



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

Ruoyu Qin, *Moonshot AI and Tsinghua University*; Zheming Li, Weiran He, and Jialei Cui, *Moonshot AI*; Feng Ren, Mingxing Zhang, Yongwei Wu, and Weimin Zheng, *Tsinghua University*; Xinran Xu, *Moonshot AI*

<https://www.usenix.org/conference/fast25/presentation/qin>

This paper is included in the Proceedings of the  
23rd USENIX Conference on File and Storage Technologies.

February 25–27, 2025 • Santa Clara, CA, USA

ISBN 978-1-939133-45-8

Open access to the Proceedings  
of the 23rd USENIX Conference on  
File and Storage Technologies  
is sponsored by

 **NetApp®**

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

Ruoyu Qin<sup>♠1</sup> Zheming Li<sup>♠1</sup> Weiran He<sup>♠</sup> Jialei Cui<sup>♠</sup> Feng Ren<sup>♡</sup>  
Mingxing Zhang<sup>♡2</sup> Yongwei Wu<sup>♡</sup> Weimin Zheng<sup>♡</sup> Xinran Xu<sup>♠2</sup>  
<sup>♠</sup>Moonshot AI <sup>♡</sup>Tsinghua University

## Abstract

MOONCAKE is the serving platform for **Kimi**, an LLM chatbot service developed by Moonshot AI. This platform features a KVCache-centric disaggregated architecture that not only separates prefill and decoding clusters but also efficiently utilizes the underexploited CPU, DRAM, SSD and NIC resources of the GPU cluster to establish a disaggregated KV-Cache. At the core of MOONCAKE is its KVCache-centric global cache and a scheduler designed to maximize throughput while adhering to stringent latency-related Service Level Objectives (SLOs).

Our experiments demonstrate that MOONCAKE excels in scenarios involving long-context inputs. In tests using **real traces**, MOONCAKE increases the effective request capacity by 59%~498% when compared to baseline methods, all while complying with SLOs. Currently, MOONCAKE is operational across thousands of nodes, processing over 100 billion tokens daily. In practical deployments, MOONCAKE’s innovative architecture enables Kimi to handle **115% and 107%** more requests on NVIDIA A800 and H800 clusters, respectively, compared to previous systems.

## 1 Introduction

With the rapid adoption of large language models (LLMs) in various scenarios [1–4], the workloads for LLM serving have become significantly diversified. These workloads differ in input/output length, distribution of arrival, and, most importantly, **demand different kinds** of Service Level Objectives (SLOs). As a Model as a Service (MaaS) provider, one of the primary goals of Kimi [5] is to solve an optimization problem with multiple complex constraints. The optimization goal is to maximize overall effective throughput, which directly impacts revenue, while the constraints reflect varying levels of SLOs. These SLOs typically involve meeting latency-related requirements, mainly the time to first token (TTFT) and the time between tokens (TBT).

To achieve this goal, a prerequisite is to make the best use of the various kinds of resources available in the GPU cluster. Specifically, although GPU servers are currently provided as highly integrated nodes (e.g., DGX/HGX supercomputers [6]),

<sup>1</sup> Ruoyu Qin’s part of work done as an intern at Moonshot AI, contributed equally with Zheming Li.

<sup>2</sup> Corresponding to zhang\_mingxing@mail.tsinghua.edu.cn, xuxinran@moonshot.cn.

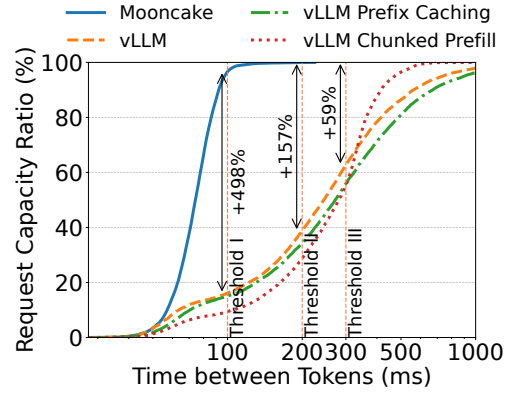


Figure 1: The experiment of the effective request capacity of MOONCAKE under the real-world conversation workload and different TBT SLOs. In this experiment, MOONCAKE and three baseline systems utilize 16  $8\times A800$  nodes each. More on §5.2.

it is necessary to decouple and restructure them into several disaggregated resource pools, each optimized for different but collaborative goals. For example, many other researchers [7–9] have suggested separating prefill servers from decoding servers, because these two stages of LLM serving have very different computational characteristics.

Further advancing this disaggregation strategy, we have engineered a disaggregated KVCache by pooling CPU, DRAM, SSD and RDMA resources of the GPU cluster, referred to as **MOONCAKE Store**. This novel architecture harnesses underutilized resources to enable efficient near-GPU prefix caching, significantly enhancing the global cache capacity and inter-node transfer bandwidth. The resultant distributed KVCache system embodies the principle of **trading more storage for less computation**. Thus, as demonstrated in Figure 1, it substantially boosts Kimi’s maximum throughput capacity in meeting the required SLOs for many important real-world scenarios. Later in this paper, we will first delve into a mathematical analysis of this strategy’s benefits for LLM serving and empirically assess its efficacy using real-world data (§2.2). Then, we will detail the design choices made in implementing this petabyte-level disaggregated cache, which is interconnected via an RDMA network up to  $8\times 400$  Gbps (§3.2).

Building on this idea, we also found that the scheduling of KVCache is central to LLM serving, and hence propose a corresponding disaggregated architecture. Figure 2 presents

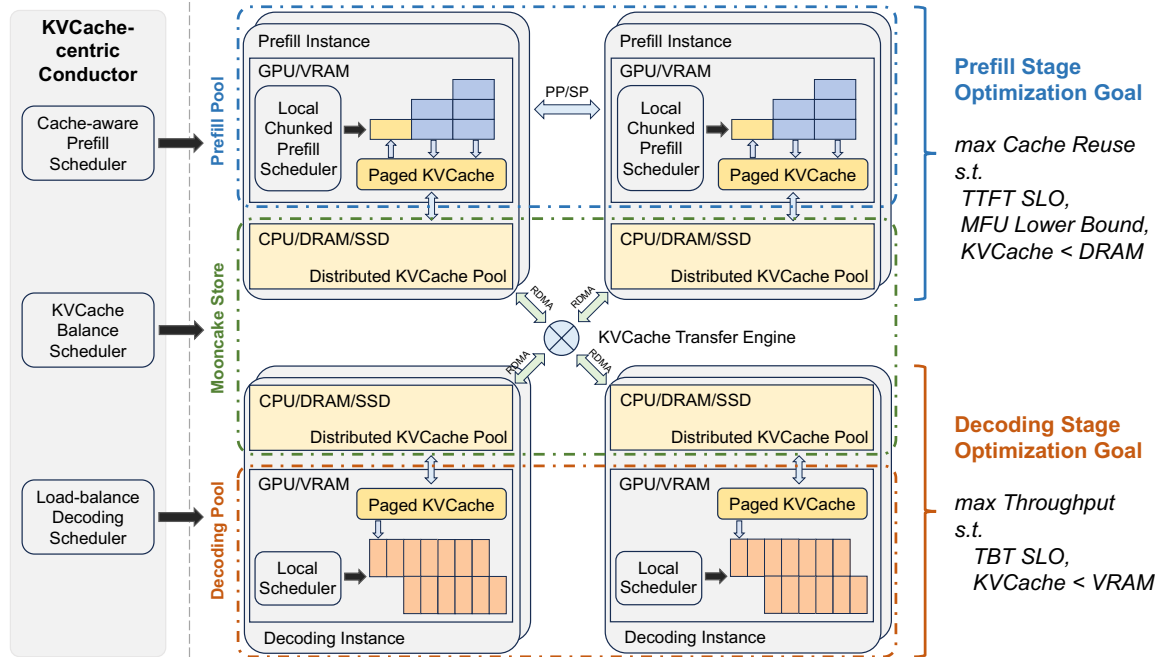


Figure 2: MOONCAKE Architecture.

our current **KVCache-centric disaggregated architecture** for LLM serving, named MOONCAKE. For each request, the global scheduler (Conductor) will select a pair of prefill and decoding instances and schedule the request in the following steps: 1) transfer as much reusable KVCache as possible to the selected prefill instance; 2) complete the prefill stage in chunks/layers and continuously stream the output KVCache to the corresponding decoding instance; 3) load the KVCache and add the request to the continuous batching process at the decoding instance for generating request outputs.

Although this process seems straightforward, the selection policy is complex due to many restrictions. In the prefill stage, the main objective is to reuse the KVCache as much as possible to avoid redundant computation. However, the distributed KVCache pool faces challenges in terms of both capacity and access latency. Thus Conductor is responsible for scheduling requests with KVCache-awareness and executing scheduling operations such as swapping and replication accordingly. The hottest blocks should be replicated to multiple nodes to avoid fetching congestion, while the coldest ones should be swapped out to reduce reserving costs. In contrast, the decoding stage has different optimization goals and constraints. The aim is to aggregate as many tokens as possible in a decoding batch to improve the Model FLOPs Utilization (MFU). However, this objective is restricted not only by the TBT SLO but also by the total size of aggregated KVCache that can be contained in the VRAM.

