



DistServe: Disaggregating Prefill and Decoding for Goodput-optimized Large Language Model Serving

OSDI'24

Authors: Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu, Yibo Zhu,
Xuanzhe Liu, Xin Jin, Hao Zhang

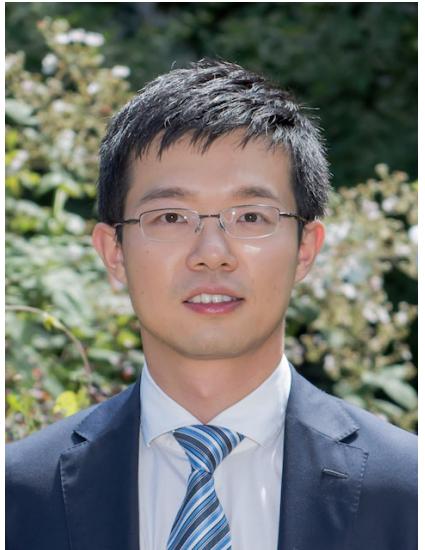


北京大学
PEKING UNIVERSITY



Presenter: Yitao Wang
2025.3.7

Author / Team Introduce



PKU

<https://xinjin.github.io/>

Projects:

- Serverless Computing
- Training and Serving Large Language Models
- Network Testing and Verification
- Disaggregated Storage with RDMA and DPUs

Papers:

- [SOSP'24] LoongServe: Efficiently Serving Long-Context Large Language Models with Elastic Sequence Parallelism
- [OSDI'24] dLoRA: Dynamically Orchestrating Requests and Adapters for LoRA LLM Serving
- [SIGCOMM'23] Ditto: Efficient Serverless Analytics with Elastic Parallelism
- [NSDI'23] Transparent GPU Sharing in Container Clouds for Deep Learning Workloads

Author / Team Introduce



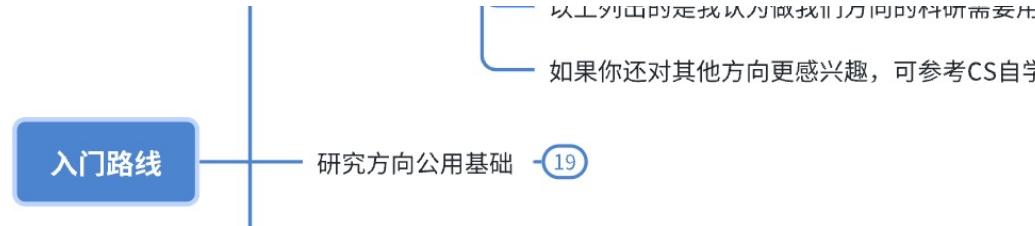
PKU

<https://yinminzhong.com/>

Papers:

- [NSDI'24]MegaScale: Scaling Large Language Model Training to More Than 10,000 GPUs
- [OSDI'24]AlphaServe: Statistical Multiplexing with Model Parallelism for Deep Learning Serving
- [OSDI'24]DistServe: Disaggregating Prefill and Decoding for Goodput-optimized Large Language Model Serving
- [arXiv'23]Fast Distributed Inference Serving for Large Language Models

Author / Team Introduce



如果他还对其他方向更感兴趣，可参考CS自学指南，挑选感兴趣的学习

<https://csdiy.wiki/>



Yinmin Zhong

PKUFlyingPig

Unfollow

CS Ph.D @ PKU; Machine Learning Systems; Distributed Systems;

8k followers · 33 following

PKUFlyingPig / README.md

Hi there 🤚

Yinmin Zhong's GitHub Stats

☆ Total Stars Earned:	69.9k
⌚ Total Commits (2025):	56
🍴 Total PRs:	67
ⓘ Total Issues:	44
✍ Contributed to (last year):	3

A-

Pinned

- cs-self-learning (Public)
- MIT6.S081-2020fall (Public)

CS自学指南

cs-self-learning v1.1.0 60.8k 7.1k

目录

- 梦开始的地方 -- CS61A
- 为什么写这本书
- 自学的好处
- 自学的坏处
- 这本书适合谁
- 特别鸣谢
- 你也想加入到贡献者的行列
- 关于交流群的建立
- 请作者喝杯下午茶

CS自学指南

前言

最近更新: Release v1.1.0 已发布!

这是一本计算机的自学指南, 也是对自己大学三年自学生涯的一个纪念。

这同时也是一份献给北大信科学弟学妹们的礼物。如果这本书能对你们的信科生涯有哪怕一丝一毫的帮助, 都是对我极大的鼓励和慰藉。

本书目前包括了以下部分(如果你有其他好的建议, 或者想加入贡献者的行列, 欢迎邮件 zhongyinmin@pku.edu.cn 或者在 issue 里提问):

- 本书使用指南: 由于书内涵盖资源众多, 我根据不同人群的空闲时间和学习目标制定了对应的使用指南。

Contents

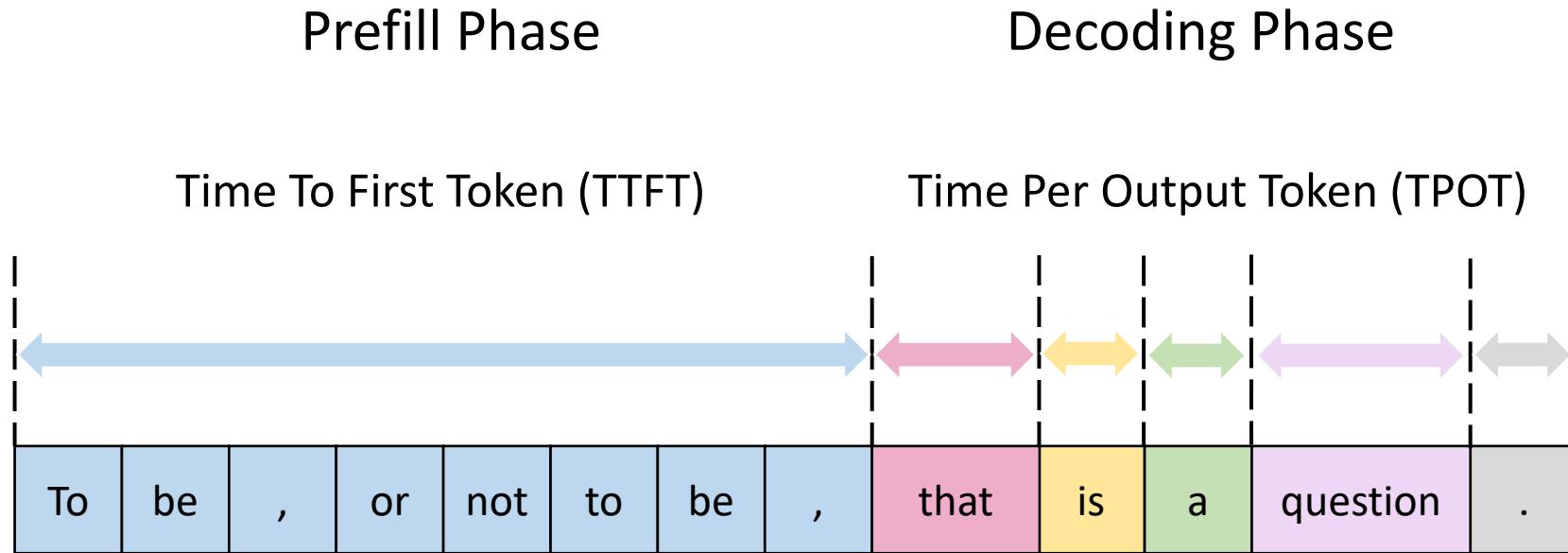


- **Background**
 - Problem
 - Opportunity
 - Challenge
 - Method
 - Evaluation
 - Thinking

Background



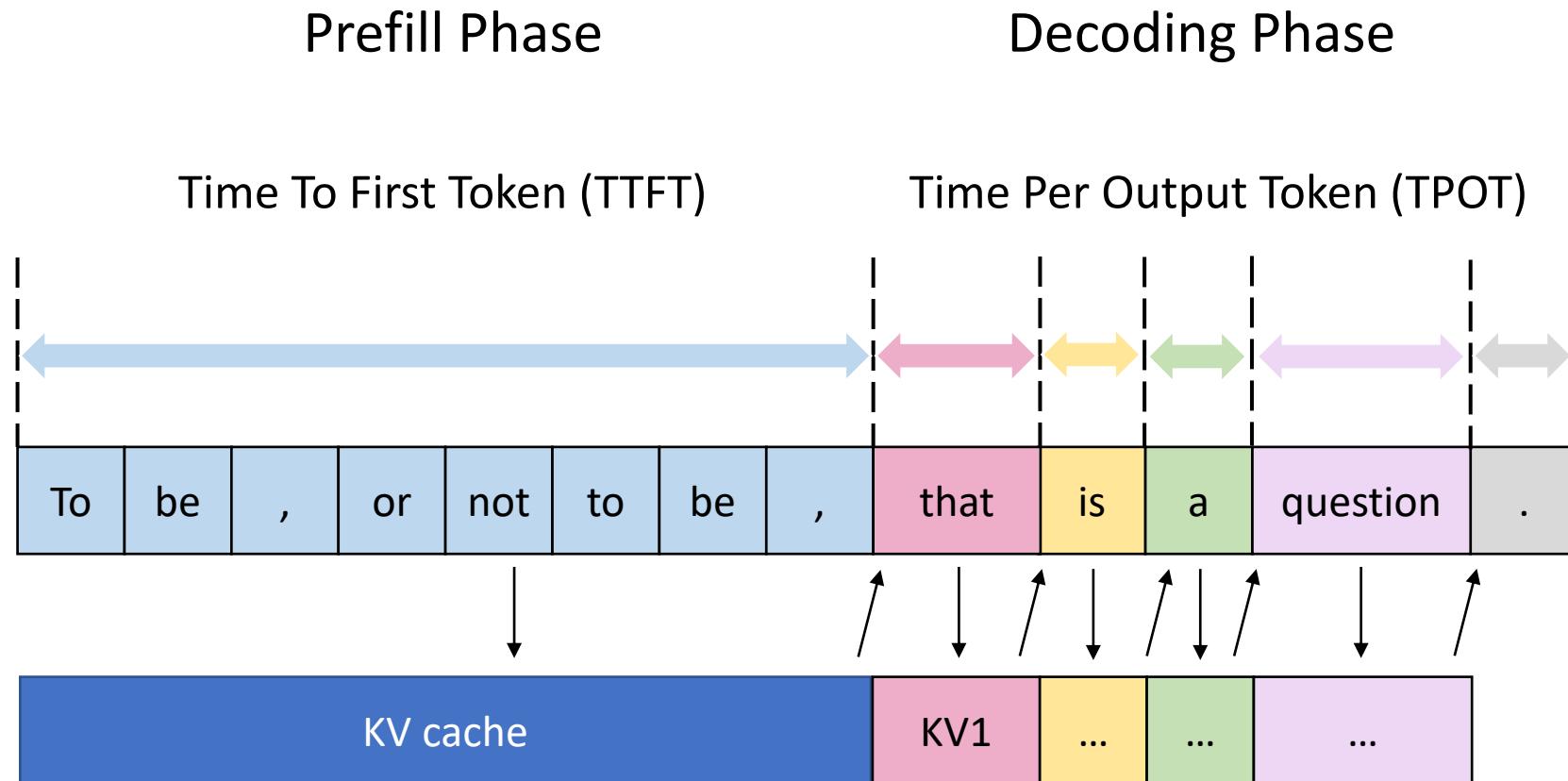
LLM Inference



Background



LLM Inference



Background



LLM Inference

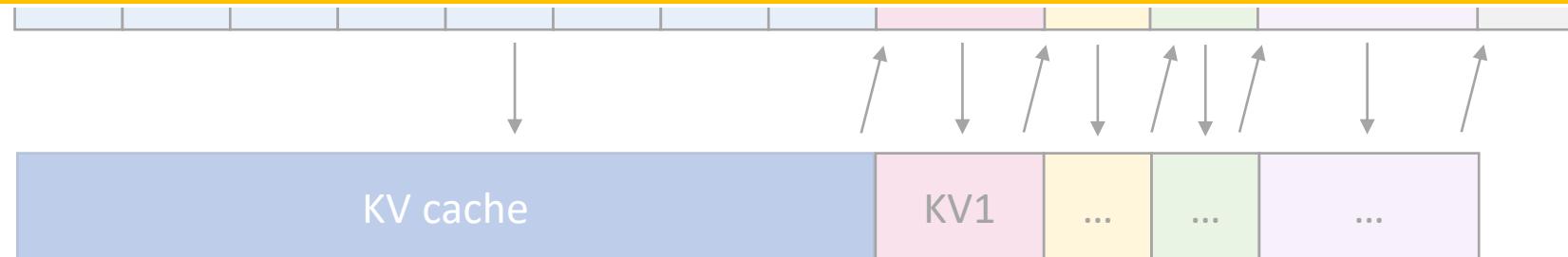
Prefill Phase

Decoding Phase

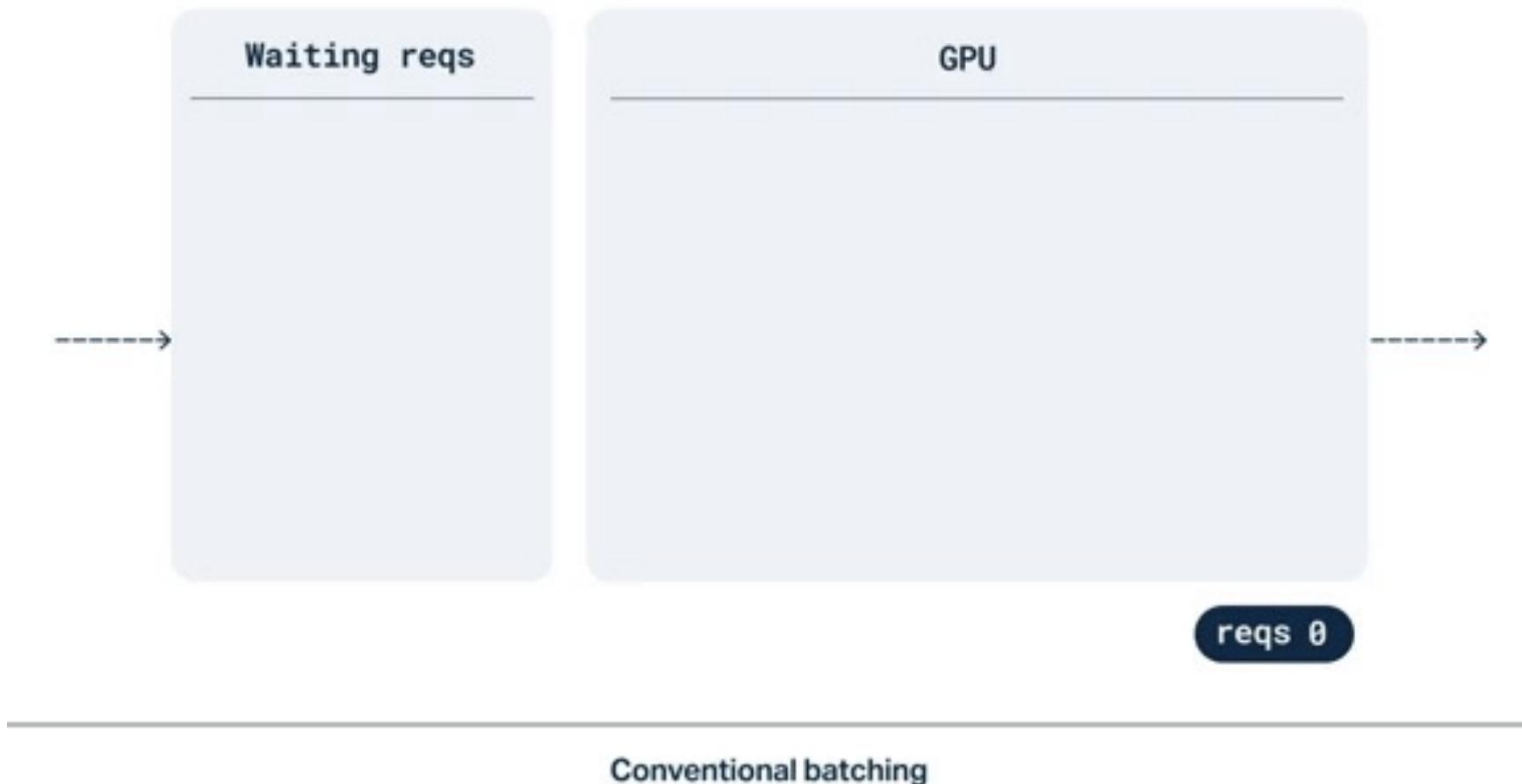
Time To First Token (TTFT)

