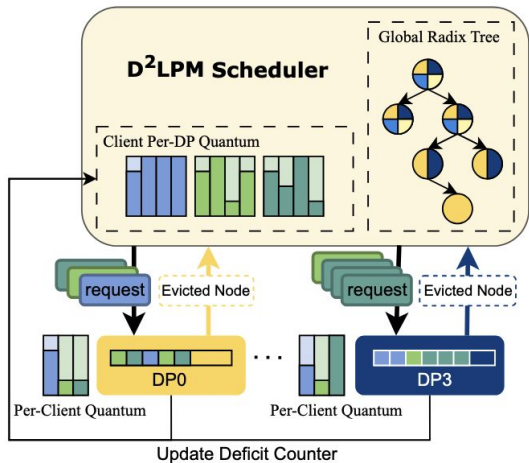
1: let  $l$  denotes the client list
2: let  $B$  denotes current running batch
3: function CHECKREFILL( $l$ ,  $Queue$ )
4:   for all  $i \in \{client(r) \mid r \in Queue\}$  do
5:     if  $q_i > 0$  then return
6:   for all  $i \in l$  do
7:     if  $q_i \leq 0$  then  $q_i \leftarrow q_i + Q^u$ 
8: end function
9:  $\triangleright$  with monitoring stream:
10: while True do
11:   if new request  $r$  from client  $i$  arrived then
12:     if  $i \notin l$  then  $q_i \leftarrow 0$ ,  $l \leftarrow l + u$ 
13:      $Queue \leftarrow Queue + r$ 
14:  $\triangleright$  with execution stream 1:
15: while True do
16:    $Queue \leftarrow \text{SORTBYPREFIX}(Queue)$ 
17:   while not  $Queue.empty()$  do
18:     for each  $r \in Queue$  do
19:        $i \leftarrow client(r)$ 
20:       if  $q_i \leq 0$  then CHECKREFILL( $l$ ,  $Queue$ )
21:       if  $q_i > 0$  and CANADD( $r$ ) then
22:          $B \leftarrow B + r$ 
23:          $q_i \leftarrow q_i - w_e \cdot extend\_length(r)$ 
24:          $Queue \leftarrow Queue - r$ 
25:   FORWARDSTEP( $B$ )
26:    $q_i \leftarrow q_i - w_q \cdot |\{r \mid client(r) = i, r \in B\}|$ 
27:    $B \leftarrow filter\_finished\_requests(B)$ 

```

---



# D2LPM



## Algorithm 2 D<sup>2</sup>LPM Scheduling

---