In §4, we will detail our KVCache-centric request scheduling algorithm, which balances instance loads and user experience as measured by TTFT and TBT SLOs. This includes a heuristic-based automated hotspot migration scheme that

replicates hot KVCache blocks without requiring precise predictions of future KVCache usage. Experimental results show that our KVCache-centric scheduling can significantly lower TTFT in real-world scenarios.

We will also describe the main design choices made during its implementation, especially those not covered in current research. For example, regarding P/D disaggregation, there are currently debates on its feasibility in large-scale practice due to bandwidth requirements and trade-offs associated with chunked prefill (e.g., Sarathi-Serve [10]). We demonstrate, through comparison with vLLM, that with a highly optimized transfer engine, the communication challenges can be managed, and P/D disaggregation is preferable for scenarios with stringent SLO limits (§5.2). Additionally, we discuss how to implement a separate prefill node pool that seamlessly handles the dynamic distribution of context length. We employ a chunked pipeline parallelism (CPP) mechanism to scale the processing of a single request across multiple nodes, which is necessary for reducing the TTFT of long-context inputs. Compared to traditional sequence parallelism (SP) based solutions, CPP reduces network consumption and simplifies the reliance on frequent elastic scaling (§3.3).

MOONCAKE is currently the serving platform of Kimi and has successfully handled exponential workload growth (more than 100 billion tokens a day). According to our historical statistics, the innovative architecture of MOONCAKE enables Kimi to handle 115% and 107% more requests on the A800 and H800 clusters, respectively, compared to previous systems.

To ensure the reproducibility of our results while safeguarding proprietary information, we also provide detailed



experimental outcomes using a dummy model mirroring the architecture of LLaMA3-70B, based on replayed traces of actual workloads. These traces, along with the KVCache transfer infrastructure of MOONCAKE, are open-sourced at <https://github.com/kvcache-ai/Mooncake>.

In end-to-end experiments using public datasets and real workloads, MOONCAKE excels in long-context scenarios. Compared to the baseline method, MOONCAKE can achieve up to a 498% increase in the effective request capacity while meeting SLOs. In §5.3, we compare MOONCAKE Store with the local cache design and find that the global cache design of MOONCAKE Store significantly improves the cache hit rate. In our experiments, the cache hit rate is up to  $2.36\times$  higher than that of the local cache, resulting in up to 48% savings in prefill computation time. MOONCAKE, to the best of our knowledge, is the first system to demonstrate the significant benefits of using a distributed KVCache pool to share KV-Cache across different chat sessions and queries in large-scale deployment scenarios. We also evaluate the performance of the transfer engine that supports high-speed RDMA transfers in MOONCAKE, which shows it is approximately  $2.4\times$  and  $4.6\times$  faster than existing solutions (§5.4).

## 2 Preliminary and Problem Definition

### 2.1 Service Level Objectives of LLM Serving

Modern large language models (LLMs) are based on the Transformer architecture, which utilizes attention mechanisms and multilayer perceptrons (MLP) to process input. Popular Transformer-based models, such as GPT [11] and LLaMA [12], employ a decoder-only structure. Each inference request is logically divided into two stages: the prefill stage and the decoding stage.

During the prefill stage, all input tokens are processed in parallel, and hence it is typically computationally intensive. This stage generates the first output token while storing intermediate results of computed keys and values, referred to as the KVCache. The decoding stage then uses this KVCache to autoregressively generate new tokens. It processes only one token at a time per batch due to the limitation of autoregressive generation, which makes it memory-constrained and causes computation time to increase sublinearly with batch size. Thus, a widely used optimization in the decoding stage is continuous batching [13, 14]. Before each iteration, the scheduler checks the status and adds newly arrived requests to the batch while removing completed requests.

Due to the distinct characteristics of the prefill and decoding stages, MaaS providers set different metrics to measure their corresponding Service Level Objectives (SLOs). Specifically, the prefill stage is mainly concerned with the latency between the request arrival and the generation of the first token, known as the time to first token (TTFT). On the other hand, the decoding stage focuses on the latency between suc-

Notation	Description	Value
$l$	Num layers	80
$d$	Model dimension	8192
$a, b$	Constant coefficients in Equation 1	4, 22
$gqa$	Num $q$ heads / Num $kv$ heads	8
$s$	Tensor element size	2 B (BFloat16)
$G$	GPU computation throughput	$8\times 312$ TFLOPS
$B_{h2d}$	Host-to-device bandwidth	128 GB/s
$B_{nic}$	NIC bandwidth	800 Gbps
$n, p$	Prompt and matched prefix length, respectively	

Table 1: Notations and parameters. Model and machine parameters are set according to LLaMA3-70B and  $8\times A800$ .

cessive token generations for the same request, referred to as the time between tokens (TBT).

In real deployments, if the monitor detects unmet SLOs, we need to either add inference resources or reject some incoming requests. However, due to the current contingent supply of GPUs, elastically scaling out the inference cluster is typically unfeasible. Therefore, we proactively reject requests that are predicted not to meet the SLOs to alleviate the cluster’s load. Our main objective is to maximize overall throughput while adhering to SLOs, a concept referred to as goodput in other research [8, 15].

### 2.2 More Storage for Less Computation

To meet the stringent SLOs described above, a commonly adopted solution is to cache previously generated KVCache and reuse it upon finding a prefix match. However, existing approaches [16–18] typically restrict caching to local HBM and DRAM, assuming that the transfer bandwidth required for global scheduling would be prohibitively high. But, as we will described later in §5.3, the capacity of local DRAM supports only up to 50% of the theoretical cache hit rate, making the design of a global cache essential. In this section, we present a mathematical analysis of the actual bandwidth necessary to benefit from this strategy, explaining why distributed caching is advantageous, especially for larger models like LLaMA3-70B. More experimental results will be given later in §5.4.2.

We base our analysis on the model using notations described in Table 1 and incorporate specific parameters of LLaMA3-70B. Essentially, current popular LLMs are autoregressive language models where each token’s KVCache depends only on itself and preceding tokens. Therefore, KVCache corresponding to the same input prefix can be reused without affecting output accuracy. If a current request’s prompt of length  $n$  shares a common prefix of length  $p$  with previously cached KVCache, its prefill process can be optimized as follows:

$$\begin{aligned}
q[p:n], k[p:n], v[p:n] &= \text{MLP}(\text{hidden}[p:n]) \\
k[1:n], v[1:n] &\leftarrow \text{KVCache} + (k[p:n], v[p:n]) \\
o[p:n] &= \text{Attention}(q[p:n], k[1:n], v[1:n]) \\
\text{KVCache} &\leftarrow (k[1:n], v[1:n])
\end{aligned}$$

Given input length  $n$ , the FLOPS of the prefill stage can be calculated as:

$$flops(n) = l \times (an^2d + bnd^2) \quad (1)$$

Thus, reusing KVCache approximately reduces the computation cost of prefill by  $l \times (ap^2d + bpd^2)$ . However, this requires transferring the cached KVCache into the prefill GPU's HBM, with a size of  $p \times l \times (2 \times d / gqa) \times s$ . Assuming the average computation throughput is  $G$  and the average KVCache loading speed is  $B$  (where  $B$  is determined by the minimum of  $B_{h2d}$  and  $B_{nic}$ ), the reuse of KVCache is beneficial in terms of TTFT if:

$$\frac{B}{G} > \frac{2ds}{gqa \times (apd + bd^2)} \quad (2)$$

In such scenarios, reusing the KVCache not only reduces GPU time and costs but also enhances the user experience by improving TTFT. The criteria for bandwidth  $B$  relative to computation throughput  $G$  are more readily met with larger values of  $d$ , which is proportional to the model size. For example, when running LLaMA3-70B on a machine with  $8 \times A800$  GPUs and assuming a prefix length of 8192, Equation 2 yields a minimum required  $B$  of 6 GB/s. The requirement for  $B$  will be enlarged to 19 GB/s for an  $8 \times H800$  machine. Moreover, in practical scenarios, because the transfer stages cannot be perfectly overlapped with each other, the actual bandwidth requirement is even higher. However, as we will demonstrate in §5.4.2, a fully utilized 100 Gbps NIC per NVIDIA A800 HGX network is sufficient to meet these criteria.

### 3 Design of MOONCAKE

#### 3.1 Overview

As depicted in Figure 2, MOONCAKE employs a disaggregated architecture that not only separates prefill from decoding nodes but also groups the CPU, DRAM, SSD, and RDMA resources of the GPU cluster to implement a disaggregated KVCache. To schedule all these disaggregated components, at its center, MOONCAKE implements a global scheduler named Conductor. Conductor is responsible for dispatching requests based on the current distribution of the KVCache and workload characteristics. MOONCAKE Store, detailed in §3.2, manages the storage and transfer of these KVCache blocks.