Time Per Output Token (TPOT)

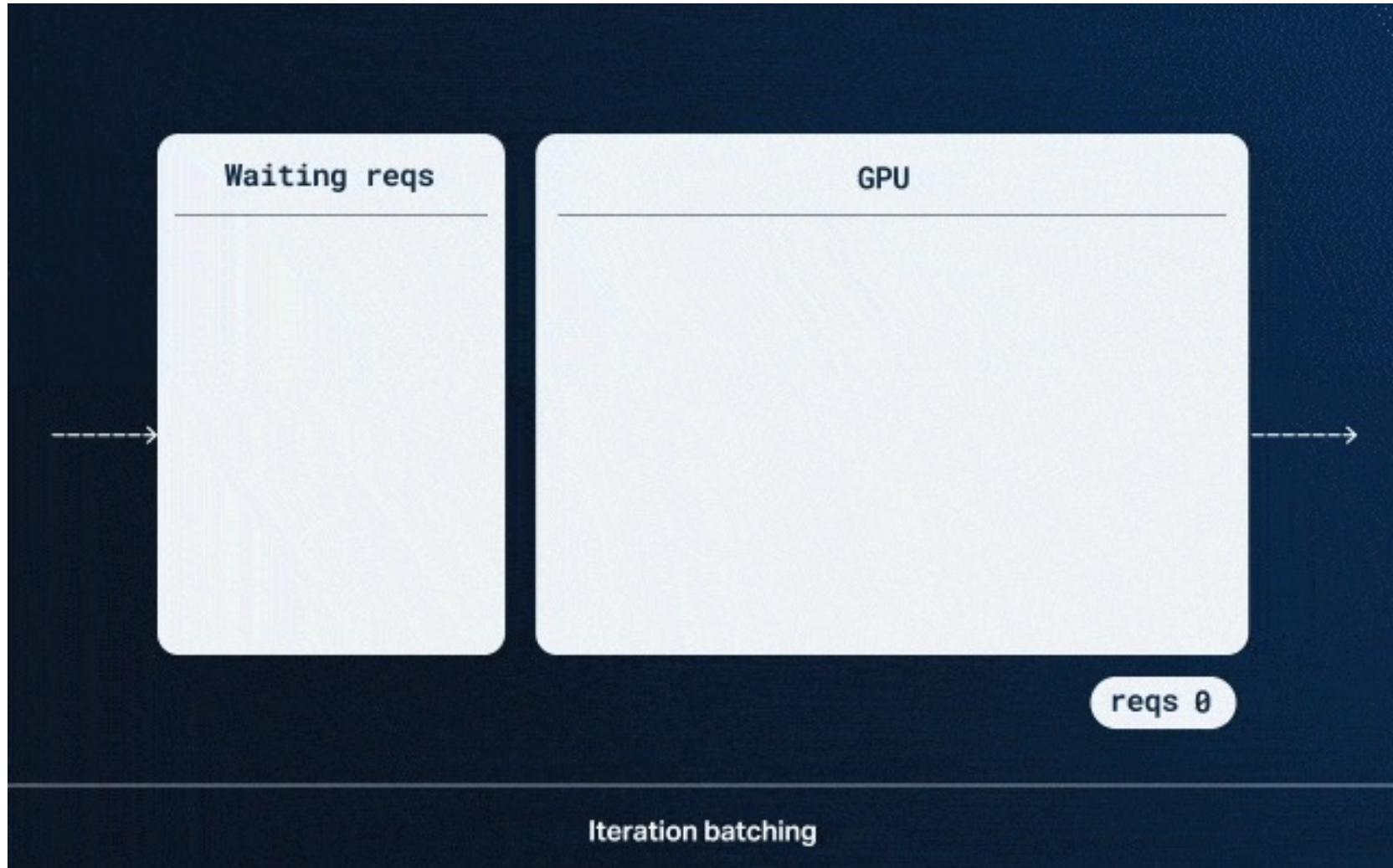
**Prefill phase is compute bounded;
While Decode phase is memory bounded!**



Naïve / Static Batching



Continuous Batching



Orca: A Distributed Serving System for Transformer-Based Generative Models [OSDI'22]



Applications have Diverse SLO

- **TTFT**

Time to first token

Initial response time



Chatbot

Fast initial response



Summarization



User can tolerate longer initial response

- **TPOT**

Time per output token

Average time between two subsequent generated tokens



Human reading speed (P99 latency = 250ms)



Data output generation (P99 latency = 35ms)

SLO for TTFT & TPOT



Different apps have various latency requirements

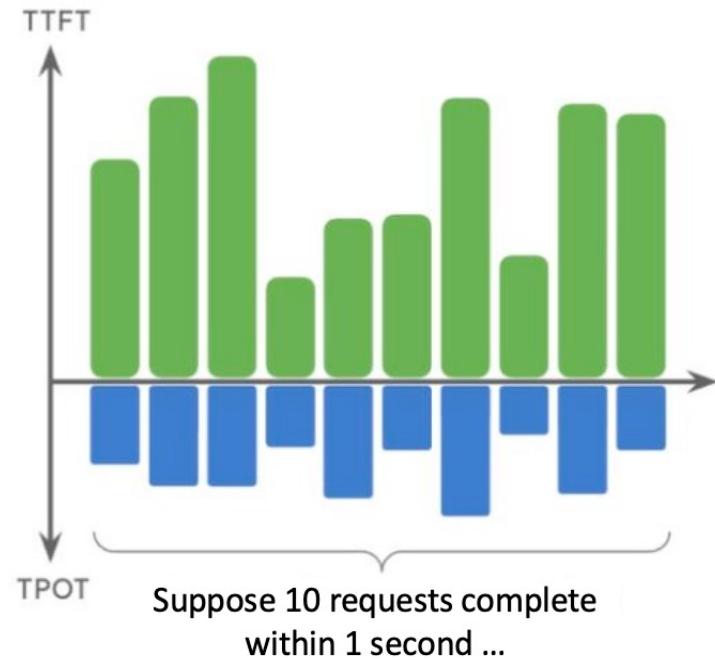
	TTFT	TPOT
Chat	Low (< 1s)	Match read speed (~100ms)
Search	Very Low (~200ms)	Match read speed (~100ms)
Program	Very Low (~200ms)	Very Low (~50ms)

Throughput vs. Goodput

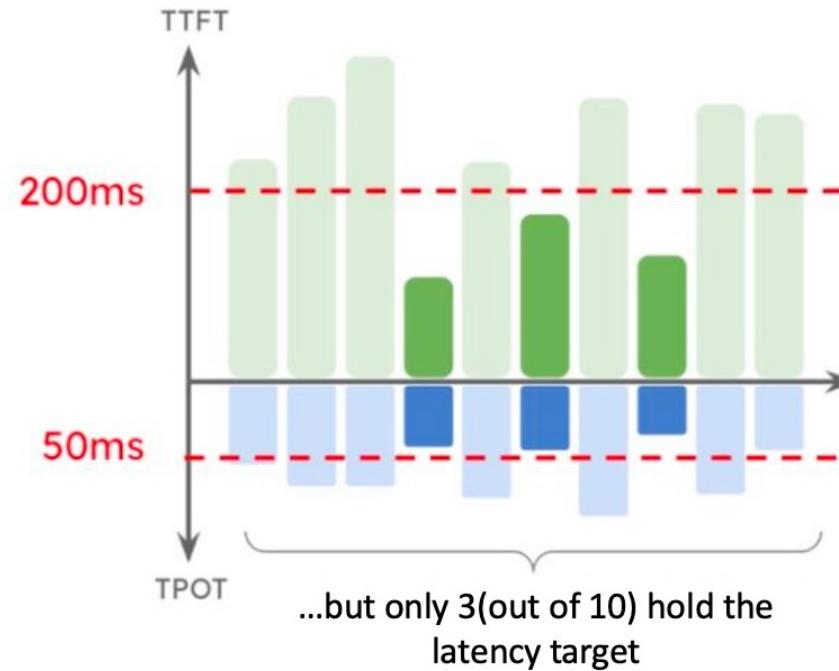


High Throughput \neq High Goodput

Throughput = completed request / time
= 10 req / s



Goodput = completed requests **within SLO** / time
= 3 req / s

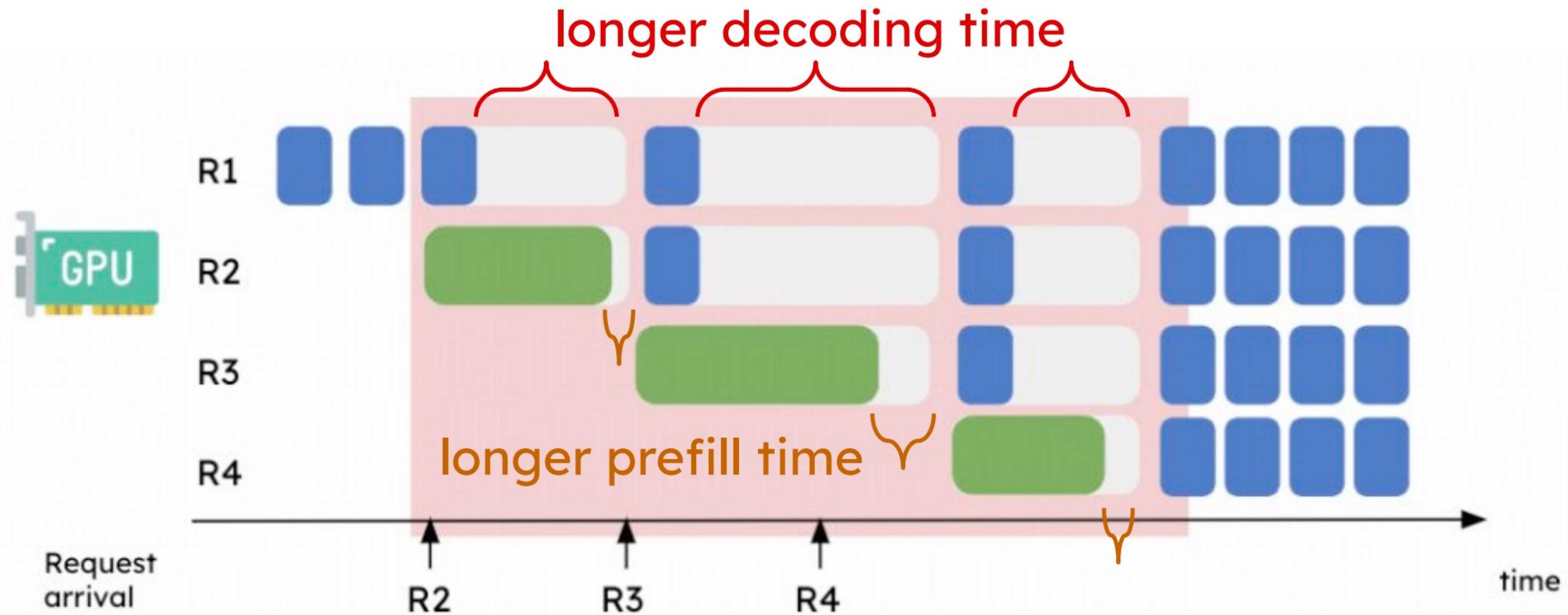


Contents

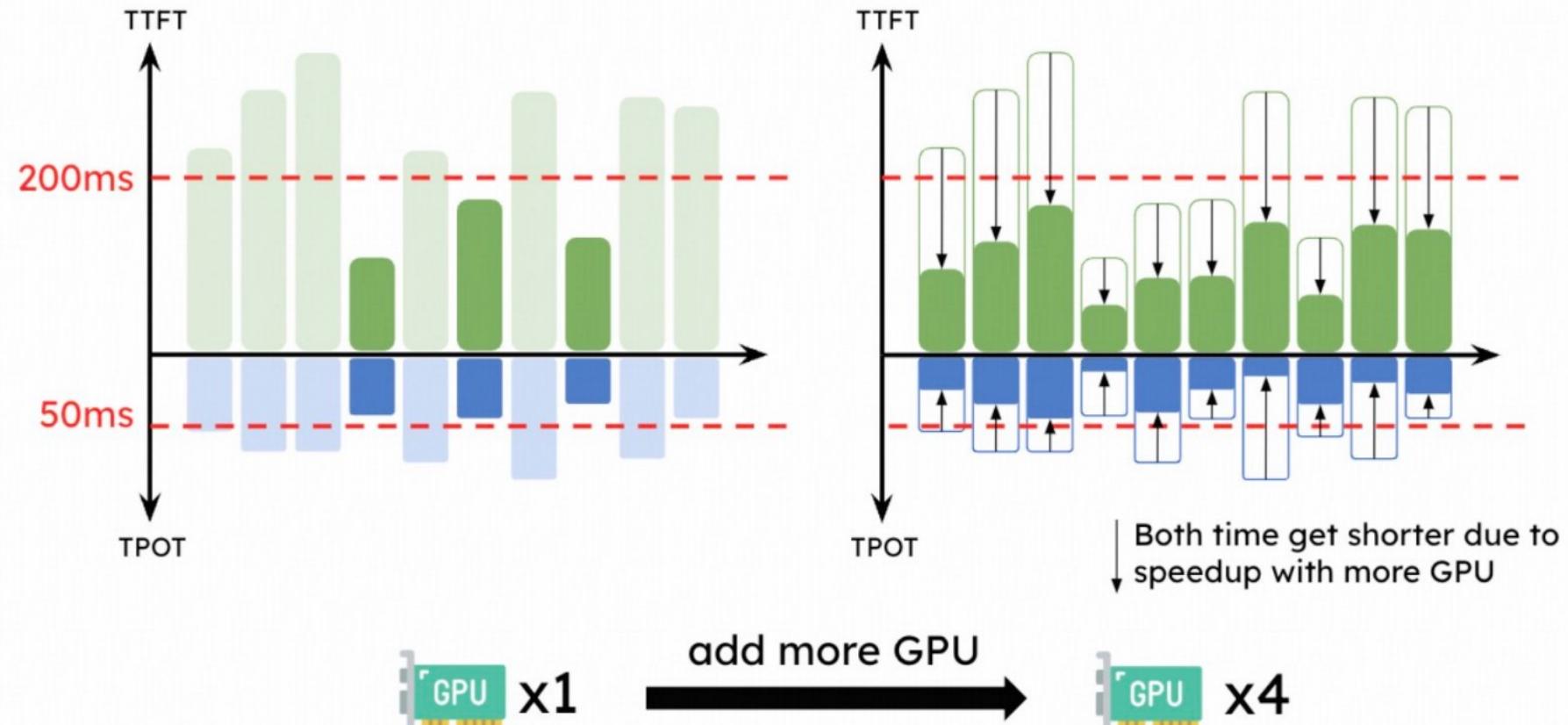


- Background
- Problem
- Opportunity
- Challenge
- Method
- Evaluation
- Thinking

