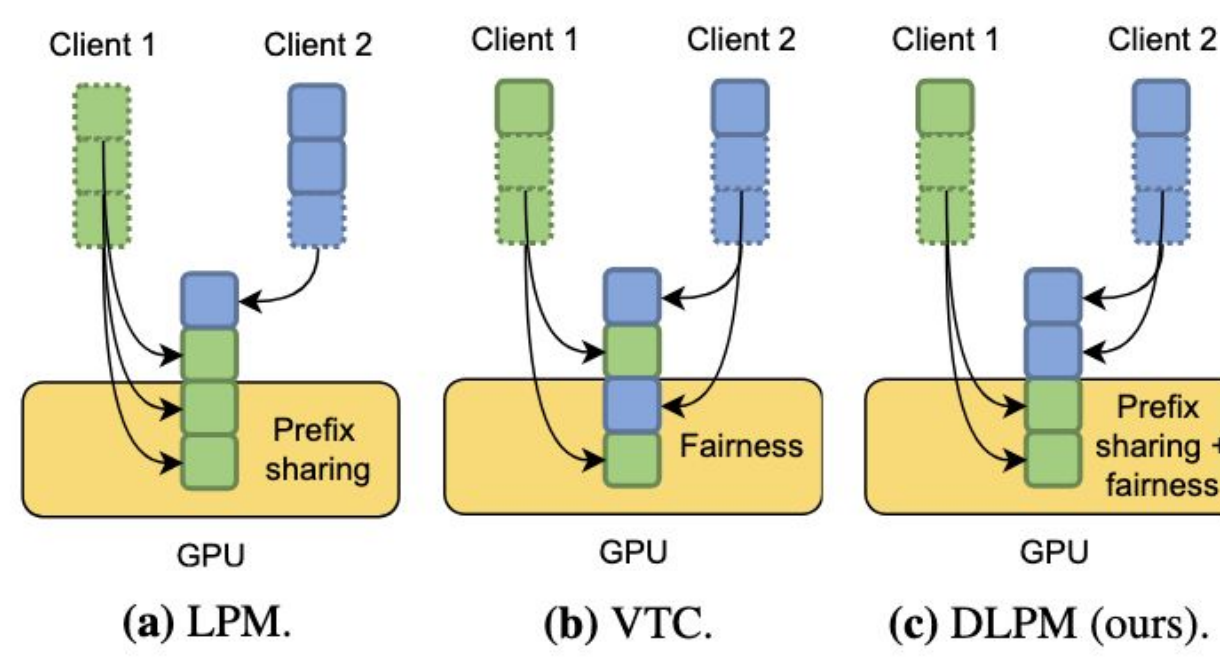
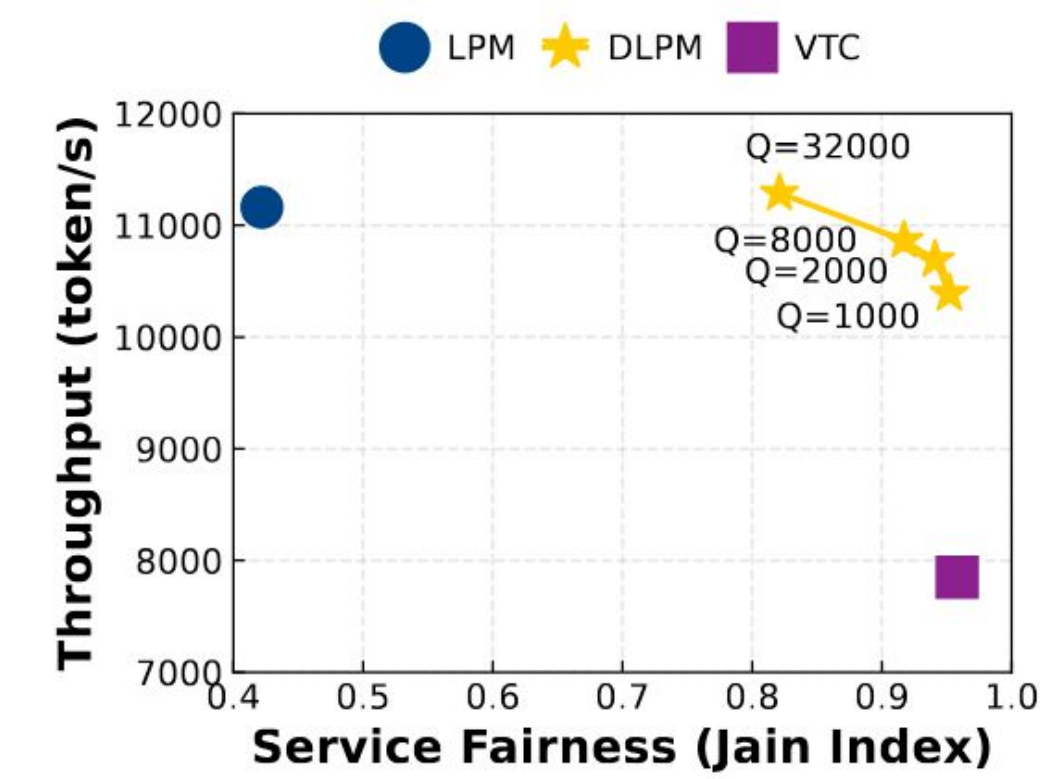