Specifically, Figure 3 demonstrates the typical workflow of a request. Once tokenizing is finished, the conductor selects a pair of prefill nodes and a decoding node, and starts a workflow comprising four steps:

1) *KVCache Reuse*: The selected prefill node (group) receives a request that includes the raw input, the block keys of the prefix cache that can be reused, and the block keys of the full cache allocated to the request. It loads the prefix cache from remote CPU memory into GPU memory based on the prefix cache block keys to bootstrap the request. This step is

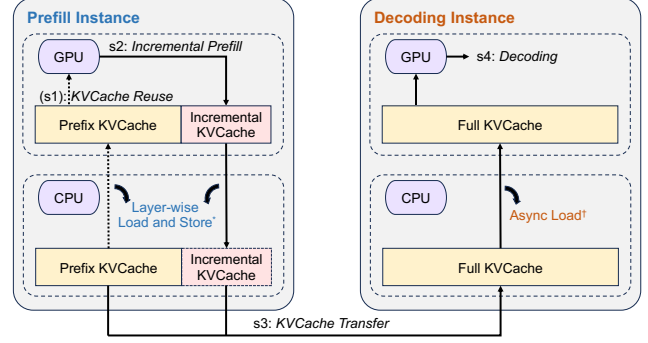


Figure 3: Workflow of inference instances.

skipped if no prefix cache exists. This selection balances three objectives: reusing as much KVCache as possible, balancing the workloads of different prefill nodes, and guaranteeing the TTFT SLO. It leads to a KVCache-centric scheduling that will be further discussed in §4.

2) *Incremental Prefill*: The prefill node completes the prefill stage using prefix cache and stores the newly generated incremental KVCache back into CPU memory. If the number of uncached input tokens exceeds a certain threshold, the prefill stage is split into multiple chunks and executed in a pipeline manner. This threshold is selected to fully utilize the corresponding GPU's computational power and is typically larger than 1000 tokens. The reason for using chunked but still disaggregated prefill nodes is explained in §3.3.

3) *KVCache Transfer*: MOONCAKE Store is deployed in each node to manage and transfer these caches. This step is asynchronously executed and overlapped with the above incremental prefill step, streaming the KVCache generated by each model layer to the destination decoding node's CPU memory to reduce waiting time.

4) *Decoding*: After all the KVCache is received in the CPU memory of the decoding node, the request joins the next batch in a continuous batching manner. The decoding node is pre-selected by Conductor based on its current load to ensure it does not violate the TBT SLO.

#### 3.2 MOONCAKE Store: Cache of KVCache

Central to MOONCAKE is its efficient implementation of a distributed global cache of KVCache, referred to as MOONCAKE Store. As described in §2.2, reusing cached KVCache not only cuts computation costs but also improves user experience by reducing the TTFT, particularly when the aggregated bandwidth is fully utilized. However, achieving full utilization is challenging because the bandwidth can reach up to  $8 \times 400$  Gbps, comparable to DRAM bandwidth.

We first introduce how MOONCAKE Store manages KVCache in §3.2.1, including its storage scheme and eviction policy. In §3.2.2, we describe the object-based APIs and memory transfer APIs of MOONCAKE Store. In §3.2.3, we will detail the design of MOONCAKE Store's transfer engine, a

high-performance, zero-copy KVCache transfer system designed to maximize the benefits of using multiple RDMA NICs per machine. It enhances execution efficiency and reliability through techniques such as topology-aware path selection and endpoint pooling.

### 3.2.1 KVCache Management

In MOONCAKE Store, all KVCache is stored as paged blocks within a distributed cache pool. The block size, i.e., the number of tokens contained in each block, is determined by the model size and the optimal network transmission size, typically ranging from 16 to 512 tokens. Each block is attached with a hash key determined by both its own hash and its prefix for deduplication. The same hash key may have multiple replicas across different nodes to mitigate hot-cache access latency, controlled by our cache-load-balancing policy described in §4.2.

MOONCAKE Store allocates space for each cache block in the cache pool and logs metadata such as the block key and its address. When the cache pool is full, MOONCAKE Store employs a LRU (Least Recently Used) strategy to evict an existing cache block—unless the block is currently being accessed by an ongoing request—and overwrites the evicted block's space with the new block.

### 3.2.2 Interface

At the higher layer, MOONCAKE Store offers object-based APIs such as `put`, `get`, and `change_replica`. These facilitate the caching of KVCache in a disaggregated manner, organizing mini blocks of KVCache as memory objects and enabling Conductor to adjust the number of replicas for each KVCache block to achieve higher bandwidth aggregation. These functions are supported by a set of synchronous batch transfer APIs, detailed in Listings 1.

Transfer operations are available for both DRAM and GPU VRAM and will utilize GPU Direct RDMA when optimal, provided that the specified memory region has been pre-registered. The completion of these operations can be monitored asynchronously via the `getTransferStatus` API, which reports whether transfers are ongoing or have encountered errors.

Listing 1: Memory transfer APIs in MOONCAKE Store.

```
int registerLocalMemory(void *vaddr, size_t len,
                        const string &type);
BatchID allocateBatchID(size_t batch_size);
int submitTransfer(BatchID batch_id,
                  const vector<Request> &entries);
int getTransferStatus(BatchID batch_id,
                     int request_index,
                     Status &status);
int freeBatchID(BatchID batch_id);
```

### 3.2.3 Transfer Engine

To efficiently implement the above APIs, a transfer engine is designed to achieve several key objectives: 1) Effectively distribute transfer tasks across multiple RDMA NIC devices; 2) Abstract the complexities of RDMA connection management from the APIs; and 3) Appropriately handle temporary network failures. This transfer engine has been meticulously engineered to fulfill each of these goals.

**Network Setup** The benefits of MOONCAKE rely on a high-bandwidth network interconnect. Currently, we use standard HGX machines where each A800 GPU is paired with a 100/200 Gbps NIC, and each H800 GPU is paired with a 200/400 Gbps NIC, which is comparable to memory bandwidth and existing libraries (other than NCCL) fail to fully utilize this capacity. As for NCCL, it cannot gracefully handle dynamic topology changes due to the addition or removal of nodes/NICs and does not support DRAM-to-DRAM pathes. In contrast, the transfer engine endeavors to find alternative paths upon failure.

To address congestion, the network utilizes RoCEv2 tuned by cloud providers. In the scheduler, we mitigate congestion by increasing the number of replicas for hot KVCaches (§4.2).

**Topology-aware path selection.** Modern inference servers often consist of multiple CPU sockets, DRAM, GPUs, and RDMA NIC devices. Although it's technically possible to transfer data from local DRAM or VRAM to a remote location using any RDMA NIC, these transfers can be limited by the bandwidth constraints of the Ultra Path Interconnect (UPI) or PCIe Switch. To overcome these limitations, MOONCAKE Store implements a topology-aware path selection algorithm.

Before processing requests, each server generates a *topology matrix* and broadcasts it across the cluster. This matrix categorizes network interface cards (NICs) into "preferred" and "secondary" lists for various *types* of memory, which types are specified during memory registration. Under normal conditions, a NIC from the preferred list is selected for transfers, facilitating RDMA operations within the local NUMA or GPU Direct RDMA through the local PCIe switch only. In case of failures, NICs from both lists may be utilized. The process involves identifying the appropriate local and target NICs based on the memory addresses, establishing a connection, and executing the data transfer.

For instance, as illustrated in Figure 4, to transfer data from buffer 0 (assigned to `cpu:0`) in the local node to buffer 1 (assigned to `cpu:1`) in the target node, the engine first identifies the preferred NICs for `cpu:0` using the local server's topology matrix and selects one, such as `mlx5_1`, as the local NIC. Similarly, the target NIC, such as `mlx5_3`, is selected based on the target memory address. This setup enables establishing an RDMA connection from `mlx5_1@local` to `mlx5_3@target` to carry out RDMA read and write operations.

To further maximize bandwidth utilization, a single request's transfer is internally divided into multiple slices at a

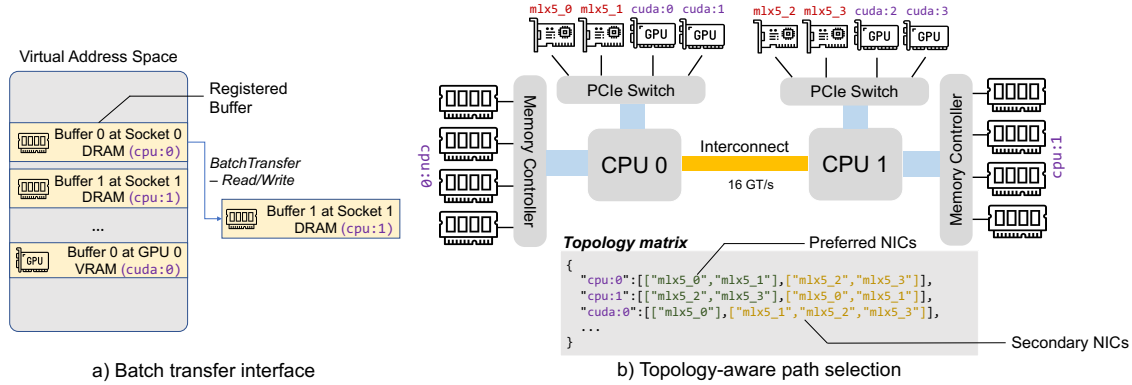


Figure 4: Transfer engine of MOONCAKE Store.

granularity of 16 KB. Each slice might use a different path, enabling collaborative work among all RDMA NICs.