Problem1: Conflict between P and D



Problem2.1: GPU Strategy



Cost	✓	✗
Goodput	✗	✓

Problem2.2: Parallelism Strategy

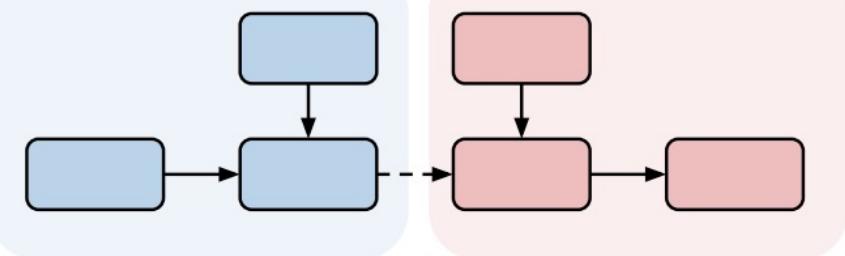


Batching the two phases make them share the same parallel strategy

Pipeline Parallelism

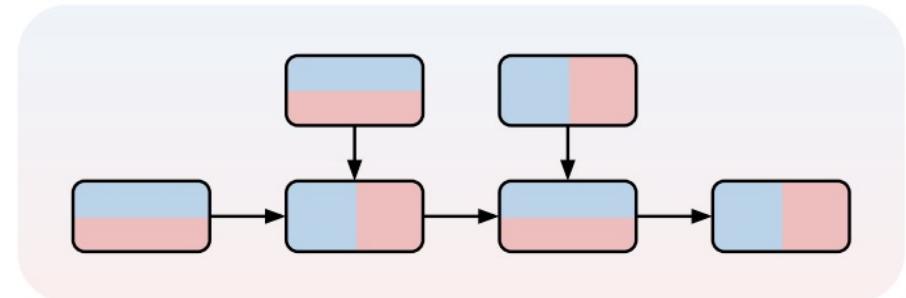


Inter-op Parallelism

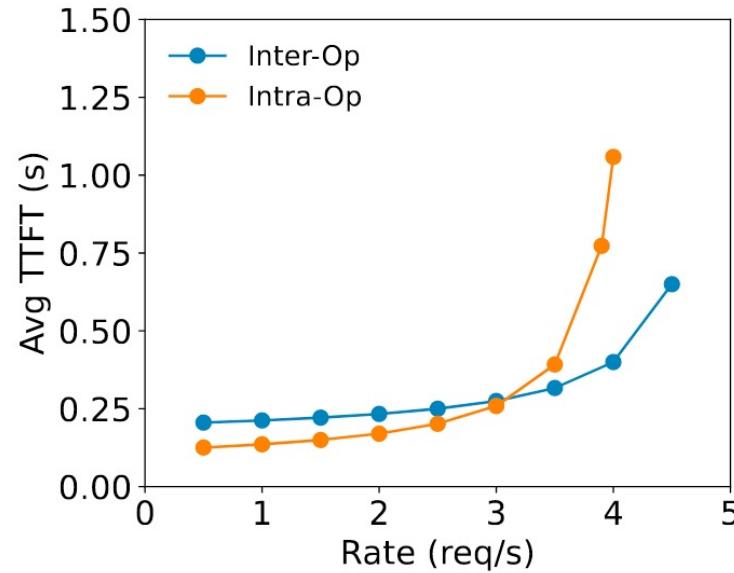


Tensor Parallelism

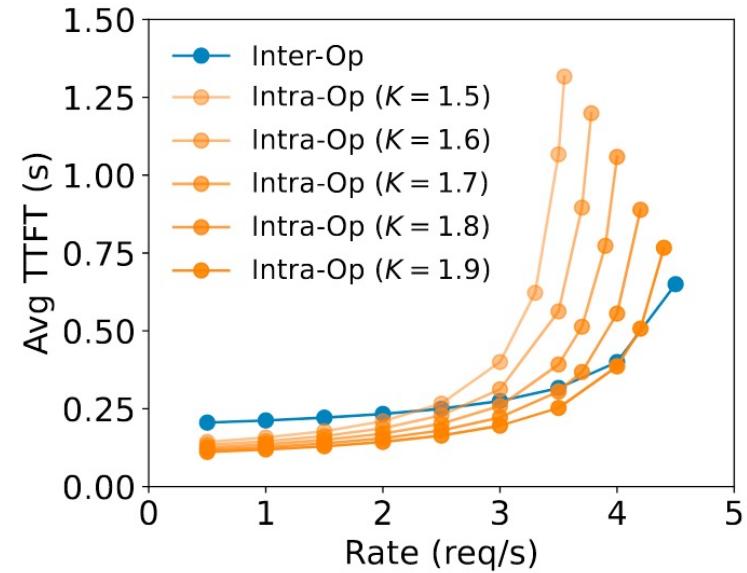
Intra-op Parallelism



Problem2.2: Parallelism Strategy



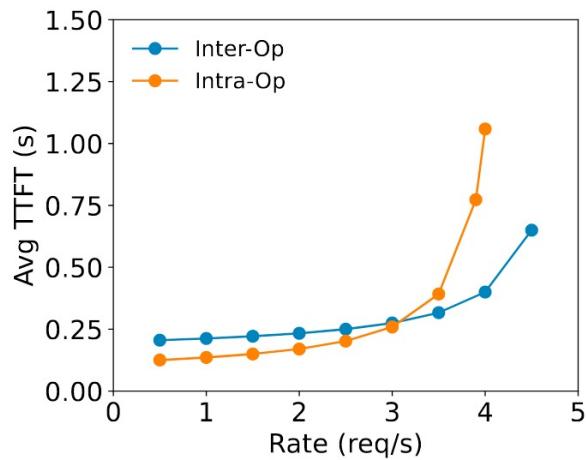
(a) Real experiment results



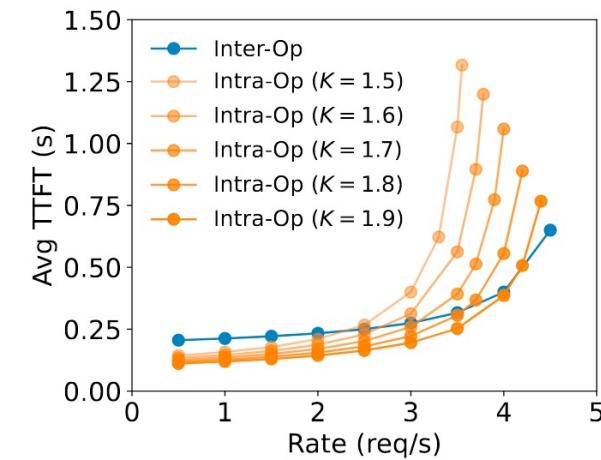
(b) Changing intra-op speedup

Figure 4: Average TTFT when serving an LLM with 66B parameters using different parallelism on two A100 GPUs.

Problem2.2: Parallelism Strategy



(a) Real experiment results



(b) Changing intra-op speedup

Figure 4: Average TTFT when serving an LLM with 66B parameters using different parallelism on two A100 GPUs.

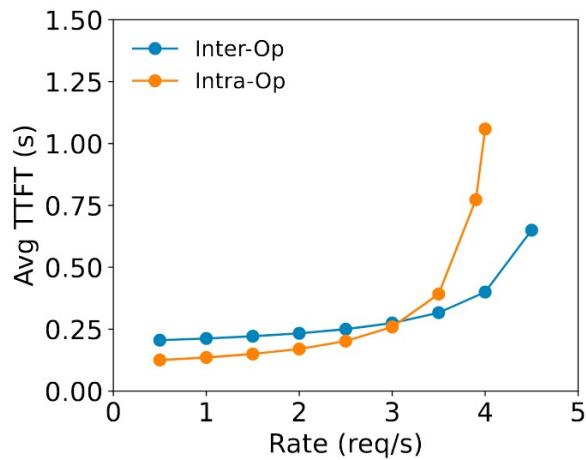


$$Avg_TTFT = D + \frac{RD^2}{2(1-RD)},$$

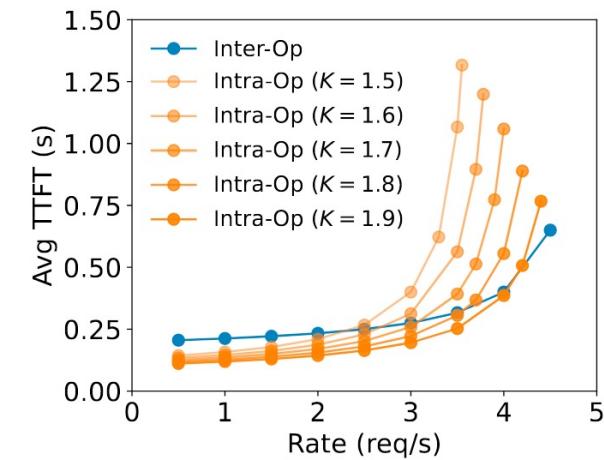
$$Avg_TTFT_{intra} = \frac{D}{K} + \frac{RD^2}{2K(K-RD)}.$$

$$Avg_TTFT_{inter} = D_s + \frac{RD_m^2}{2(1-RD_m)} = D + \frac{RD^2}{4(2-RD)}.$$

Problem2.2: Parallelism Strategy



(a) Real experiment results



(b) Changing intra-op speedup

Figure 4: Average TTFT when serving an LLM with 66B parameters using different parallelism on two A100 GPUs.

$$Avg_TTFT = D + \frac{RD^2}{2(1-RD)},$$

Queueing Theory



$$Avg_TTFT_{inter} = D_s + \frac{RD_m^2}{2(1-RD_m)} = D + \frac{RD^2}{4(2-RD)}.$$

Problem2.2: Parallelism Strategy

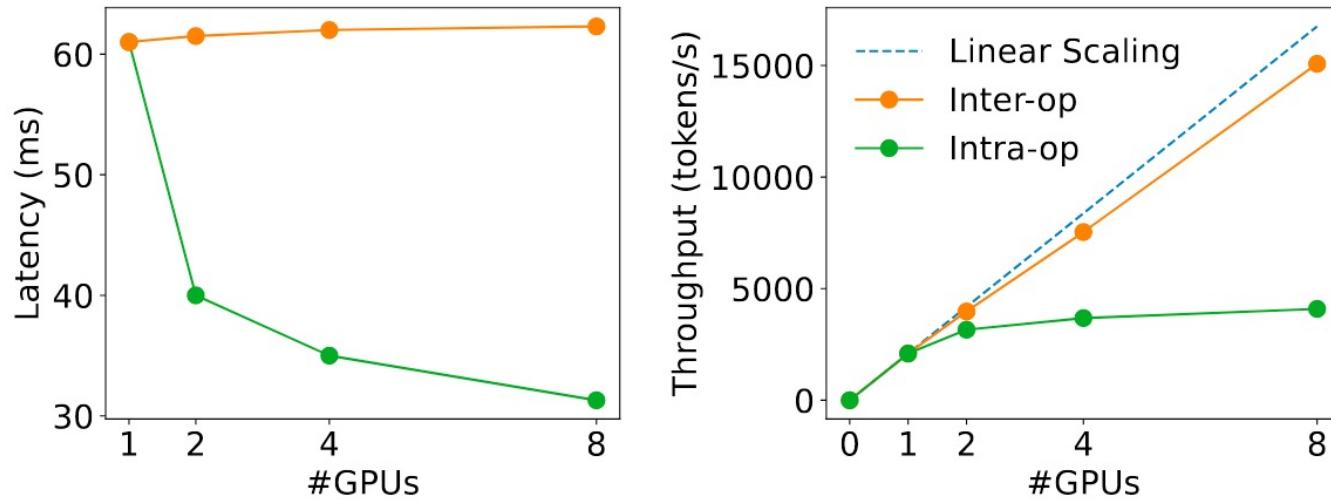
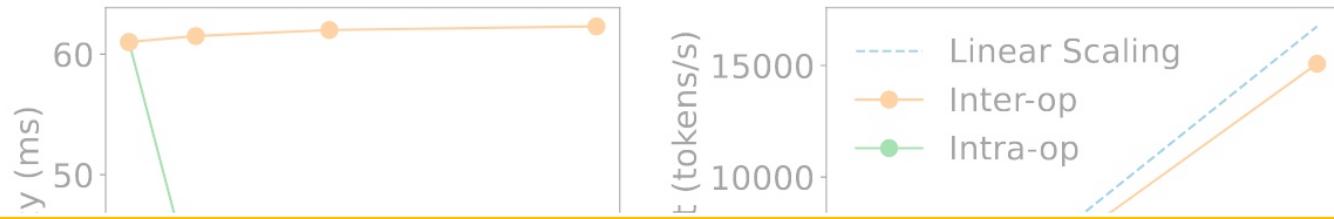


Figure 5: Decoding phase latency and throughput when serving a 13B LLM with batch size = 128 and input length = 256 under different parallel degrees.

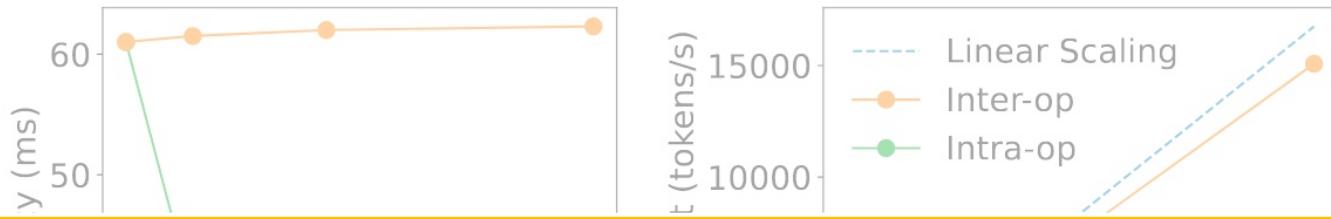
Problem2.2: Parallelism Strategy



Different phase need to meet the target SLO, so their parallelism strategies are also different.

