# Mooncake: Trading More Storage for Less Computation — A KVCache-centric Architecture for Serving LLM Chatbot

Ruoyu Qin, Zheming Li, Weiran He, Jialei Cui, Feng Ren, Mingxing Zhang, Yongwei Wu, Weimin Zheng, Xinran Xu
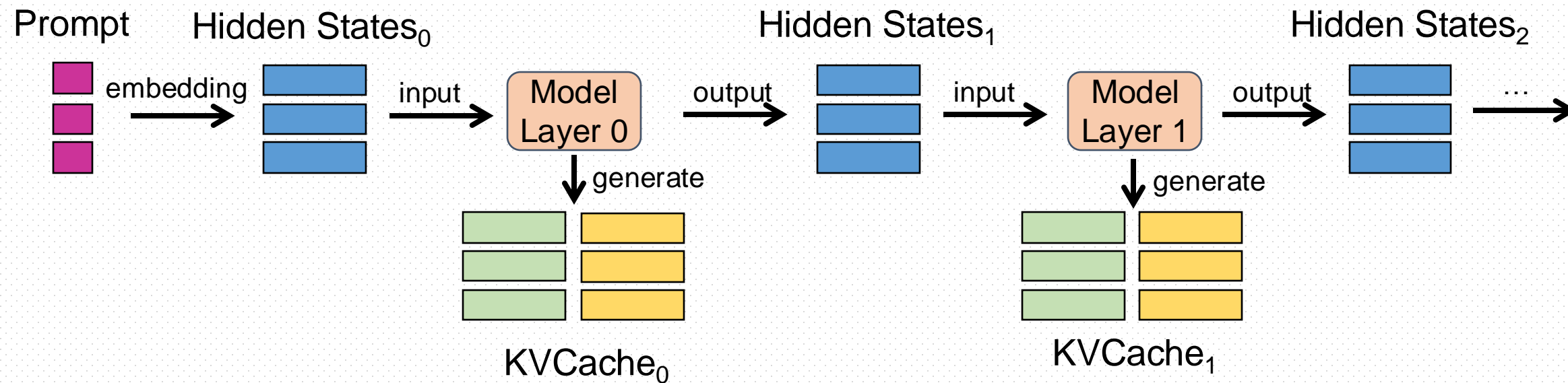
Shared by Juncheng Zhang

# Outline

- ❑ **Prefix Caching**

- ❑ **PD Disaggregation**

- ❑ **Evaluation**

- ❑ **Discussion**

# LLM Inference

❑ LLM inference process can be divided into two phases
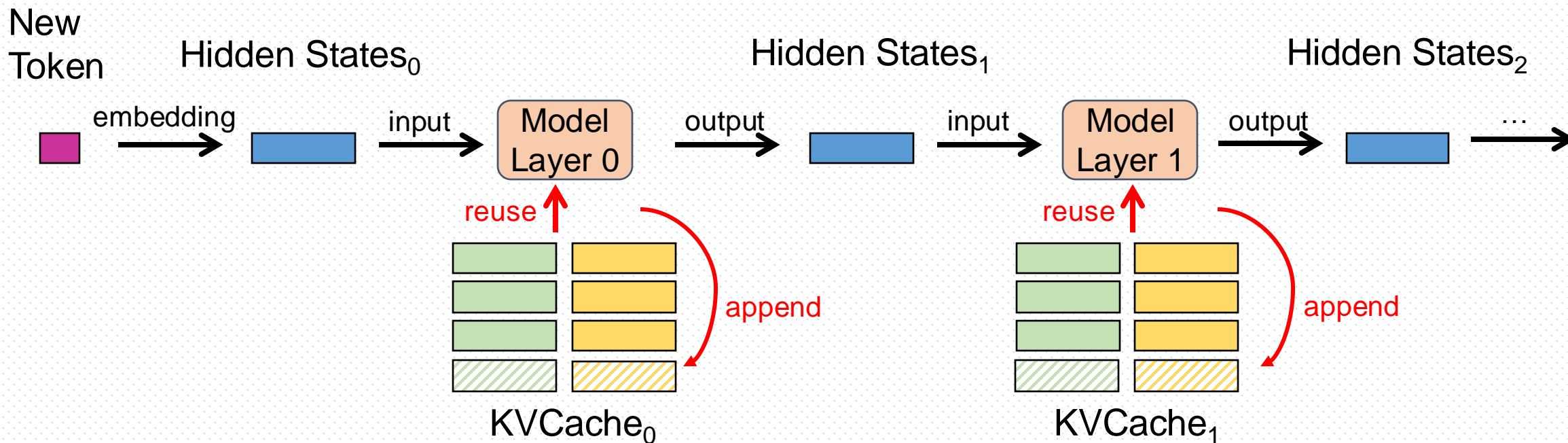
❖ **Prefill Phase**: generate KVCache and output first token

Prompt    Hidden States$_0$                Hidden States$_1$                Hidden States$_2$

embedding → input → | Model Layer 0 | → output → input → | Model Layer 1 | → output → …

↓ generate                                    ↓ generate

KVCache$_0$                                    KVCache$_1$

# LLM Inference

❑ LLM inference process can be divided into two phases

❖ **Decode Phase**: generate next token

New Token

Hidden States$_0$

Hidden States$_1$

Hidden States$_2$

embedding → input → Model Layer 0 → output → input → Model Layer 1 → output → …

reuse

append

KVCache$_0$

reuse

append

KVCache$_1$

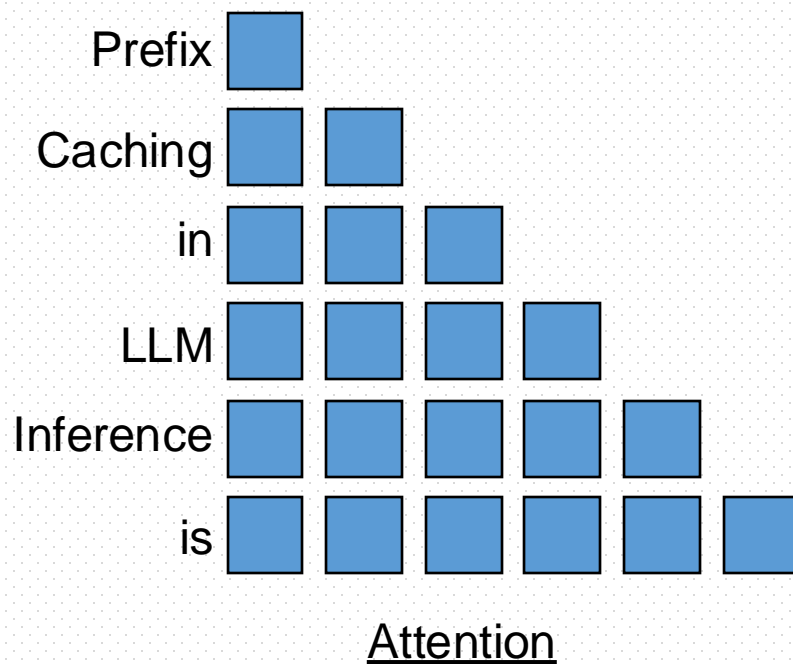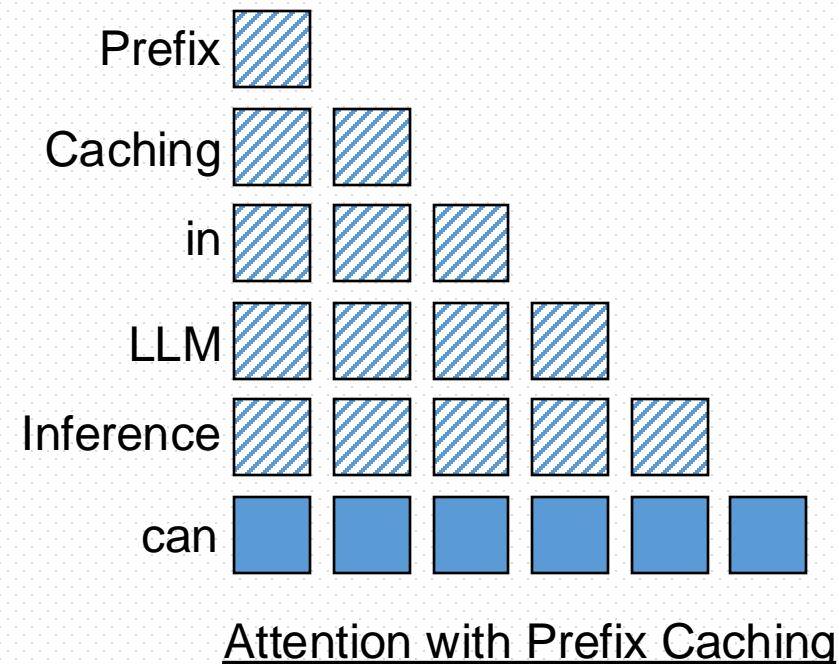# Prefix Caching in LLM Inference

❑Requests with the same prefix can shared the same KVCache

❖Computation reduction in prefill phase
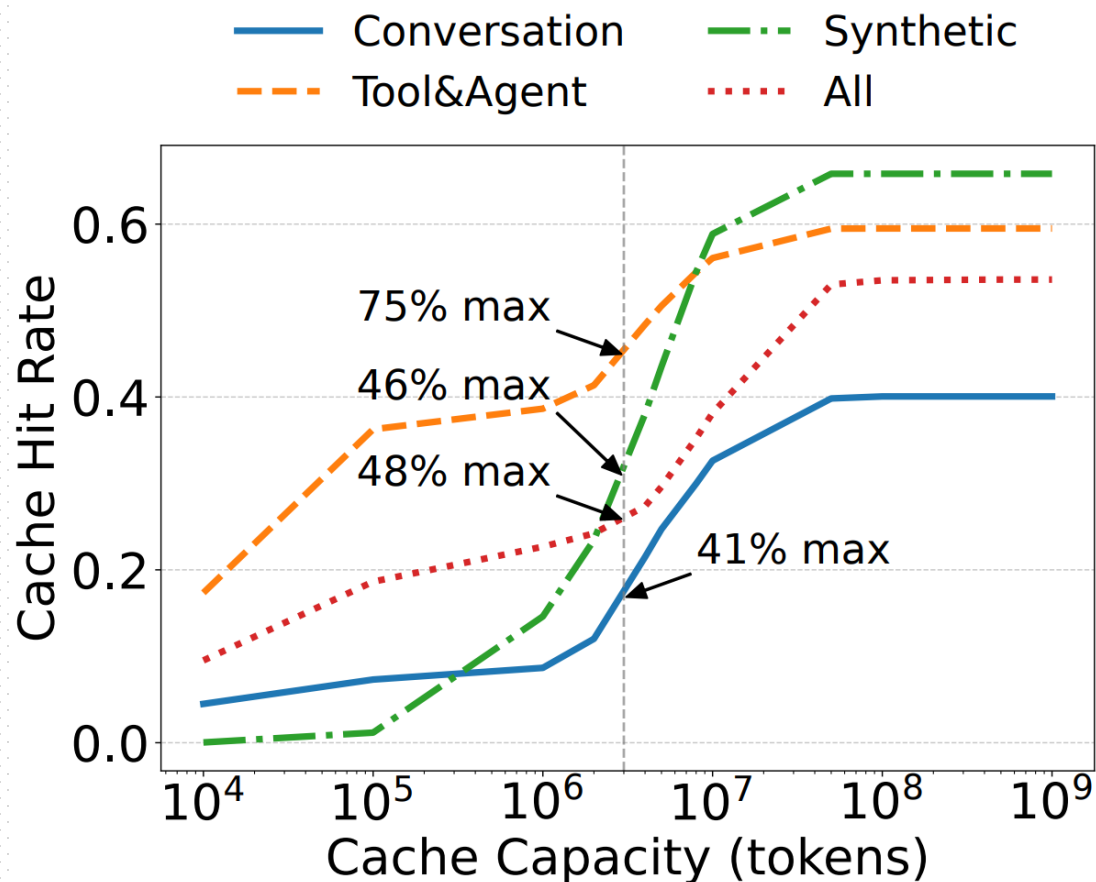


Reuse!

Attention

Attention with Prefix Caching

# Prefix Caching in LLM Inference

❑ Not easy in Kimi's real system deployment !!!



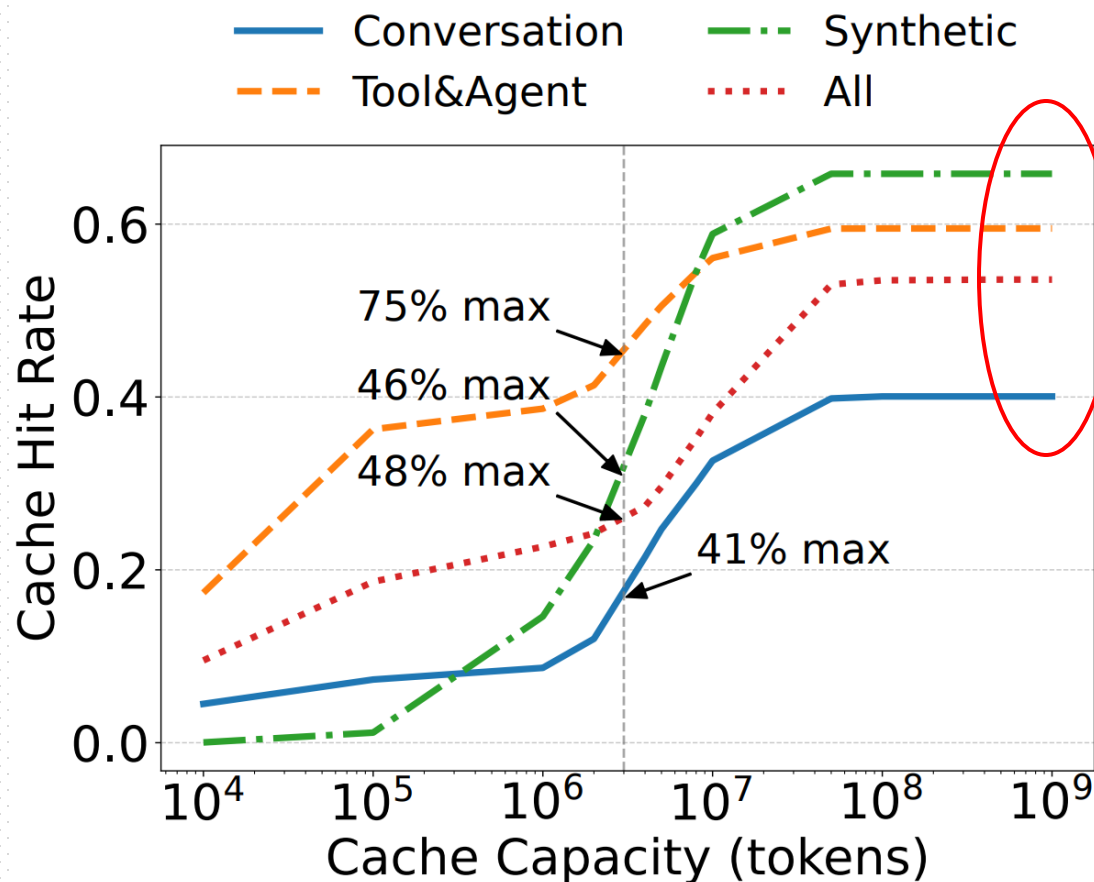**Conversation**: collected from real-world online conversation requests