**Endpoint management.** MOONCAKE Store employs a pair of *endpoints* to represent the connection between a local RDMA NIC and a remote RDMA NIC. In practice, each endpoint includes one or more RDMA queue pair objects. Connections in MOONCAKE Store are established in an on demand manner; endpoints remain unpaired until the first request is made.

To prevent a large number of endpoints from slowing down request processing, MOONCAKE Store employs endpoint pooling, which caps the maximum number of active connections. We use the SIEVE [19] algorithm to manage endpoint eviction. If a connection fails due to link errors, it is removed from the endpoint pools on both sides and re-established during the next data transfer attempt.

**Failure handling.** In a multi-NIC environment, one common failure scenario is the temporary unavailability of a specific NIC, while other routes may still connect two nodes. MOONCAKE Store is designed to adeptly manage such temporary failures effectively. If a connection is identified as unavailable, MOONCAKE Store automatically identifies an alternative, reachable path and resubmits the request to a different RDMA NIC device. Furthermore, MOONCAKE Store is capable of detecting problems with other RDMA resources, including RDMA contexts and completion queues. It temporarily avoids using these resources until the issue, such as a downed link, is resolved.

### 3.3 MOONCAKE’s Prefill Pool

Unlike the inviolable decoding nodes, the necessity and best practices for designing a separate and elastic prefill pool remain under debate. For example, although many researchers [7–9] share our intuition to use a disaggregated architecture, it is worth discussing whether this separation is still necessary with the introduction of chunked prefill [10].

However, after careful consideration, we decided to maintain MOONCAKE’s disaggregated architecture. This decision

is primarily driven by the fact that online services typically have more stringent SLOs. While chunked prefill reduces decoding interference, it remains challenging to simultaneously maximize MFU during the prefill stage and meet the TBT SLO during the decoding stage. We will demonstrate this in the end-to-end experiments in §5.2. Another important reason is that we think prefill nodes require different cross-node parallelism settings to handle long contexts as the available context length of recent LLMs is increasing rapidly, from 8k to 128k and even up to 1 million tokens [20]. Typically, for such long context requests, the input tokens can be 10 to 100 times larger than the output tokens, making optimizing the TTFT crucial. Due to the abundant parallelism in long context prefill, using more than a single  $8 \times \text{GPU}$  node to process them in parallel is desirable. However, extending tensor parallelism (TP) across more than one node requires two expensive RDMA-based all-reduce operations per layer, significantly reducing the MFU of prefill nodes.

Recently, many works have proposed sequence parallelism (SP) [21–27]. SP partitions the input sequences of requests across different nodes to achieve acceleration, allowing even long requests to meet the TTFT SLO. However, when applied to shorter input requests, SP results in a lower MFU compared to using single-node TP only. Recent research [15] proposes elastic sequence parallelism to dynamically scale up or down the SP group. Although possible, this adds complexity to our architecture. Additionally, SP still requires frequent cross-node communication, which lowers the MFU and competes with network resources for transferring KVCache across nodes.

To address this, MOONCAKE leverages the autoregressive property of decoder-only transformers and implements chunked pipeline parallelism (CPP) for long context prefill. We group every  $X$  nodes in the prefill cluster into a pipelined prefill node group. For each request, its input tokens are partitioned into chunks, each no longer than the *prefill\_chunk*. Different chunks of the same request can be processed simultaneously by different nodes, thus parallelizing the processing and reducing TTFT.

CPP offers two main benefits: 1) Similar to pipeline par-



allelism in training, it requires cross-node communication only at the boundaries of each pipeline stage, which can be easily overlapped with computation. This leads to better MFU and less network resource contention with KVCache transfer. 2) It naturally fits both short and long contexts, bringing no significant overhead for short context prefill and avoiding frequent dynamic adjustment of node partitioning. This pipeline-based acceleration method has been explored in training systems [28], but to our knowledge, this is the first application in the inference stage, as long context inference has only recently emerged.

## 4 Scheduling

### 4.1 Prefill Global Scheduling

Previous research on LLM serving typically uses a load-balancing strategy that evaluates the load on each instance based on the number of assigned requests. In MOONCAKE, however, the selection of prefill instances considers additional factors—not just load but also the prefix cache hit length and the distribution of reusable KVCache blocks. While there is a preference to route requests to prefill instances with longer prefix cache lengths to reduce computation costs, it may be beneficial to schedule them to other nodes to ensure overall system balance and meet TTFT SLOs. To address these complexities, we propose a cache-aware global scheduling algorithm that accounts for both the prefill time due to the prefix cache and the local queuing time.

Algorithm 1 details the mechanism for our KVCache-centric prefill scheduling. For every new request, block keys are then compared one by one against each prefill instance’s cache keys to identify the prefix match length (*prefix\_len*). With this matching information, Conductor estimates the corresponding execution time based on the request length and *prefix\_len* (which varies by instance), using a polynomial regression model fitted with offline data. It then adds the estimated waiting time for that request to get the TTFT on that instance. Finally, Conductor assigns the request to the instance with the shortest TTFT and updates the cache and queue times for that instance accordingly. If the SLO is not achievable, Conductor directly returns the HTTP 429 Too Many Requests response status code to the upper layers.

The backbone of this scheduling framework is straightforward, but complexities are hidden in the engineering implementation of various components. For example, to predict the computation time of the prefill stage for a request, we employ a predictive model derived from offline test data. This model estimates the prefill duration based on the request’s length and prefix cache hit length. Thanks to the regular computation pattern of Transformers, the error bound of this prediction is small as long as enough offline data is available. The queuing time for a request is calculated by aggregating the prefill times of all queued requests. In practical implementations,

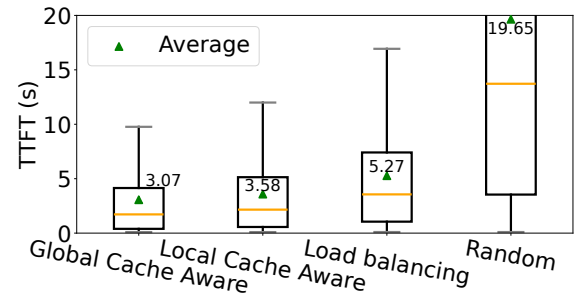


Figure 5: The prefill scheduling experiment.

TTFTs are computed in parallel, rendering the processing time negligible compared to the inference time.

More difficulty lies in predicting the transfer time because it is determined not only by the size of the transferred data but also by the current network status, especially whether the sending node is under congestion. This also necessitates the replication of hot KVCache blocks, which will be discussed in §4.2.

### 4.2 Cache Load Balancing

In MOONCAKE, each prefill instance has its own set of local prefix caches. The usage frequency of these caches varies significantly. For example, system prompts are accessed by almost every request, whereas caches storing content from a local long document may be used by only one user. As discussed in §4.1, Conductor’s role is crucial in achieving an optimal balance between cache matching and instance load. Thus, from the perspective of the distributed cache system, load balancing also plays an important role. Specifically, it involves strategizing on how to back up caches to ensure that global prefill scheduling can achieve both high cache hits and low load.

A straw-man solution to this KVCache scheduling problem could be collecting the global usages of each block, using a prediction model to forecast their future usages, and making scheduling decisions accordingly. However, unlike the estimation of prefill time, workloads are highly dynamic and change significantly over time. Especially for a MaaS provider experiencing rapid growth in its user base, it is impossible to accurately predict future usages. Thus, we propose a heuristic-based automated hotspot migration scheme to enhance cache load balancing.

As previously noted, requests may not always be directed to the prefill instance with the longest prefix cache length due to high instance load. In such cases, Conductor forwards the cache’s location and the request to an alternative instance if the estimated additional prefill time is shorter than the transfer time. This instance proactively retrieves the KVCache from the holder and stores it locally. More importantly, we prefer to compute the input tokens if the best remote prefix match length is no larger than the current local reusable prefix



---

**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:  $TTFT, p \leftarrow \text{inf}, \emptyset$ 
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  $(p, d)$ 
```

---

multiplied by a threshold<sup>1</sup>. Both strategies not only reduce the prefill time for requests but also facilitate the automatic replication of hotspot caches, allowing for their broader distribution across multiple instances.

To validate the effectiveness of our strategy, we conduct a scheduling experiment that compares random scheduling and load-balancing scheduling with our strategy. We further compare the local cache-aware scheduling described in §4.1 and the global cache-aware scheduling described in this section that considers cache load balancing. In random scheduling, a prefill instance is selected arbitrarily for each request. In load-balancing scheduling, the instance with the lightest load is chosen. Specifically, we build a MOONCAKE cluster consisting of 16  $8 \times A800$  nodes, and replay the conversation trace detailed in §5.2.1 for the experiment. We assess the performance of each scheduling algorithm based on the TTFTs. The experimental results, depicted in Figure 5, demonstrate that our KVCache-centric scheduling algorithms outperform random and load-balancing scheduling. By incorporating cache load balancing, the global cache-aware algorithm reduces the average TTFT by an additional 14% compared to the local cache-aware algorithm.

---

<sup>1</sup>This threshold is currently adjusted manually but can be adaptively adjusted by an algorithm in the future.

## 5 Evaluation

As described before, according to historical statistics of Kimi, MOONCAKE enables Kimi to handle 115% and 107% more requests on the A800 and H800 clusters, respectively, compared to our previous systems based on vLLM. To further validate this results and ensure reproducibility, in this section, we conduct a series of end-to-end and ablation experiments on MOONCAKE with a dummy LLaMA3-70B model to address the following questions: 1) Does MOONCAKE outperform existing LLM inference systems in real-world scenarios? 2) Compared to conventional prefix caching methods, does the design of MOONCAKE Store significantly improve MOONCAKE's performance?

### 5.1 Setup

**Testbed.** During the reproducing experiments, the system was deployed on a high-performance computing node cluster to evaluate its performance. Each node in the cluster is configured with eight NVIDIA-A800-SXM4-80GB GPUs and four 200 Gbps RDMA NICs. The KVCache block size in MOONCAKE Store is set to 256. For deploying MOONCAKE, each node operates as either a prefill instance or a decoding instance based on the startup parameters. For deploying other systems, each node hosts a single instance.

**Metric.** Specifically, we measure the TTFT and TBT of each request, where the TBT is calculated as the average of the longest 10% of the token arrival intervals. As mentioned in §2, the threshold for TTFT is set to 30 s, and TBT thresholds are set to 100 ms, 200 ms, and 300 ms, depending on the scenario. We consider requests with both TTFT and TBT below their respective thresholds as effective requests, and the proportion of effective requests among all requests as the effective request capacity. For brevity, the subsequent experiments not mentioning TTFT are assumed to meet the TTFT threshold. To more intricately compare the caching performance, we also measure the GPU time during the prefill stage and the cache hit rate for each request.

**Baseline.** We employ vLLM [14], one of the state-of-the-art open-source LLM serving systems, as our experimental baseline. vLLM features continuous batching and PagedAttention technologies, significantly enhancing inference throughput. Despite its strengths, vLLM's architecture, which couples the prefill and decoding stages, can disrupt decoding especially in scenarios involving long contexts. Recent updates to vLLM have integrated features like prefix caching and chunked prefill to improve performance metrics in long-context scenarios, such as TTFT and TBT. In our experiments, we also compare these features of vLLM. In our experiments, we utilize the latest release (v0.5.1) of vLLM. Due to limitations in the current implementation, we test the prefix cache and chunked prefill features of this version separately.

	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

Table 2: Workload Statistics.

## 5.2 End-to-end Performance

In our end-to-end experiments, we evaluate the request handling capabilities of MOONCAKE and baseline systems under various workloads. Specifically, we measure the maximum throughput that remains within the defined SLO thresholds. We employ three types of workloads in our tests: two real-world traces sampled from Kimi that represent online conversations and tool&agent interactions, respectively, and a synthetic workload to cover different inference scenarios. We will first describe the unique characteristics of these workloads and then discuss the results. Lastly, we analyze the GPU computation time during the prefill stage, further demonstrating the advantages of MOONCAKE Store in enhancing cache utilization and reducing computation costs.

### 5.2.1 Workload

**Conversation workload.** Chatbots [1, 5] represent one of the most prevalent applications of LLMs, making conversational requests a highly representative workload for LLM inference. As shown in Table 2, the conversation workload contains a significant portion of long-context requests—reaching up to 128k tokens and averaging around 12k tokens—which is comparable to the data lengths found in current long-context datasets [29, 30]. Moreover, the workload has an average of approximately 40% prefix caching ratio brought about by multi-turn conversations. We sampled 1 hour of conversation traces from an online inference cluster, where each record includes the input and output lengths along with timestamps of arrival. Requests are dispatched according to these timestamps and are preemptively terminated once the model output reaches the predetermined length.

**Tool&Agent workload.** Recent studies [31] involving LLMs deployed as tools or agents to perform tasks have been increasing. These tasks are typically characterized by the incorporation of pre-designed, often lengthy, system prompts that are fully repetitive. We collected traces of the tool&agent workload, also sampled over a 1-hour period. As indicated in Table 2, this workload exhibits a high proportion of prefix caching, with shorter input and output lengths.

**Synthetic workload.** The synthetic workload was constructed from a combination of publicly available datasets. We categorized the requests in the real trace into three types: short conversations, tool and agent calls, and long text summarization and QA. For each category, we selected the follow-

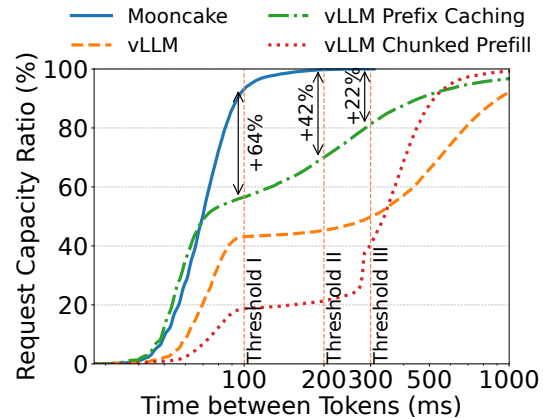


Figure 6: The experiment of the effective request capacity of MOONCAKE under the tool&agent workload.

ing datasets: ShareGPT [32], Leval [29], and LooGLE [30]. ShareGPT comprises multi-turn conversations with short input lengths. Leval serves as a benchmark for evaluating model performance over long contexts, simulating scenarios where requests involve lengthy system prompts typical of tool and agent interactions. LooGLE is tailored for long-context QA and summarization tasks, featuring input lengths of up to 100k tokens and including both multi-turn QA and single-turn summarizations, making it well-suited for long text summarization and QA scenarios. Overall, the synthetic workload has the longest average input length. Despite having the highest proportion of prefix caching, its cache hits are quite dispersed, thus requiring a substantial cache capacity.

During preprocessing, each conversation turn was mapped into a separate request, incorporating both the input and outputs from previous interactions. For datasets featuring multiple questions with the same lengthy prompt, each question and its preceding prompt were treated as a single request. We combined the processed datasets in a 1:1:1 ratio, preserving the sequential relationships within the multi-turn dialogue requests while randomly shuffling them. Since the datasets do not specify arrival times, we simulated realistic conditions by dispatching requests at a defined rate using a Poisson process.

### 5.2.2 Effective Request Capacity

To assess the maximum number of requests that can adhere to the SLOs under different workloads, we test four system configurations: MOONCAKE, vLLM, vLLM with the prefix caching feature, and vLLM with the chunked prefill feature, each utilizing 16 nodes.

**Conversation workload.** The results for this workload are presented in Figure 1. This workload, characterized by varying input lengths and longer output lengths, causes significant fluctuations in TBT for the vLLM system due to the lengthy contexts in the prefill stage. While chunked prefill reduces decoding interference, balancing the enhancement of MFU

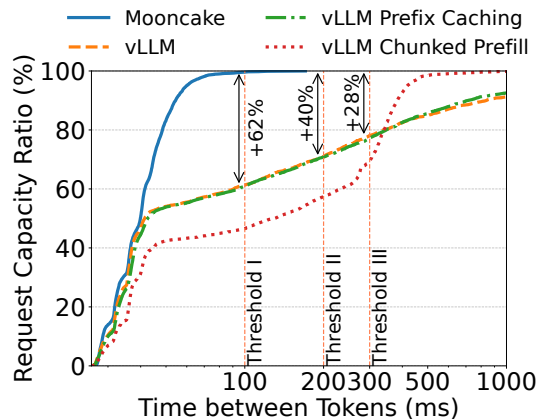


Figure 7: The experiment of the effective request capacity of MOONCAKE under the synthetic workload.

in the prefill stage with the TBT constraints in the decoding stage remains challenging. Despite meeting the TTFT SLO, its effective request capacity is still suboptimal. Compared to vLLM, MOONCAKE achieves a very significant increase in effective request capacity.

**Tool&Agent workload.** In contrast, the tool&agent workload has a high proportion of prefix caching and shorter output lengths, favoring the vLLM system as the short prefill time minimally impacts output. However, as illustrated in Figure 6, vLLM and vLLM with chunked prefill experience more severe disruptions in decoding due to longer prefill processing times, resulting in a lower effective caching capacity than vLLM with prefix caching. MOONCAKE uses a global cache pool to significantly increase caching capacity and optimize cache utilization through internode transfers, excelling in scenarios with high prefix caching. As a result, it enhances effective caching capacity by 42% compared to vLLM with prefix caching under the 200 ms threshold.

**Synthetic workload.** The synthetic workload features the longest average input lengths and dispersed cache hotspots which leads to poor cache utilization under smaller cache capacities. As depicted in Figure 7, most requests processed by MOONCAKE maintain a TBT within 100 ms, whereas about 20% of requests handled by vLLM exceed 300 ms. The performance of systems with prefix caching and chunked prefill is similar to vLLM, as they fail to mitigate the impact of long contexts on the decoding stage. Compared to vLLM, MOONCAKE increases effective request capacity by 40% under the 200 ms threshold.

### 5.2.3 Prefill GPU Time

Prefill GPU time is positively correlated with requests' TTFT and serving cost and is determined by requests' input lengths and cache hit rates. We analyze the average GPU time during the prefill stage under different workloads, as shown in Figure 8. For MOONCAKE, the conversation workload incurs the

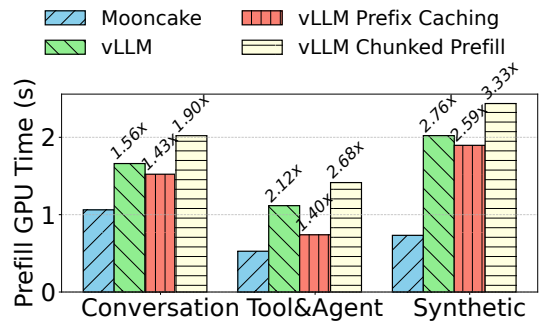


Figure 8: Average GPU time of each request during the prefill stage under different workloads.

longest prefill GPU time due to its longer input lengths and lower prefix cache ratio. The synthetic workload, featuring the highest prefix cache ratio and dispersed cache hotspots, achieves optimal cache hit rates within MOONCAKE's global cache pool. Consequently, despite having the longest average input lengths, it requires less prefill GPU time than the conversation workload. Finally, the tool&agent workload exhibits the shortest prefill GPU time because it has the shortest average input length and a relatively high prefix cache ratio.

Across different systems, MOONCAKE significantly reduces GPU time by fully utilizing global cache for prefix caching, achieving reductions of 36%, 53%, and 64% for conversation, tool&agent, and synthetic workloads, respectively, compared to vLLM. vLLM featuring prefix caching uses local cache stored on HBM, where the cache capacity is far lower than that of MOONCAKE. Its prefill GPU time is  $1.43\times$  and  $1.40\times$  higher than MOONCAKE for the conversation and tool&agent workloads, respectively. However, in the synthetic workload, where cache hotspots are more dispersed, the prefill GPU time of vLLM with prefix caching is nearly equivalent to vLLM, and is  $2.59\times$  that of MOONCAKE. vLLM with chunked prefill sacrifices some prefill efficiency to maintain lower TBT during the decoding stage, resulting in the longest prefill GPU times, which are  $1.90\times$ ,  $2.68\times$ , and  $3.33\times$  that of MOONCAKE for the three workloads.

## 5.3 MOONCAKE Store

To address Question 2, we examine the effects of MOONCAKE Store's global cache pool on system performance. Our analysis reveals that although using local DRAM to construct KVCache memory increases cache capacity than HBM only, restricting the cache to a single node still leads to suboptimal cache utilization. We will first conduct a quantitative analysis of cache capacity requirements and then showcase the benefits through practical workload experiments.

### 5.3.1 Quantitative Analysis of Cache Capacity

Considering the LLaMA3-70B model, the KVCache size required for a single token is 320 KB. Despite the possibility



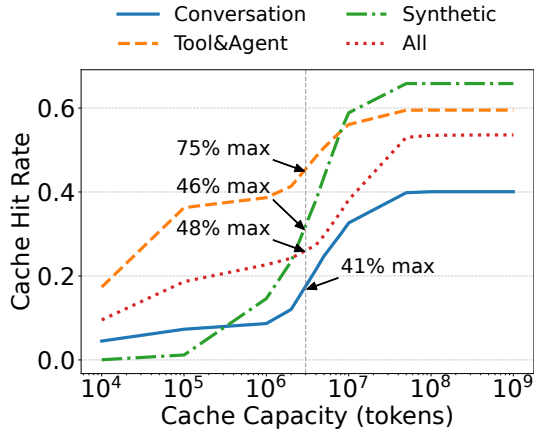


Figure 9: Quantitative analysis of prefix cache hit rates with varying cache capacities. We consider only the sequence of requests and do not account for factors such as prefill computation time or the replication of hotspots in the cache. The dashed line for 3M tokens capacity represents the local cache capacity, with the intersection points indicating the ratio of the cache hit rate to the theoretical maximum hit rate.

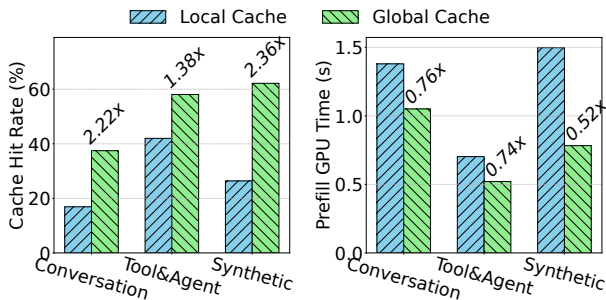


Figure 10: Cache hit rates and average GPU computation time for prefill in global and local caches.

of reserving approximately 1 TB of DRAM for local caching, this setup only supports storage for about 3 million tokens, which proves insufficient. Figure 9 displays theoretical cache hit rates under various workloads and their combinations. The findings indicate that a local cache with a 3M token capacity does not achieve 50% of the theoretical maximum hit rate in most scenarios. We also determine that, in these workloads, a cache capacity of 50M tokens nearly reaches the theoretical maximum hit rate of 100%, which require to pool at least 20 nodes’ DRAM. The results highlight that a global cache significantly enhances capacity over local caches, thus improving cache hit rates and reducing GPU times.

### 5.3.2 Practical Workload Experiment

To evaluate the effectiveness of global versus local caching mechanisms, we focus on two metrics: cache hit rate and average GPU computation time for prefill. We configure a cluster with 10 prefill nodes and restrict all request outputs

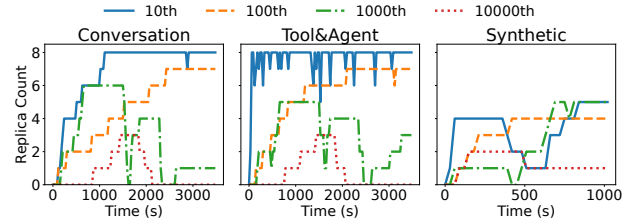


Figure 11: Replication count of cache keys across various workloads. We continuously monitor and record the keys and counts of all cache blocks every 30 seconds, subsequently ranking the cache keys by the cumulative counts from all samples. This figure depicts the temporal variation in replication numbers for cache keys ranked at the 10th, 100th, 1000th, and 10,000th positions.

to 1 to isolate the impact of the decoding stage. Each node in the local cache setup has a 3M token capacity but can only access its own cache. The global scheduler is programmed to direct requests to nodes with higher prefix match ratios to maximize cache utilization. Conversely, in the global cache setup, each node also has a 3M token capacity but can share caches across all nodes, supported by proactive inter-node cache migration. The experimental data, shown in Figure 10, indicates that the global cache achieves higher cache hit rates and shorter average prefill GPU computation times across all tested workloads. Compared to the local cache, the global cache exhibits a maximum increase of 136% in cache hit rate and a reduction of up to 48% in prefill computation time.

### 5.3.3 Cache Replica

Building upon the cache load balancing scheduling strategy discussed in §4.2, the cache keys in MOONCAKE Store may have replicas distributed across different machines, thereby reducing access latency for hot caches. To further investigate the system’s dynamic behavior, we count the number of cache replicas for keys across three workloads, as shown in Figure 11.

It can be observed that in the conversation and tool&agent workloads, there are highly concentrated hot caches (e.g., the top 100 keys), which, after the system stabilizes, have replicas on almost every instance in the prefill pool. In contrast, the synthetic workload has fewer shared prefix caches, resulting in fewer replicas and potential fluctuations, even for the top 10 blocks. This demonstrates that our scheduling strategy in §4.2 effectively provides replicas for hot caches, particularly in scenarios with highly concentrated prefix caches.