Figure 5: Decoding phase latency and throughput when serving a 13B LLM with batch size = 128 and input length = 256 under different parallel degrees.

Problem2.2: Parallelism Strategy



**Our goal: Meet both SLOs while minimizing the cost of GPUs.
In other word, *maximizing per-GPU goodput*.**

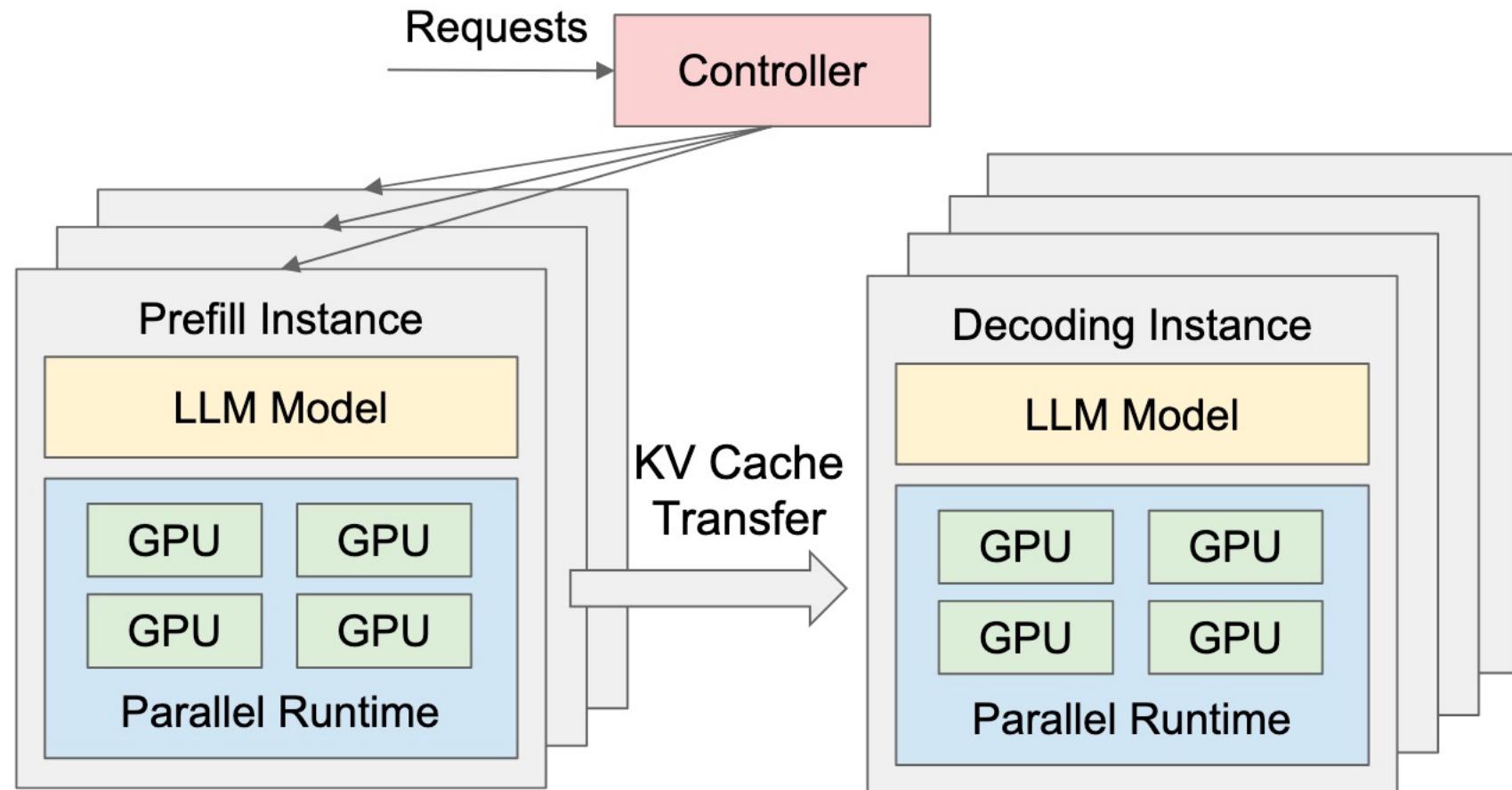
Figure 5: Decoding phase latency and throughput when serving a 13B LLM with batch size = 128 and input length = 256 under different parallel degrees.

Contents

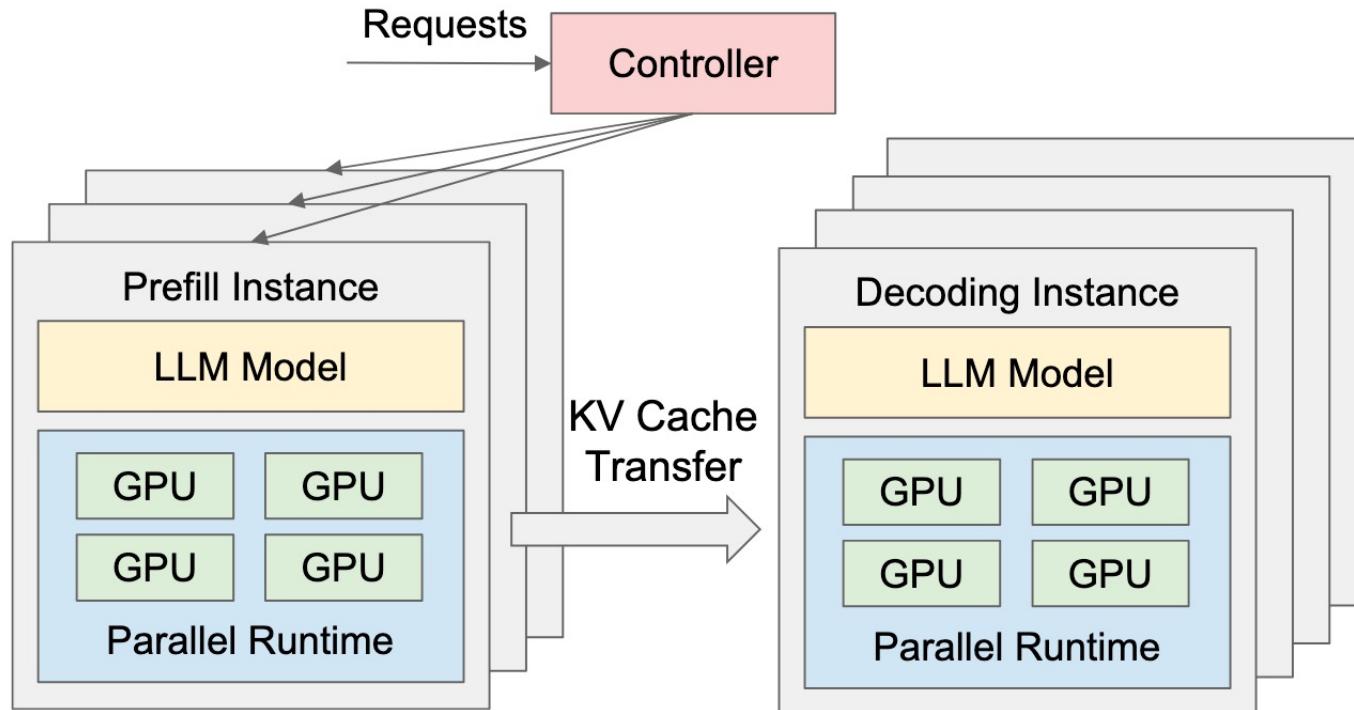


- Background
- Problem
- **Opportunity**
- Challenge
- Method
- Evaluation
- Thinking

Opportunity



Opportunity



1. TTFT and TPOT are independent.
2. Decoupling GPU / Parallelism Strategy.

Contents



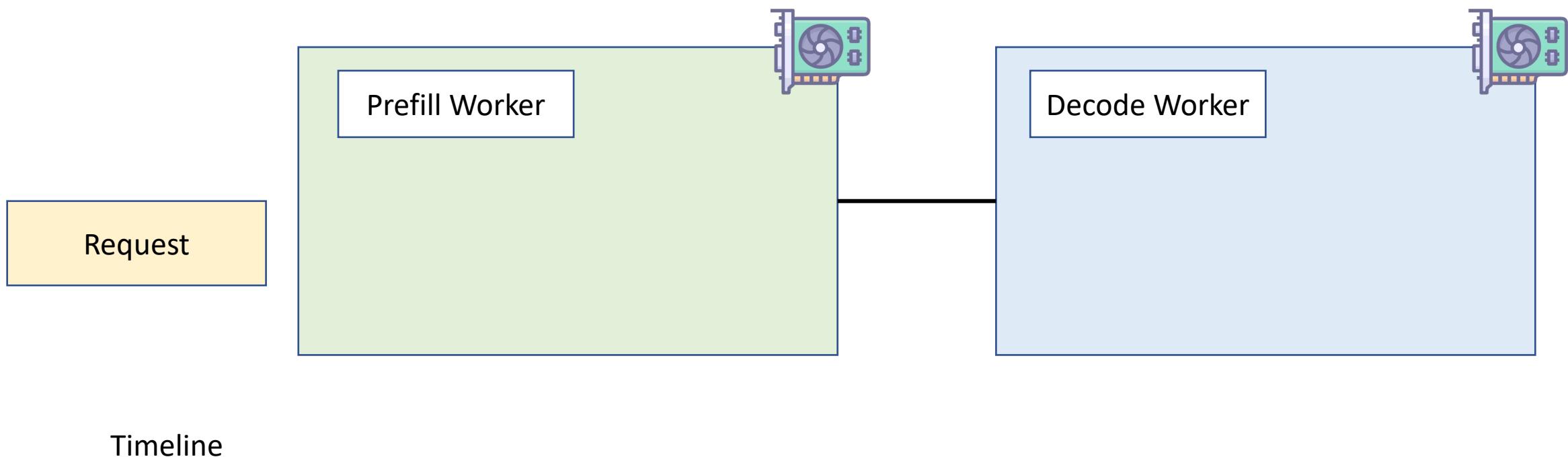
- Background
- Problem
- Opportunity
- **Challenge**
- Method
- Evaluation
- Thinking

Challenge

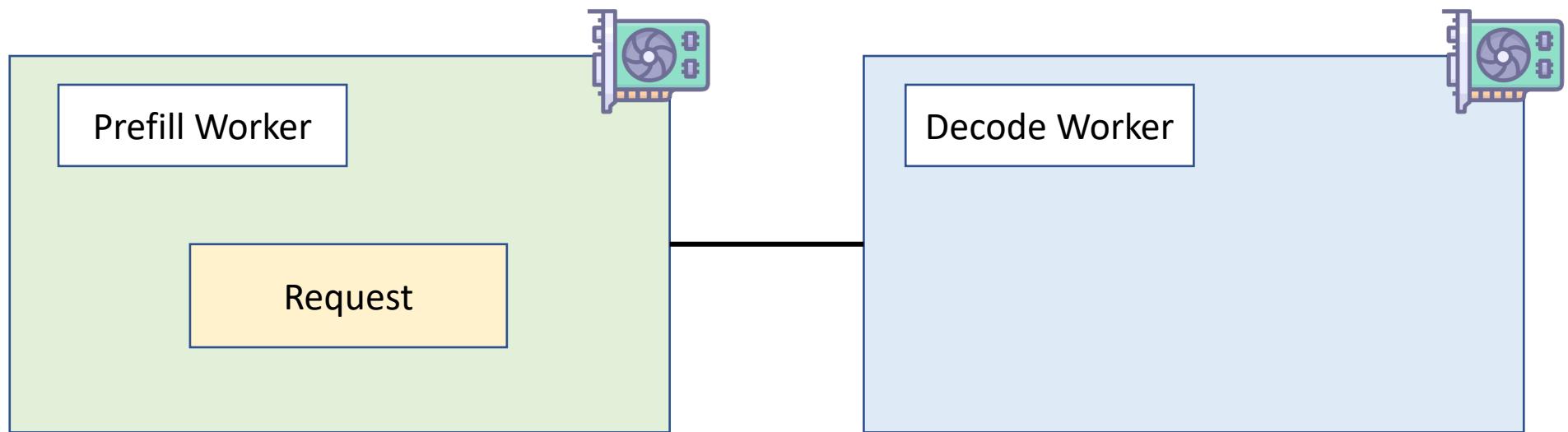


- **Communication overhead for KV cache migration**
- The optimization target – *per-GPU goodput*, is difficult to optimize:
 - The workload pattern
 - SLO requirements
 - Parallelism strategies
 - Resource allocation
 - Network bandwidth

