

Chameleon: Adaptive Caching and Scheduling for Many-Adapter LLM Inference Environments

Nikoleta Iliakopoulou, Jovan Stojkovic, Chloe Alverti, Tianyin Xu, Hubertus Franke^{*}, Josep Torrellas
University of Illinois Urbana-Champaign^{*} IBM Research
{nmi4,jovans2,xalverti,tyxu,torrella}@illinois.edu, frankeh@us.ibm.com

Abstract

The widespread adoption of LLMs has driven an exponential rise in their deployment, imposing substantial demands on inference clusters. These clusters must handle numerous concurrent queries for different LLM downstream tasks. To handle multi-task settings with vast LLM parameter counts, methods like Low-Rank Adaptation (LoRA) enable task-specific fine-tuning while sharing most of the base LLM model across tasks. Hence, they allow concurrent task serving with minimal memory requirements. However, existing LLM serving systems face inefficiencies: they overlook workload heterogeneity, impose high link bandwidth from frequent adapter loading, and suffer from head-of-line blocking in their schedulers.

To address these challenges, we present *Chameleon*, a novel LLM serving system optimized for many-adapter environments, that relies on two core ideas: adapter caching and adapter-aware scheduling. First, Chameleon caches popular adapters in GPU memory, minimizing the adapter loading times. Importantly, it uses the otherwise idle GPU memory, avoiding extra memory costs. Second, Chameleon uses a non-preemptive multi-queue scheduling to efficiently account for workload heterogeneity. In this way, Chameleon simultaneously prevents head of line blocking and starvation. We implement Chameleon on top of a state-of-the-art LLM serving platform and evaluate it with real-world production traces and open-source LLMs. Under high loads, Chameleon reduces P99 and P50 TTFT latency by 80.7% and 48.1%, respectively, while improving throughput by 1.5× compared to state-of-the-art baselines.

1 Introduction

Generative Large Language Models (LLMs) have seen an exponential growth in recent years, triggering a long line of research [2, 7, 20, 31, 39, 50, 59, 63, 64]. They have become integral to numerous technological advancements and applications of different domains, e.g. healthcare [40], coding [12], data analytics [55], or education [3]. As their popularity increases, the number of online queries received by datacenter inference clusters is getting substantially larger [21]. These queries typically target a variety of downstream tasks, e.g. chat-bot conversation, coding or text summarization. These different tasks require different or special-purpose fine-tuned LLMs to achieve their highest accuracy.

This imposes a huge hardware [39] and energy [49] tax to datacenters, as each of these models typically requires large memory to store the multi-billion parameters. Adapter-based techniques, such as Low-Rank Adaptation (*LoRA*) [16, 56], have been explored to alleviate this problem. These methods fine-tune only a small (low-rank) subset of a base model’s parameters for every task, and have been originally proposed to speed up LLM training. Recent studies on inference serving environments [4, 47] capitalize on this. They decouple each model’s base and fine-tuned adapter parameters,

allowing different colocated LLMs to share a large fraction of their memory, i.e., their base model. This enables serving of potentially hundreds of LoRA fine-tuned LLMs at much lower memory cost.

Despite the memory footprint reduction in multi-task settings, this method imposes two previously unseen challenges. First, inference clusters have to additionally orchestrate the adapters required by the incoming requests as they are being scheduled. Punica [4] and S-LoRA [47] keep the base model stored in GPU memory and the adapters in the host memory. They consider storing adapters in the GPU memory too expensive as it can interfere with the memory allocations of incoming requests, i.e. for the key-value caches of their intermediate state. Instead, they fetch on-demand the adapters required by the running requests and discard them as soon as the requests terminate. S-Lora [47] and dLora [58] further fetch in advance the adapters for the requests waiting in the system’s queue to hide some of the loading overheads. Our study reveals that even *the asynchronous adapter fetching increases the time-to-first-token (TTFT) latency, especially when the system is heavily loaded, as it increases the contention in the CPU-GPU PCIe link bandwidth.*

Second, apart from the loading costs, adapters introduce inference overheads. Decoupled computations over the base model and adapter increase the execution time of a single query [58]. Moreover, prior art [4, 25, 47, 58] executes batches of heterogeneous requests, i.e. for different tasks and LoRA adapters, increasing resource utilization and throughput but penalizing tail latency. Requests for adapters of higher rank take longer to execute and, thus, they stall the execution of smaller requests within the same batch [25]. To reduce this bottleneck, one can cluster the requests for the same adapter [58] or for adapters of the same rank [25] within the same node, i.e., server in a cluster that stores a given model and adapter replicas, potentially causing load imbalance and frequent request migration [58]. While focusing on multi-node distributed scheduling, most systems fail to account for the inevitable workload heterogeneity within a replica, i.e., at *server-level*.