**Tool & Agent**: collected from real-world online requests that include tool use

**Synthetic**: synthesized from publicly available long context datasets

# Prefix Caching in LLM Inference
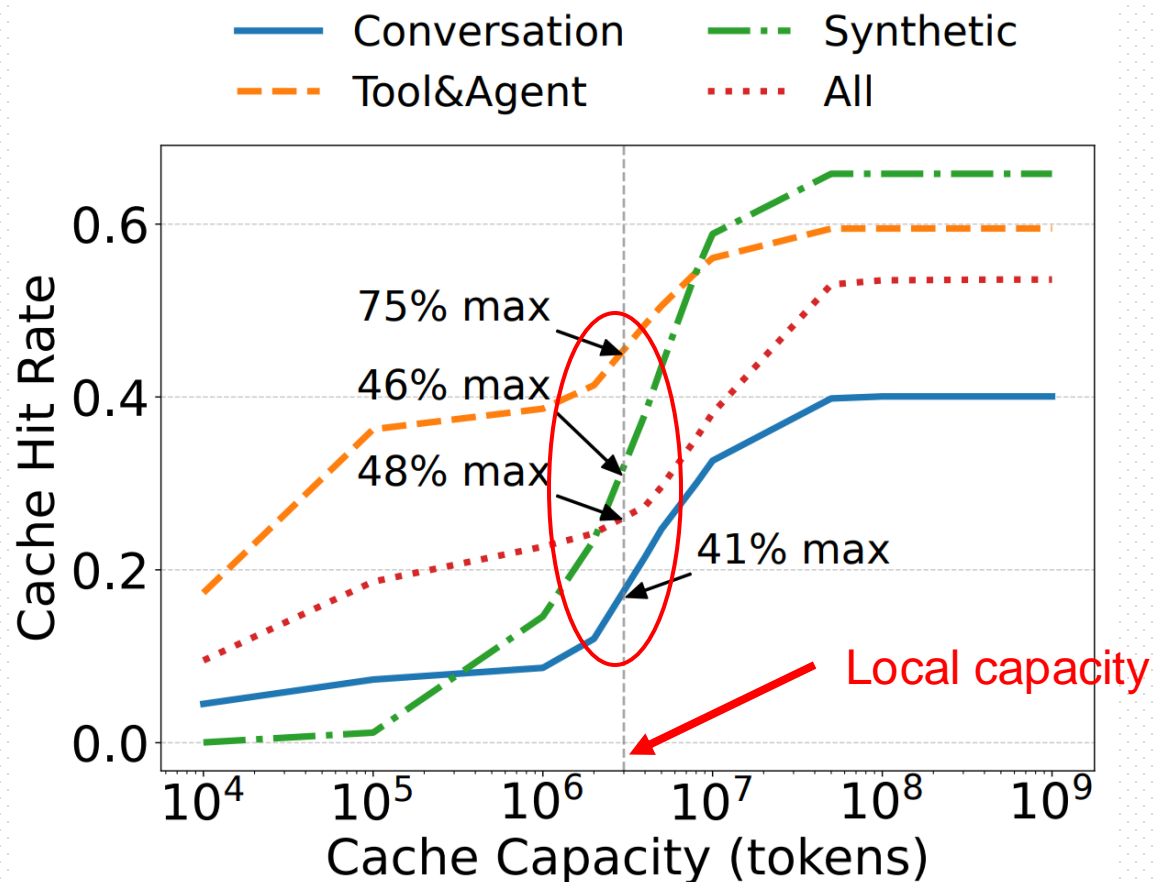
❑Not easy in Kimi's real system deployment !!!



In ideal case, around 50% of the token's KVCache can be reused
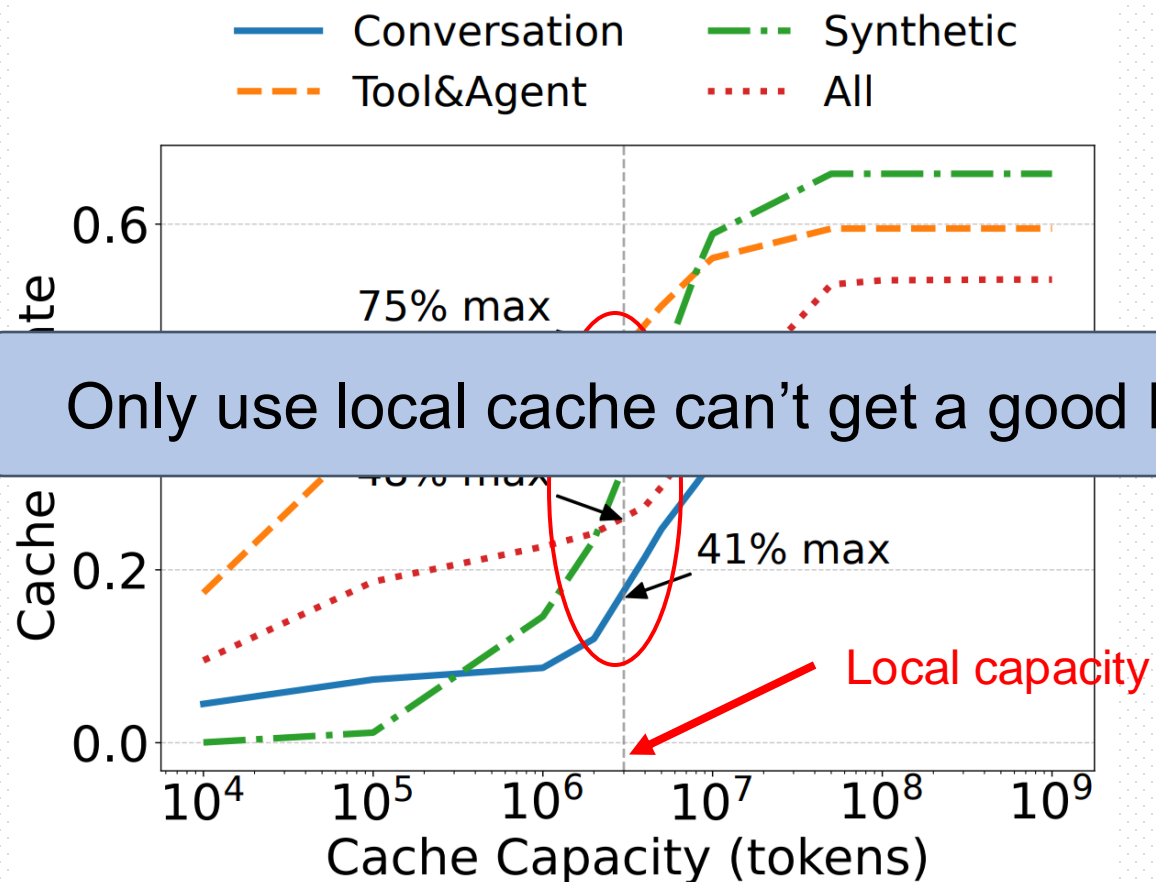
# Prefix Caching in LLM Inference
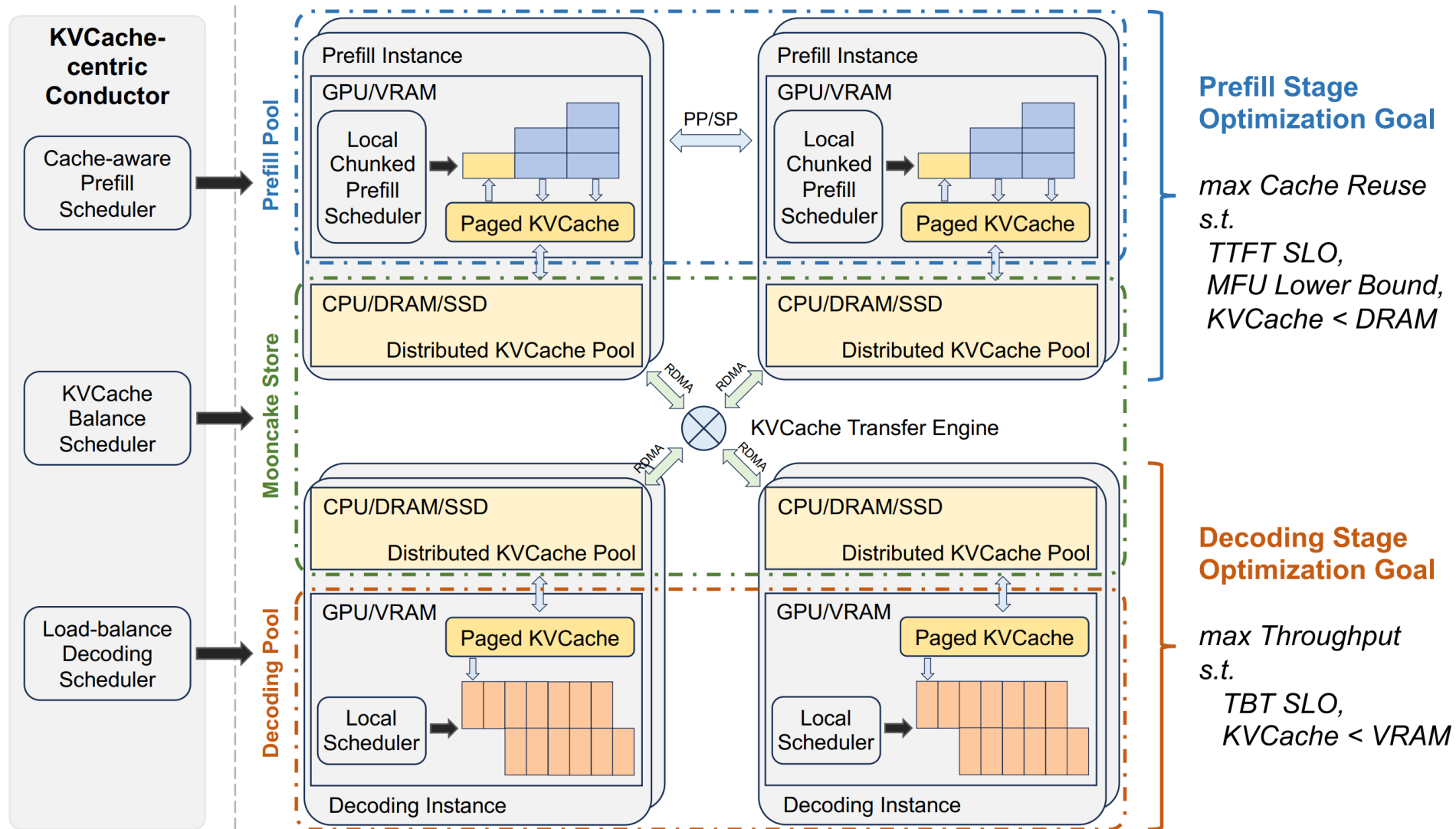
❏Not easy in Kimi's real system deployment !!!

# Prefix Caching in LLM Inference

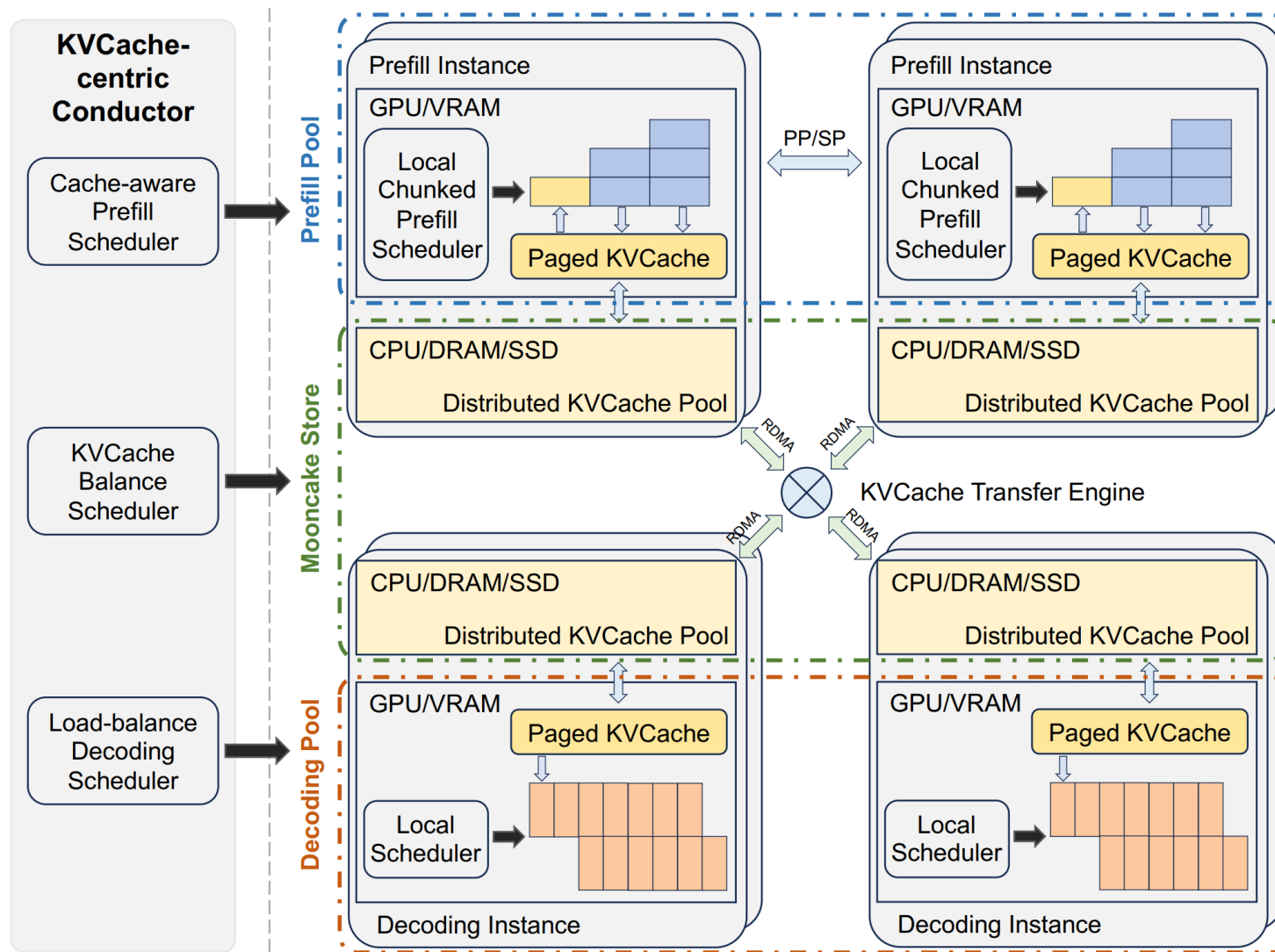❏ Not easy in Kimi's real system deployment !!!