Challenge: Migrate Overhead



Challenge: Migrate Overhead

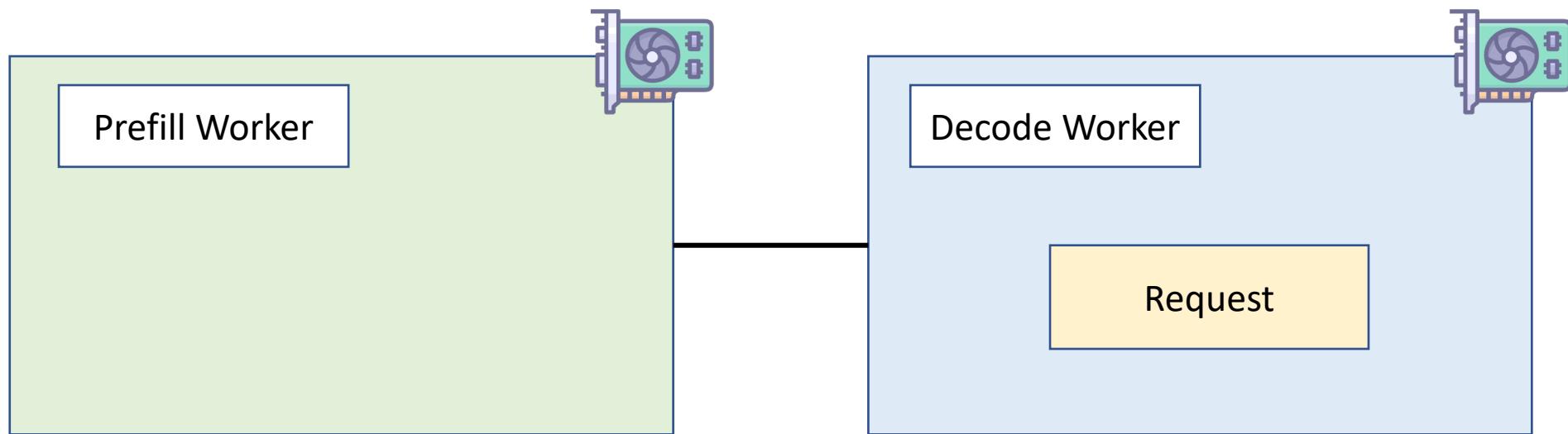


Timeline



Prefill

Challenge: Migrate Overhead

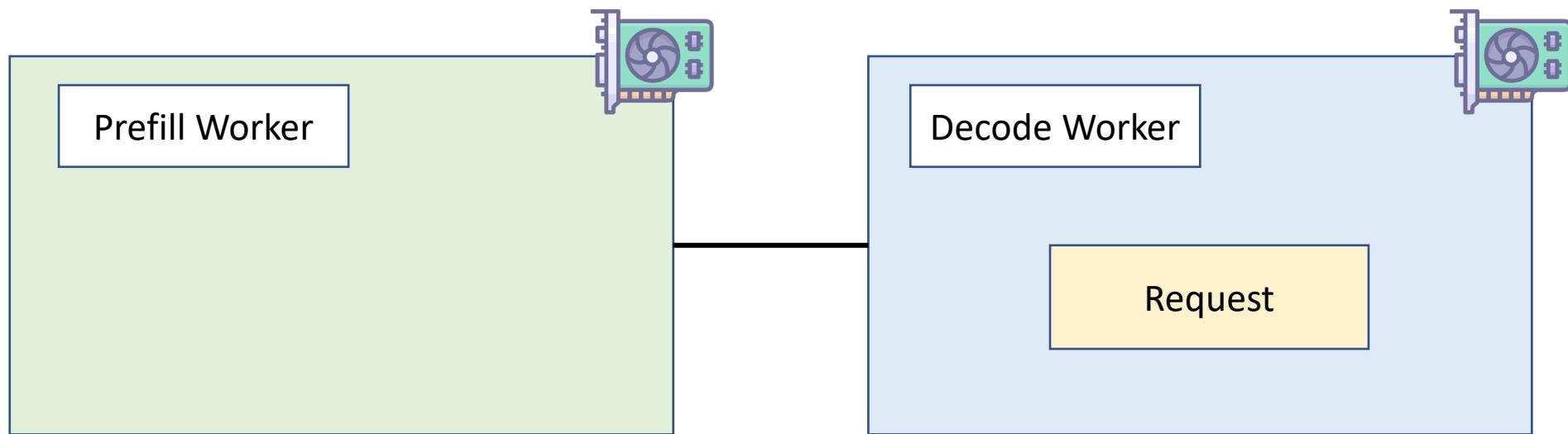


Timeline

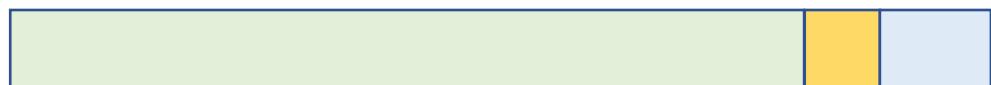


Migrate

Challenge: Migrate Overhead

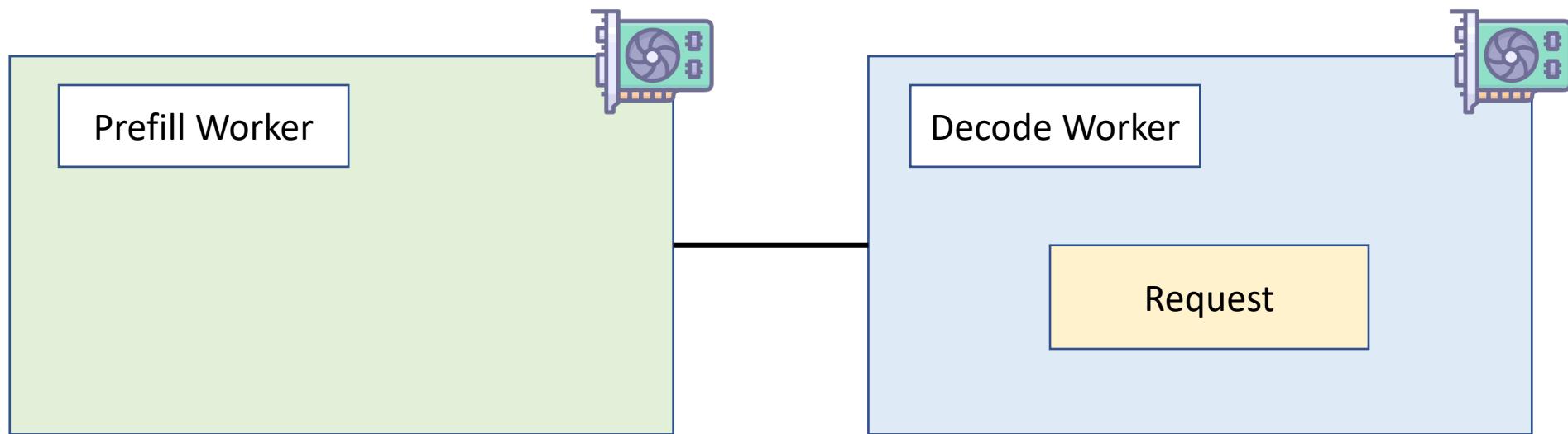


Timeline

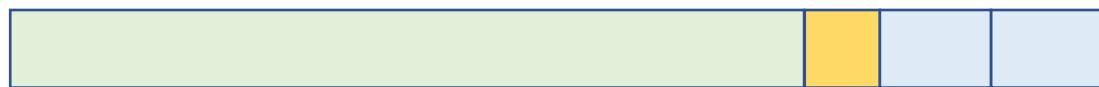


Decode

Challenge: Migrate Overhead

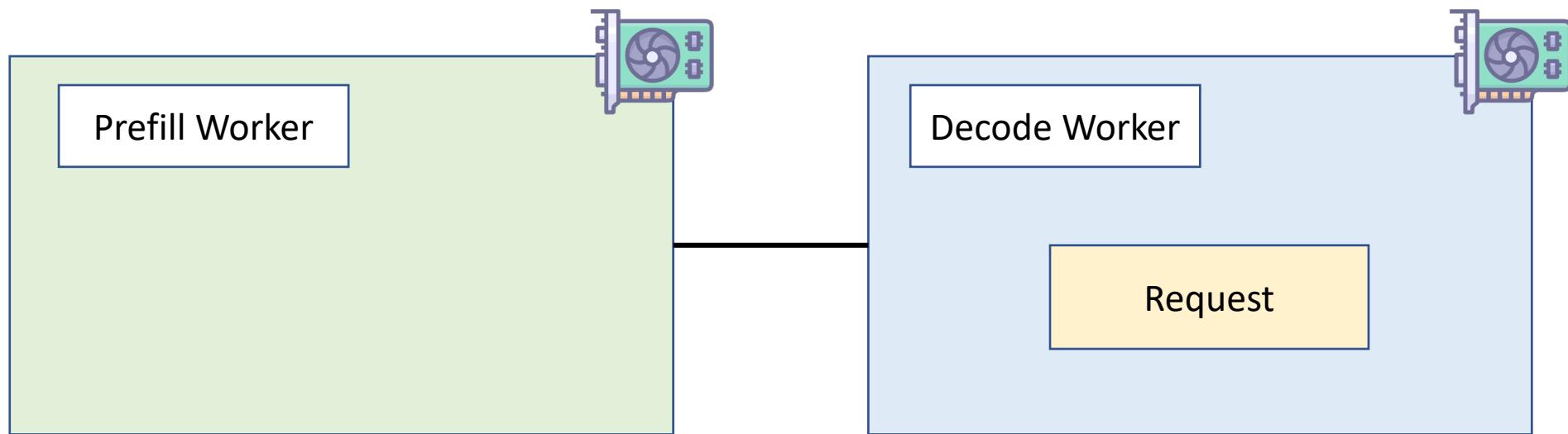


Timeline

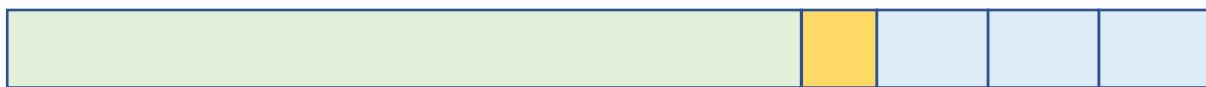


Decode

Challenge: Migrate Overhead

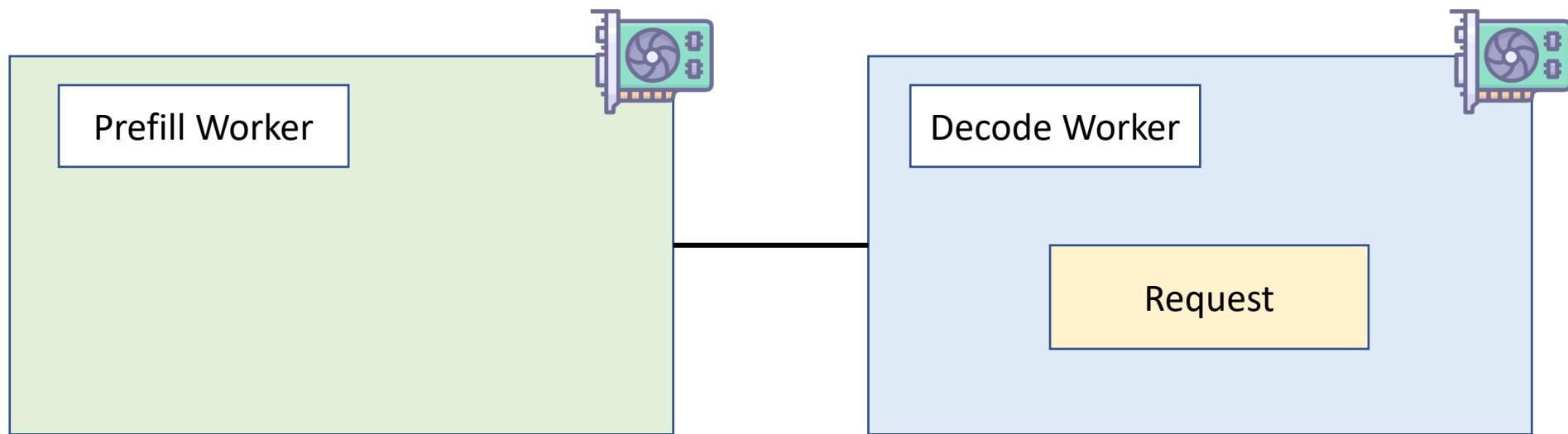


Timeline

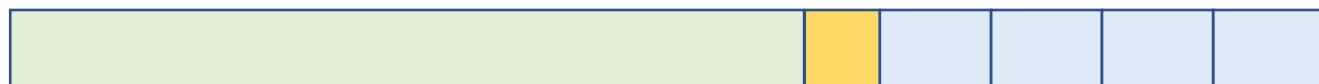


Decode

Challenge: Migrate Overhead



Timeline



Decode

Challenge: Migrate Overhead

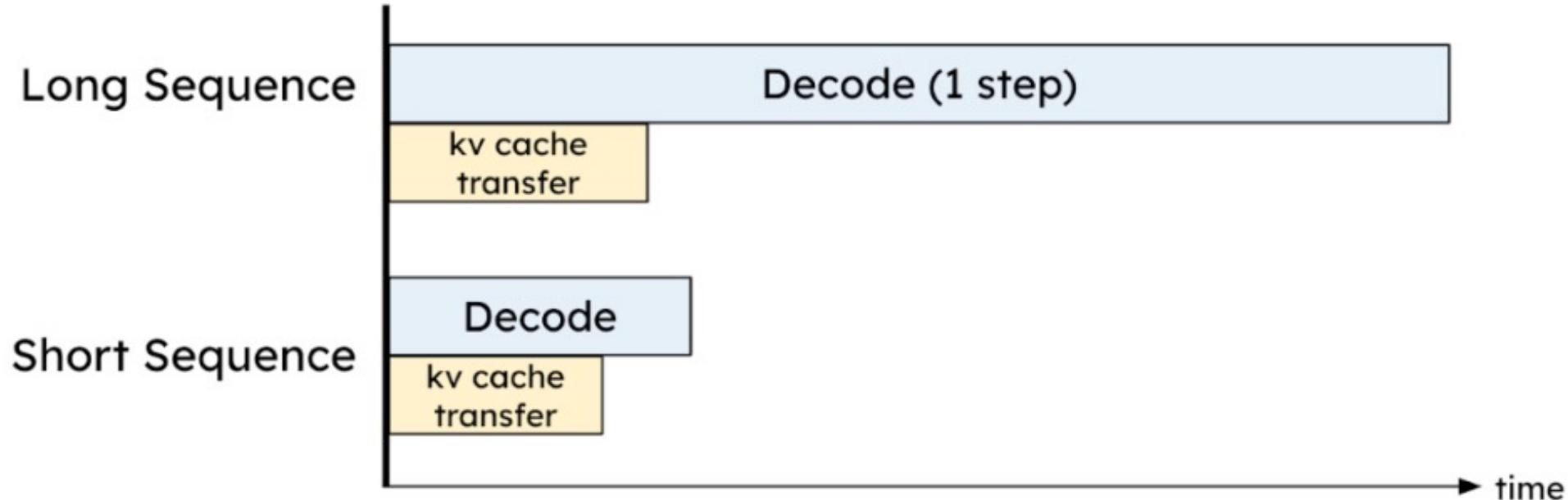


- 8-channel Pcle 5.0 * 16 (64GB/s per link)
- Seq Length = 2048
- OPT-175B
- Latency = $2048 * (4.5\text{MB/token}) / (64\text{GB/s} * 8) = 17.6\text{ms}$
- Decode latency = 30-50ms

Challenge: Migrate Overhead



KV cache transfer is not the bottleneck



Challenge: Goodput-optimization



- Communication overhead for KV cache migration
- The optimization target – *per-GPU goodput*, is difficult to optimize:

- The workload pattern
- SLO requirements
- Resource allocation
- Network bandwidth

DistServe

- Parallelism strategies for P/D
- The #P,#D to deploy
- How to place them onto cluster

Contents



- Background
- Problem
- Opportunity
- Challenge
- **Method**
- Evaluation
- Thinking

Method



Inter-machine Search

Algorithm 1 High Node-Affinity Placement Algorithm

Input: LLM G , #node limit per-instance N , #GPU per-node M , GPU memory capacity C , workload W , traffic rate R .