Our study corroborates previous observations that LLM inference requests are prone to head-of-line blocking [44, 57]. We analyze real-world production workloads [39] and observe that these requests follow a heavy-tailed distribution: most are completed in a very short time, while a small fraction experiences significantly longer execution durations. While prior work has largely attributed this heterogeneity to differences in input [57] and output [44] request sizes, our study is the first to shed light on how the variability in adapters rank [47] and popularity [9, 50, 58] affect the requests at the tail, underlying the necessity to take the adapter size into account. Moreover, we find that simply prioritizing short requests is insufficient to address the issue. For instance, the speculative Shortest-Job-First (SJF) scheduler employed by μ Serve [44],

along with its aging mechanism to mitigate starvation, inadvertently increases the tail latency of longer requests missing their Service Level Objectives (SLOs). Overall, our findings emphasize the need for a more nuanced scheduling strategy—one that addresses adapter-level heterogeneity, offers expedited processing for short requests, and ensures that longer requests still meet their SLOs.

We use these insights to design *Chameleon*, an LLM serving system optimized for datacenter environments with many colocated LLM tasks. Tasks share their base LLM, which occupies a large fraction of GPU memory, while each task uses its own specific adapter. Chameleon attains high efficiency through two main principles.

First, Chameleon designs a *transparent, adaptive and interference-free cache for adapters*. Our study reveals that contrary to common wisdom [4, 47], there is enough idle GPU memory, even during high load, that can be repurposed to cache the adapters that are likely to be reused and expensive to reload. However, as this idle memory fluctuates over time, the size of the Chameleon cache has to dynamically adapt to the incoming load to avoid interfering with the key-value cache and other memory allocations of the incoming requests. To that end, Chameleon integrates its cache manager with the system scheduler to judiciously up- and downsize its resources at runtime while monitoring the request traffic. Furthermore, our study shows that simple eviction policies, e.g., LRU, are not optimal for this set-up as different ranks (sizes) of the adapters affect cache miss penalties. Instead, Chameleon implements a cost-aware eviction policy that combines the recency and frequency in adapters use with their reloading cost to decide which adapters to discard when the cache is full or has to shrink.

Second, Chameleon incorporates a *non-preemptive adapter-aware multi-level queue (MLQ) scheduler* to minimize the head-of-line blocking and guarantee SLOs for all types of requests. Prior art considers feedback queues (MLFQ) to schedule LLM inference requests, and relies on preemption to deal with their non-deterministic execution times [57]. This requires the orchestration of the intermediate states of preempted requests, introducing non-negligible complexities. Chameleon instead speculates a *weighted request size (WRS)* and uses it to assign priorities, i.e. admit requests to specific queues. WRS takes into account the number of input tokens, an estimated output length [44] and the adapter rank for every request. Unlike preemptive solutions, Chameleon admits requests from all queues to every batch but partitions resources in proportion to the queues priority. Specifically, Chameleon assigns a different resource quota to each queue [33], i.e. a maximum number of tokens it can admit to a batch, and gradually decreases it for the queues used by larger requests. This *enables a faster lane for smaller requests* but also *eliminates starvation*. To guarantee SLOs for all requests in every queue and maximize throughput, Chameleon dynamically adjusts the number of queues and their resource quotas at runtime.

We implement Chameleon on top of the open-source S-LoRA [47] LLM serving platform. Chameleon does not require any hardware or operating system (OS) support, or changes to CUDA kernels. We evaluate Chameleon with open-source LLMs using real-world production traces [39] and show that Chameleon is very effective. Compared to state-of-the-art baselines [44, 47], Chameleon reduces the P99 time-to-first-token (TTFT) latency by 80.7%, reduces the P50 TTFT latency by 48.1%, and improves the throughput by 1.5 \times .

This paper makes the following contributions:

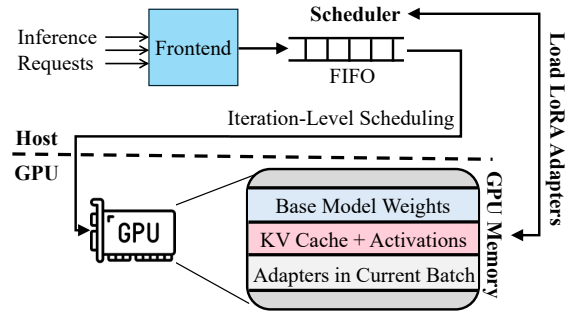


Figure 1: Conventional LoRA online serving system

- A characterization of state-of-the-art LLM serving systems in many LoRA adapter environments.
- The Chameleon LLM serving platform that introduces the first cache design for LoRA adapters and a novel adapter-aware multi-queue scheduler that eliminates head-of-line blocking.
- An implementation and evaluation of Chameleon.

2 Background on LLMs

LLM inference. Generative LLMs [28, 34, 45, 54] execute inference in two stages. First, they process the entire input at once (*prefill phase*) and then they generate output tokens auto-regressively, i.e. one by one (*decode phase*). During prefill, all input tokens are processed in parallel making this phase compute-bound and its performance depends on the input size which is *known in advance*. Conversely, during decode, the output tokens are generated sequentially in iterations. Each iteration generates a token based on the input prompt and all previously generated tokens, typically cached on the GPU memory in *key-value (KV) caches*. Hence, decode phase is memory-intensive and its performance depends on the output size, i.e. the number of decode iterations, which is *determined on the fly and thus is unknown at the time a request is admitted to execute*.