Only use local cache can't get a good hit ratio, Kimi needs a *global cache*!

# Mooncake Architecture



**KVCache-centric Conductor**

Cache-aware Prefill Scheduler

KVCache Balance Scheduler

Load-balance Decoding Scheduler

**Prefill Pool**

Prefill Instance

GPU/VRAM

Local Chunked Prefill Scheduler

Paged KVCache

PP/SP

Prefill Instance

GPU/VRAM

Local Chunked Prefill Scheduler

Paged KVCache

CPU/DRAM/SSD

Distributed KVCache Pool

CPU/DRAM/SSD

Distributed KVCache Pool

**Mooncake Store**

RDMA  RDMA

KVCache Transfer Engine

RDMA  RDMA

**Decoding Pool**

CPU/DRAM/SSD

Distributed KVCache Pool

CPU/DRAM/SSD

Distributed KVCache Pool

GPU/VRAM

Paged KVCache

Local Scheduler

Decoding Instance

GPU/VRAM

Paged KVCache

Local Scheduler

Decoding Instance

**Prefill Stage Optimization Goal**

*max Cache Reuse s.t.*
  *TTFT SLO,*
  *MFU Lower Bound,*
  *KVCache < DRAM*

**Decoding Stage Optimization Goal**

*max Throughput s.t.*
  *TBT SLO,*
  *KVCache < VRAM*

# Mooncake Architecture



For simplicity, we leave the PD disaggregation in later discussion

# Mooncake Architecture



1. User request

2. Select *Prefill* and *Decode* instances

# Mooncake Architecture



3. Prefill phase utilize *Prefix Caching*

4. Transfer *KVCache* from Prefill instance to Decode instance
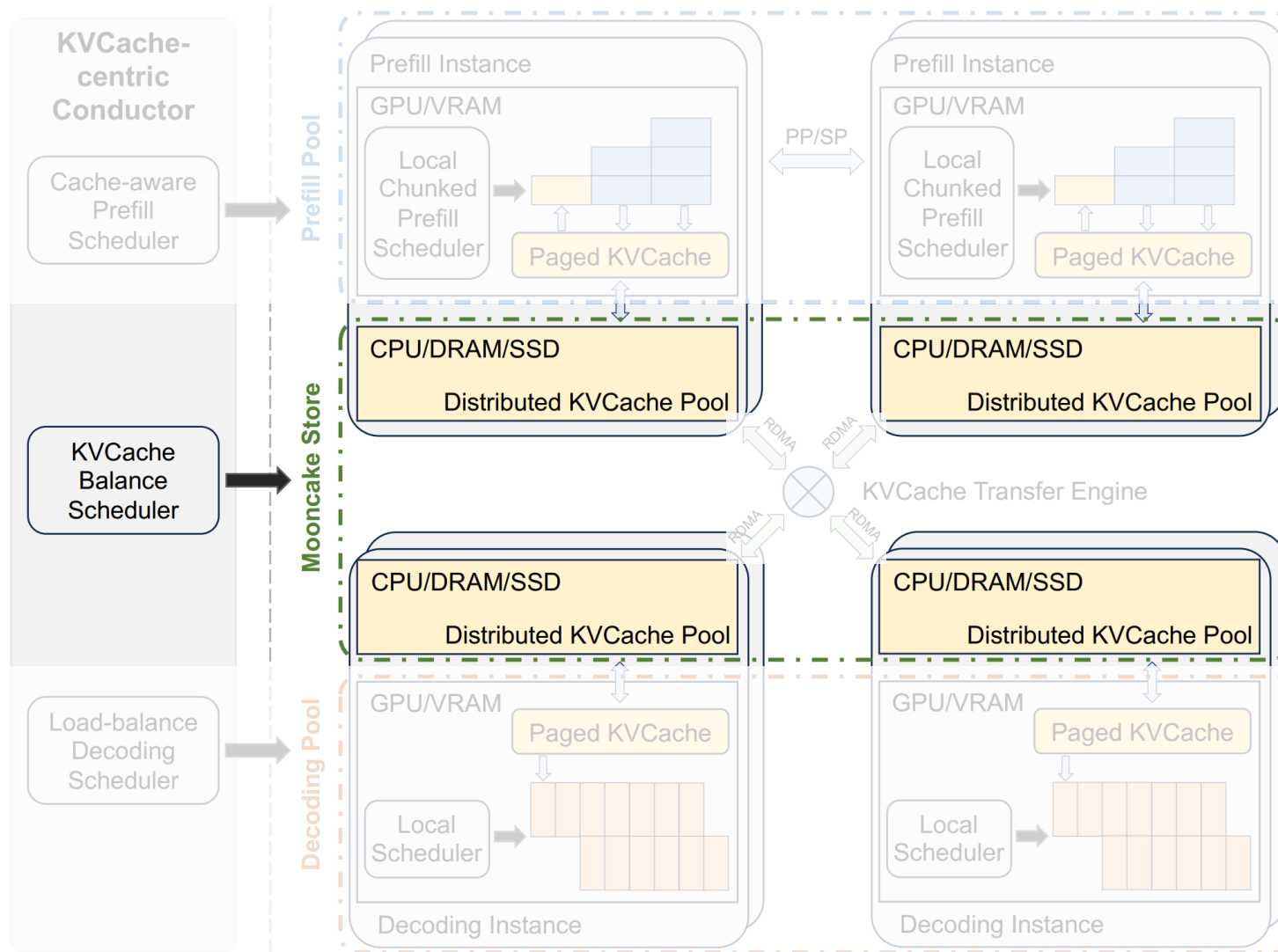
5. Finish inference, return to user

# Challenge in global cache

❑ Aggregate all the usable resource

   ❖ A distributed multi-layer *KVCache pool*

❑ KVCache needs to be transferred between different machines

   ❖ A low latency, high bandwidth *transfer engine*

❑ User request dispatch should consider prefix cache

   ❖ *KVCache-aware scheduling*

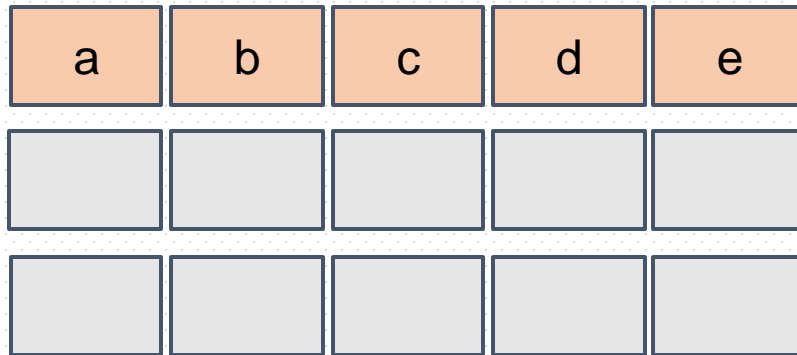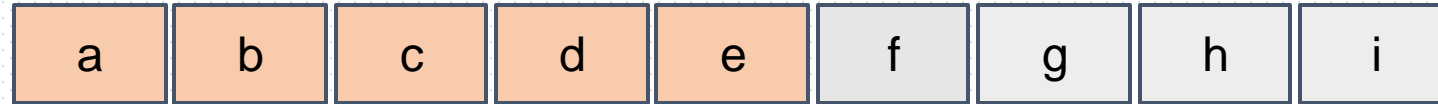# Mooncake Architecture

# Mooncake Store

❑ Block-level KVCache management
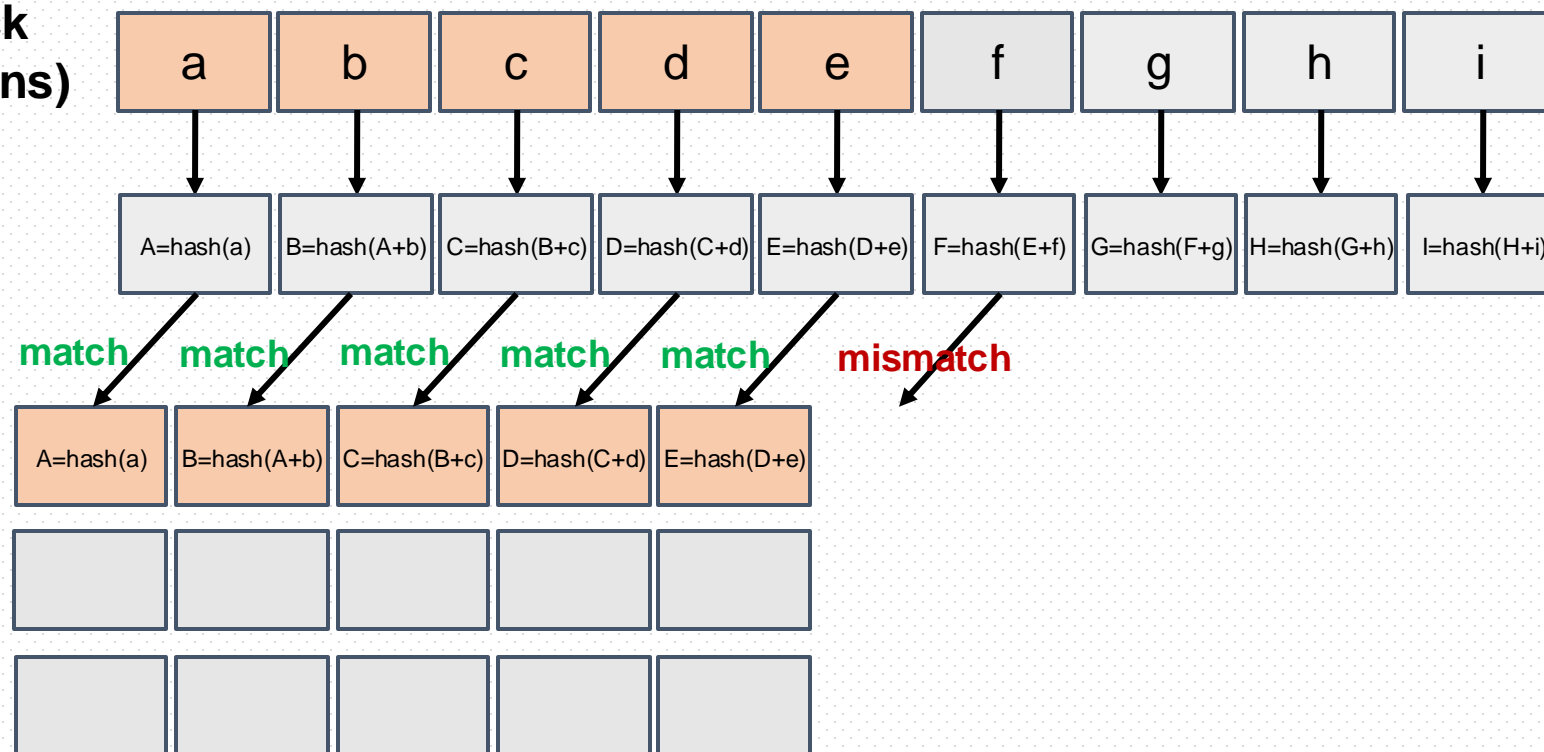
**Token block (16~512 tokens)**

| a | b | c | d | e | f | g | h | i |
|---|---|---|---|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|
| | | | | |
| | | | | |

# Mooncake Store

❏ Block-level KVCache management

❏ Prefix-hashed for fast match and deduplication

**Token block (16~512 tokens)**

| a | b | c | d | e | f | g | h | i |

| A=hash(a) | B=hash(A+b) | C=hash(B+c) | D=hash(C+d) | E=hash(D+e) | F=hash(E+f) | G=hash(F+g) | H=hash(G+h) | I=hash(H+i) |