Output: the placement $best_plm$.

```
configp, configd ← 0, 0
for intra_op ∈ {1, 2, ..., M} do
    for inter_op ∈ {1, 2, ..., ⌈ N × M / intra_op ⌉} do
        if G.size / (inter_op × intra_op) < C then
            config ← (inter_op, intra_op)
             $\hat{G}$  ← parallel(G, config)
            config.goodput ← simu_prefill( $\hat{G}$ , W)
            if configp.goodput / configp.num_gpus < config.goodput / config.num_gpus then
                configp ← config
            config.goodput ← simu_decode( $\hat{G}$ , W)
            if configd.goodput / configd.num_gpus < config.goodput / config.num_gpus then
                configd ← config
        n, m ← ⌈ R / configp.goodput ⌉, ⌈ R / configd.goodput ⌉
        best_plm ← (n, configp, m, configd)
return best_plm
```

Method



Inter-machine Search

TP ← **for** $intra_op \in \{1, 2, \dots, M\}$ **do**
PP ← **for** $inter_op \in \{1, 2, \dots, \frac{N \times M}{intra_op}\}$ **do**

$config_p, config_d \leftarrow \emptyset, \emptyset$
if $\frac{G.size}{inter_op \times intra_op} < C$ **then**
 $config \leftarrow (inter_op, intra_op)$
 $\hat{G} \leftarrow parallel(G, config)$
 $config.goodput \leftarrow simu_prefill(\hat{G}, W)$
 if $\frac{config_p.goodput}{config_p.num_gpus} < \frac{config.goodput}{config.num_gpus}$ **then**
 $config_p \leftarrow config$
 $config.goodput \leftarrow simu_decode(\hat{G}, W)$
 if $\frac{config_d.goodput}{config_d.num_gpus} < \frac{config.goodput}{config.num_gpus}$ **then**
 $config_d \leftarrow config$
 $n, m \leftarrow \lceil \frac{R}{config_p.goodput} \rceil, \lceil \frac{R}{config_d.goodput} \rceil$
 $best_plm \leftarrow (n, config_p, m, config_d)$
return $best_plm$

Method



Inter-machine Search

Algorithm 1 High Node-Affinity Placement Algorithm

Input: LLM G , #node limit per-instance N , #GPU per-node M , GPU memory capacity C , workload W , traffic rate R .

Output: the placement $best_plm$.

```
configp, configd ← 0, 0
for intra_op ∈ {1, 2, ..., M} do
    TP ← ...
    PP ← ...
    for inter_op ∈ {1, 2, ..., ⌈ N × M / intra_op ⌉} do
        if  $\frac{G.size}{inter\_op \times intra\_op} < C$  then
            config ← (inter_op, intra_op)
             $\hat{G} \leftarrow parallel(G, config)$ 
            config.goodput ← simu_prefill( $\hat{G}, W$ )
            if  $\frac{config_p.goodput}{config_p.num_gpus} < \frac{config.goodput}{config.num_gpus}$  then
                configp ← config
            config.goodput ← simu_decode( $\hat{G}, W$ )
            if  $\frac{config_d.goodput}{config_d.num_gpus} < \frac{config.goodput}{config.num_gpus}$  then
                configd ← config
        n, m ← ⌈  $\frac{R}{config_p.goodput}$  ⌉, ⌈  $\frac{R}{config_d.goodput}$  ⌉
        best_plm ← (n, configp, m, configd)
return best_plm
```

Simulate

TP

PP

Method



Inter-machine Search

Algorithm 1 High Node-Affinity Placement Algorithm

Input: LLM G , #node limit per-instance N , #GPU per-node M , GPU memory capacity C , workload W , traffic rate R .

Output: the placement $best_plm$.

```
configp, configd ← 0, 0
for intra_op ∈ {1, 2, ..., M} do
    TP ← ...
    PP ← ...
    for inter_op ∈ {1, 2, ..., ⌈ N × M / intra_op ⌉} do
        if  $\frac{G.size}{inter\_op \times intra\_op} < C$  then
            config ← (inter_op, intra_op)
             $\hat{G} \leftarrow parallel(G, config)$ 
            config.goodput ← simu_prefill( $\hat{G}, W$ )
            if  $\frac{config_p.goodput}{config_p.num_gpus} < \frac{config.goodput}{config.num_gpus}$  then
                configp ← config
            config.goodput ← simu_decode( $\hat{G}, W$ )
            if  $\frac{config_d.goodput}{config_d.num_gpus} < \frac{config.goodput}{config.num_gpus}$  then
                configd ← config
n, m ← ⌈  $\frac{R}{config_p.goodput}$  ⌉, ⌈  $\frac{R}{config_d.goodput}$  ⌉
best_plm ← (n, configp, m, configd)
return best_plm
```

Solution ←

Simulate ←

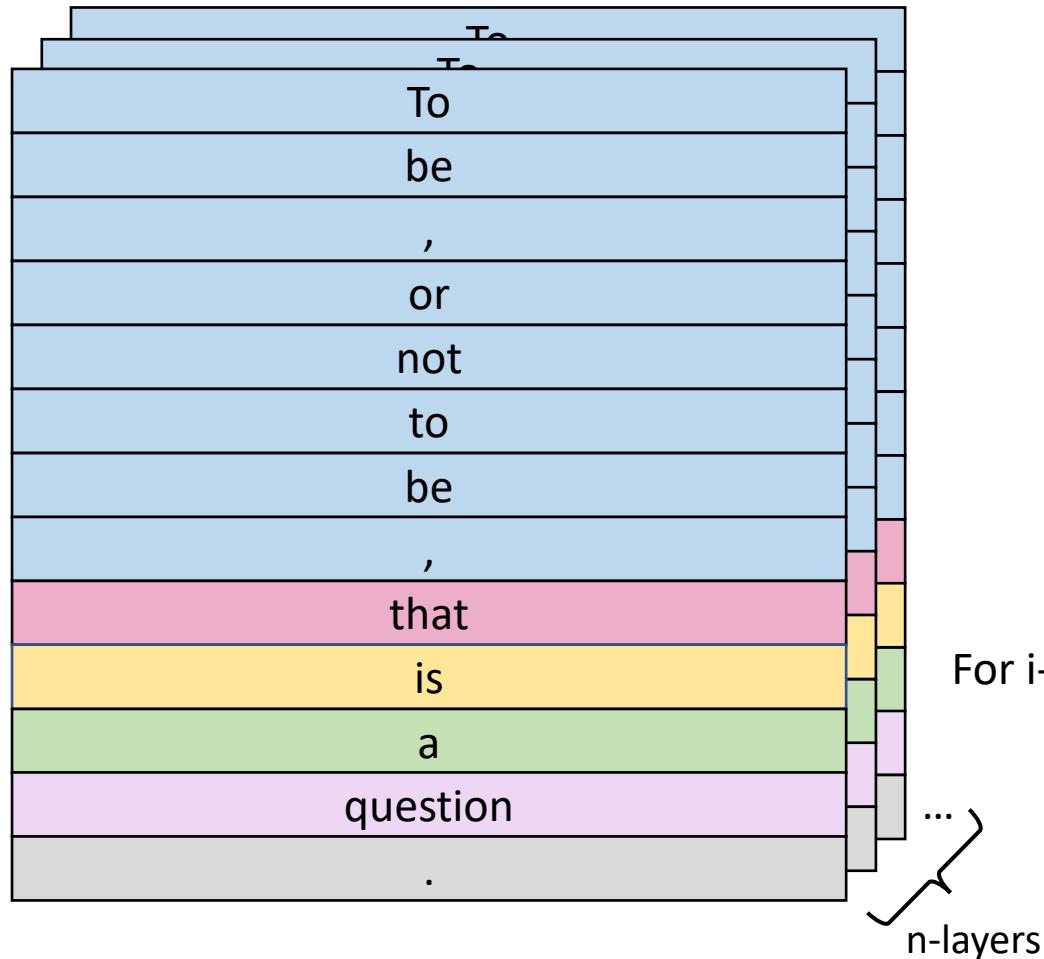
TP ←

PP ←

Method



Intra-machine Search



If LLM G (such as OPT-175B) is too large,
single machine cannot afford.

Cross-machine transmission is required,
increases the migration time 😭

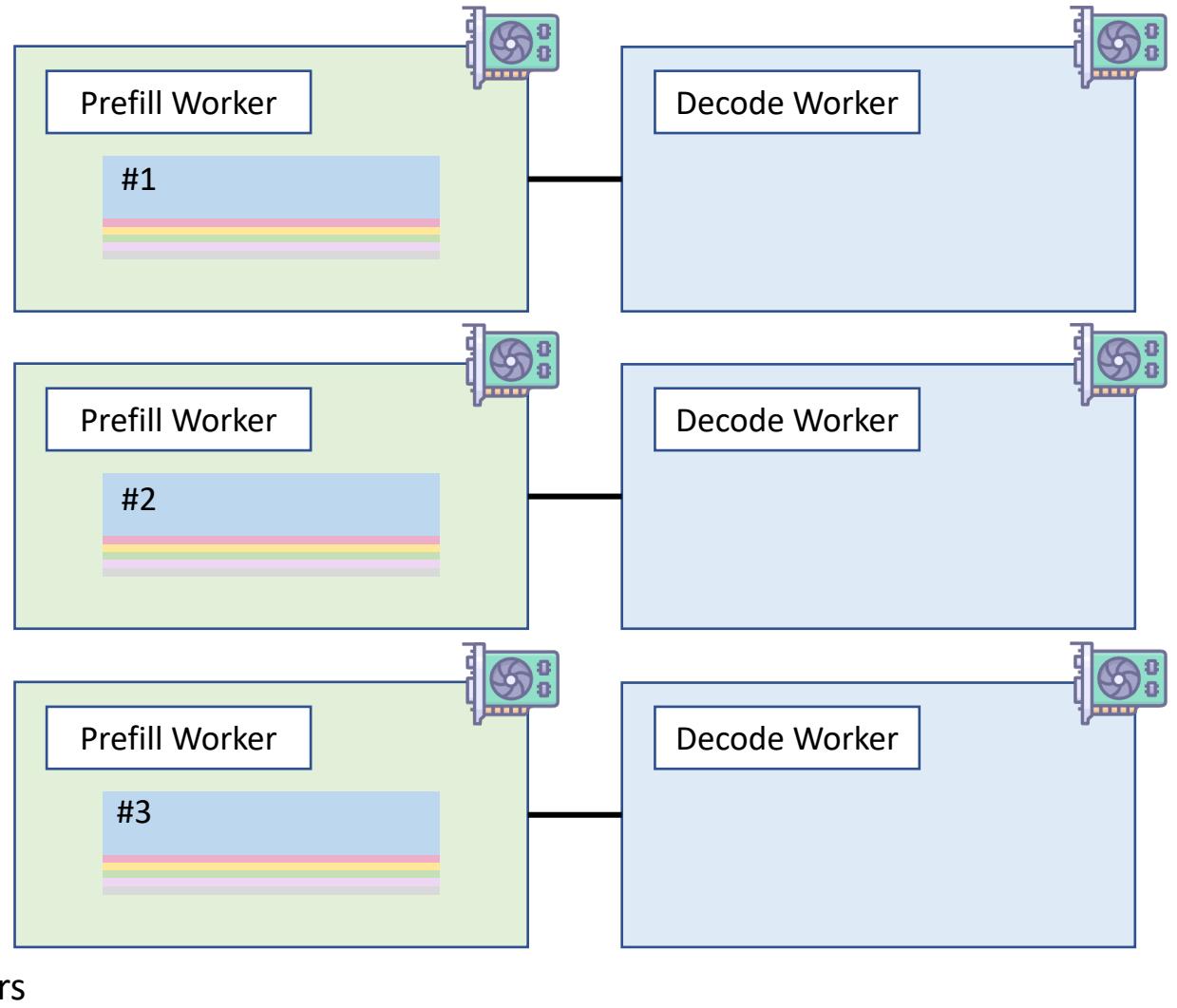
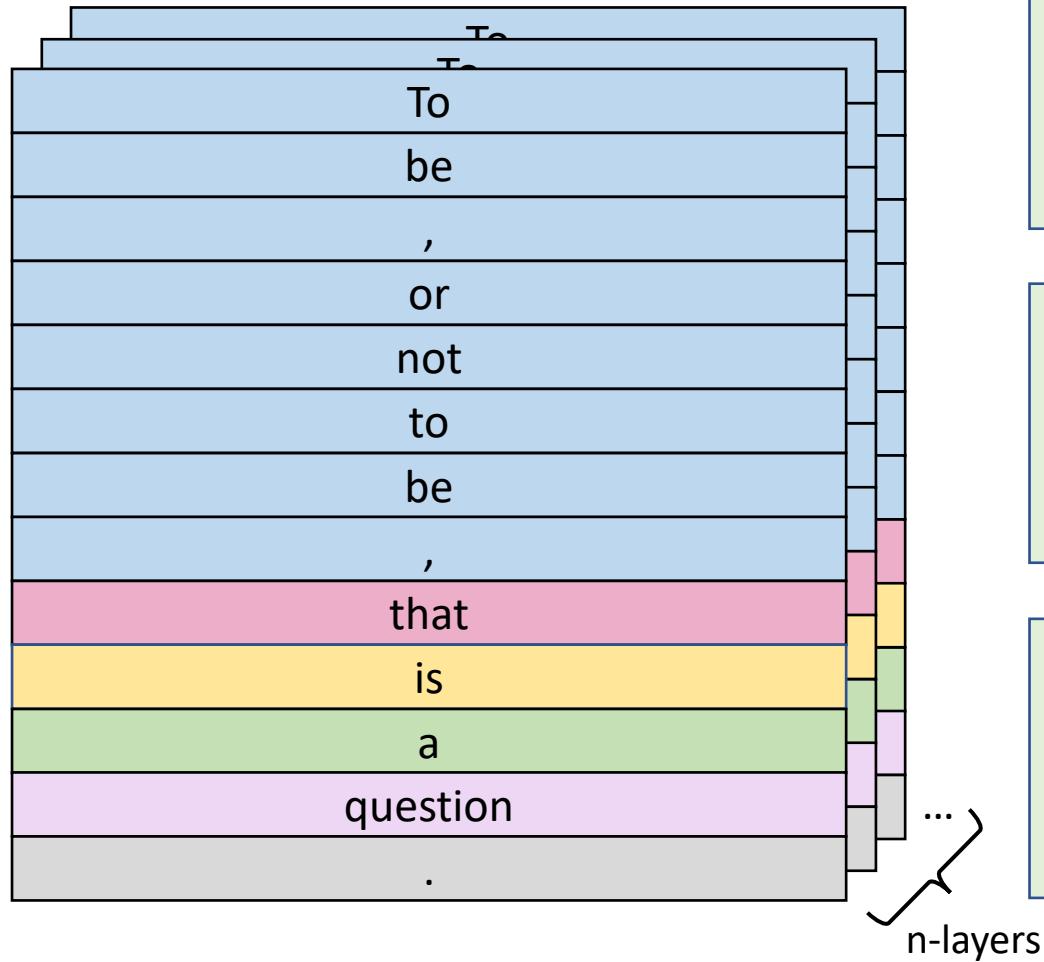
Need to place workers reasonably:

For i -th layer, make the Prefill worker and Decode worker with in the same
machine.

Method



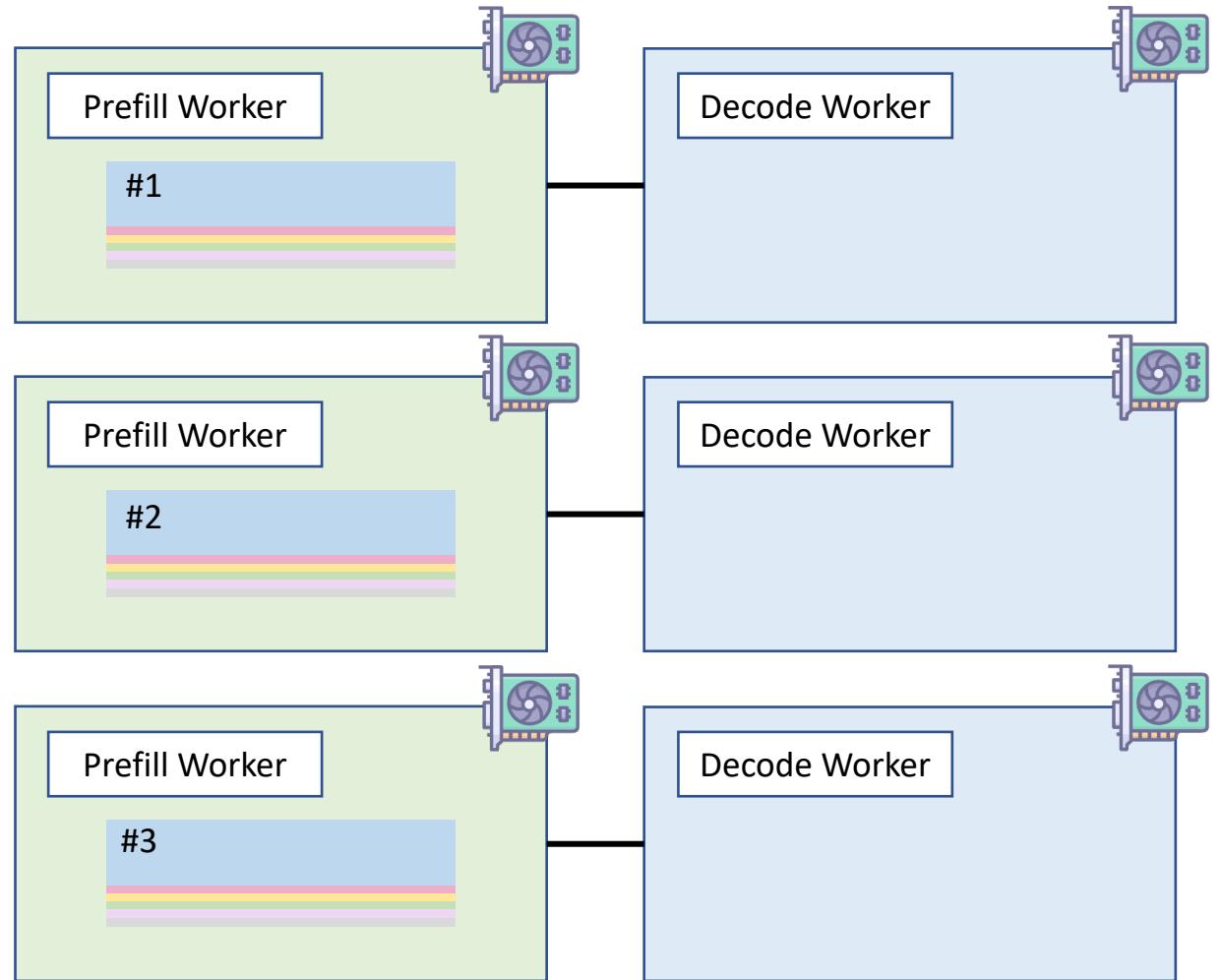
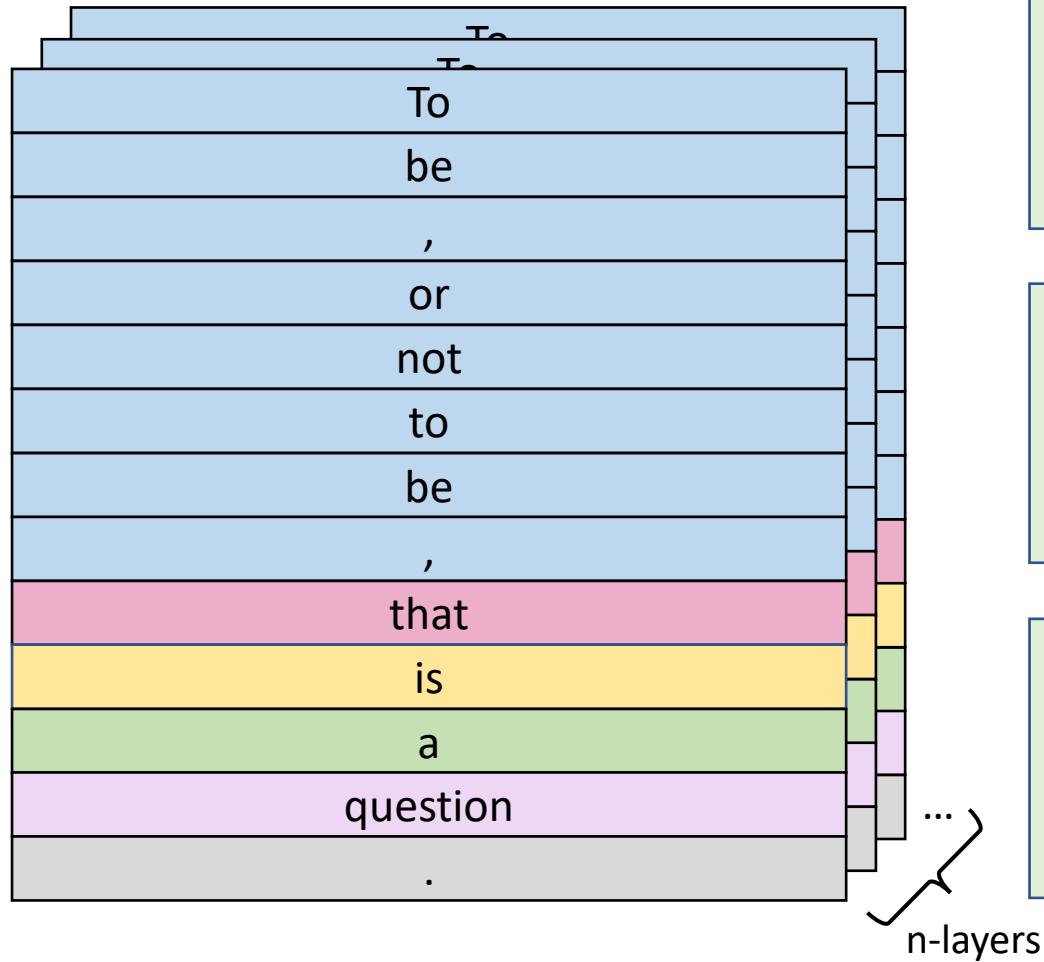
Intra-machine Search



Method



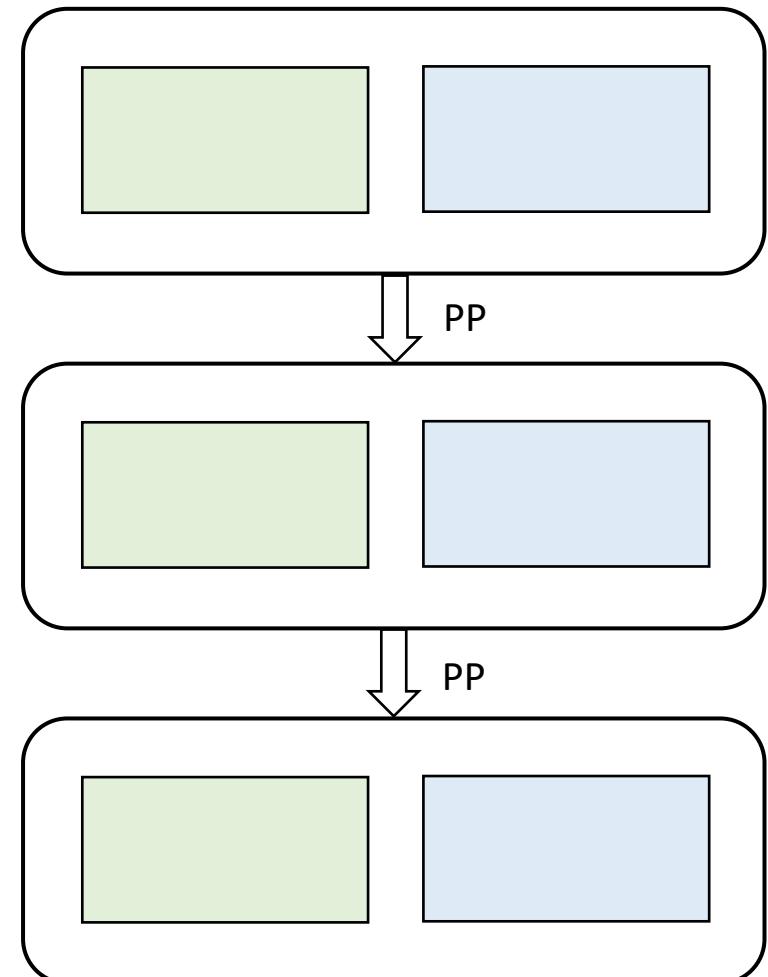
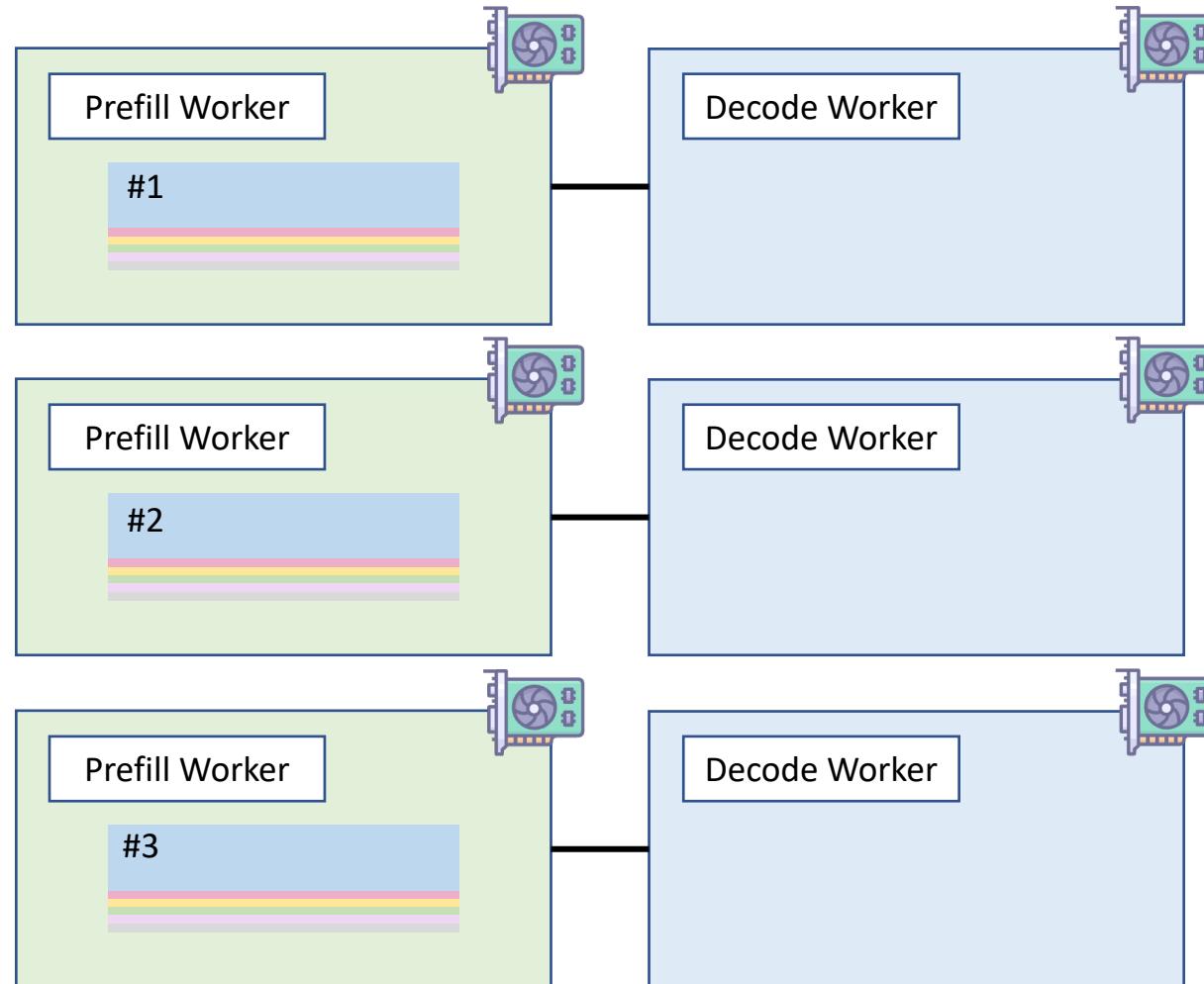
Intra-machine Search



Method



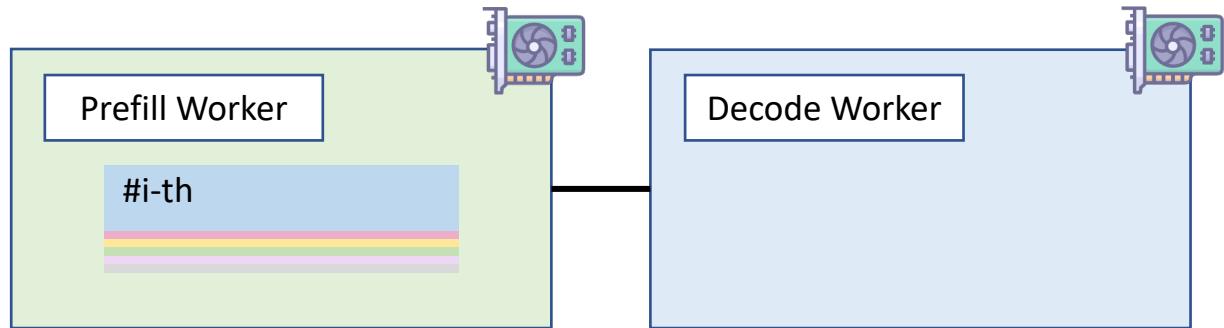
Intra-machine Search



Method



Intra-machine Search



Algorithm 2 Low Node-Affinity Placement Algorithm

Input: LLM G , #node limit per-instance N , #GPU per-node M , GPU memory capacity C , workload W , traffic rate R .

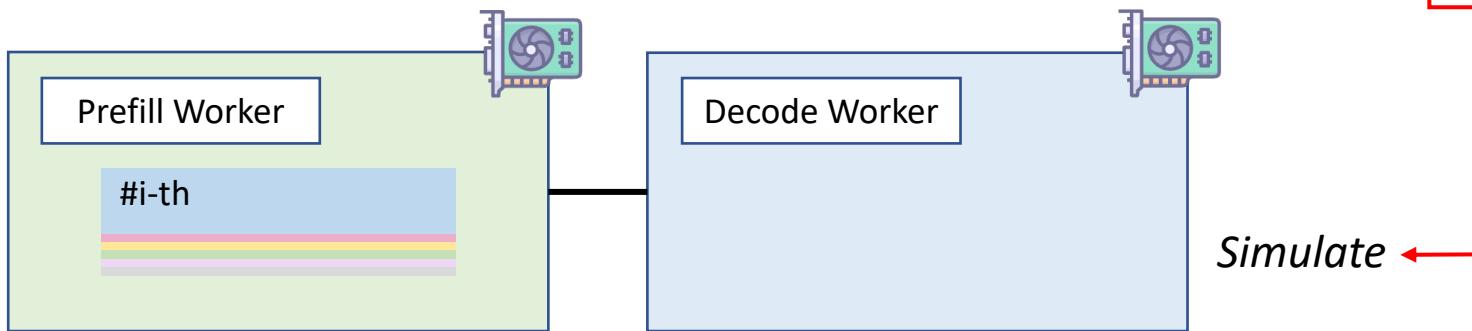
Output: the placement $best_plm$.

```
config* ← ∅  
for inter_op ∈ {1, 2, ..., N} do  
    P ← get_intra_node_configs(G, M, C, inter_op)  
    for Pp ∈ P do  
        for Pd ∈ P do  
            if Pp.num_gpus + Pd.num_gpus ≤ M then  
                config ← (inter_op, Pp, Pd)  
                Gp, Gd ← parallel(G, config)  
                config.goodput ← simulate(Gp, Gd, W)  
                if config.*goodput < config.goodput then  
                    config* ← config  
n ← ⌈ R / config.*goodput ⌉  
best_plm ← (n, config*)  
return best_plm
```

Method



Intra-machine Search



Algorithm 2 Low Node-Affinity Placement Algorithm

Input: LLM G , #node limit per-instance N , #GPU per-node M , GPU memory capacity C , workload W , traffic rate R .

Output: the placement $best_plm$.

$config^* \leftarrow \emptyset$

for $inter_op \in \{1, 2, \dots, N\}$ **do**

$\mathcal{P} \leftarrow get_intra_node_configs(G, M, C, inter_op)$

for $P_p \in \mathcal{P}$ **do**

for $P_d \in \mathcal{P}$ **do**

if $P_p.num_gpus + P_d.num_gpus \leq M$ **then**

$config \leftarrow (inter_op, P_p, P_d)$

$\hat{G}_p, \hat{G}_d \leftarrow parallel(G, config)$

$config.goodput \leftarrow simulate(\hat{G}_p, \hat{G}_d, W)$

if $\frac{config.^*goodput}{config.^*num_gpus} < \frac{config.goodput}{config.num_gpus}$ **then**
 $config^* \leftarrow config$

$n \leftarrow \lceil \frac{R}{config.^*goodput} \rceil$

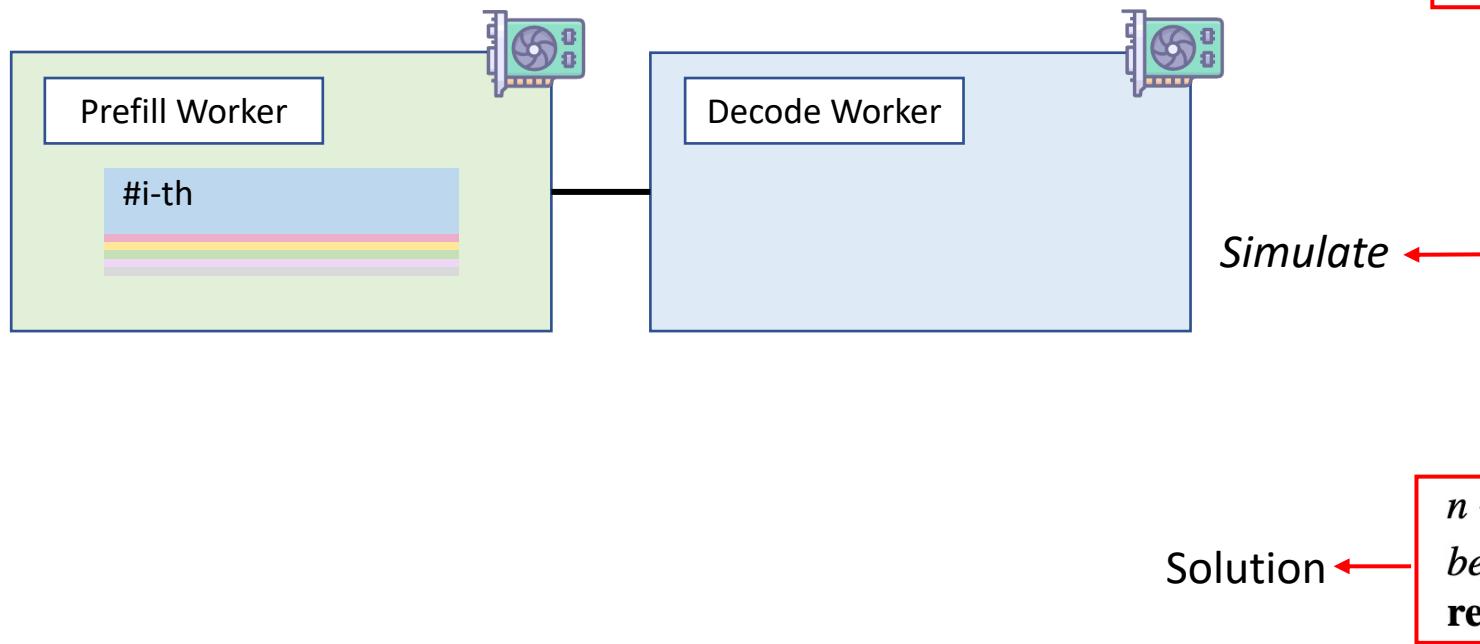
$best_plm \leftarrow (n, config^*)$

return $best_plm$

Method



Intra-machine Search



Algorithm 2 Low Node-Affinity Placement Algorithm

Input: LLM G , #node limit per-instance N , #GPU per-node M , GPU memory capacity C , workload W , traffic rate R .

Output: the placement $best_plm$.

```
config* ← Ø
for inter_op ∈ {1, 2, ..., N} do
    P ← get_intra_node_configs(G, M, C, inter_op)
    for Pp ∈ P do
        for Pd ∈ P do
            if Pp.num_gpus + Pd.num_gpus ≤ M then
                config ← (inter_op, Pp, Pd)
                Gp, Gd ← parallel(G, config)
                config.goodput ← simulate(Gp, Gd, W)
                if config.*goodput < config.goodput
                    config* ← config
n ← ⌈ R / config.*goodput ⌉
best_plm ← (n, config*)
return best_plm
```

Contents



- Background
- Problem
- Opportunity
- Challenge
- Method
- **Evaluation**
- Thinking

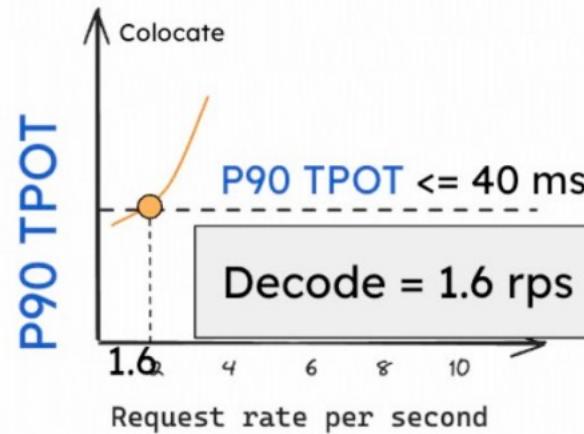
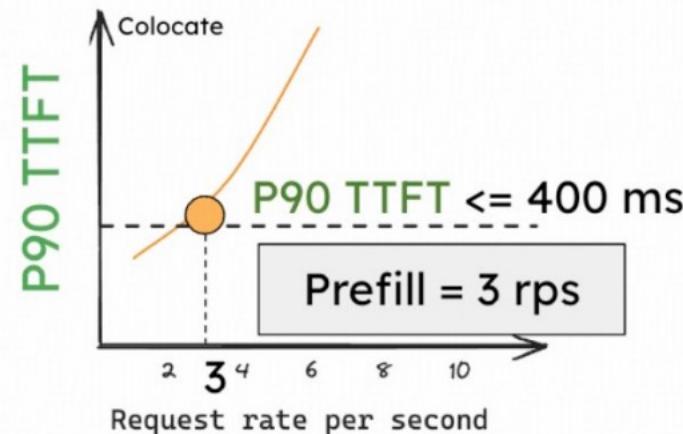
Evaluation



LLM App	Dataset	TTFT	TPOT
Chatbot	shareGPT	Tight	Medium
Code completion	HumanEval	Tight	Tight
Summarization	LongBench	Loose	Medium

Model	Dataset	Prefill		Decoding	
		TP	PP	TP	PP
OPT-13B	ShareGPT	2	1	1	1
OPT-66B	ShareGPT	4	1	2	2
OPT-66B	LongBench	4	1	2	2
OPT-66B	HumanEval	4	1	2	2
OPT-175B	ShareGPT	3	3	4	3

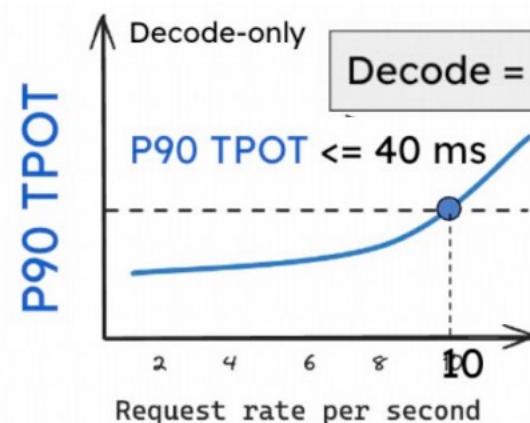
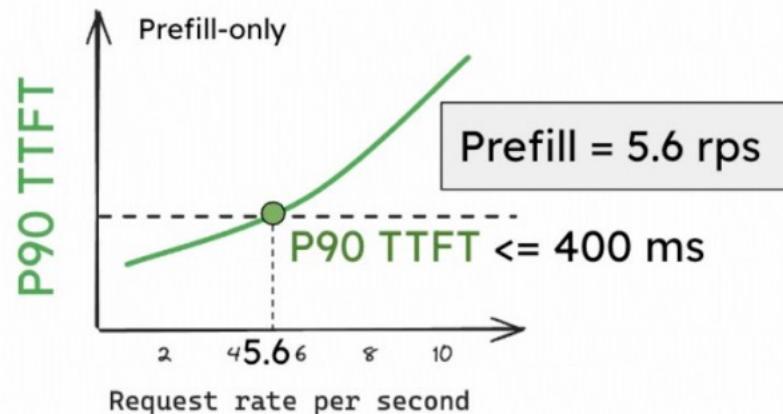
Evaluation



Colocation



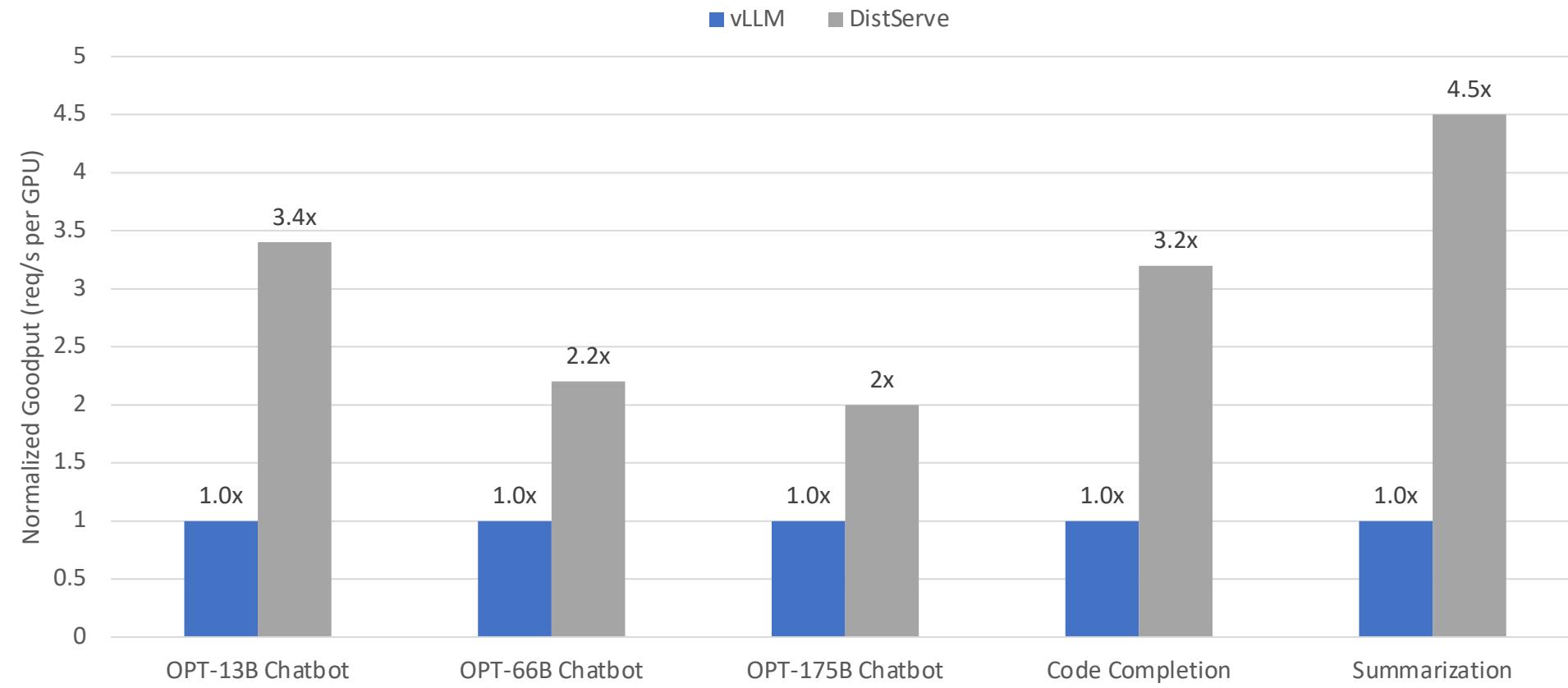
Max System rps
= Min(Prefill, Decode)
= 1.6 rps / GPU



Disaggregation (2P1D)

Max System rps
= Min (5.6 x 2, 10) rps / 3 GPU
= 3.3 rps / GPU

Evaluation



Contents



- Background
- Problem
- Opportunity
- Challenge
- Method
- Evaluation
- Thinking



- 文章有什么问题，基于这个 paper 还能做什么
 - P/D分离服务架构，会导致可维护性下降，整个引擎耦合性变高
 - KV transfer 可以增量更新，但实际上传输开销并不大，这种方式是否有益还需要进行 trade-off 判断
- Paper 的 idea 能不能应用在自己的工作上面
 - 重新考虑不同应用的 “goodput”，针对不同的 SLO 侧重不同的策略
- 这个 paper 能不能泛化
 - 根据模型各层/各阶段计算性质不同，对资源的偏好不同，应用不同的推理手段。例如端侧的 NPU 适合做 prefill 性质的计算，需要重新回顾下 powerinfer2（个人认为是在端侧做PD分离的一种尝试）



Q & A

**Presenter: Yitao Wang
2024.3.7**