```

1: let  $s_w$  denotes the current queue size of worker  $w$ .
2:  $W \leftarrow \text{GETWORKERS}()$ ,  $R \leftarrow \text{INITRADIXTREE}(|W|)$ 
3: function SELECTWORKER( $G, i$ )
4:    $G_{\text{avail}} \leftarrow \{w \mid q_{i,w} > 0\}$ 
5:   while  $G_{\text{avail}} == \emptyset$  do
6:     for all  $w \in W$  do  $q_{i,w} \leftarrow q_{i,w} + Q^w$ 
7:    $G_{\text{cand}} \leftarrow G \cap G_{\text{avail}}$ 
8:   if  $G_{\text{cand}} == \emptyset$  then return  $\arg \min_{w \in G_{\text{avail}}} s_w$ 
9:   return  $\arg \min_{w \in G_{\text{cand}}} s_w$ 
10: end function
11:  $\triangleright$  with concurrent stream 1:
12: while True do
13:   if new request  $r$  from client  $i$  arrived then
14:      $G \leftarrow R.\text{LONGESTMATCHWORKERS}(r)$ 
15:      $w \leftarrow \text{SELECTWORKER}(G, \text{client}(r))$ 
16:     DISPATCH( $w, r$ )
17:      $q_{i,w} \leftarrow q_{i,w} - w_e \cdot r.\text{input\_tokens}$ 
18:      $s_w \leftarrow s_w + 1$ 
19:      $R.\text{INSERT}(r.\text{input\_tokens}, w)$ 
20:  $\triangleright$  with concurrent stream 2:
21: while True do
22:   if request  $r$  from client  $i$  has finished at worker  $w$  then
23:      $q_{i,w} \leftarrow q_{i,w} - w_q \cdot r.\text{output\_tokens}$ 
24:      $s_w \leftarrow s_w - 1$ 