**match**  **match**  **match**  **match**  **match**  **mismatch**

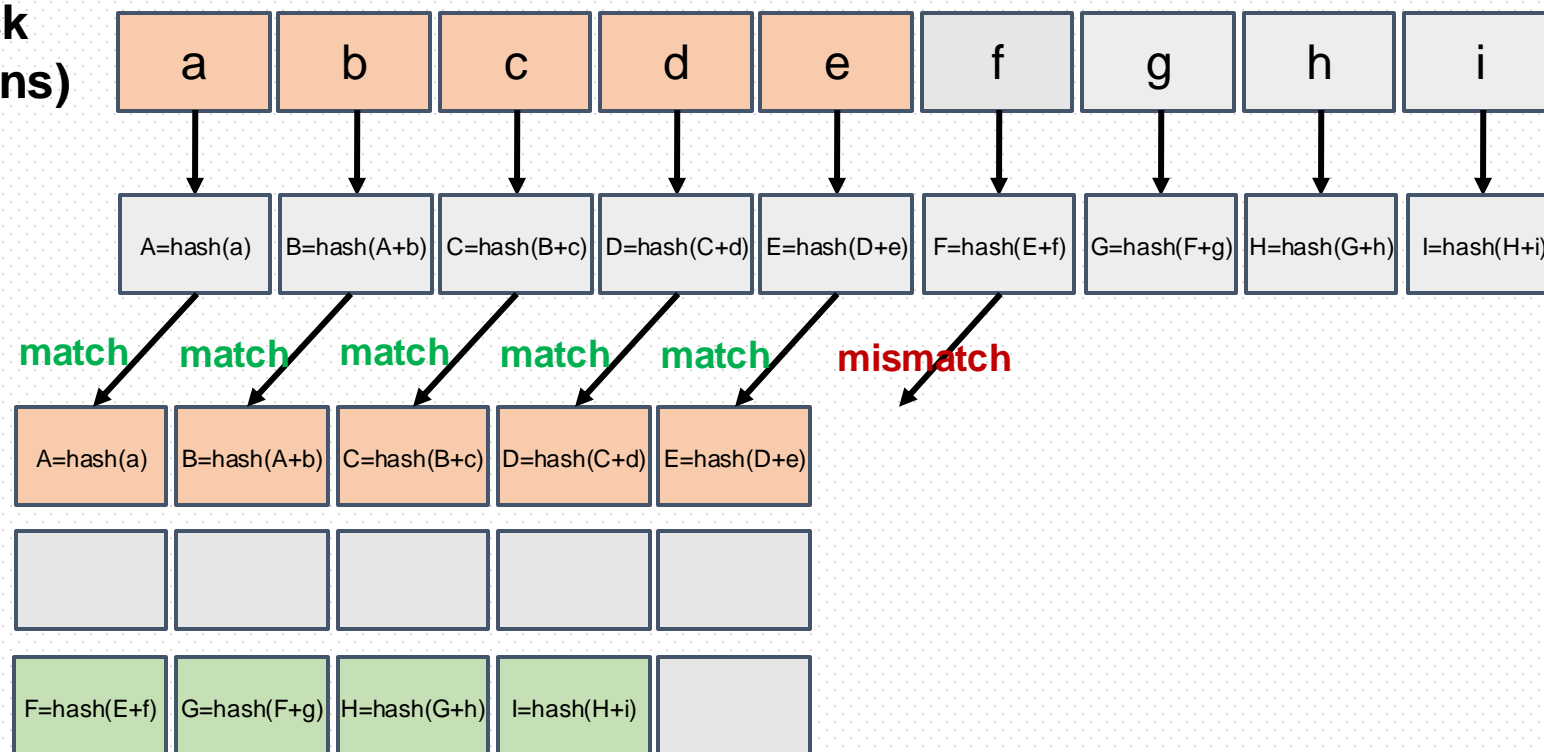| A=hash(a) | B=hash(A+b) | C=hash(B+c) | D=hash(C+d) | E=hash(D+e) |

# Mooncake Store

❑ Block-level KVCache management

❑ Prefix-hashed for fast match and deduplication

**Token block (16~512 tokens)**

| a | b | c | d | e | f | g | h | i |
|---|---|---|---|---|---|---|---|---|

| A=hash(a) | B=hash(A+b) | C=hash(B+c) | D=hash(C+d) | E=hash(D+e) | F=hash(E+f) | G=hash(F+g) | H=hash(G+h) | I=hash(H+i) |
|---|---|---|---|---|---|---|---|---|

**match** **match** **match** **match** **match** **mismatch**

| A=hash(a) | B=hash(A+b) | C=hash(B+c) | D=hash(C+d) | E=hash(D+e) |
|---|---|---|---|---|

| F=hash(E+f) | G=hash(F+g) | H=hash(G+h) | I=hash(H+i) | |
|---|---|---|---|---|

# Mooncake Store

❑Block-level KVCache management

❑Prefix-hashed for fast match and deduplication

**Token block (16~512 tokens)**

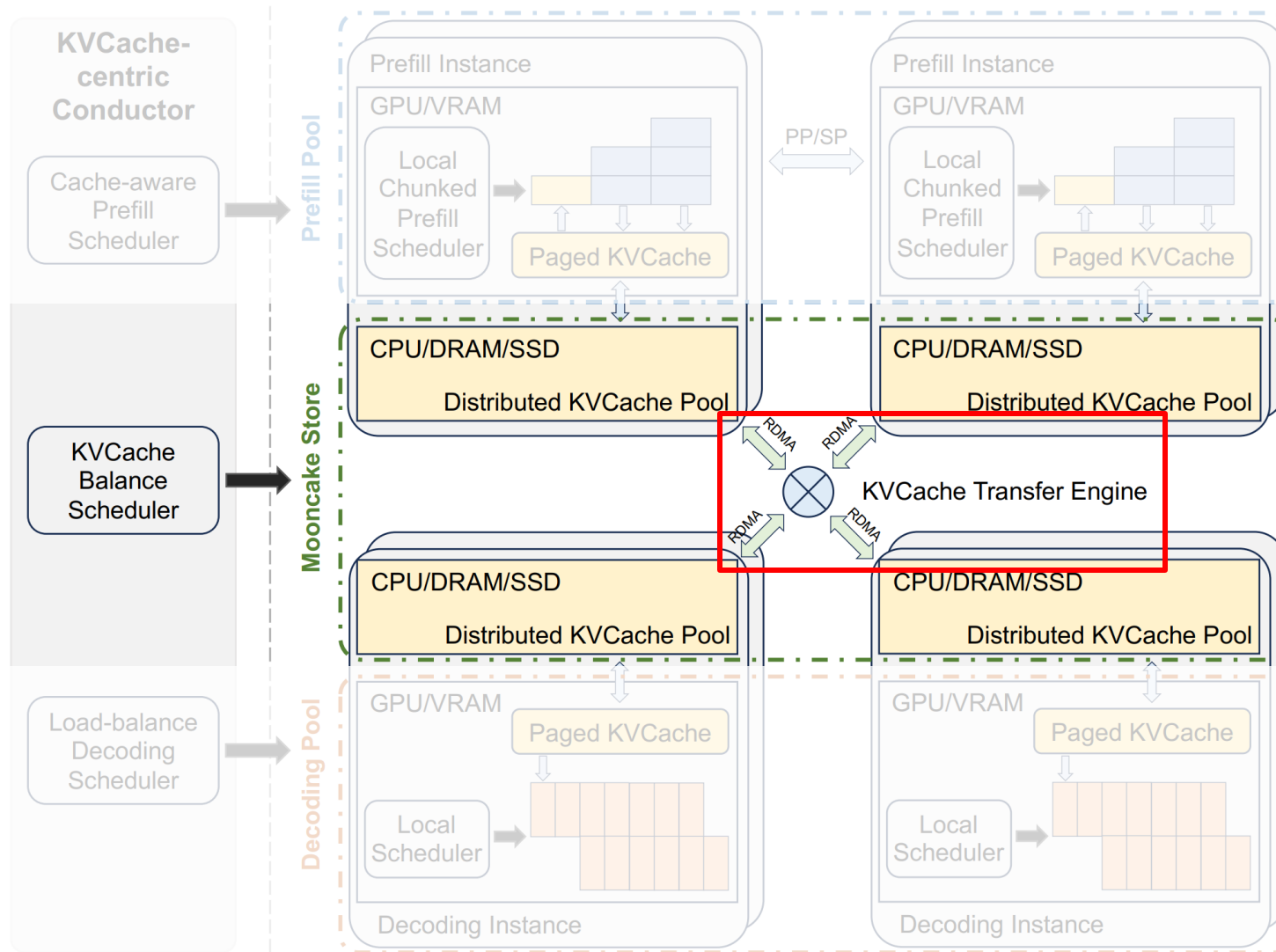| a | b | c | d | e | f | g | h | i |
|---|---|---|---|---|---|---|---|---|

| A=hash(a) | B=hash(A+b) | C=hash(B+c) | D=hash(C+d) | E=hash(D+e) | F=hash(E+f) | G=hash(F+g) | H=hash(G+h) | I=hash(H+i) |

match  match  match  match  match  **mismatch**

| A=hash(a) | B=hash(A+b) | C=hash(B+c) | D=hash(C+d) | E=hash(D+e) |

**Transfer()**

| F=hash(E+f) | G=hash(F+g) | H=hash(G+h) | I=hash(H+i) | |

**Get()**

Client

**Put()**

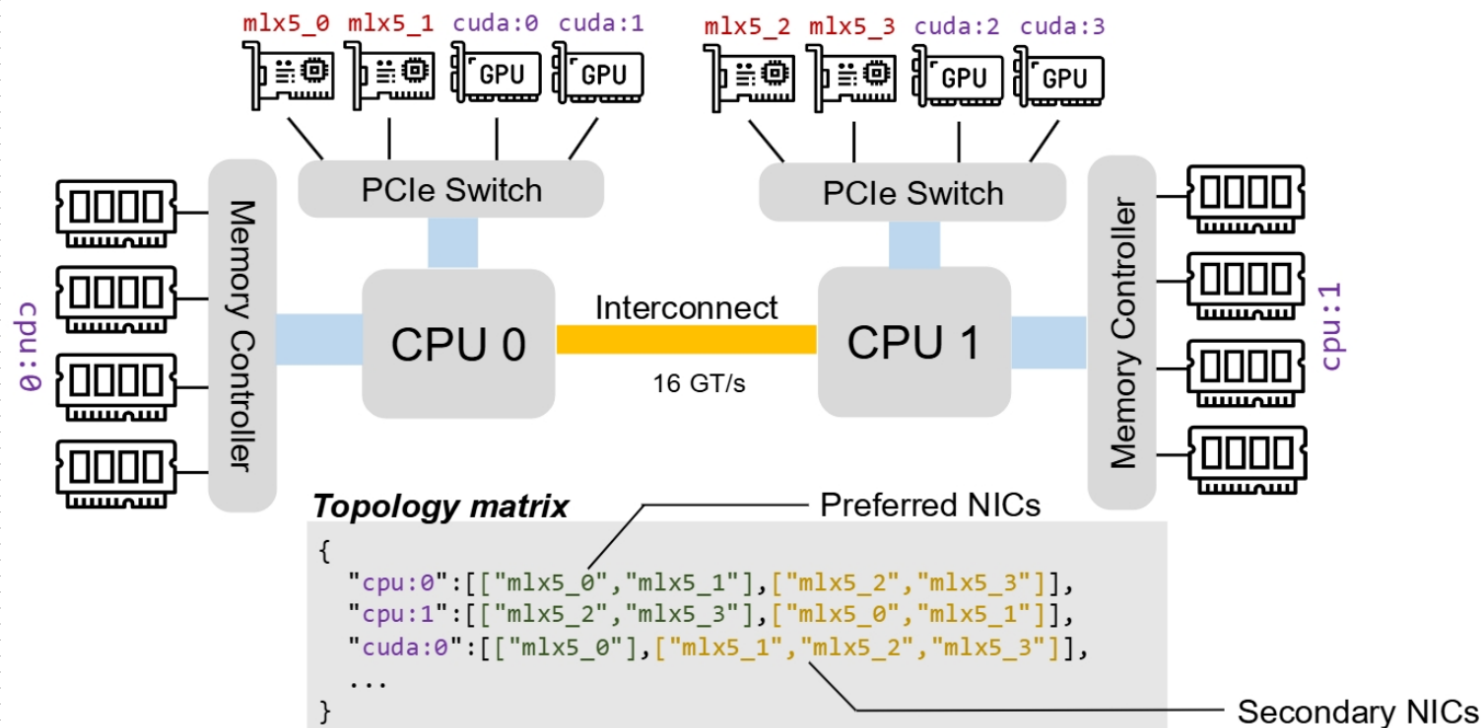**Mooncake Store offers object-based API**

# MoonCake Architecture

# Transfer Engine

❑Topology Aware Path Selection

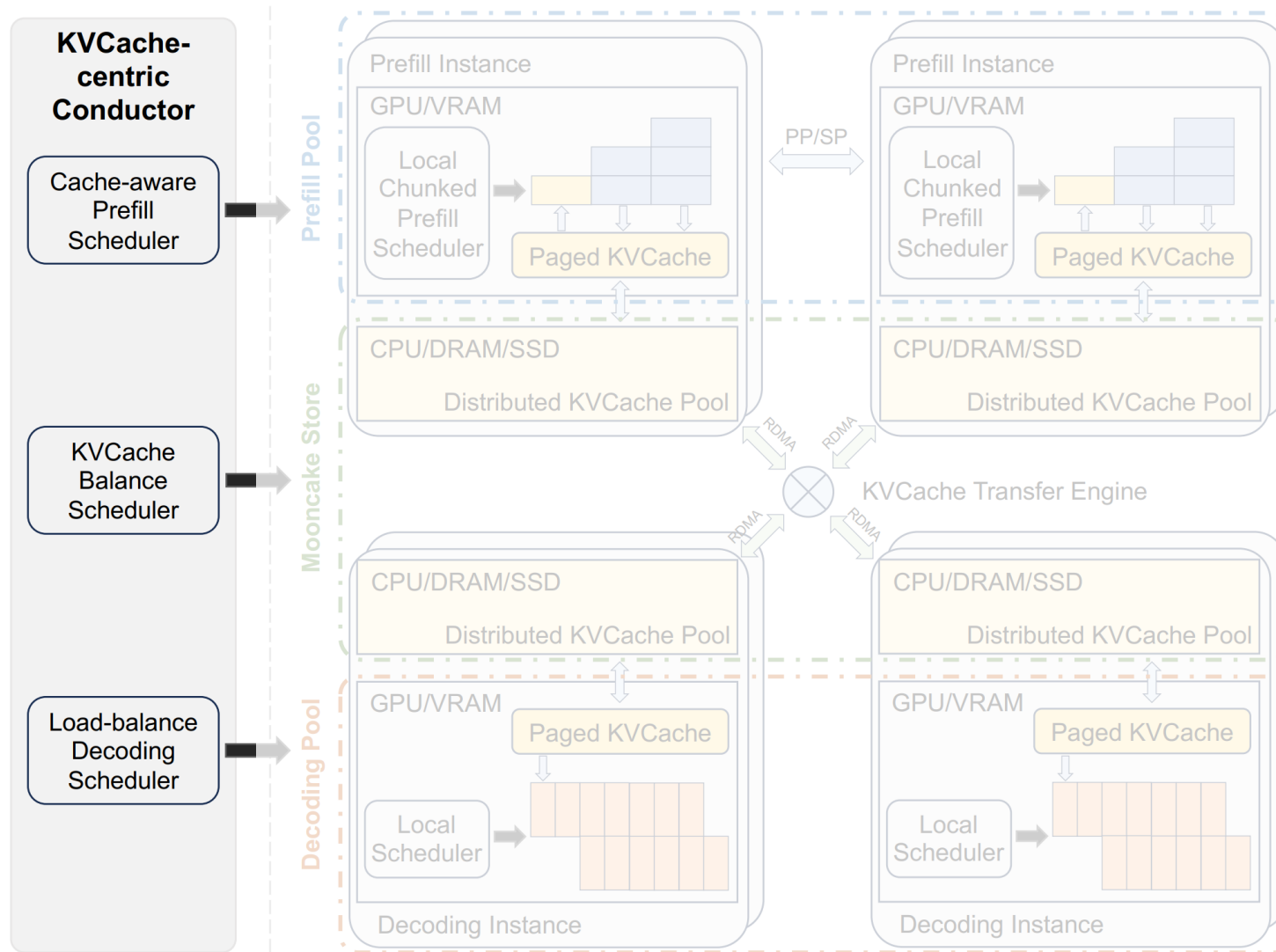❖ Broadcast the topology matrix across the cluster