Locality-aware Fair Scheduling in LLM Serving

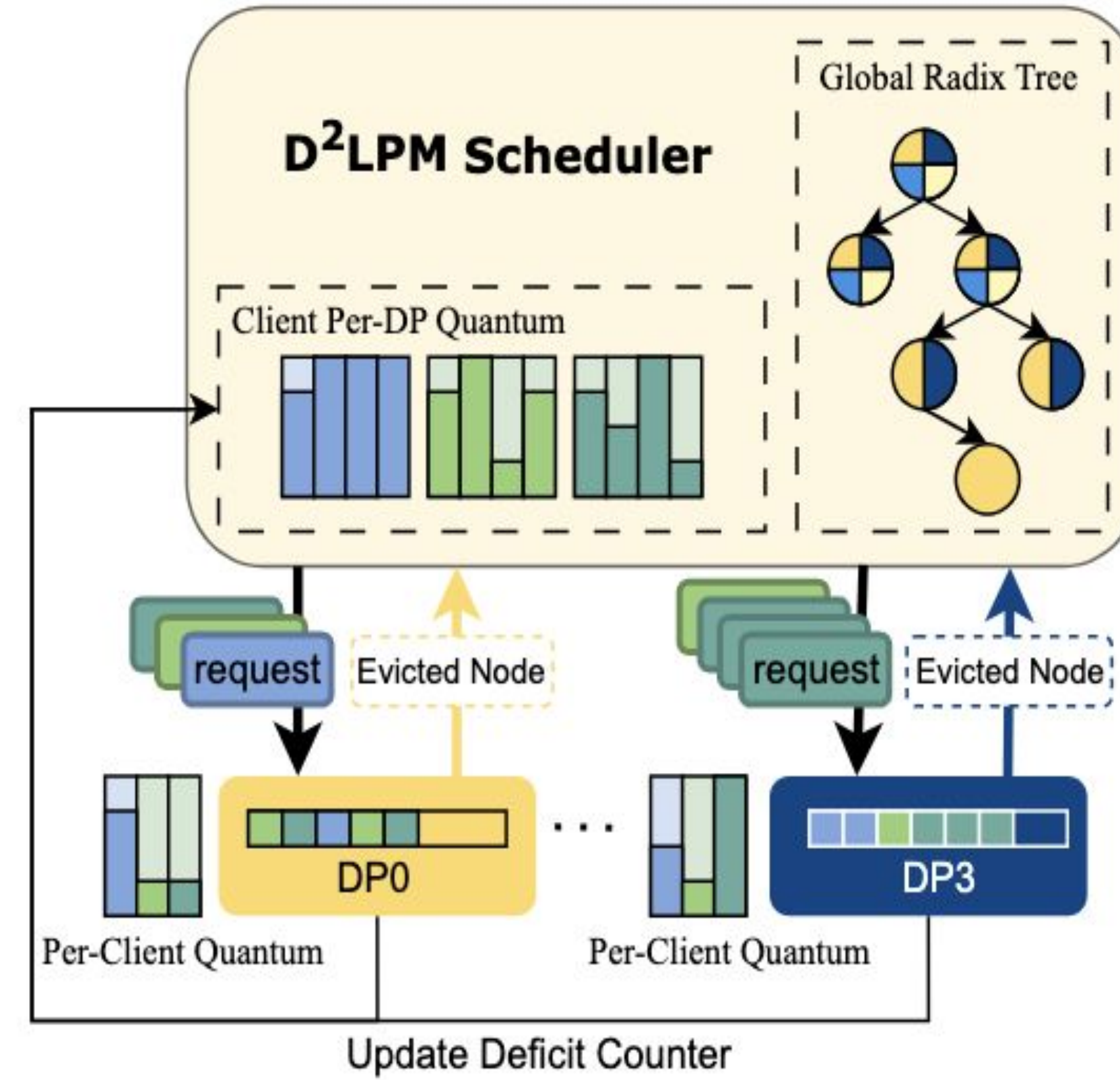
Shiyi Cao^{*1 2}, Yichuan Wang^{*1 2}, Ziming Mao¹, Pin-Lun Hsu², Liangsheng Yin^{1 2}, Tian Xia¹, Dacheng Li^{1 2}, Shu Liu¹, Yineng Zhang², Yang Zhou¹, Ying Sheng², Joseph E. Gonzalez¹, Ion Stoica¹
¹UC Berkeley ²LMSYS *indicates equal contribution.