## 5.4 KVCache Transfer Performance

### 5.4.1 Transfer Engine

MOONCAKE’s transfer engine is designed to facilitate efficient cache transfers between nodes. We compare its latency

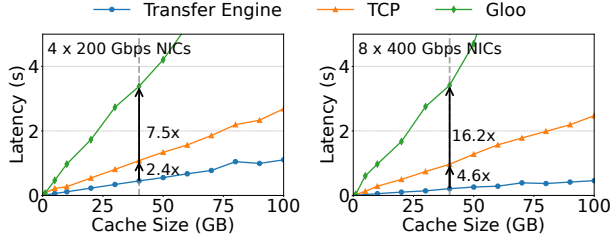
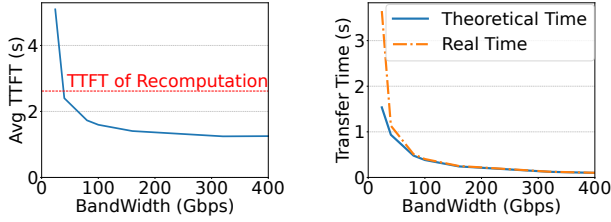


Figure 12: Latency of inter-node cache transfer.



(a) Average TTFT.

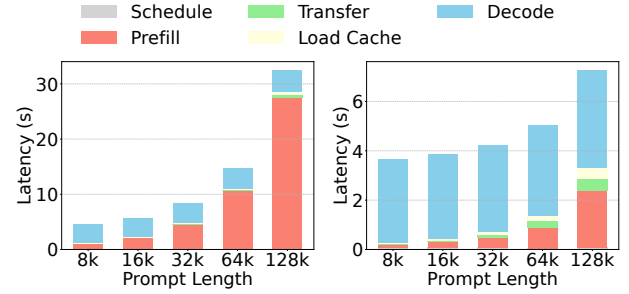
(b) Transfer time.

Figure 13: The synthetic workload experiment with varying network bandwidths.

with other popular schemes, considering two alternative baselines: `torch.distributed` with a Gloo backend and TCP-based transfers. All schemes are tested with a concurrency level of 64 and a minimum transfer granularity of 128 KB. As depicted in Figure 12, the transfer engine consistently exhibits significantly lower latency than the alternative methods. In the scenario of transferring 40 GB of data, corresponding to the cache size for LLaMA3-70B with 128k tokens, the transfer engine achieves bandwidth of 87 GB/s and 190 GB/s under network configurations of 4×200 Gbps and 8×400 Gbps, respectively. These rates are approximately 2.4× and 4.6× faster than those achieved using the TCP protocol. The code of this transfer engine will also be open sourced later as it is a decoupled and basic tool that can be used in many scenarios (e.g., it is also used in the checkpoint transfer service of Moonshot AI).

#### 5.4.2 Bandwidth Demand by MOONCAKE

MOONCAKE’s global cache pool relies on efficient inter-node cache transfer to hide cache transfer times within GPU computation times. We evaluate the impact of network bandwidth on the system’s performance by simulating a range of bandwidths from 24 Gbps to 400 Gbps and measuring the transfer time and TTFT under the synthetic workload described in §5.2.1. Figure 13a shows that the average TTFT of requests decreases as bandwidth increases. When the total communication bandwidth exceeds 100 Gbps, the average TTFT remains below 2 s, significantly less than the TTFT of the recomputation baseline. However, when bandwidth falls below 100 Gbps, system performance is significantly compromised. This is marked by a sharp increase in TTFT and evident



(a) Prefix cache ratio 0%.

(b) Prefix cache ratio 95%.

Figure 14: End-to-end latency breakdown of MOONCAKE. In the figure, *Prefill* represents the time for layer-wise prefill that integrates cache loading and storing, and *Decode* represents the time to decode 128 tokens. All processes with diagonal stripes can proceed asynchronously with model inference and do not affect MOONCAKE’s throughput.

network congestion, as demonstrated by the substantial divergence between actual and theoretical transfer times illustrated in Figure 13b. Consequently, we recommend a minimum network bandwidth of 100 Gbps to ensure optimal system performance.

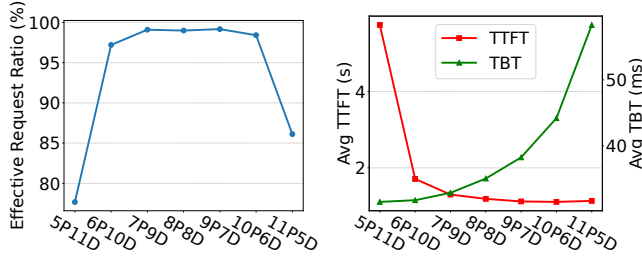
#### 5.4.3 E2E Latency Breakdown

The latency of a single inference request in MOONCAKE can be decomposed into five components: 1) scheduling and queuing time; 2) layer-wise prefill time; 3) cache transfer time; 4) time required for the decoding node to load cache from DRAM to HBM; and 5) decoding time. We experimentally analyze the proportion of these five components under settings with prefix cache ratios of 0% and 95%, as shown in Figure 14.

First, it is evident from the figure that the introduction of prefix caching significantly reduces the prefill time. Specifically, with an input length of 128k tokens, prefix caching reduces the prefill time by 92%. Furthermore, the overhead introduced by MOONCAKE has minimal impact on the system’s performance. The *Schedule*, *Transfer*, and *Load Cache* components can proceed asynchronously with model inference and therefore do not affect MOONCAKE’s throughput. Moreover, the increase in TTFT due to these overheads is smaller than the reduction achieved by prefix caching. Even when accounting for the overhead, prefix caching in MOONCAKE can reduce TTFT by 86% with an input length of 128k tokens.

#### 5.5 P/D Ratio

As a deployed P/D disaggregation system, in this section, we explore the impact of different P/D ratios on system performance. We define the P/D ratio as the number of prefill nodes to decoding nodes. Using the clusters comprising 16 nodes



(a) Effective request capacity with varying P/D ratio. (b) Latency with varying P/D ratio.

Figure 15: The impact of the P/D ratio on the system performance. P is short for prefill nodes and D is short for decoding nodes.

but with varying P/D ratios, we measure the average TTFT and TBT under the synthetic workload described in §5.2.1. We then calculate the effective request capacity as introduced in §5.2.2, setting the thresholds for TTFT and TBT to 10 seconds and 100 milliseconds, respectively. Increasing the number of prefill nodes reduces TTFT but increases TBT, and vice versa (Figure 15b). Therefore, we need to find a balance between TTFT and TBT. Figure 15a demonstrates that when the P/D ratio is approximately 1:1, MOONCAKE achieves its highest effective request capacity, indicating that the loads on the prefill and decoding clusters are relatively balanced.

We also note that some prior work [7, 9] has proposed dynamically switching the roles of nodes between prefill and decoding. However, in practical deployments, we find that the statistical characteristics of online traffic are generally stable. Therefore, we choose to fix the P/D ratio while continuously monitoring the loads of the prefill and decoding clusters, only switching node roles when significant load fluctuations occur.

## 6 Related Work

Significant efforts have been dedicated to enhancing the efficiency of LLM serving systems through scheduling, memory management, and resource disaggregation. Production-grade systems like FasterTransformer [33], TensorRT-LLM [34], and DeepSpeed Inference [35] are designed to significantly boost throughput. Orca [13] employs iteration-level scheduling to facilitate concurrent processing at various stages, while vLLM [14] leverages dynamic KVCache management to optimize memory. FlexGen [36], Sarathi-Serve [10], and Fast-Serve [37] incorporate innovative scheduling and swapping strategies to distribute workloads effectively across limited hardware, often complementing each other’s optimizations. Further optimizations [7–9] lead to the separation of prefill and decoding stages, leading to the disaggregated architecture of MOONCAKE. Our design of MOONCAKE builds on these developments, particularly drawing from the open-source community of vLLM, for which we are deeply appreciative.

Prefix caching is also widely adopted to enable the reuse of KVCache across multiple requests, reducing computational overhead in LLM inference systems [14, 34]. Prompt Cache [16] precomputes and stores frequently used text KV-Cache on inference servers, facilitating their reuse and significantly reducing inference latency. SGLang [17] leverages RadixAttention, which uses a least recently used (LRU) cache within a radix tree structure to efficiently enable automatic sharing across various reuse patterns.

Among these approaches, CachedAttention [38], a concurrent work with us, proposes a hierarchical KV caching system that utilizes cost-effective memory and storage media to accommodate KVCache for all requests. The architecture of MOONCAKE shares many design choices with CachedAttention. However, in long-context inference, the KVCache becomes extremely large, requiring high capacity and efficient data transfer along with KVCache-centric global scheduling. Additionally, MOONCAKE is not a standalone cache service, it incorporates both a memory-efficient cache storage mechanism and a cache-aware scheduling strategy, further improving prefix caching efficiency.

The benefits of a distributed KVCache pool depend on the cache hit rate, which increases as the per-token cache size decreases under a fixed capacity. Consequently, orthogonal techniques such as KVCache compression [39–41] and KVCache-friendly attention architectures [42, 43] can further enhance our approach.

## 7 Conclusion

This paper presents MOONCAKE, a KVCache-centric disaggregated architecture designed for efficiently serving LLMs, particularly in long-context scenarios. We discuss the necessity, challenges, and design choices involved in balancing the goal of maximizing overall effective throughput while meeting latency-related SLO requirements.

## Acknowledgments

We thank the anonymous reviewers and our shepherd, Mr. Kan Wu, for their valuable feedback. The authors affiliated with Tsinghua University are all in the Department of Computer Science and Technology, Beijing National Research Center for Information Science and Technology (BN-Rist), Tsinghua University, China. This work is supported by National Key Research & Development Program of China (2022YFB4502004), Natural Science Foundation of China (62141216) and Tsinghua University Initiative Scientific Research Program, Young Elite Scientists Sponsorship Program by CAST (2022QNRC001), and Beijing HaiZhi XingTu Technology Co., Ltd.



## References

- [1] OpenAI. Introducing chatgpt. <https://openai.com/blog/chatgpt>, 2022.
- [2] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [3] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [4] Charles Packer, Vivian Fang, Shishir G Patil, Kevin Lin, Sarah Wooders, and Joseph E Gonzalez. Memgpt: Towards llms as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.
- [5] Moonshot AI. Kimi. <https://kimi.moonshot.cn>, 2023.
- [6] NVIDIA. Nvidia h100 tensor core gpu architecture. <https://resources.nvidia.com/en-us-tensor-core>, 2022.
- [7] Pratyush Patel, Esha Choukse, Chaojie Zhang, Aashaka Shah, Íñigo Goiri, Saeed Maleki, and Ricardo Bianchini. Splitwise: Efficient generative llm inference using phase splitting. In *2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)*, pages 118–132. IEEE, 2024.
- [8] Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu, Yibo Zhu, Xuanzhe Liu, Xin Jin, and Hao Zhang. Dist-serve: Disaggregating prefill and decoding for goodput-optimized large language model serving. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*, pages 193–210, 2024.
- [9] Cunchen Hu, Heyang Huang, Liangliang Xu, Xusheng Chen, Jiang Xu, Shuang Chen, Hao Feng, Chenxi Wang, Sa Wang, Yungang Bao, et al. Inference without interference: Disaggregate llm inference for mixed downstream workloads. *arXiv preprint arXiv:2401.11181*, 2024.
- [10] Amey Agrawal, Nitin Kedia, Ashish Panwar, Jayashree Mohan, Nipun Kwatra, Bhargav Gulavani, Alexey Tumanov, and Ramachandran Ramjee. Taming throughput-latency tradeoff in llm inference with sarathi-serve. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*, pages 117–134, 2024.
- [11] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [12] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [13] Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. Orca: A distributed serving system for transformer-based generative models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, pages 521–538, 2022.
- [14] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with page-dattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626, 2023.
- [15] Bingyang Wu, Shengyu Liu, Yinmin Zhong, Peng Sun, Xuanzhe Liu, and Xin Jin. Loongserve: Efficiently serving long-context large language models with elastic sequence parallelism. In *Proceedings of the ACM SIGOPS 30th Symposium on Operating Systems Principles*, pages 640–654, 2024.
- [16] In Gim, Guojun Chen, Seung-seob Lee, Nikhil Sarda, Anurag Khandelwal, and Lin Zhong. Prompt cache: Modular attention reuse for low-latency inference. *Proceedings of Machine Learning and Systems*, 6:325–338, 2024.
- [17] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Jeff Huang, Chuyue Sun, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, et al. Efficiently programming large language models using sglang. *arXiv preprint arXiv:2312.07104*, 2023.
- [18] Vikranth Srivatsa, Zijian He, Reyna Abhyankar, Dongming Li, and Yiyang Zhang. Preble: Efficient distributed prompt scheduling for llm serving. 2024.
- [19] Yazhuo Zhang, Juncheng Yang, Yao Yue, Ymir Vigfusson, and KV Rashmi. Sieve is simpler than lru: an efficient turn-key eviction algorithm for web caches. In *21st USENIX Symposium on Networked Systems Design and Implementation (NSDI 24)*, pages 1229–1246, 2024.
- [20] Google. Our next-generation model: Gemini 1.5. <https://blog.google/technology/ai/google-gemini-next-generation-model-february-2024>, 2024.

- [21] Sam Ade Jacobs, Masahiro Tanaka, Chengming Zhang, Minjia Zhang, Leon Song, Samyam Rajbhandari, and Yuxiong He. DeepSpeed Ulysses: System optimizations for enabling training of extreme long sequence transformer models. *arXiv preprint arXiv:2309.14509*, 2023.
- [22] Hao Liu, Matei Zaharia, and Pieter Abbeel. Ring attention with blockwise transformers for near-infinite context. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
- [23] William Brandon, Aniruddha Nrusimha, Kevin Qian, Zachary Ankner, Tian Jin, Zhiye Song, and Jonathan Ragan-Kelley. Striped attention: Faster ring attention for causal transformers. *arXiv preprint arXiv:2311.09431*, 2023.
- [24] Dacheng Li, Rulin Shao, Anze Xie, Eric P Xing, Joseph E Gonzalez, Ion Stoica, Xuezhe Ma, and Hao Zhang. Lightseq:: Sequence level parallelism for distributed training of long context transformers. In *Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ NeurIPS 2023)*, 2023.
- [25] Vijay Anand Korthikanti, Jared Casper, Sangkug Lym, Lawrence McAfee, Michael Andersch, Mohammad Shoeybi, and Bryan Catanzaro. Reducing activation recomputation in large transformer models. *Proceedings of Machine Learning and Systems*, 5:341–353, 2023.
- [26] Shenggui Li, Fuzhao Xue, Chaitanya Baranwal, Yongbin Li, and Yang You. Sequence parallelism: Long sequence training from system perspective. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2391–2404, 2023.
- [27] Jiarui Fang and Shangchun Zhao. Usp: A unified sequence parallelism approach for long context generative ai. *arXiv preprint arXiv:2405.07719*, 2024.
- [28] Zhuohan Li, Siyuan Zhuang, Shiyuan Guo, Danyang Zhuo, Hao Zhang, Dawn Song, and Ion Stoica. Terapipe: Token-level pipeline parallelism for training large-scale language models. In *International Conference on Machine Learning*, pages 6543–6552. PMLR, 2021.
- [29] Chenxin An, Shansan Gong, Ming Zhong, Xingjian Zhao, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. L-eval: Instituting standardized evaluation for long context language models. *arXiv preprint arXiv:2307.11088*, 2023.
- [30] Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. Loogle: Can long-context language models understand long contexts? *arXiv preprint arXiv:2311.04939*, 2023.
- [31] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22, 2023.
- [32] Sharegpt teams. <https://sharegpt.com/>.
- [33] NVIDIA Corporation. Fastertransformer. <https://github.com/NVIDIA/FasterTransformer>, 2019.
- [34] NVIDIA Corporation. TensorRT-LLM. <https://github.com/NVIDIA/TensorRT-LLM>, 2023.
- [35] Reza Yazdani Aminabadi, Samyam Rajbhandari, Ammar Ahmad Awan, Cheng Li, Du Li, Elton Zheng, Olatunji Ruwase, Shaden Smith, Minjia Zhang, Jeff Rasley, et al. DeepSpeed-Inference: enabling efficient inference of transformer models at unprecedented scale. In *SC22: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15. IEEE, 2022.
- [36] Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, Beidi Chen, Percy Liang, Christopher Ré, Ion Stoica, and Ce Zhang. Flexgen: High-throughput generative inference of large language models with a single gpu. In *International Conference on Machine Learning*, pages 31094–31116. PMLR, 2023.
- [37] Bingyang Wu, Yinmin Zhong, Zili Zhang, Gang Huang, Xuanzhe Liu, and Xin Jin. Fast distributed inference serving for large language models. *arXiv preprint arXiv:2305.05920*, 2023.
- [38] Bin Gao, Zhuomin He, Puru Sharma, Qingxuan Kang, Djordje Jevdjic, Junbo Deng, Xingkun Yang, Zhou Yu, and Pengfei Zuo. Cost-efficient large language model serving for multi-turn conversations with cached attention. In *2024 USENIX Annual Technical Conference (USENIX ATC 24)*, pages 111–126, 2024.
- [39] Yuhan Liu, Hanchen Li, Yihua Cheng, Siddhant Ray, Yuyang Huang, Qizheng Zhang, Kuntai Du, Jiayi Yao, Shan Lu, Ganesh Ananthanarayanan, et al. CacheGen: Kv cache compression and streaming for fast large language model serving. In *Proceedings of the ACM SIGCOMM 2024 Conference*, pages 38–56, 2024.
- [40] Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Yakun Sophia Shao, Kurt Keutzer, and Amir Gholami. Kvquant: Towards 10 million context length llm inference with kv cache quantization. *arXiv preprint arXiv:2401.18079*, 2024.

- [41] Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint arXiv:2402.02750*, 2024.
- [42] Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, et al. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434*, 2024.
- [43] William Brandon, Mayank Mishra, Aniruddha Nrusimha, Rameswar Panda, and Jonathan Ragan Kelly. Reducing transformer key-value cache size with cross-layer attention. *arXiv preprint arXiv:2405.12981*, 2024.