❑Endpoint pooling

❖ Use SIEVE for eviction

# Mooncake Architecture

# KVCache-centric Scheduling

---

**Algorithm 1** KVCache-centric Scheduling Algorithm

---

**Input:** prefill instance pool $P$, decoding instance pool $D$, request $R$,
cache block size $B$.

**Output:** the prefill and decoding instances $(p, d)$ to process $R$.

1: $block\_keys \leftarrow \text{PrefixHash}(R.prompt\_tokens, B)$
2: $TTFT, p \leftarrow \inf, \emptyset$
3: $best\_len, best\_instance \leftarrow \text{FindBestPrefixMatch}(P, block\_keys)$
4: **for** $instance \in P$ **do**
5:     **if** $\frac{best\_len}{instance.prefix\_len} > \textbf{kvcache\_balancing\_threshold}$ **then**
6:         $prefix\_len \leftarrow best\_len$
7:         $transfer\_len \leftarrow best\_len - instance.prefix\_len$
8:         $T_{transfer} \leftarrow \text{EstimateKVCacheTransferTime}(transfer\_len)$
9:     **else**
10:         $prefix\_len \leftarrow instance.prefix\_len$
11:         $T_{transfer} \leftarrow 0$
12:     $T_{queue} \leftarrow \text{EstimatePrefillQueueTime}(instance)$
13:     $T_{prefill} \leftarrow \text{EstimatePrefillExecutionTime}($
            $\text{len}(R.prompt\_tokens), prefix\_len)$
14:     **if** $TTFT > T_{transfer} + T_{queue} + T_{prefill}$ **then**
15:         $TTFT \leftarrow T_{transfer} + T_{queue} + T_{prefill}$
16:         $p \leftarrow instance$
17: $d, TBT \leftarrow \text{SelectDecodingInstance}(D)$
18: **if** $TTFT > TTFT\_SLO$ **or** $TBT > TBT\_SLO$ **then**
19:     **reject** $R$; **return**
20: **if** $\frac{best\_len}{p.prefix\_len} > \textbf{kvcache\_balancing\_threshold}$ **then**
21:     $\text{TransferKVCache}(best\_instance, p)$
22: **return** $(p, d)$

---

# KVCache-centric Scheduling

**Algorithm 1** KVCache-centric Scheduling Algorithm

**Input:** prefill instance pool $P$, decoding instance pool $D$, request $R$, cache block size $B$.

**Output:** the prefill and decoding instances $(p, d)$ to process $R$.

1: $block\_keys \leftarrow$ PrefixHash($R.prompt\_tokens, B$)
2: $TTFT, p \leftarrow \inf, \emptyset$
3: $best\_len, best\_instance \leftarrow$ FindBestPrefixMatch($P, block\_keys$)
4: **for** $instance \in P$ **do**
5:    **if** $\frac{best\_len}{instance.prefix\_len} >$ **kvcache_balancing_threshold then**

**Estimate Transfer Time**

6:      $pre$
7:      $transfer\_len \leftarrow best\_len - instance.prefix\_len$
8:      $T_{transfer} \leftarrow$ EstimateKVCacheTransferTime($transfer\_len$)
9:    **else**
10:      $prefix\_len \leftarrow instance.prefix\_len$
11:      $T_{transfer} \leftarrow 0$

12:    $T_{queue} \leftarrow$ EstimatePrefillQueueTime($instance$)

**Estimate Queue and Prefill Time**

13:    $len(R.prompt\_tokens), prefix\_len)$
14:    **if** $TTFT > T_{transfer} + T_{queue} + T_{prefill}$ **then**
15:      $TTFT \leftarrow T_{transfer} + T_{queue} + T_{prefill}$

**Choose Prefill and Decode Instance**

17: $d, TBT \leftarrow$ SelectDecodingInstance($D$)
18: **if** $TTFT > TTFT\_SLO$ **or** $TBT > TBT\_SLO$ **then**
19:    **reject** $R$; **return**
20: **if** $\frac{best\_len}{p.prefix\_len} >$ **kvcache_balancing_threshold then**

**Other**

21:    TransferKVCache($best\_instance, p$)
22: **return** $(p, d)$

# KVCache-centric Scheduling

**Algorithm 1** KVCache-centric Scheduling Algorithm

**Input:** prefill instance pool $P$, decoding instance pool $D$, request $R$, cache block size $B$.

**Output:** the prefill and decoding instances $(p, d)$ to process $R$.

1: $block\_keys \leftarrow \text{PrefixHash}(R.prompt\_tokens, B)$
2: $TTFT, p \leftarrow \inf, \emptyset$
3: $best\_len, best\_instance \leftarrow \text{FindBestPrefixMatch}(P, block\_keys)$
4: **for** $instance \in P$ **do**
5:    **if** $\frac{best\_len}{instance.prefix\_len} > \textbf{kvcache\_balancing\_threshold}$ **then**
6:       $prefix\_len \leftarrow best\_len$
7:       $transfer\_len \leftarrow best\_len - instance.prefix\_len$
8:       $T_{transfer} \leftarrow \text{EstimateKVCacheTransferTime}(transfer\_len)$
9:    **else**
10:       $prefix\_len \leftarrow instance.prefix\_len$
11:       $T_{transfer} \leftarrow 0$
12:    $T_{queue} \leftarrow \text{EstimatePrefillQueueTime}(instance)$
13:    $T_{prefill} \leftarrow \text{EstimatePrefillExecutionTime}($
            $len(R.prompt\_tokens), prefix\_len)$
14:    **if** $TTFT > T_{transfer} + T_{queue} + T_{prefill}$ **then**
15:       $TTFT \leftarrow T_{transfer} + T_{queue} + T_{prefill}$
16:       $p \leftarrow instance$
17: $d, TBT \leftarrow \text{SelectDecodingInstance}(D)$
18: **if** $TTFT > TTFT\_SLO$ **or** $TBT > TBT\_SLO$ **then**
19:    reject $R$; **return**
20: **if** $\frac{best\_len}{p.prefix\_len} > \textbf{kvcache\_balancing\_threshold}$ **then**
21:    $\text{TransferKVCache}(best\_instance, p)$
22: **return** $(p, d)$

← **Design1: TTFT prioritized scheduling**

# KVCache-centric Scheduling

**Algorithm 1** KVCache-centric Scheduling Algorithm

**Input:** prefill instance pool $P$, decoding instance pool $D$, request $R$, cache block size $B$.

**Output:** the prefill and decoding instances $(p,d)$ to process $R$.

1: $block\_keys \leftarrow \text{PrefixHash}(R.prompt\_tokens, B)$
2: $TTFT, p \leftarrow \inf, \emptyset$
3: $best\_len, best\_instance \leftarrow \text{FindBestPrefixMatch}(P, block\_keys)$
4: **for** $instance \in P$ **do**
5:     **if** $\frac{best\_len}{instance.prefix\_len} >$ **kvcache_balancing_threshold then**
6:         $prefix\_len \leftarrow best\_len$
7:         $transfer\_len \leftarrow best\_len - instance.prefix\_len$
8:         $T_{transfer} \leftarrow \text{EstimateKVCacheTransferTime}(transfer\_len)$
9:     **else**
10:         $prefix\_len \leftarrow instance.prefix\_len$
11:         $T_{transfer} \leftarrow 0$
12:     $T_{queue} \leftarrow \text{EstimatePrefillQueueTime}(instance)$
13:     $T_{prefill} \leftarrow \text{EstimatePrefillExecutionTime}($
                $\text{len}(R.prompt\_tokens), prefix\_len)$
14:     **if** $TTFT > T_{transfer} + T_{queue} + T_{prefill}$ **then**
15:         $TTFT \leftarrow T_{transfer} + T_{queue} + T_{prefill}$
16:         $p \leftarrow instance$
17: $d, TBT \leftarrow \text{SelectDecodingInstance}(D)$
18: **if** $TTFT > TTFT\_SLO$ **or** $TBT > TBT\_SLO$ **then**
19:     **reject** $R$; **return**
20: **if** $\frac{best\_len}{p.prefix\_len} >$ **kvcache_balancing_threshold then**
21:     $\text{TransferKVCache}(best\_instance, p)$
22: **return** $(p,d)$

**Design2: KVCache load balancing**

# Summary of Prefix Caching

❑ Mooncake Store

  ❖ Prefix-hashed, Object-based API

❑ Transfer Engine

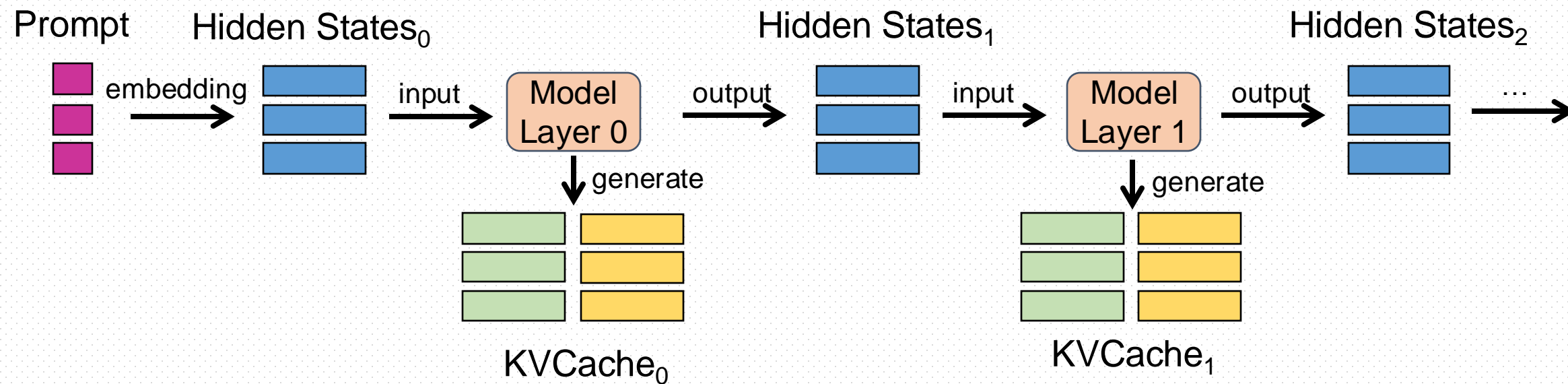  ❖ Topology-aware path selection, endpoint pooling

❑ KVCache-centric Scheduling

  ❖ TTFT prioritized scheduling, KVCache load balance

# LLM Inference

❑ LLM inference process can be divided into two phases

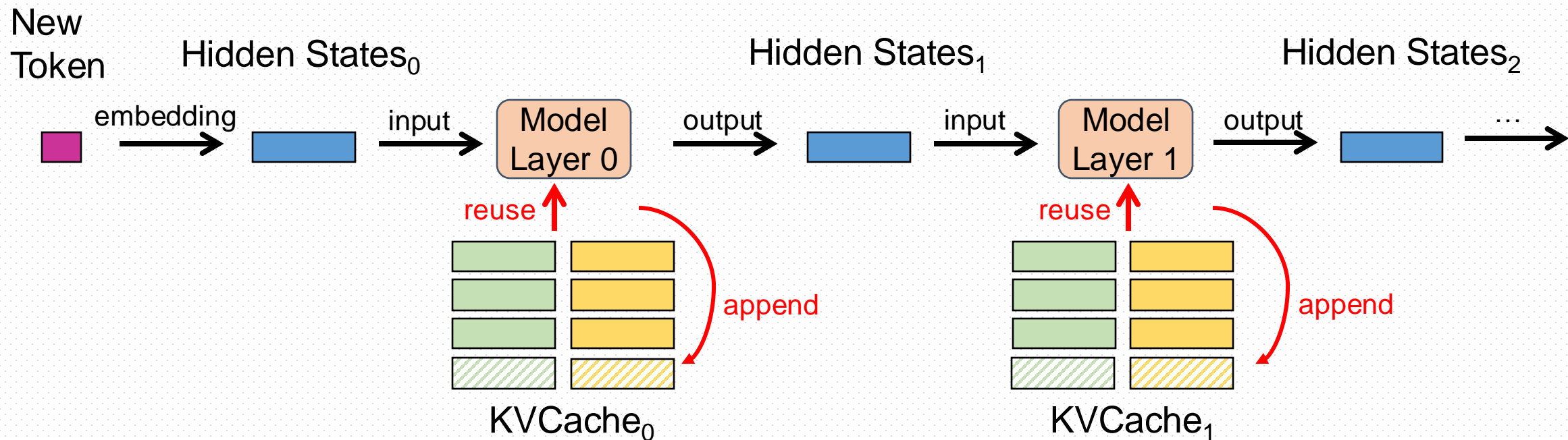❖ **Prefill Phase**: generate KVCache and output first token

# LLM Inference

❑ LLM inference process can be divided into two phases
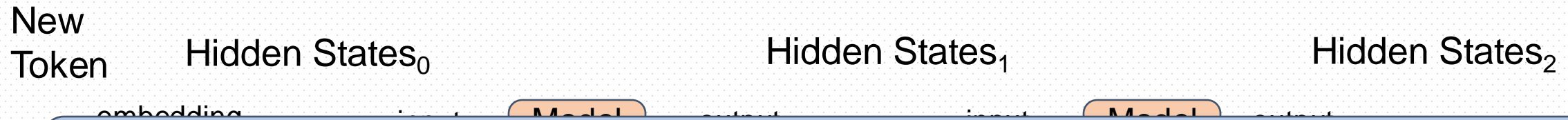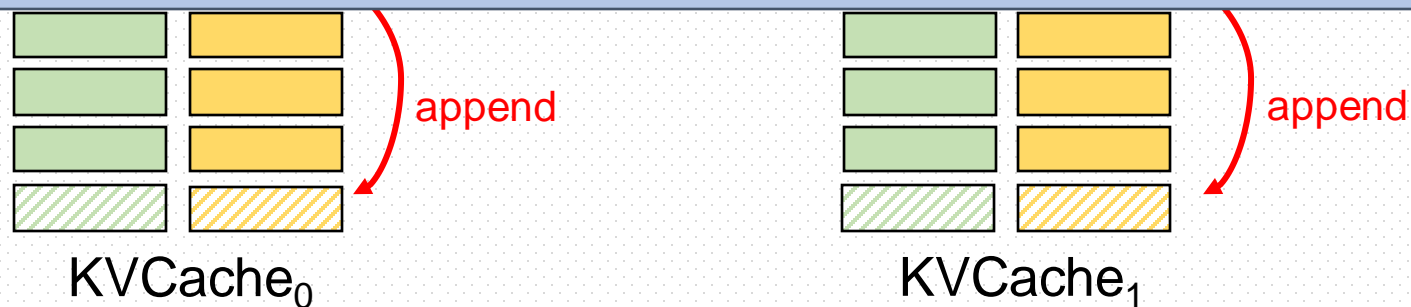
❖ **Decode Phase**: generate next token

New
Token

Hidden States$_0$            Hidden States$_1$            Hidden States$_2$

embedding    input    Model Layer 0    output    input    Model Layer 1    output    ...

reuse        append        reuse        append

KVCache$_0$            KVCache$_1$

# LLM Inference

❑ LLM inference process can be divided into two phases

&#10022; **Decode Phase**: generate next token

New
Token

Hidden States$_0$

Hidden States$_1$

Hidden States$_2$

The most significant difference between Prefill and Decode is the ***token shape***!!!
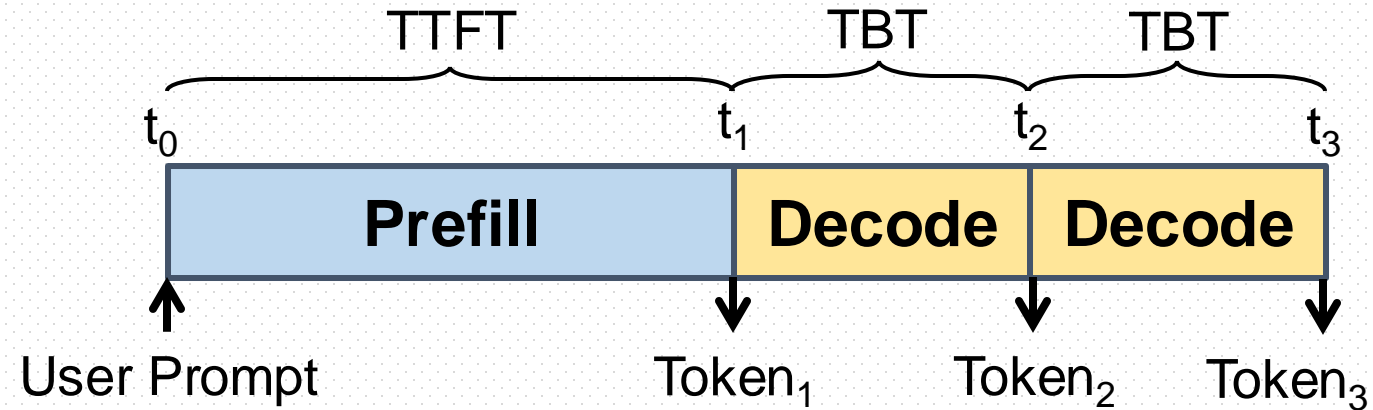
append

append

KVCache$_0$

KVCache$_1$

# Performance Metrics

❑ Latency

  ❖ Time-to-first-token (_TTFT_)

  ➢ Latency metric for _Prefill_

  ❖ Time-between-token (_TBT_)

  ➢ Latency metric for _Decode_
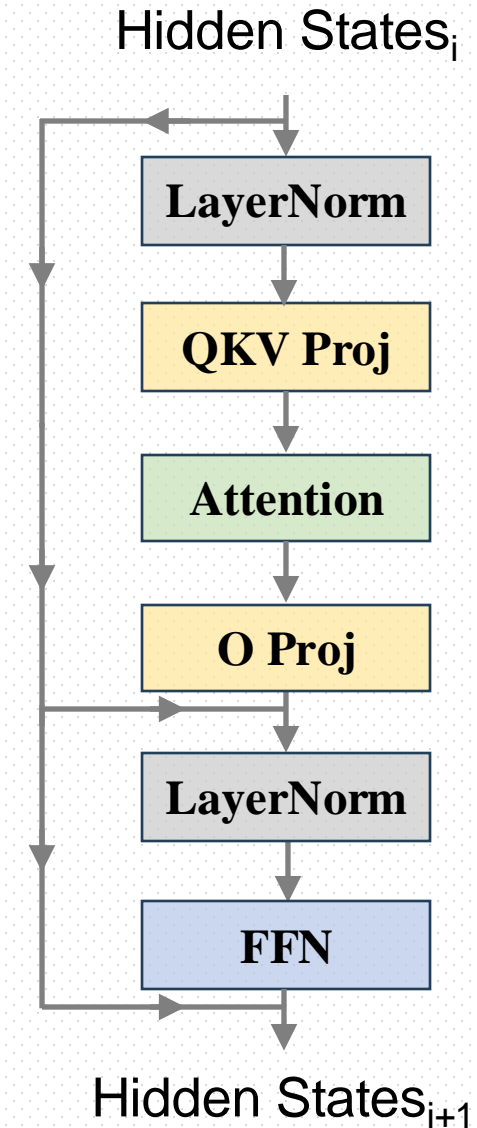
❑ Throughput

  ❖ Model-Flops-Utilization (_MFU_)

  ➢ Measured Flops / Theoretical Upper Bound

TTFT  TBT  TBT

$t_0$  $t_1$  $t_2$  $t_3$

**Prefill**  **Decode**  **Decode**

User Prompt  Token$_1$  Token$_2$  Token$_3$

# Layer Computation Details

❑There are four operations in single layer

❖Proj

❖Attention

❖FFN

❖Layer Norm

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Layer Computation Details

❑ Some notations on model parameters

| Parameter | Notation |
| --- | --- |
| Batch size | $b$ |
| Num head | $n$ |
| Sequence length | $s$ |
| Head dimension | $h$ |
| Hidden dimension | $d$ |

# Layer Computation Details

## ❑ Proj

❖ Shape: [$d$, $d$]

❖ Input

➢ Prefill: [$b, s, d$]

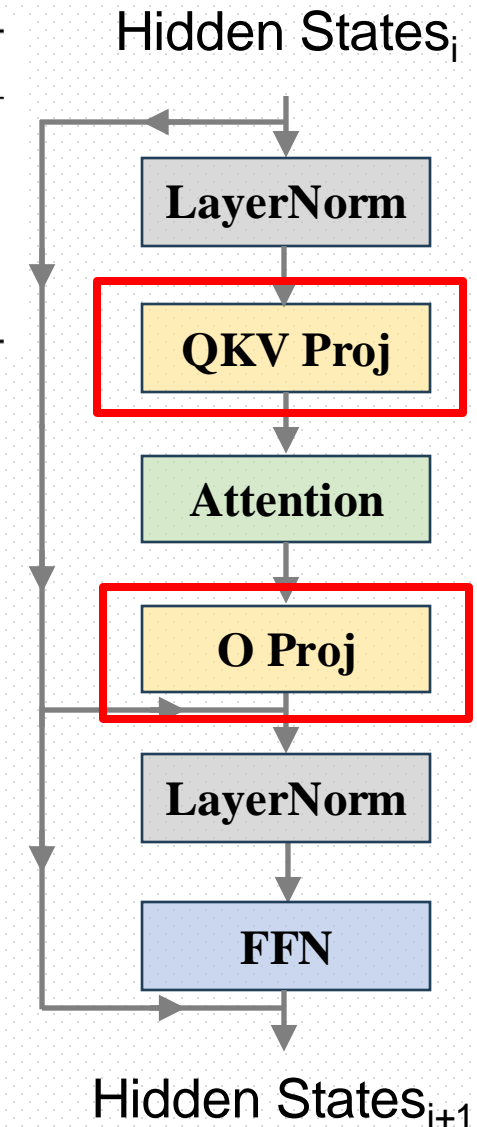➢ Decode: [$b, 1, d$]

❖ HBM load num:

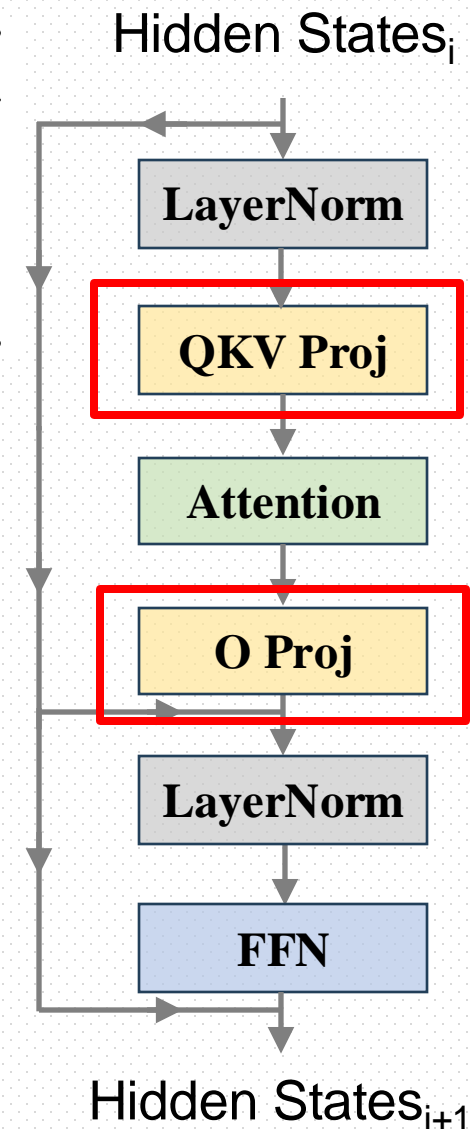➢ Prefill: $d^2+bsd$

➢ Decode: $d^2+bd$

❖ Computation Ops:

➢ Prefill: $bsd^2$

➢ Decode: $bd^2$

| Parameter | Notation |
|---|---|
| Batch size | $b$ |
| Num head | $n$ |
| Sequence length | $s$ |
| Head dimension | $h$ |
| Hidden dimension | $d$ |

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Layer Computation Details

☐ **Proj**

❖ Shape: [$d$, $d$]

❖ Input

➢ Prefill: [$b, s, d$]

➢ Decode: [$b, 1, d$]

❖ HBM load num:

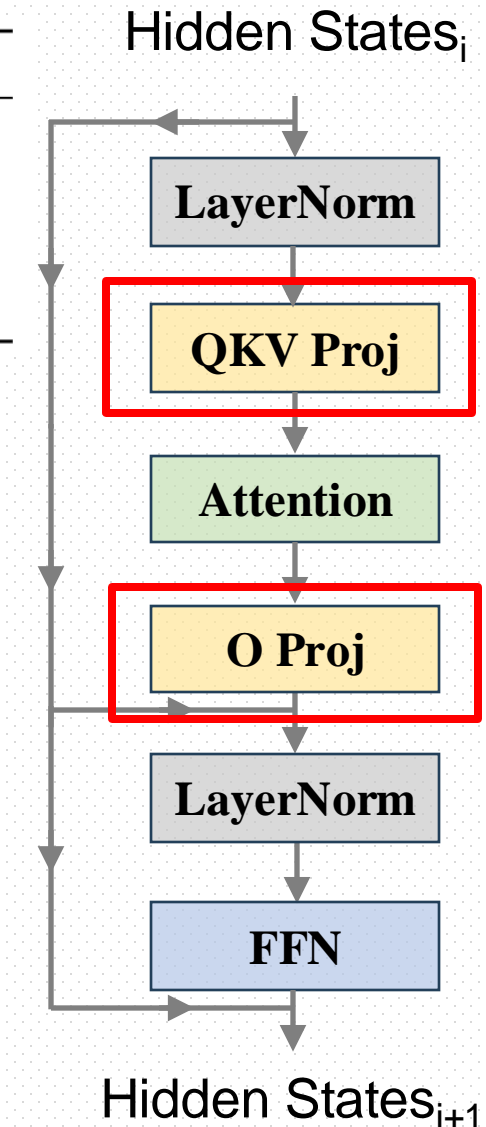➢ Prefill: $d^2 + bsd$

➢ Decode: $d^2 + bd$

❖ Computation Ops:

➢ Prefill: $bsd^2$

➢ Decode: $bd^2$

| Parameter | Notation |
|---|---|
| Batch size | $b$ |
| Num head | $n$ |
| Sequence length | $s$ |
| Head dimension | $h$ |
| Hidden dimension | $d$ |

As $b$ and $s$ increases, *proj* quickly becomes *comp-bound* task

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Layer Computation Details

## ❑Proj

❖Shape: [$d$, $d$]

❖Input

➢Prefill: [$b, s, d$]

➢Decode: [$b, 1, d$]

❖HBM load num:

➢Prefill: $d^2+bsd$

➢Decode: $d^2+bd$

❖Computation Ops:

➢Prefill: $bsd^2$

➢Decode: $bd^2$

| Parameter | Notation |
|-----------|----------|
| Batch size | $b$ |
| Num head | $n$ |
| Sequence length | $s$ |
| Head dimension | $h$ |
| Hidden dimension | $d$ |

As $b$ and $s$ increases, *proj* quickly becomes ***comp-bound*** task

In other words, even with small $b$, *proj* can use up comp resource

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Layer Computation Details

## ❑ Proj

❖ Shape: [$d$, $d$]

❖ Input

➤ Prefill: [$b, s, d$]

➤ Decode: [$b, 1, d$]

❖ HBM load num:

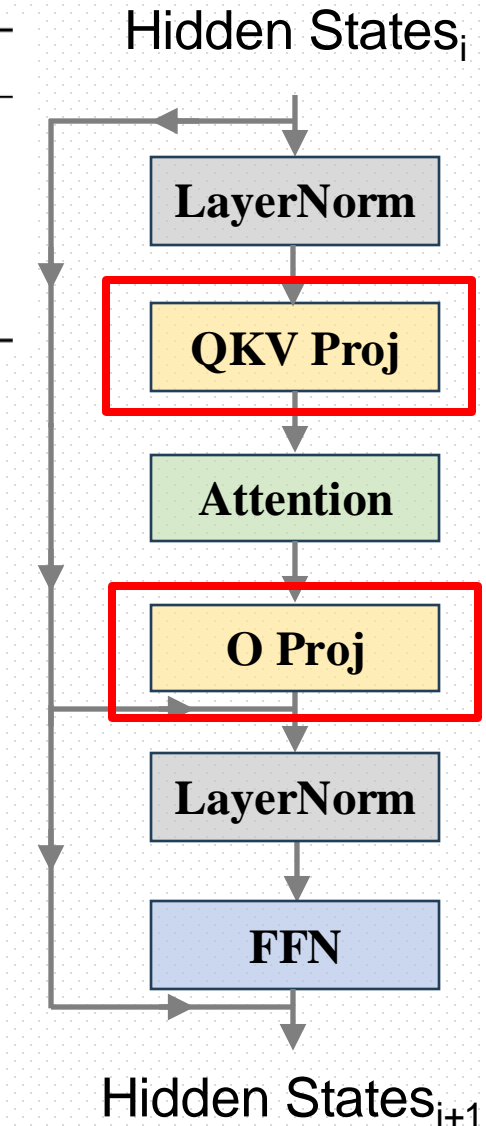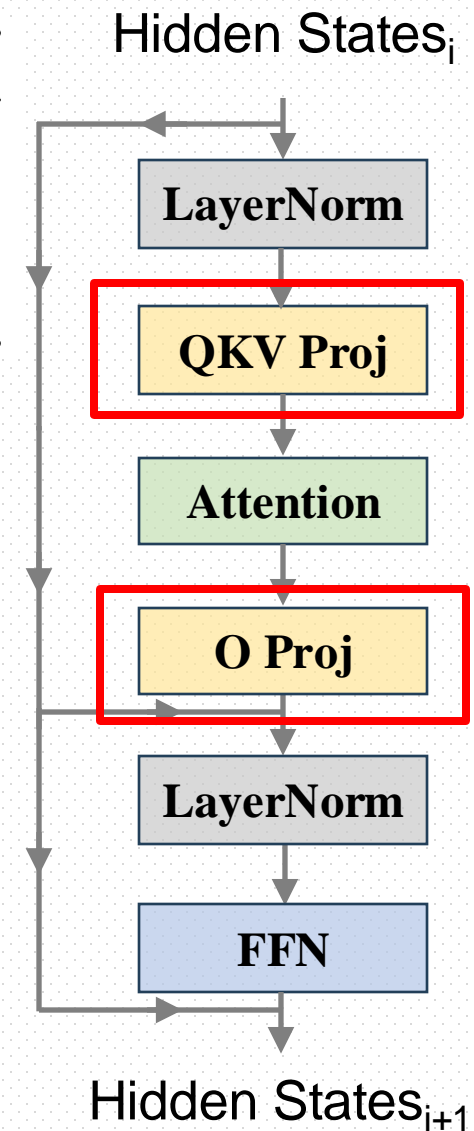➤ Prefill: $d^2 + bsd$

➤ Decode: $d^2 + bd$

❖ Computation Ops:

➤ Prefill: $bsd^2$

➤ Decode: $bd^2$

| Parameter | Notation |
|---|---|
| Batch size | $b$ |
| Num head | $n$ |
| Sequence length | $s$ |
| Head dimension | $h$ |
| Hidden dimension | $d$ |

But if we give a small $b$ in decode phase, **_proj_** wound become an **_IO-bound_** task

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Layer Computation Details

## ❑Proj

❖Shape: [$d$, $d$]

❖Input

➢Prefill: [$b, s, d$]

➢Decode: [$b, 1, d$]

❖HBM load num:

➢Prefill: $d^2+bsd$

➢Decode: $d^2+bd$

❖Computation Ops:

➢Prefill: $bsd^2$

➢Decode: $bd^2$

| Parameter | Notation |
|---|---|
| Batch size | $b$ |
| Num head | $n$ |
| Sequence length | $s$ |
| Head dimension | $h$ |
| Hidden dimension | $d$ |

But if we give a small $b$ in decode phase, *proj* wound become an *IO-bound* task

On the other hand, *proj* wound become more *comp-intensive* if we give a bigger $b$

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Layer Computation Details

## ❑FFN

❖Shape: [$d$, $4d$] and [$4d,d$]

❖Input

➢Prefill: [$b, s, d$]
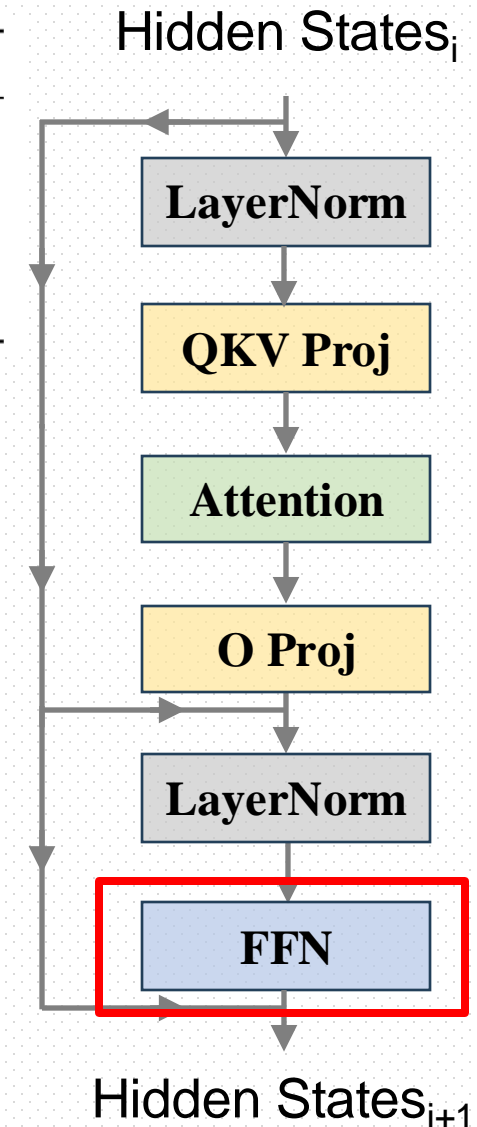
➢Decode: [$b, 1, d$]
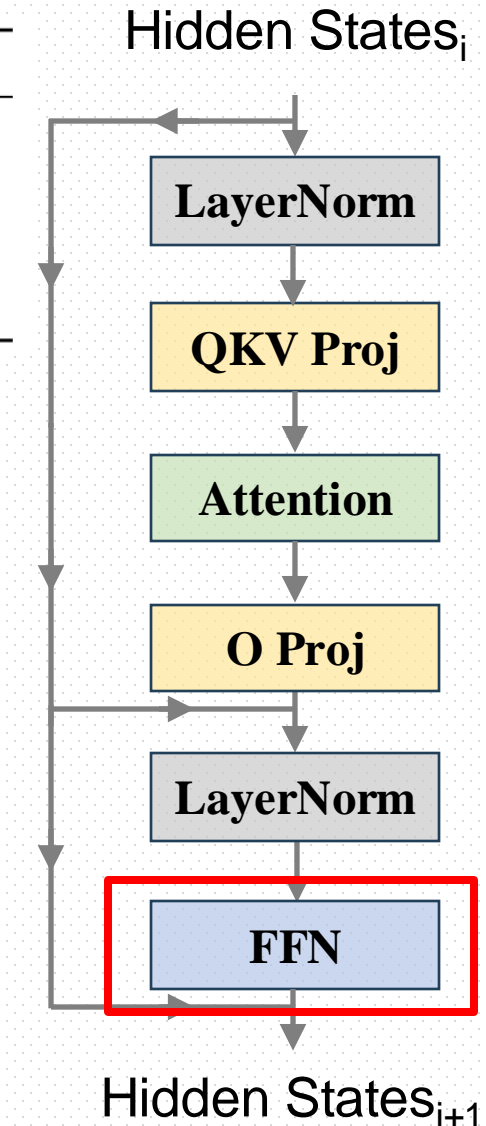
❖HBM load num:

➢Prefill: $8d^2+bsd$

➢Decode: $8d^2+bd$

❖Computation Ops:

➢Prefill: $8bsd^2$

➢Decode: $8bd^2$

| Parameter | Notation |
| --- | --- |
| Batch size | $b$ |
| Num head | $n$ |
| Sequence length | $s$ |
| Head dimension | $h$ |
| Hidden dimension | $d$ |

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Layer Computation Details

❑**FFN**

❖Shape: [$d$, $4d$] and [$4d,d$]

❖Input

➢Prefill: [$b, s, d$]

➢Decode: [$b, 1, d$]

❖HBM load num:

➢Prefill: $8d^2+bsd$

➢Decode: $8d^2+bd$

❖Computation Ops:

➢Prefill: $8bsd^2$

➢Decode: $8bd^2$

| Parameter | Notation |
|---|---|
| Batch size | $b$ |
| Num head | $n$ |
| Sequence length | $s$ |
| Head dimension | $h$ |
| Hidden dimension | $d$ |

Similar to *proj*, *FFN* is *comp-bound* in prefill, *IO-bound* in decode

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Layer Computation Details

## ❑ Attention

❖ Input

➢ Prefill: [$b, n, s, h$] for QKV

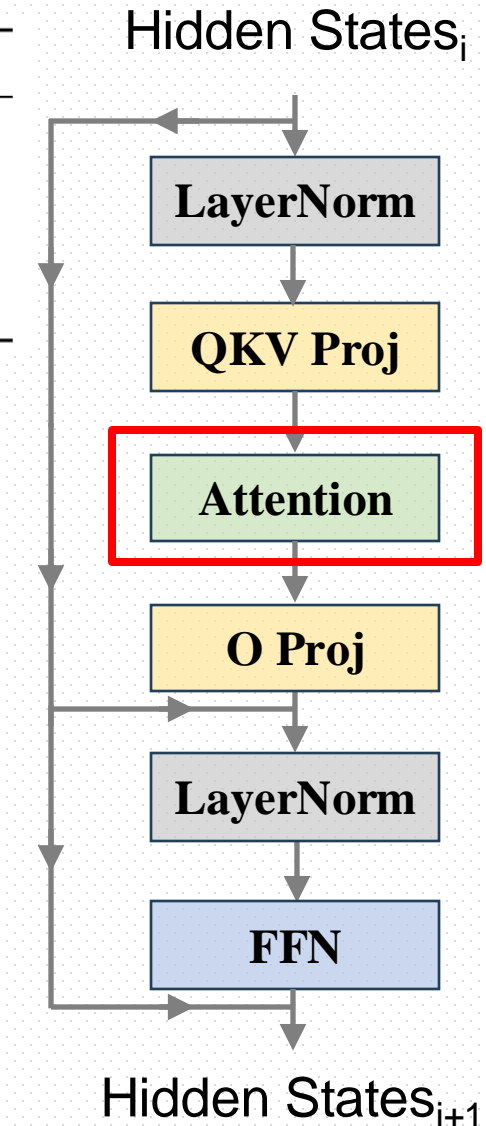➢ Decode: [$b, n, 1, h$] for Q, [$b, n, s, h$] for KV

❖ HBM load num:

➢ Prefill: **3bnsh**

➢ Decode: **2bnsh+bnh**

❖ Computation Ops:

➢ Prefill: **2bns²h**

➢ Decode: **2bnsh**

| Parameter | Notation |
|---|---|
| Batch size | $b$ |
| Num head | $n$ |
| Sequence length | $s$ |
| Head dimension | $h$ |
| Hidden dimension | $d$ |

$$Attn = Softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Layer Computation Details

❑ **Attention**

❖ Input

➤ Prefill: [*b, n, s, h*] for QKV

➤ Decode: [*b, n, 1, h*] for Q, [*b, n, s, h*] for KV

❖ HBM load num:
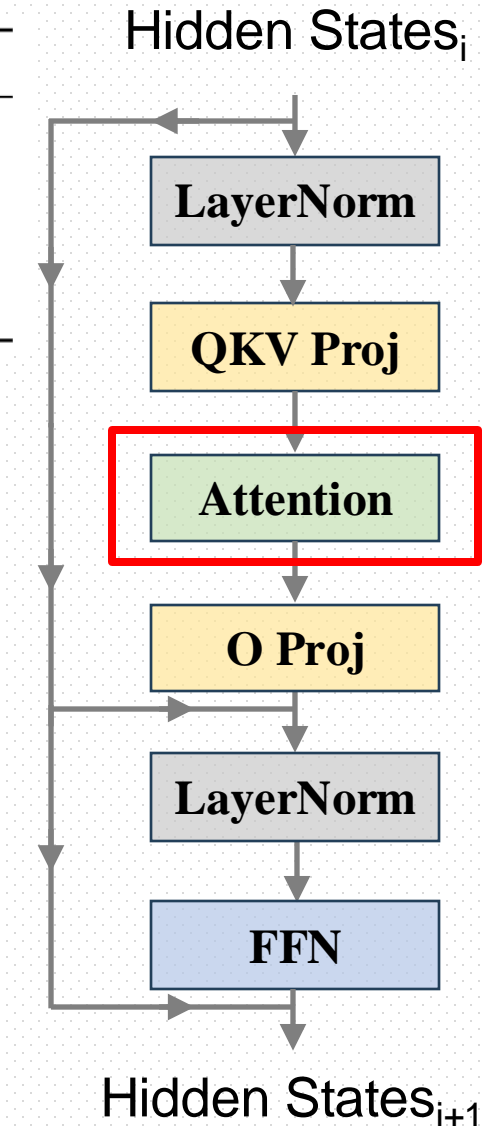
*comp-bound*

➤ Prefill: **3bnsh**

➤ Decode: **2bnsh+bnh**

❖ Computation Ops:

➤ Prefill: **2bns²h**

➤ Decode: **2bnsh**

*IO-bound*

| Parameter | Notation |
|---|---|
| Batch size | *b* |
| Num head | *n* |
| Sequence length | *s* |
| Head dimension | *h* |
| Hidden dimension | *d* |

$$Attn = Softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Hidden States$_i$

LayerNorm

QKV Proj

Attention

O Proj

LayerNorm

FFN

Hidden States$_{i+1}$

# Decode Need Bigger Batch Size

❑A40, KV tokens=1000

| Batch Size | QKV Proj | Attn | O Proj | FFN |
|------------|----------|-------|--------|--------|
| 1 | 0.35% | 0.29% | 0.35% | 0.37% |
| 2 | 0.68% | 0.35% | 0.67% | 0.72% |
| 4 | 1.34% | 0.37% | 1.34% | 1.43% |
| 8 | 2.69% | 0.36% | 2.67% | 2.85% |
| 16 | 5.35% | 0.36% | 5.33% | 5.62% |
| 32 | 10.72% | 0.39% | 10.72% | 11.34% |
| 64 | 21.13% | 0.40% | 21.00% | 21.10% |
| 128 | 36.39% | 0.40% | 35.63% | 37.50% |