Background & Motivation

- **Background:** Achieving efficient online LLM inference with SLO guarantees necessitates **isolation** among **different clients** (introduced in VTC OSDI24)
- However, existing LLM serving scheduler does not carefully consider both **fairness** and **locality** (in addition, **load balance** in distributed scenario).
- LPM(longest prefix match) prioritizes locality but will result in severe fairness issue, VTC (virtual token counter) ensures fairness but leads to poor performance



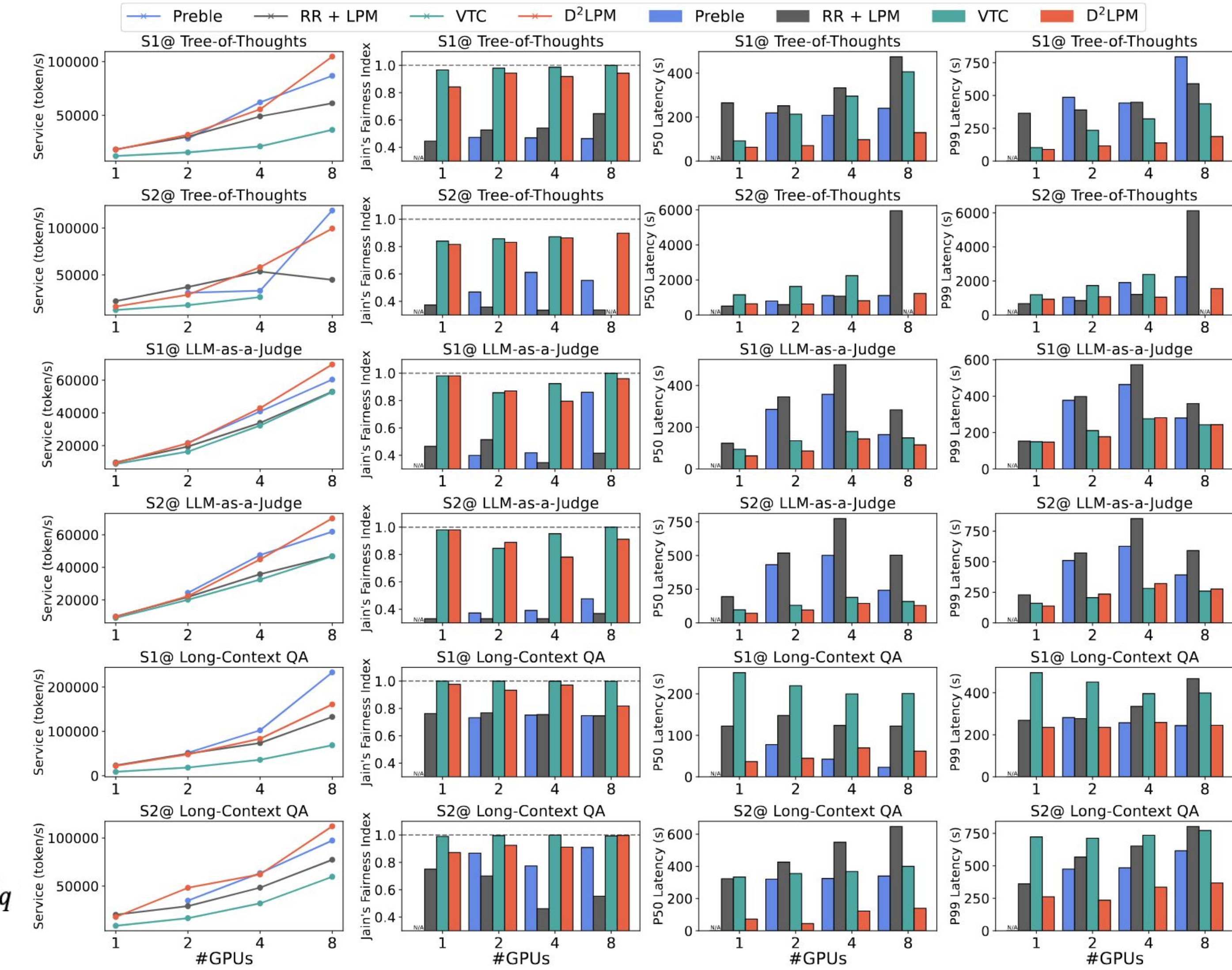
2. D²LPM for Distributed scenario:



3. Strict fairness bound

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

2. Results:



* S1 means malicious client using high request rate, S2 means malicious client using longer prefix.

Main Algorithms

1. DLPM for Single GPU:

Algorithm 1 Deficit Longest Prefix Match (DLPM)

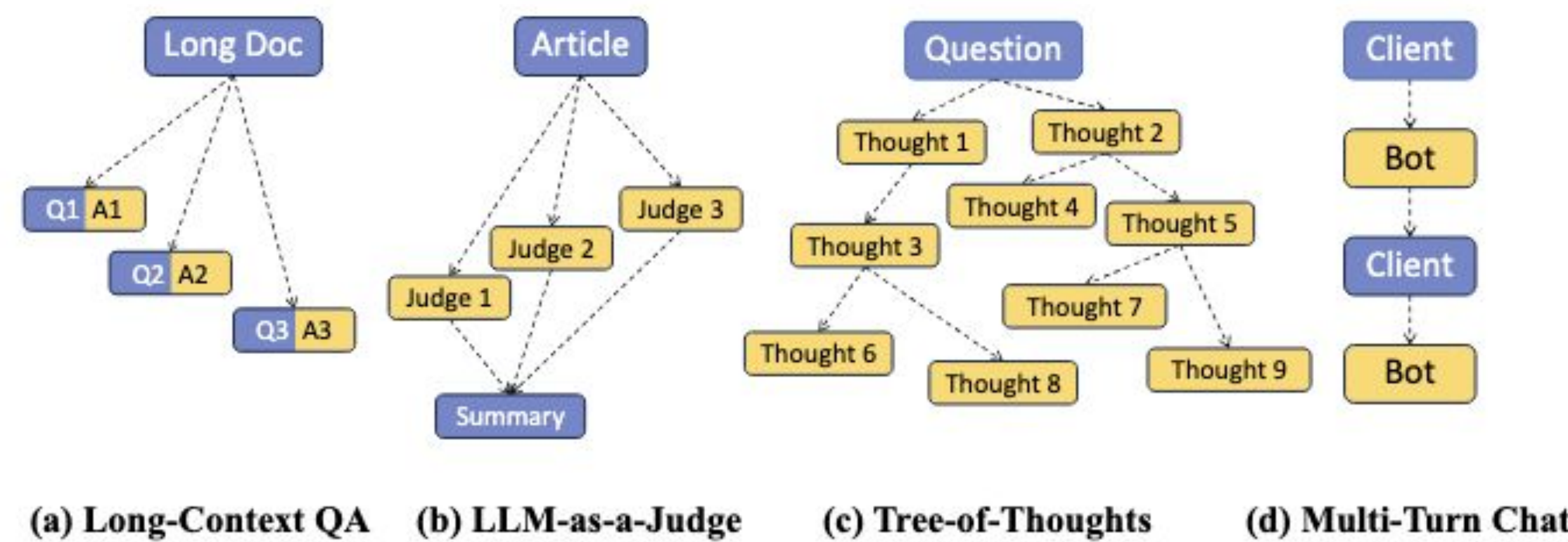
```

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:    $B \leftarrow filter\_finished\_requests(B)$ 

```

Evaluation

1. Workload:



3. Case study:

