# Mooncake: Kimi's KVCache-centric Architecture for LLM Serving

## (W6) Paper Reading

Nan Lin

Shanghai Jiao Tong University

2024-08-10

# Mooncake: Kimi's KVCache-centric Architecture for LLM Serving

Ruoyu Qin[♠♡1]   Zheming Li[♠1]   Weiran He[♠]

Mingxing Zhang[♡2]   Yongwei Wu[♡]   Weimin Zheng[♡]   Xinran Xu[♠2]

[♠]Moonshot AI   [♡]Tsinghua University

- Click the link to open: https://arxiv.org/html/2407.00079v1

# Outline

**Review**

**Motivation**

**Architecture**

**Evaluations**

**References**

# Outline

## Review

## Motivation

## Architecture

## Evaluations

## References

# Review 1: Transformer-based Inference Workflow

Generating from an input text sequences (**prompts**) to text (**completion**) consists of:

1. Loading weights to GPU
2. *Tokenizing* **prompts** on CPU
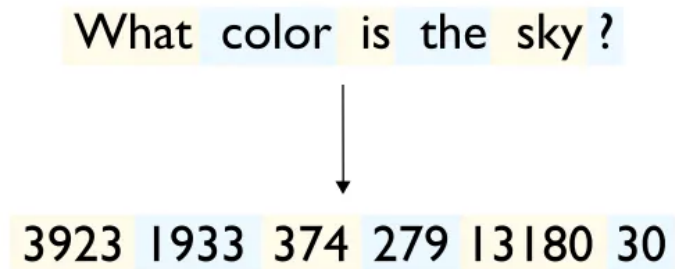3. Transfering the tokenized prompts to GPU

Figure 2: Tokenization

# Review 1: Transformer-based Inference Workflow

(cont.)

4. **Prefill stage** (or **Initialization stage**):
   Generate the first token

5. **Decoding stage** (or **Generation stage**,
   **Auto-regression stage**, **Incremental
   stage**): Repeat …
   - Append the generated token to the
     sequence of input tokens
   - Using it as a new input to generate the
     second token

   End repeat if EOS token generated or
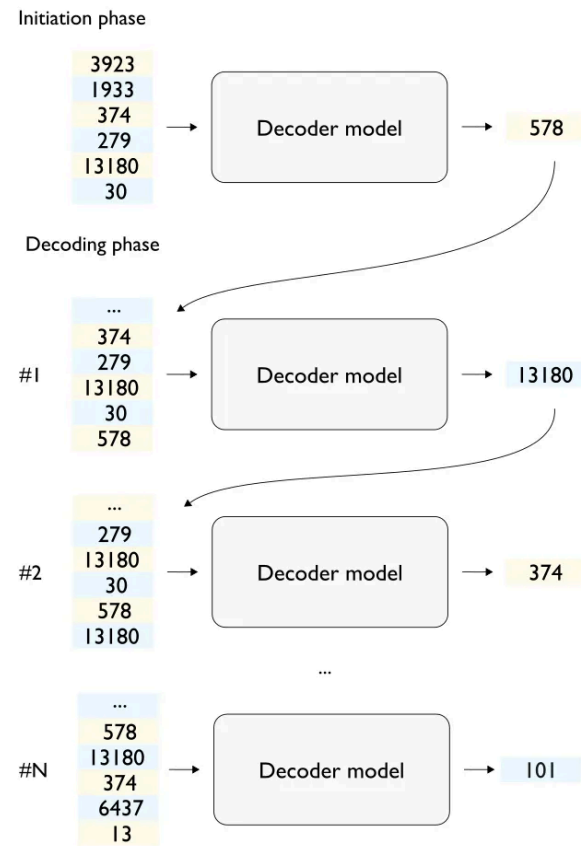   maximum token sequence length reached.



Figure 3: Prefill Stage and Decoding Stage

# Review 1: Transformer-based Inference Workflow

(cont.)

6. *Detokenizing* the tokenized incremented
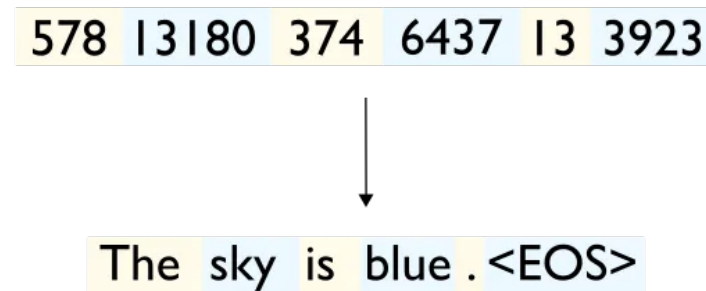   token to get your generated text.

Figure 4: Detokenization

# Review 2: KVCache

Motivation:

- **Masking** technique: *Causality*, we cannot use token generated of later words to influence the former ones.
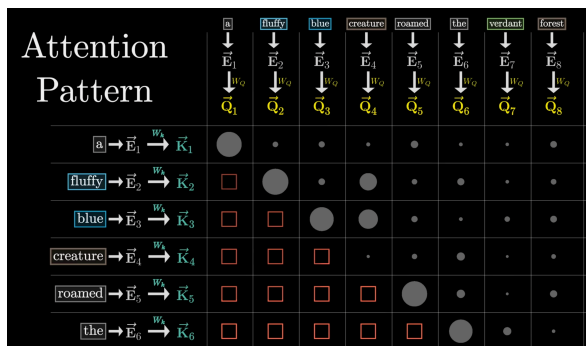


Figure 5: Masking

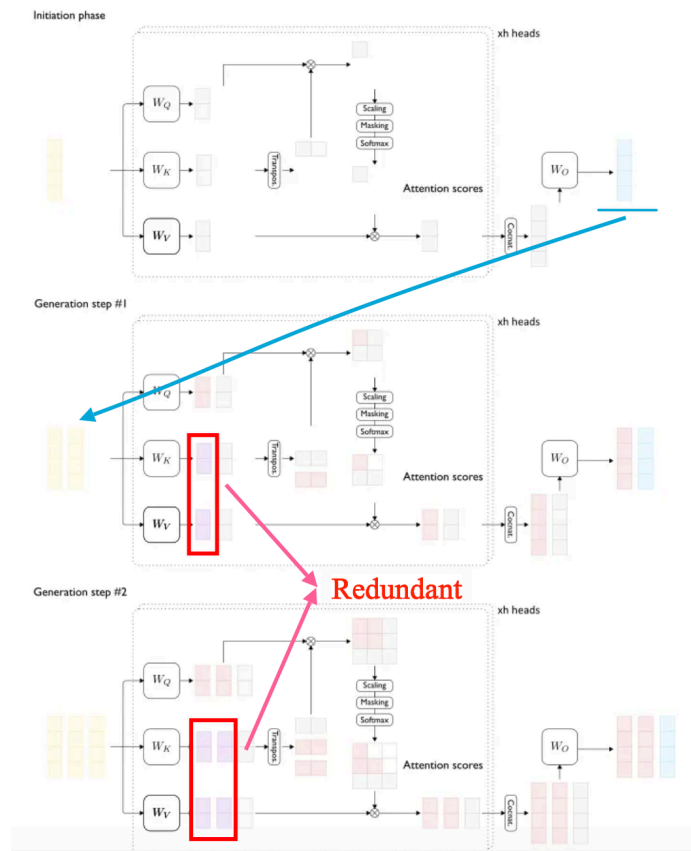- Previous tokens are identical across all the stages.



Figure 6: Excessive Datas in Prefill Stage and Decoding Stage

# Review 2: KVCache

Changes to the original workflow:

- No longer need tokens but their key and value vectors.
- Take the KVCache and the last generated token as input.
- Initilization phase could be renamed as *pre-fill* stage since it prepares all the key and value vectors of input text sequence.

Issues of KVCache:

- *Large memory consumption in GPU*
- *Low GPU Utilization problem*: data transmission time > calculation time



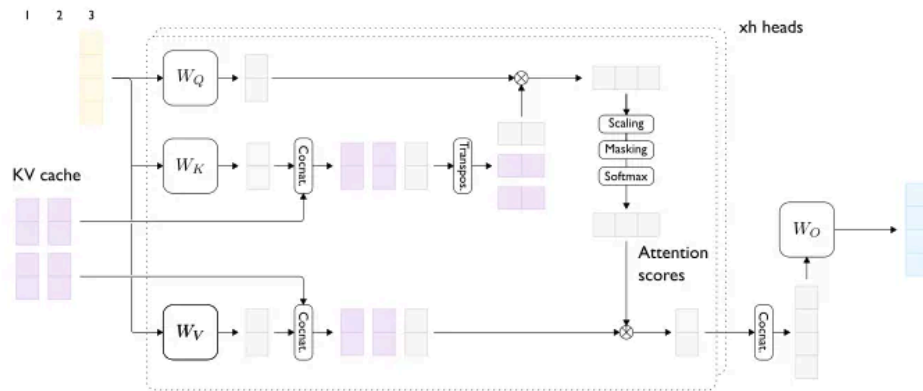Figure 7: Generating Steps with KVCache Enabled

# Outline

Review

## Motivation

Architecture

Evaluations

References

# Modelization

## Primary Goal

*Diversified workload* of LLM models $\rightarrow$ *Optimization problem* with multiple constraints:
- Maximize overall throughput
- Constraints of latency-related SLOs:

Workload difference in:
- I/O length
- Frequency
- Distribution of arrival
- SLOs

# Two Types of SLOs



Figure 2: Normalized throughput and latency of prefill and decoding stages with different sequence lengths or batch sizes for the dummy LLaMA2-70B model.

Computation time increases:

- At *prefill stage*, superlinearly with input length due to parallel processing of tokens → **Time to first token (TTFT)**
- At *decoding stage*, sublinearly with batch size due to auto-regressive processing on one token at a time per batch → **Time between tokens (TBT)**

# Going More Concrete

## Disaggregation Idea

Make best use of resources $\rightarrow$

1. Scheduling of KVCache is central to LLM serving scheduling;
2. Decouple and reconstruct nodes for different but collabrative goals

Current approaches:

- Reuse KVCache as much as possible
- Maximize the number of tokens in each batch (to improve **Model FLOPs Utilization (MFU)**)

However, corresponding issues:

- If tokens stored in remote places $\rightarrow$ TTFT problem
- Large batch size $\rightarrow$ TBT problem

# Another Objective

Limited GPU/accelerator supply while excessive demand.

- Early reject policy if no available slots could be use
- However, fluctuation in the workloads

Their work:

## Overload-Oriented Scheduling

- Prediction of generation length
- Load prediction in the short-tern future
- Classify request priorities.

# Outline

# Overview: KVCache-Centric Disaggregated Architecture



Figure 1: Mooncake Architecture.

Quick-reminder: reuse KVCache / max num of tokens in each batch.
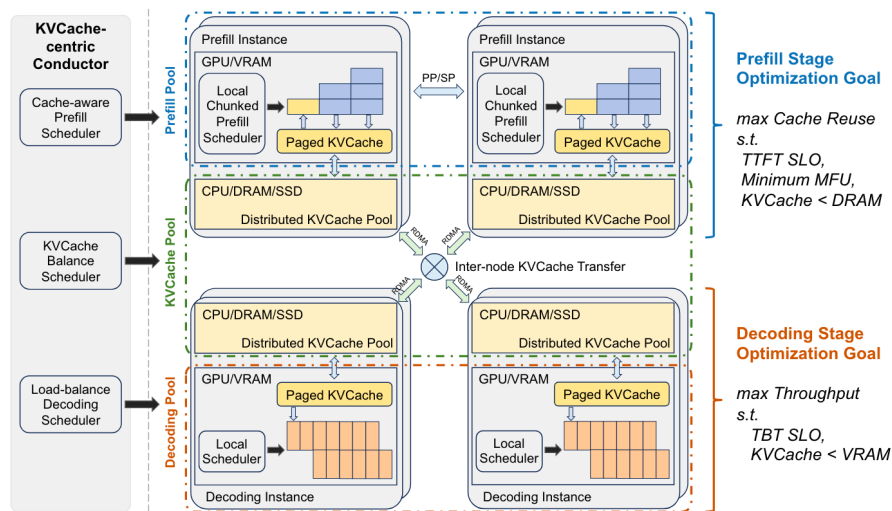
## Conductor: The Global Scheduler

Selects a pair of prefill and decoding instances and schedule the request to:
1. Transfer reusable KVCache
2. (Prefill Stage) Complete in chunks/layers, stream the output to the corresponding decoding instance
3. (Decoding Stage) Load KVCache and add it to the batching process, generate output

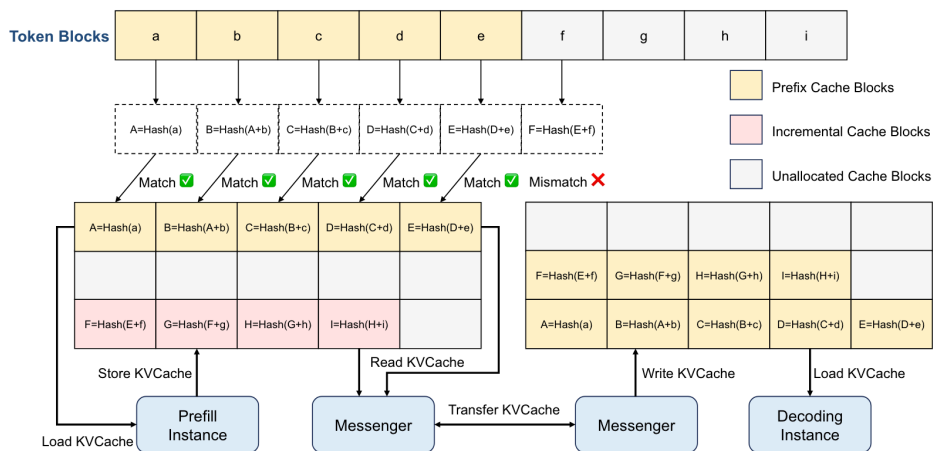# Storage and Transfer Logic of the KVCache Blocks



Figure 3: The KVCache pool in CPU memory. Each block is attached with a hash value determined by both its own hash and its prefix for deduplication.

- In CPU memory, KVCache stored as page blocks.
- Cache eviction algorithms used
- Transfer across CPUs and GPUs handled by a RDMA-based component called *Messenger*
- Global scheduler: *Conductor*

# Workflow of a Request (C3)



Figure 4: Workflow of inference instances. (∗) For prefill instances, the load and store operations of the KVCache layer are performed layer-by-layer and in parallel with the prefill computation to mitigate transmission overhead (see §4.2). (†) For decoding instances, asynchronous loading is performed concurrently with GPU decoding to prevent GPU idle time.

After tokenizing is finished, conductor selects a pair of prefill nodes and a decoding node:

1. **KVCache reuse**: loads the prefix cache from remote CPU memory into GPU memory
2. **Incremental Prefill**: Prefill node complete prefill stage using the prefix cache, afterwards it stores the *newly generated and incremental* KVCache back into CPU memory
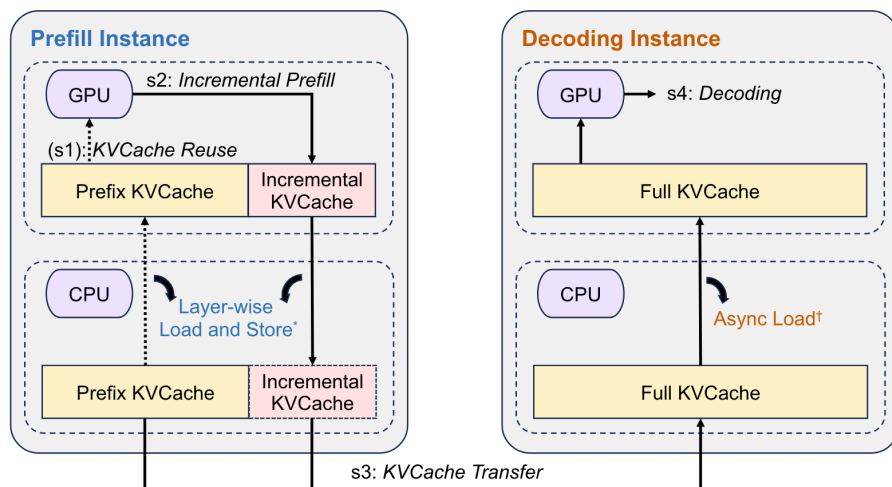
# Workflow of a Request (C3)



Figure 4: Workflow of inference instances. (∗) For prefill instances, the load and store operations of the KVCache layer are performed layer-by-layer and in parallel with the prefill computation to mitigate transmission overhead (see §4.2). (†) For decoding instances, asynchronous loading is performed concurrently with GPU decoding to prevent GPU idle time.

3. **KVCache Transfer**: Operated by *Messenger*, cross-machine KVCache transformer, asynchronously exectued and streaming the KVCache

4. **Decoding**: KVCache is received in the CPU DRAM of the decoding node, request joins the next batch in a continuous batching manner

# KVCache-centric Scheduling

## Prefill Global Scheduling

Selection of prefill instances depend not only on load, but also *prefix cache hit length* and *distribution of reusable KVCache blocks.*



**Prefill Instance**

GPU

s2: *Incremental Prefill*

(s1): *KVCache Reuse*

Prefix KVCache | Incremental KVCache

CPU

Layer-wise Load and Store

Prefix KVCache | Incremental KVCache

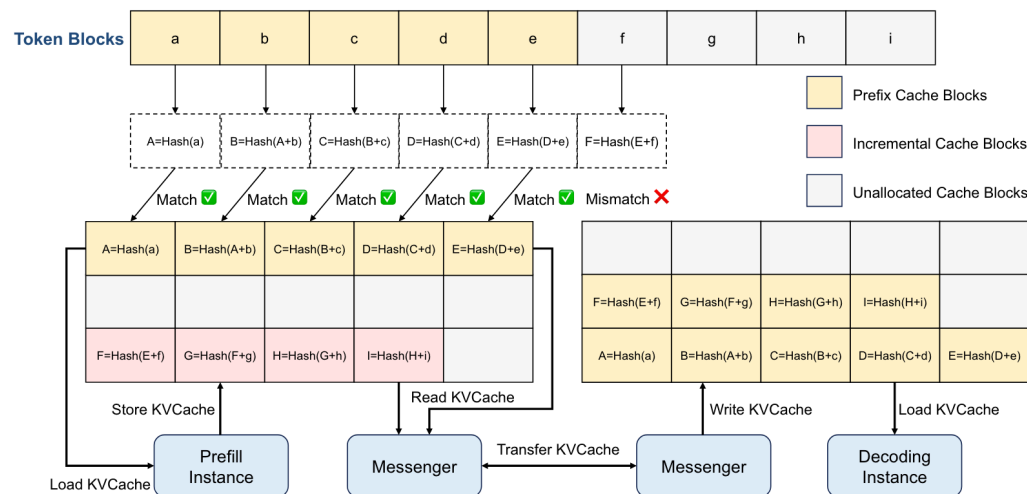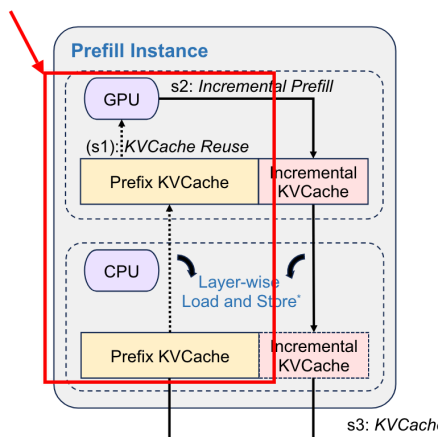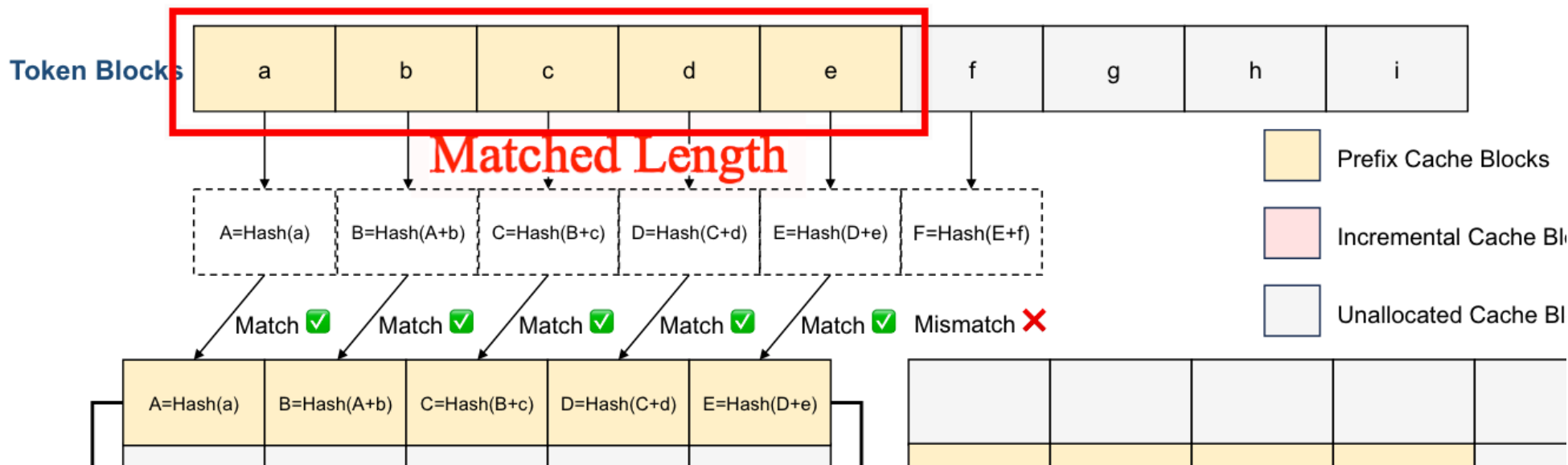Load KVCache

s3: *KVCache*



Figure 3: The KVCache pool in CPU memory. Each block is attached with a hash value determined by both its own hash and its prefix for deduplication.

# KVCache-centric Scheduling

- Input token for each request divided into blocks, assign hash keys one by one.
- Compare one by one to identify the prefix match length, named `prefix_len`

# 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:  $TTFT \leftarrow \inf$
3:  $p \leftarrow \emptyset$
4:  $best\_prefix\_len, best\_matched\_instance \leftarrow \text{FindBestPrefixMatch}(P, block\_keys)$
5:  **for** $instance \in P$ **do**
6:      $prefix\_len \leftarrow instance.prefix\_len$  ◀ **Matched Length**
7:      $T_{queue} \leftarrow \text{EstimatePrefillQueueTime}(instance)$  ◀ **Waiting Time**
8:      **if** $\frac{best\_prefix\_len}{prefix\_len} < $ **kvcache_balancing_threshold then**  ▷ Cache-aware prefill scheduling
9:          $T_{prefill} \leftarrow \text{EstimatePrefillExecutionTime}(\text{len}(R.prompt\_tokens), prefix\_len)$
10:         **if** $TTFT > T_{queue} + T_{prefill}$ **then**
11:             $TTFT \leftarrow T_{queue} + T_{prefill}$  ◀ **Prefill Exe Time**
12:             $p \leftarrow instance$
13:         **end if**
14:     **else**  ▷ Cache-aware and -balancing prefill scheduling
15:         $transfer\_len \leftarrow best\_prefix\_len - prefix\_len$
16:         $T_{transfer} \leftarrow \text{EstimateKVCacheTransferTime}(instance, best\_matched\_instance, transfer\_len)$
17:         $T_{prefill} \leftarrow \text{EstimatePrefillExecutionTime}(\text{len}(R.prompt\_tokens), best\_prefix\_len)$
18:         **if** $TTFT > T_{transfer} + T_{queue} + T_{prefill}$ **then**
19:             $TTFT \leftarrow T_{transfer} + T_{queue} + T_{prefill}$
20:             $p \leftarrow instance$  ◀ **Transfer Cache Time**
21:         **end if**
22:     **end if**
23: **end for**
24: $d, TBT \leftarrow \text{SelectDecodingInstance}(D)$  ▷ Load-balancing decoding scheduling
25: **if** $TTFT > TTFT\_SLO$ **or** $TBT > TBT\_SLO$ **then**
26:     **reject** $R$; **return**
27: **end if**
28: **if** $\frac{best\_prefix\_len}{p.prefix\_len} > $ **kvcache_balancing_threshold then**
29:     $\text{TransferKVCache}(best\_matched\_instance, p)$  ▷ KVCache hot-spot migration
30: **end if**
31: **return** $(p, d)$

---

- Find shortest TTFT
- Update $T_{\text{cache}}$ and $T_{\text{queue}}$
- If SLO not achievable, return `HTTP 429 TOO MANY REQUESTS`.

# KVCache-centric Scheduling

- Each prefill machine mangages its own set of prefix caches.
- Goal: achieve balance between cache matching and instance load (e.g. system prompts vs long document)
- Collect global usages of each block, use for forecasting ? (No way, requests fluctuation)

**Cache Load Balancing**

Requests may not always be directed to the prefill instance with the longest prefix cache length due to high instance load.

*Conductor* forwards the cache's location and the request to an alternative instance.

- Improve replication of hot-spot caches.

# Overload-Oriented Scheduling (C6)

**Overload**: Unrealistic to process every incoming request.

System should process as many requests as possibles until the system load reaches a predefined threshold.

**Load** = Ratio of num of requests being processed to the system's max capacity

Decision making:
- Whether to accept the prefill stage based on the *prefill instance's load*
- Whether to proceed with the decoding stage based on the *decoding instance's load*

# Overload-Oriented Scheduling (C6)

**Time-lag** between two instances for a single request: If a request is rejected by the decoding instance due to high load after the prefill stage has been completed $\rightarrow$ Waste of resources

**Early Rejection**

Advance the load assessement of the decoding instance to precede the beginning of the prefill stage.

# Overload-Oriented Scheduling (C6)

Problems of Early Rejection: Fluctuation



(a) Early Rejection.

# Overload-Oriented Scheduling (C6)

**Early Rejection based on Prediction**

**Request level**

Predict Output Length, thus estimate accurately TTFT and TBT.

**System level**

Estimate the overall batch count or the TBT status for instances after a specified time.

Their work: System level, request level left for future work.

# Outline

# Evaluations

Mooncake demonstrates significantly higher throughput, with enhancements ranging from 50% to 525%, while adhering to the same TTFT and TBT SLO constraints compared to vLLM.

Mooncake can mitigate load fluctuations, increasing the request handling capacity.

# Outline

Review

Motivation

Architecture

Evaluations

**References**

# References

- Review: LLM Inference
  - https://medium.com/@plienhar/llm-inference-series-2-the-two-phase-process-behind-llms-responses-1ff1ff021cd5
  - https://medium.com/@plienhar/llm-inference-series-3-kv-caching-unveiled-048152e461c8
  - https://www.bilibili.com/video/BV1TZ421j7Ke/?spm_id_from=333.337.search-card.all.click&vd_source=5ade9da381cec8d2c191f450ccd0cf57