**Bigger batch size can increate decode MFU**

# **Decode Need Bigger Batch Size**

❑A40, KV tokens=1000

| Batch Size | QKV Proj | Attn | O Proj | FFN |
|---|---|---|---|---|
| 1 | 20.82% | 3.97% | 6.97% | 68.24% |
| 2 | 20.18% | 6.29% | 6.80% | 66.73% |
| 4 | 19.19% | 11.23% | 6.43% | 63.15% |
| 8 | 16.89% | 20.78% | 5.66% | 56.67% |
| 16 | 14.15% | 33.87% | 4.72% | 47.25% |
| 32 | 10.94% | 49.27% | 3.65% | 36.15% |
| 64 | 7.43% | 64.07% | 2.48% | 26.02% |
| 128 | 5.17% | 75.57% | 1.77% | 17.50% |

**MFU increase has an upper bound**

# Workload Feature

❑ *Prefill*

  ❖ *Comp-bound*

  ❖ Small batch size can utilize most computation resource

❑ **Decode**

  ❖ *IO-bound*

  ❖ Most time in model weight loading

# Serving Scheduling

❑Tradeoff between Latency and Throughput

**C,D enter**        **vLLM [SOSP 23]**        **Decodes for A,B stalled**

| $A_d,B_d$ | $C_p$ | $D_p$ | $A_d,B_d,C_d,D_d$ | ... |

**Prefill Prioritized Schedules**

**C,D enter**        **Orca [OSDI 22]**        **Decodes for A,B stalled**

| $A_d,B_d$ | $C_p,D_p,A_d,B_d$ | $A_d,B_d,C_d,D_d$ | ... |

**Prefill finish early -> Large decode batch size -> High Throughput**
**Decode stall -> High Latency**

**C,D enter**        **FasterTransformer (NV)**        **Prefill for C,D stalled**

| $A_d,B_d$ | $A_d,B_d$ | $A_d,B_d$ | $B_d$ | $B_d$ | $C_p$ | ... |

**Decode Prioritized Schedules**

**Decode without interference -> Low Latency**
**Prefill Stall -> Small batch in decode -> Poor throughput**

# Serving Scheduling

❑Amortize weight loading overhead in decode phase

**C,D enter**       **vLLM [SOSP 23]**       **Decodes for A,B stalled**

| $A_d,B_d$ | $C_p$ | $D_p$ | $A_d,B_d,C_d,D_d$ | ... |

**C,D enter**       **Orca [OSDI 22]**       **Decodes for A,B stalled**

| $A_d,B_d$ | $C_p,D_p,A_d,B_d$ | $A_d,B_d,C_d,D_d$ | ... |

**C,D enter**       **FasterTransformer (NV)**       **Prefill for C,D stalled**

| $A_d,B_d$ | $A_d,B_d$ | $A_d,B_d$ | $B_d$ | $B_d$ | $C_p$ | ... |

**C,D enter**       **Chunk Prefill [OSDI 24]**       **No stalls**

| $A_d,B_d$ | $A_d,B_d,C_p$ | $A_d,B_d,C_p$ | $A_d,B_d,D_p$ | $A_d,B_d,D_p$ | ... |

# Serving Scheduling

❑ Prefill and Decode share the same model weight

  ❖ QKV Proj, O Proj, FFN

  ❖ Amortized model weight loading overhead

❑ Attention is processed separately

❑ Increased latency to both Prefill and Decode

  ❖ Prefill need to load more KV from HBM

  ❖ More computation in Decode phase

# Serving Scheduling

❏ Amortize weight loading overhead in decode phase

**C,D enter**      **vLLM [SOSP 23]**      **Decodes for A,B stalled**

| $A_d,B_d$ | $C_p$ | $D_p$ | $A_d,B_d,C_d,D_d$ | ... |

**C,D enter**      **Orca [OSDI 22]**      **Decodes for A,B stalled**

| $A_d,B_d$ | $C_p,D_p,A_d,B_d$ | $A_d,B_d,C_d,D_d$ | ... |

**C,D enter**      **FasterTransformer (NV)**      **Prefill for C,D stalled**

| $A_d,B_d$ | $A_d,B_d$ | $A_d,B_d$ | $B_d$ | $B_d$ | $C_p$ | ... |

**C,D enter**      **Chunk Prefill [OSDI 24]**      **No stalls**

| $A_d,B_d$ | $A_d,B_d,C_p$ | $A_d,B_d,C_p$ | $A_d,B_d,D_p$ | $A_d,B_d,D_p$ | ... |

**C,D enter**      **PD Disaggregation**      **No interference**

| $C_p$ | $D_p$ |

**Transfer**

| $A_d,B_d$ | $A_d,B_d$ | $A_d,B_d$ | $B_d$ | $B_d,C_d$ | $B_d,C_d$ |

# Serving Scheduling

❑Minimize Prefill and Decode interference

  ❖Additional KVCache transfer overhead

  ❖Still low MFU in Decode phase

# Serving Scheduling —— Summary

| | Prefill prioritized | Decode prioritized | Chunk-Prefill | PD Disaggregation |
|---|---|---|---|---|
| **TTFT** | +++ | --- | ++ | +++ |
| **TBT** | --- | +++ | ++ | +++ |
| **Prefill-MFU** | +++ | +++ | ++ | +++ |
| **Decode-MFU** | + | - | ++ | + |

# Serving Scheduling —— Summary

Suitable for *relax SLO throughput-oriented* scenario

| | Prefill prioritized | Decode prioritized | Chunk-Prefill | PD Disaggregation |
|---|---|---|---|---|
| **TTFT** | +++ | --- | ++ | +++ |
| **TBT** | --- | +++ | ++ | +++ |
| **Prefill-MFU** | +++ | +++ | ++ | +++ |
| **Decode-MFU** | ++ | - | ++ | + |

Suitable for *stringent SLO latency-oriented* scenario

# Some other discussion

❑ **Long context support**

❖ Emphasize the importance of <u>*TTFT*</u> optimization

❑ **Parallelism Choice**

❖ <u>*TP*</u> -> require two RDMA-based all-reduce in cross-node setting

❖ <u>*SP*</u> -> not suitable for short input

❖ <u>*ESP*</u> -> add complexity to architecture design

❖ <u>*CPP*</u> -> <u>*C*</u>hunk <u>*P*</u>ipeline <u>*P*</u>arallelism

  ➢ Less communication overhead

  ➢ Fit both short and long contexts

# Evaluation

❑ **Testbed**

  ❖ **16 nodes each with 8 * A800-80GB, 4 * 200 Gbps RDMA NICs**

❑ **Metric**

  ❖ **TTFT (SLO: 30s)**

  ❖ **TBT (SLO: 100ms, 200ms, 300ms)**

❑ **Baseline**

  ❖ **vLLM**

  ❖ **vLLM with prefix caching**

  ❖ **vLLM with chunk prefill**

# Evaluation

❑ **Workload**

|  | Conversation | Tool&Agent | Synthetic |
|---|---|---|---|
| Avg Input Len | 12035 | 8596 | 15325 |
| Avg Output Len | 343 | 182 | 149 |
| Cache Ratio | 40% | 59% | 66% |
| Arrival Pattern | Timestamp | Timestamp | Poisson |
| Num Requests | 12031 | 23608 | 3993 |

# Evaluation —— TBT



**Conversation**
**Input:12035**
**Output:343**
**Cache:40%**

**Synthetic**
**Input:15325**
**Output:149**
**Cache:66%**

**Tool & Agent**
**Input:8596**
**Output:182**
**Cache:59%**

# Evaluation —— TTFT

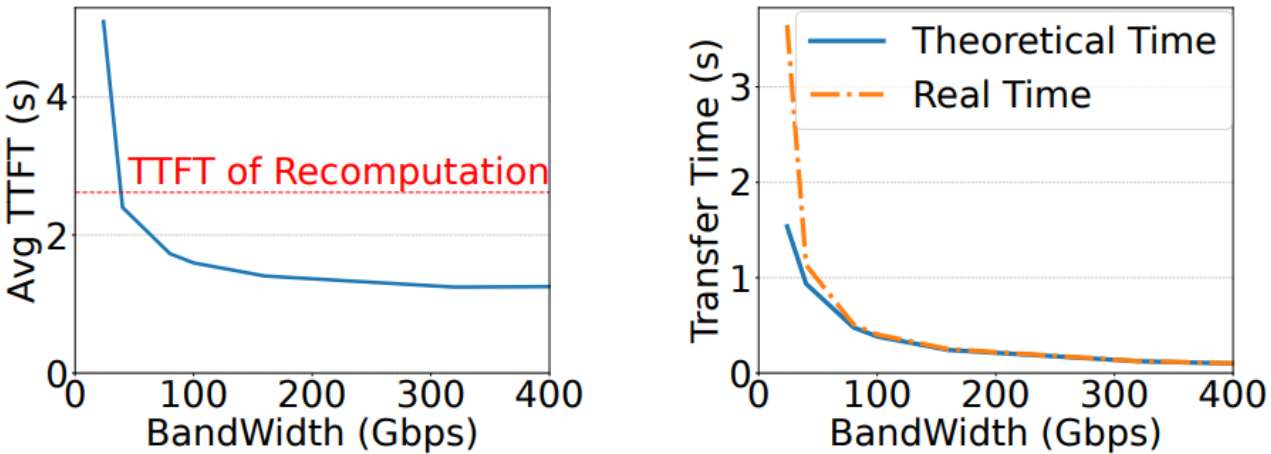# Evaluation —— Global Cache Efficiency

# Evaluation —— Transfer Engine

**Transfer Engine Performance**

**Bandwidth Analyze**

# Evaluation —— Breakdown



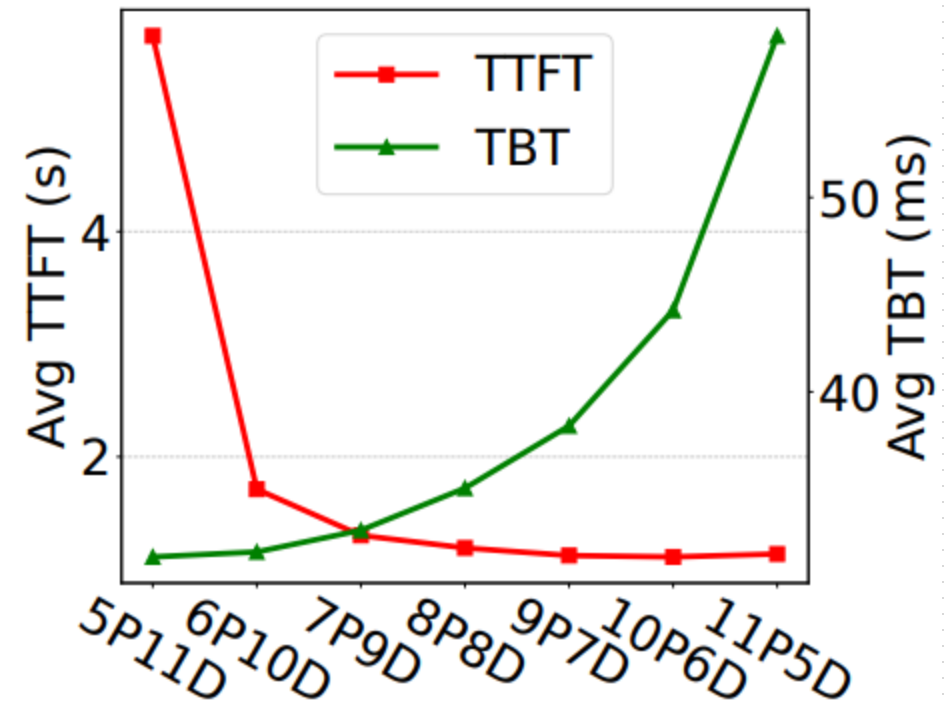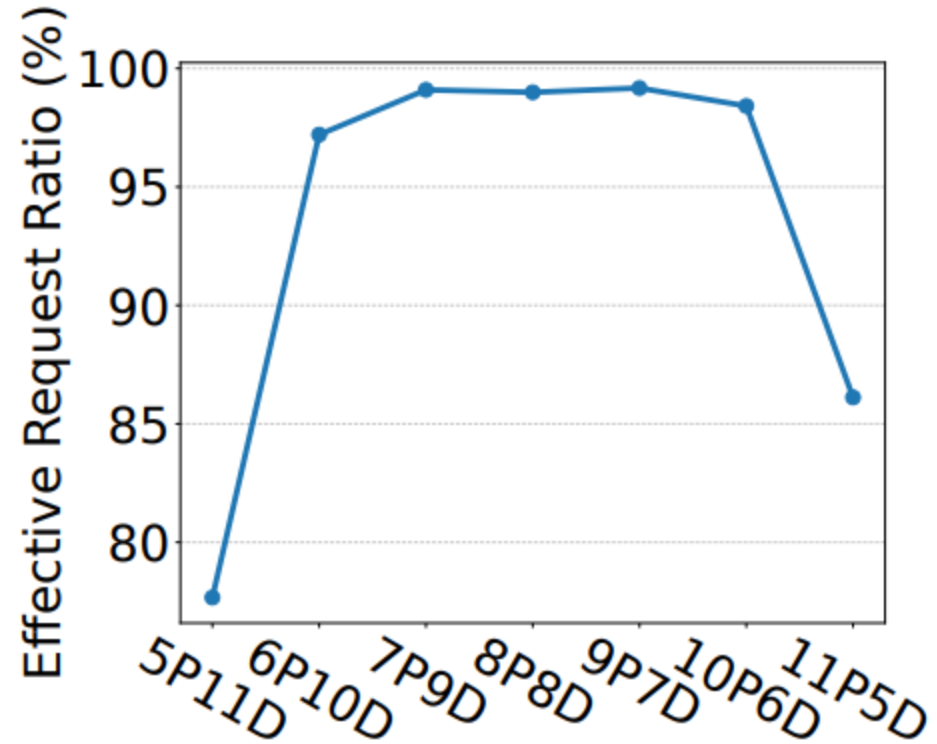**Prefix Cache Ratio 0%**                    **Prefix Cache Ratio 95%**

# Evaluation —— Breakdown

# Summary & Discussion

❑ **Prefix-cache based PD disaggregation LLM serving system**

  ❖ **Mooncake Store**

  ❖ **Transfer Engine**

  ❖ **Scheduling**

❑ **Provide an industrial view on comparison between _chunk prefill_ and _PD disaggregation_**

❑ **Open sourced LLM serving trace and codebase**