25:  $\triangleright$  with concurrent stream 3:
26: while True do
27:   if prefix  $P$  has been evicted at worker  $w$  then
28:      $R.\text{EVICT}(P, w)$ 

```

---

## ➔ Theoretical Fairness Bound

- Single GPU:

$$|W_f(t_1, t_2) - W_g(t_1, t_2)| \leq 2 \cdot (U + Q^u), \text{ where } U = w_e \cdot L_{input} + w_q$$

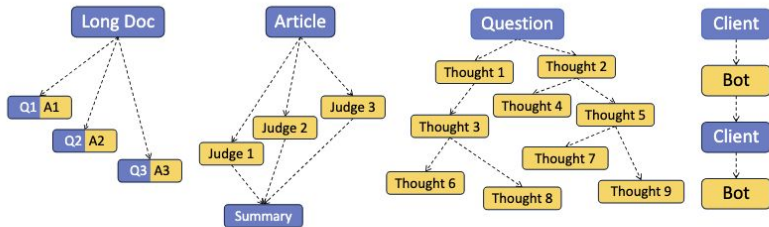
- Multiple GPU

$$|W_f(t_1, t_2) - W_g(t_1, t_2)| \leq 2 \cdot |W| \cdot (U + Q^u)$$

## **→ 3. Evaluation**

# → Evaluation Setup:

## ◆ Workload



(a) Long-Context QA (b) LLM-as-a-Judge (c) Tree-of-Thoughts (d) Multi-Turn Chat

## ◆ Eval Config

Workload	Detailed Behavior
<i>S1: More Requests</i>	
Long-Context QA	👤: Higher req rate
Tree-of-Thoughts	👤: Trees of 4 branches (340 req per tree) 👤: Trees of 2 branches (30 req per tree)
LLM-as-Judge	👤: Evaluation with 16 dimensions 👤: Evaluation with 2 dimensions
<i>S2: Longer Prefix</i>	
Long-Context QA	👤: 2× longer input documents
Tree-of-Thoughts	👤: 10× longer input questions
LLM-as-Judge	👤: Extra 600 tokens before each article

## ◆ Model

- Llama-3.2-3B on A100 80G
- Llama-3.1-8B on A10G

## ◆ Existing Scheduler

- VTC
- LPM (Round Robin+LPM in DP setting)
- Preble (locality aware router)

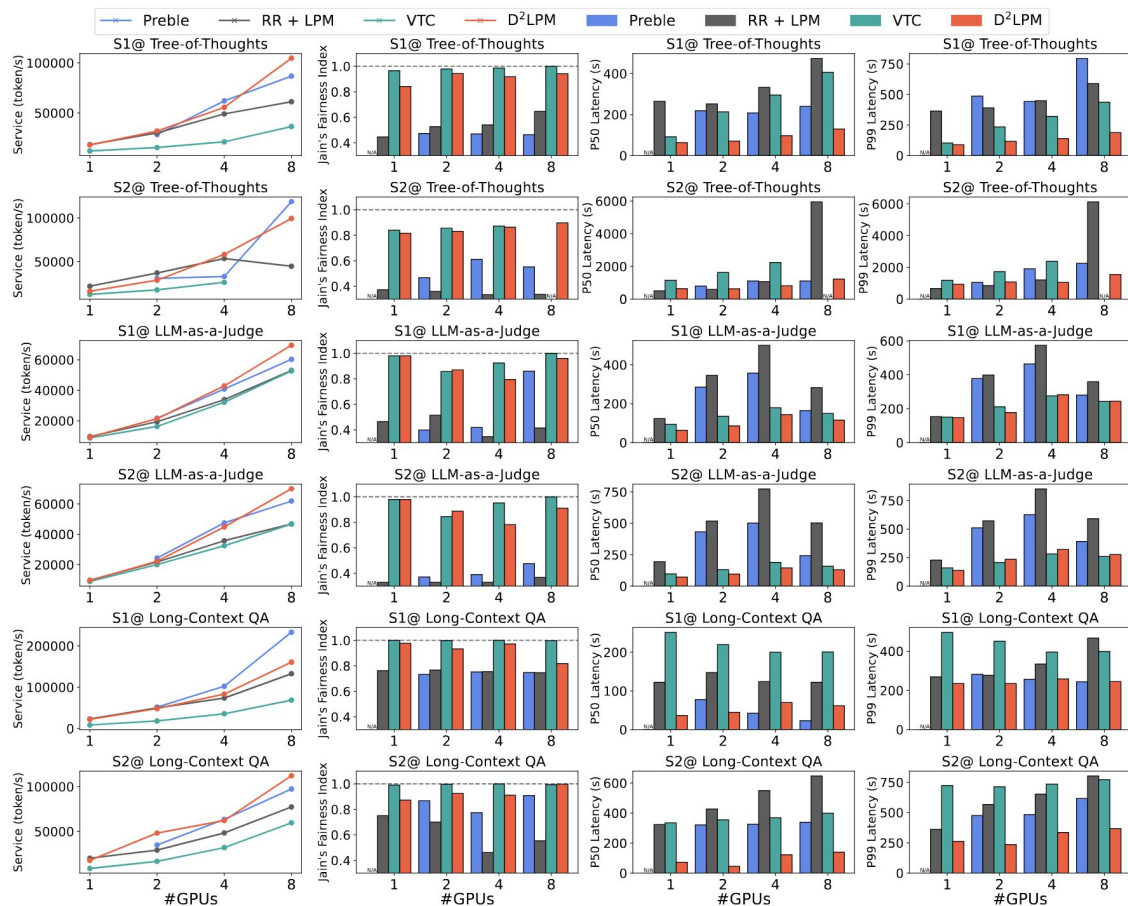
➔ Main evaluation:

◆ Throughput

◆ Fairness

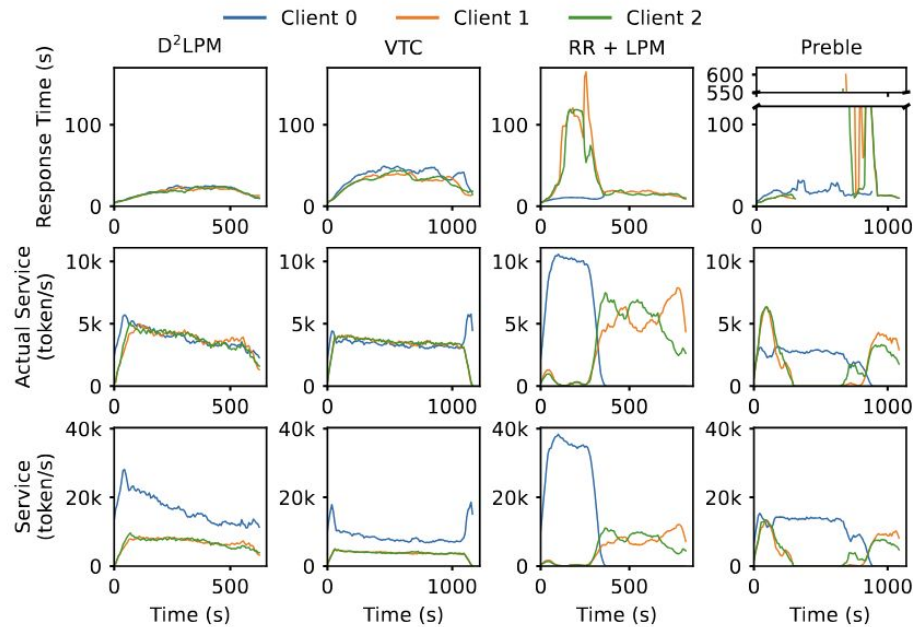
◆ Victim User

Latency



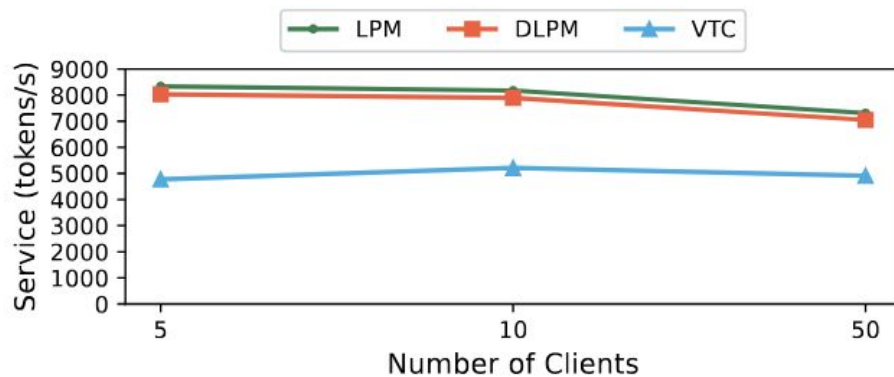


→ **Case Study:**  
◆ **4 GPU data  
parallelism**

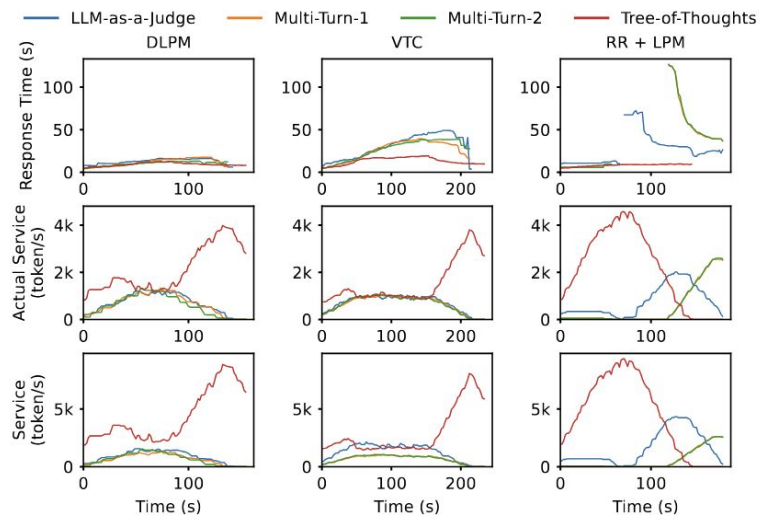


**Figure 9:** Fairness and performance visualization of different schedulers on Tree-of-Thoughts workloads with  $D = 4$  (3B model + 4 A10G GPUs). The maximum value on the X-axis represents the end-to-end completion time for each scheduler. The actual service is calculated using the cost function defined in §3, which is a weighted sum of the number of extend tokens and the number of output tokens.

➔ **Scalability:**



➔ **Mix of Workload:**



## **→ 4. Conclusion and Future Work**



## Conclusion

- ◆ We introduce the first locality-aware fair scheduling algorithm, Deficit Longest Prefix Match (DLPM), which can maintain a high degree of prefix locality with a fairness guarantee. And we can easily extend it to distributed version.



## Future Work

- ◆ Gradually upstream to SGLang codebase
- ◆ More follow up directions:
  - Share Prefix among users
  - Within user request level fairness guarantee
  - Fairness issue when considering data dependency in LLM workload

**Thank U**  
**Welcome to join our [Slack](#) and use [SGLang](#)!**