LLM serving systems. LLM serving systems typically *batch* requests for the same model to maximize hardware utilization. However, different requests generate different number of output tokens and thus requests on the same batch can have varying execution times. Systems need to dynamically update batches to avoid long requests blocking smaller ones. Specifically, state-of-the-art systems perform continuous batching [1, 59], i.e. they remove completed requests from a batch and potentially add new ready-to-run requests on every decode iteration (*iteration-level scheduling*).

Performance metrics for LLMs. The key performance metrics for LLM inference systems are Time to First Token (TTFT), Time Between Tokens (TBT), and Throughput [52]. TTFT measures the latency from the time a request is received to the generation of the first output token, capturing both the latency of the model’s prefill phase and any queuing delays introduced by the system. TBT tracks the time it takes to generate each subsequent output token during the decode phase, reflecting the latency of a decode iteration along with any queuing delays introduced by iteration-level scheduling. Throughput refers to the load that the system can handle without violating the service-level objectives (SLOs) for TTFT or TBT. In

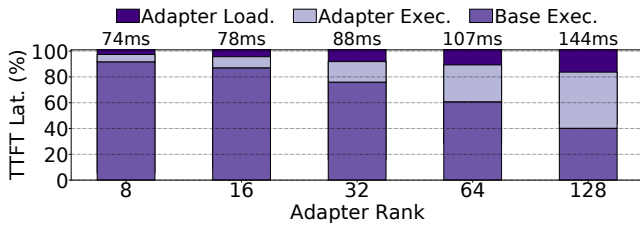


Figure 2: TTFT latency with different adapter ranks broken down into base model execution, adapter execution and adapter loading time.

this work, we set SLOs to $5\times$ the TTFT and TBT latencies achieved on a low-load system (similar to prior art [26, 50]).

2.1 LLM Low-Rank Adaptation (LoRA)

Pre-training LLMs and fine-tuning them for domain-specific tasks is a common technique to reduce training overheads. Low-Rank Adaptation (LoRA) [16, 41] is such a method that injects low-rank matrices into a base model’s layers and updates them to fine tune it. The most important performance factor is the size of a LoRA adapter, i.e. the *rank* of its matrices. Higher ranks potentially translate to better tuning and thus higher accuracy. As different tasks have different accuracy requirements, they are likely to employ adapters of different ranks over the same base model [47, 56].

LLM adaptation can be used to also reduce the memory requirements of online systems that serve such diverse tasks [64]. For example, while the straightforward way to apply adapters is to merge them with the base model and create a full size standalone specialized LLM instance [17], recent works [4, 25, 47, 58] allow sharing the weights of the base model and allocate only the specific adapters per-task. Typically, an adapter size is significantly smaller than the base model’s size thus this method reduces significantly the memory usage in multi-task environments. Moreover, such systems also enable batching of requests of different tasks, i.e. base and adapter combinations, further improving the throughput.

Figure 1 shows the organization of such a system that employs LoRA adapters [4, 25, 47, 58]. On system initialization, the base LLM model is transferred to the GPU memory from the host. A scheduler on the host manages the incoming requests, updating the batch to be executed on every iteration (iteration-level scheduling). Before it sends the batch to the inference engine on the GPU, the scheduler also loads any missing adapters required by the requests in the batch. Once there are no running requests that use a given adapter, the adapter is discarded from the GPU memory to make space for new incoming requests [4, 47]. In this work we study LLM serving in such many LoRA adapter environments.

3 Opportunities for Efficient LLM Serving in Many-Adapter Environments

In this section we examine the new challenges of many-adapter environments and why they are not efficiently handled by the conventional online LLM serving systems. We characterize the open-source Llama-7B model [54] on an NVIDIA A40 server with 48GB of GPU memory [36]. We use the S-LoRA serving platform [47], a state-of-the-art inference system for multi LoRA adapter scenarios.

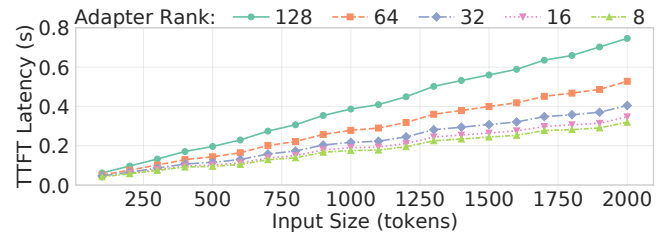


Figure 3: TTFT latency for different adapter ranks while varying the request’s input size.

3.1 Adapters Increase Workload Heterogeneity

As discussed in the previous section, the adapters employed by different tasks are expected to vary in size (*rank*), as tasks require different levels of accuracy [16, 47, 56, 58]. Figure 2 studies how this *rank heterogeneity* affects the TTFT of a single inference request, with medium input and output size [50]. We run the request over a base Llama-7B model combined with a specific adapter on an unloaded system and increase the adapter rank from 8 to 128 [47, 58]. We break down the total execution time to the time spent in the base model computation, the adapter computation, and the time spent on loading the adapter’s weights from host to the GPU memory. We observe that as the rank size increases, the adapter overheads also increase. For example, for rank 128, $\sim 60\%$ of the total TTFT latency is spent on adapter loading and computation.

We further examine the effect of the adapter rank while considering other sources of inference heterogeneity. For example, prior work observed that large inputs lead to longer prefill phases and large outputs to much longer decode phases [50]. Similarly, large batches of requests increases throughput but at the cost of longer decode iterations. Figure 3 shows the TTFT latency for different adapter ranks as we vary the input size of a request, i.e. the number of input tokens, while keeping the output size fixed. For this experiment, we keep the adapter weights in GPU memory and isolate prefill performance excluding the adapter loading. For all input sizes, TTFT latency varies significantly when we consider different adapter ranks and the rank impact is pronounced as the input size increases. Similarly, it can be shown that for large batch sizes different adapter ranks lead to diverse decode latencies for requests with similar input and output sizes. *Overall, we find the adapter rank to be an extra, equally important source of heterogeneity, next to input, output, and batch size.*

Apart from the different ranks, adapters are expected to have skewed popularity as well, following the skewed popularity of different tasks. LLM inference is a user-facing service where some tasks receive a substantially larger amount of requests than others and these requests typically arrive in bursts [9, 50, 58]. Next, we will show how this heterogeneity affects various system design decisions, such as which adapter to keep in GPU memory or how to schedule inference requests for different adapters.

Insight #1: Adapters are an additional source of heterogeneity in LLM inference that must be managed dynamically.

3.2 Adapters are Expensive to Load

When an LLM inference request for a specific adapter arrives at an online serving system, the adapter’s weights must be loaded into

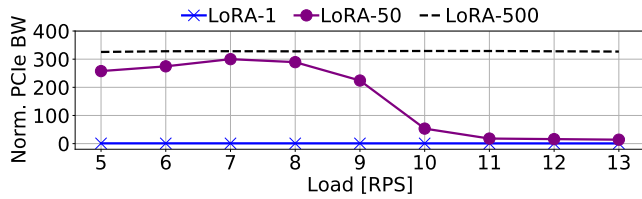


Figure 4: PCIe bandwidth usage under different loads for environments with: 1 adapter (*LoRA-1*), 50 different adapters (*LoRA-50*), and 500 different adapters (*LoRA-500*).

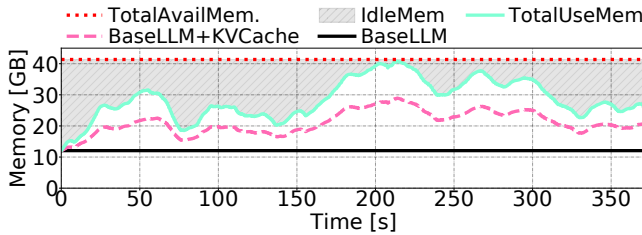


Figure 5: Memory usage over time for different parts of the computation: base LLM model, KV cache, and adapters.

GPU memory for the request to be processed. Thus, loading the adapter weights lies on the critical path of inference execution and Figure 2 shows it can cost up to 17.5% of the total TTFT latency when larger adapters are used on an unloaded system.

These overheads are pronounced when the system is loaded, i.e. the TTFT latency of LLM inference requests deteriorates as the number of adapters used by the requests increases. The main reason is the contention on the PCIe link between the host and the GPU. Figure 4 shows the PCIe bandwidth consumption at different loads when we increase the number of rank-32 adapters used by the requests. Specifically, we evaluate systems with a single adapter used by all requests (*LoRA-1*), and with 50 or 500 different adapters uniformly distributed across requests (*LoRA-50* and *LoRA-500*). We normalize the bandwidth to that of *LoRA-1* environment at 5 RPS. We observe that as the number of adapters increases, the PCIe bandwidth usage grows significantly. Interestingly, at higher loads, *LoRA-50* becomes compute-bound and due to the queuing effects, the PCIe bandwidth drops. In contrary, *LoRA-500* is constantly bound by the PCIe bandwidth: *LoRA-500* has 328 \times higher PCIe bandwidth consumption than *LoRA-1*.

An intuitive way to reduce these overheads would be to leverage idle GPU memory to cache adapters. However, LLM inference has substantial load fluctuations [50]. Figure 5 shows the GPU memory usage, i.e. base model, KV cache and requests input and output, over time when we run the Llama-7B model using production traces of requests from Azure [39]. While most of the time there is abundant idle memory, we observe that it drastically drops during load spikes. Thus, while typically there is enough idle memory to cache adapter weights, the system needs to carefully and dynamically resize the cache resources based on the incoming load employing a *cost-aware eviction policy* taking into account per-adapter reloading cost.

Our findings challenge the common design decision to discard the adapters from GPU memory if none of the currently running requests use them [47, 58]. We find that keeping them in GPU memory can significantly improve performance, especially in high

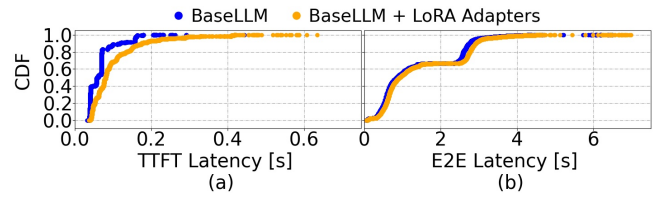


Figure 6: CDF of (a) TTFT and (b) E2E latency of requests for a real LLM trace [39]. Requests are executed one by one.

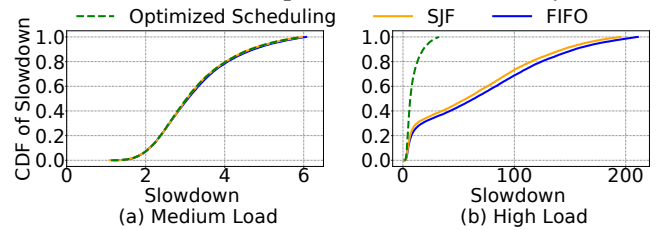


Figure 7: CDF of the per-request slowdown for different scheduling policies under (a) medium and (b) high load.

load scenarios, and, that there is a substantial fraction of otherwise idle GPU memory that can be used for this.

Insight #2: Frequent loading of adapters from host to GPU memory creates bandwidth contention, degrading the system performance. Idle GPU memory can be repurposed to cache adapters and mitigate some of these overheads. However, dynamic resizing of the cache is essential as idle memory fluctuates heavily based on request traffic.

3.3 Adapters Affect Requests at the Tail

In Section 3.1, we observed that there is a high degree of heterogeneity in the performance of LLM inference requests, based on their input, output, and adapter size. Now, we analyze how this heterogeneity impacts the effectiveness of scheduling decisions.

First, we take the open-source production traces of LLM inference requests for a conversation service [39] and execute one request at a time. We run only with a BaseLLM, and augmented with LoRA adapters. We consider a pool of 100 different adapters with rank sizes uniformly distributed among 8, 16, 32, 64, and 128 [47]. We associate every request in the trace with one of these adapters, following the power-law distribution for the rank popularity [47]. Adapters with smaller ranks are the ones more frequently requested and we distribute uniformly the requests among the adapters with the same rank. Figure 6 shows the CDF of (a) TTFT and (b) end-to-end latency of all requests. For this experiment, the latency includes both the prefill phase and the time it takes to load the adapter. This figure shows that, in a real production environment, requests remain highly heterogeneous and their execution time follows the heavy-tail pattern: the majority of requests have very short execution times, but there are a few very long requests. Importantly, adding LoRA adapters significantly affects requests at the tail.

Heterogeneity in execution times typically requires handling at the scheduling level. LLM engines schedule the requests at iteration-level [59], thus, at each iteration the scheduler decides which requests will execute in a batch. Unfortunately, the majority of conventional systems use a FIFO approach due to its simplicity [20, 47]. However, FIFO is inefficient for heterogeneous requests as it introduces head-of-line (HoL) blocking, leading to increased tail latency.

For that reason, researchers have proposed to schedule the requests in a Shortest-Job-First (SJF) manner. Specifically, the existing systems [44] predict the request’s output length and prioritize the requests with the shortest predicted outputs. However, continuously prioritizing short requests leads to the starvation of long requests, again, negatively impacting the overall tail latency. Moreover, using the output length as the only scheduling knob is insufficient, as inputs and adapters also impact the total latency (Figure 3).

To show the inefficiencies of these two popular scheduling policies, we execute the production trace [39] using the Llama-7B model. We record the slowdown of each request: how longer was the request’s response time compared to an isolated environment where the request executes alone. Figure 7 shows the CDF of slowdown per request with FIFO and SJF scheduling policies, along with the optimized scheduling policy that we will introduce in Section 4.

While scheduling policy does not play a significant role in medium loads, in high loads, conventional policies create huge slowdowns for the requests at the tail. For FIFO, these are the short requests blocked by the long requests, while for SJF, these are the long requests starved due to the prioritization of short requests. An optimized scheduling policy can efficiently handle this inter-request heterogeneity and significantly improve the performance.

Insight #3: Conventional scheduling policies, such as FIFO and SJF, are ineffective for highly heterogeneous LLM inference requests. Instead, there is a need for a scheduling policy that can efficiently manage request heterogeneity while, at the same time, taking into account all knobs that affect the execution time.

4 Chameleon Design

Based on the insights from our characterization (Section 3), we design *Chameleon*, an LLM inference serving system optimized for many-adapter environments. Chameleon is designed to address the unique challenges posed by the heterogeneity found in many-adapter serving systems, i.e. (1) *the overheads and side-effects of adapters’ weights loading* and (2) *the increased tail latency in inference requests* due to inefficient request scheduling.

To address the first issue, Chameleon leverages underutilized GPU memory to implement a software-managed adapter cache (*Chameleon Adapter Cache*). The cache stores the adapter’s weights in GPU memory, removing the adapter loading off the inference request’s critical path, and reducing the PCIe bandwidth usage. To improve the requests’ tail latency, Chameleon uses a multi-queue non-preemptive scheduler that provides an express lane for short requests, while ensuring longer requests do not starve (*Chameleon Scheduler*). The scheduler and the cache work synergistically, maximizing the system’s throughput.

Architecture Overview. Figure 8 overviews the Chameleon’s architecture. To handle workload heterogeneity, Chameleon classifies all incoming requests based on their total size. Every request first goes through an output length predictor that estimates the number of output tokens (1). We use an existing, open-source, predictor based on BERT proxy model [44]. The Chameleon scheduler combines this estimated output size with the known number of input tokens and the rank of the adapter required by the request to calculate a *weighted request size (WRS)* (2).

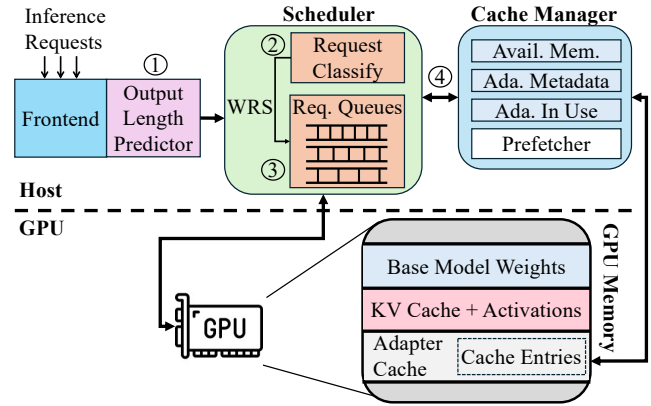


Figure 8: Chameleon architecture overview.

Chameleon uses WRS to categorize requests into classes, e.g. *small*, *medium* and *large*, and admits them to different request queues (3). Chameleon uses *iteration-level scheduling* (Section 2), thus on every decode iteration it removes and adds requests to a batch. Importantly, at every iteration, Chameleon uses requests from all queues, while respecting their assigned quotas. Specifically, each queue is assigned some amount of GPU resources that the requests from that queue can consume at each iteration.

The scheduler creates a fast lane for short requests preventing their head-of-line blocking, but still admits requests from all queues guaranteeing that no request will starve. To handle load fluctuations and changes in request properties, Chameleon dynamically adjusts both the total number of queues and the per-queue cutoffs based on the monitored WRS distribution of the incoming load.

On every scheduling decision, Chameleon also invokes its Cache Manager. Cache Manager is a software controller that manages Chameleon Adapter Cache, and it is assigned to i) (pre)fetch any necessary adapters, required by the requests to be scheduled and to ii) evict any idling adapters when the available GPU memory is not enough to store the incoming requests’ input, output and KV cache entries (4). The Cache Manager tracks all cached adapters and the necessary metadata to enforce a *cost-aware* eviction policy.

4.1 Chameleon Adapter Cache

Chameleon Adapter Cache (Chameleon Cache, for short) is a software structure that stores unused adapters, loaded by previous requests, in idling GPU memory. The goal is to eliminate the costs of fetching the adapters again in the future, on the critical path of an inference request. Chameleon maintains one cache instance per LLM instance, i.e., each LLM replica has its own local adapter cache. Each cache entry contains the adapter’s weights and some metadata used for cache management. The metadata contains:

- **Adapter ID:** A unique identifier for the adapter.
- **Adapter Rank:** The size of the adapter, which affects the amount of GPU memory it occupies.
- **Last Used Timestamp:** The last time the adapter was accessed.
- **Usage Frequency:** The total number of times the adapter has been used within a specific time frame.
- **Reference Counter (RC):** The number of active requests using this adapter. If RC is zero, the adapter is eligible for eviction.

The cache is managed by a *Cache Manager*. The manager performs the following operations: i) it *retrieves a cached adapter* for an incoming request, ii) it *loads a missing adapter* from host memory, and iii) it *employs a cost-aware eviction policy* to discard cached adapters when necessary. Recall that adapter’s weights are read only. Thus, there is no need to maintain their coherence, or to write them back to the host memory on eviction from the cache.

Dynamic Cache Sizing. Unlike hardware caches, Chameleon cache’s overall capacity dynamically changes over time. The Cache Manager adjusts its size in real-time to always meet the resource demands of incoming requests while consuming the fluctuating idle memory for adapter caching (Figure 5). For example, when the memory requirements of incoming requests, i.e. to store input activations, KV entries or missing adapters, cannot be served by the available free GPU memory, the Cache Manager downsizes accordingly the cache capacity, evicting adapters to free up the necessary space. Similarly, when a request terminates and there is enough idling GPU memory, the Cache Manager expands the cache capacity to store the departing request’s adapter.

The Chameleon Cache Manager and Scheduler work synergistically for this resizing. The Scheduler monitors the incoming requests’ memory requirements and invokes the Cache Manager. Specifically, it scans the request queues to assemble a batch of requests to be scheduled on every decode iteration (Section 4.2) and communicates the exact amount of memory required by the batch to the Cache Manager. If necessary, the manager then discards unused adapters to free up space, based on a cost-aware eviction policy, and loads any missing adapters from the CPU memory.

Cost-Aware Eviction Policy. Selecting an appropriate eviction policy is important for the cache performance, and such a policy has to consider the special characteristics of the multi-adapter environment. For example, while policies based on recency can capture the temporal locality in adapter requests, they may fail to cover the adapters skewed popularity found in LLM serving workloads, i.e. the fact that certain adapters are being accessed more frequently than others (used by a larger number of concurrent inference requests) [47]. Evicting frequently used adapters can increase miss rates and CPU-GPU link bandwidth consumption due to frequent adapter reloading. Indeed, in Section 5 we show that least recently used (LRU) policy alone fails to achieve optimal performance.

An additional reason why LRU is insufficient on its own, is that cache misses in this environment have varying costs. Unlike hardware caches, that store objects with uniform sizes (cache lines), Chameleon caches objects, i.e. adapters, with different sizes, i.e. ranks. Consequently, the latency to load an adapter on a cache miss is not fixed, i.e. larger adapters take longer to transfer from host to GPU memory. Thus, the eviction policy must be also *cost-aware* [11], i.e. prioritize the eviction of smaller adapters.

Similar to caching schemes employed by different domains [11], we show that relying solely on a single feature is insufficient to capture the complex trade-offs in our system. To address this limitation, we propose a *compound eviction algorithm* that simultaneously considers multiple factors influencing adapter importance. This scheme calculates a *score* for each adapter based on its frequency of use, recency of access, and size. The score is computed using the formula: $Score = F \times Frequency + R \times Recency + S \times Size$, where F, R, and S

are weighting coefficients. The adapter with the lowest score is considered the least critical and is evicted first. This approach balances the trade-offs between evicting adapters that are infrequently used, not recently accessed, or large in size.

The weighting coefficients F, R, and S enable adjusting the sensitivity of the eviction policy to each factor. For example, if frequency is deemed more indicative of future usage patterns than recency for the running workload, F can be assigned a higher value relative to R and S. For this study, we use static coefficients which we fine-tune by profiling offline industrial traces of inference requests [39] combined with adapter size distributions found in literature [47]: F, R, and S are set to 0.45, 0.10, and 0.45, respectively. We consider an on-line dynamic adjustment of the coefficients based on the incoming load for future work.

Chameleon never evicts adapters that are actively used by running requests. To guarantee this, the Cache Manager maintains a reference counter per adapter and considers eligible for eviction only the adapters whose counters have dropped to zero (not in use). Additionally, the manager consults the scheduler to identify adapters associated with queued requests—those not currently running but guaranteed to execute in the near future. The manager attempts to retain these adapters in the cache, provided there is sufficient memory available. However, if memory constraints arise, adapters for queued requests may be evicted to accommodate the running batch. Overall, the manager applies the eviction policy only to adapters that are not used by either currently running or queued requests. Adapters associated with queued requests are considered for eviction only when memory constraints make it necessary.

Prefetching. Chameleon builds on top of prior art optimizations, i.e. it monitors the request queues and prefetches, whenever possible, the missing adapters required by the waiting requests before they are admitted to a batch for execution [47, 58]. This approach can lead to late prefetching, i.e., the request becomes ready for execution before the prefetching is completed. Hence, we explore techniques that predict future load, such as a histogram-based approach [46], to prefetch adapters even for requests that are not currently queued. Since the effectiveness of this approach heavily depends on prediction accuracy, we do not enable this optimization by default in our evaluation and we include a separate experiment to assess the potential impact of such prefetching (Figure 15).

Prefetching has the potential to hide the loading costs of an adapter’s weights and remove it from the inference request’s critical path of execution. However, prefetching and unloading the adapters on completion still consumes high CPU-GPU link bandwidth (PCIe). Thus, caching the adapters is still necessary for high performance.

4.2 Chameleon Scheduler

Inspired by prior art on load balancing in multi-server environments with heterogeneous task size distribution [6, 14, 33], we design Chameleon Scheduler. The scheduler stores the inference requests across multiple queue lanes, each dedicated to handling requests within a specific size range. Its goal is twofold: to provide a fast lane for smaller requests, preventing their head-of-line blocking, and to ensure that requests from all lanes are scheduled in parallel, avoiding starvation for larger requests.

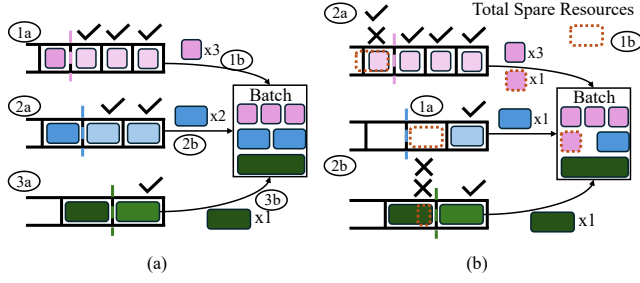


Figure 9: An example of Chameleon Scheduler operation when (a) no Spare Resources are collected (b) Spare Resources are collected and redistributed.

Each queue is assigned a resource quota, which governs the resources available for executing the requests from that queue. This quota is represented as tokens and includes input tokens, output tokens, and the memory required for the corresponding adapter, also expressed in token units. These tokens determine the resources a queue can reserve for request execution. When a request from a specific queue is admitted to the batch, the queue’s available quota is decreased by the request’s consumption. When the request finishes, it returns the borrowed quota back to the queue. Even though, in every iteration all queues have the chance to admit requests to the batch, the queues with smaller requests are accessed first, hence we will refer to them as “higher” request queues.

Admission to the Queues Requests entering the system are defined by three parameters: known input size, predicted output size, and rank of the used adapter. We calculate the weighted request size (WRS) based on these parameters using the formula:

$$\text{WRS} = A \cdot \frac{\text{InputSize}}{\text{MaxInputSize}} + B \cdot \frac{\text{OutputSize}}{\text{MaxOutputSize}} + C \cdot \frac{\text{AdapterSize}}{\text{MaxAdapterSize}}$$

Here, A, B, and C are weighting coefficients, chosen based on our sensitivity studies and profiling in Section 3. For example, the output size primarily determines the execution time, but is predictive, hence, assigning it a too high value would make the system fragile to prediction accuracy. We set these parameters to $A = 0.3$, $B=0.5$, and $C=0.2$. Using the calculated WRS and specified per-queue cut-offs, i.e., the boundaries that define the request size ranges for each queue, the scheduler places a request in the correct queue. Later, we detail how to determine per-queue cut-offs using request clustering. Note that Chameleon Scheduler uses an open-source BERT-based proxy model to predict the request’s output length [44].

Admission to the Batch. The idea behind the Chameleon scheduler is depicted in Algorithm 1. It operates in two phases: *Initial Request Admission* and *Redistribution of Spare Resources*. In the first phase, each queue attempts to admit requests up to the queue’s maximum available resources. If certain queues have few or no requests to admit, any unused resources are collected. At the end of this phase, the *total spare resources* are consolidated. In the second phase, the scheduler redistributes the spare resources to queues with pending requests, aiming to maximize resource utilization. Starting with the highest-priority queue and moving downward, the scheduler allocates as much of the reclaimed resources as possible to admit waiting requests. If requests from a given queue

still cannot be admitted due to insufficient tokens available, no additional resources are allocated to that queue.

The phases of this process are illustrated in Figure 9, which depicts three request queues, for “short”, “medium” and “large” requests. Figure 9(a) shows the *Initial Request Admission* phase. The scheduler starts with the highest-priority queue, admitting three requests, that fit within the queue’s resource quota (1a). However, the fourth request is not admitted due to insufficient resources. These admitted requests are then placed into the batch (1b). The same procedure is subsequently applied to the second-highest-priority queue (2), and the lowest-priority queue (3). At the end of this phase, there are no remaining resources to redistribute, so the process concludes without entering the second phase.

Figure 9(b) shows a case where spare resources are available for redistribution. Spare resources are left in the “medium” request queue during the Initial Request Admission phase (1a). These resources form the *Total Spare Resources* (1b). The process proceeds to the *Redistribution of Spare Resources* phase. Each queue reattempts to admit any remaining requests into the batch (2). The highest-priority queue evaluates if its pending request can be admitted (2a). Since the *Total Spare Resources* are sufficient, the scheduler allows these requests to be admitted. The second queue has no pending requests and is skipped. The queue with “large” requests attempts to admit its pending request (2b). As the available resources are insufficient, the request remains in the queue. At the conclusion of this phase, the batch is finalized and ready for execution.

Bypassing Adapter Blocking. If a request at the head of a given queue cannot be admitted to the batch because there is not enough memory for its adapter, all requests within the queue are blocked. However, it might happen that some of the later requests are either for an adapter that is already loaded in the Chameleon Adapter Cache or for an adapter that is small enough to fit the remaining space in the cache. To address this challenge, Chameleon implements a bypass mechanism. Specifically, it allows younger requests to bypass the requests at the head of the queue if their memory requirement can be satisfied with the current resources. This mechanism improves system throughput by allowing more requests to be processed without waiting for cache space to become available. However, constant bypassing might lead to starvation and unfairness. Instead, Chameleon allows a request to bypass the current head of the queue, only if the admitted request will not exceed the waiting time for the request at the head of the queue. We rely on the prediction of request length when deciding whether to allow adapter bypassing. However, if the prediction was wrong and the request does not finish within the predicted number of iterations, the request is squashed, and later re-executed. In our experiments, we see at most 5% of requests getting squashed.

Determining the Number of Queues. The efficiency of Chameleon Scheduler depends on the number of used queues. Too few queues may cause head-of-line blocking when there is a high variability in request sizes, while too many queues can result in load imbalance and underutilized queue resources when requests are homogeneous due to the resource fragmentation. To decide the optimal number of queues, Chameleon Scheduler uses a method based on *K-Means*

Algorithm 1: Generate a new batch of requests.

```

def generate_batch:
  Inputs: Queues = req. queues; PQ_Tokens = per queue tokens
  Result: Batch of requests to be sent to the GPU.
  batch ← [];
  leftover ← 0;
  for each q in Queues do // Phase 1
    consumed ← put_batch(q, PQ_Tokens[q], batch);
    if q is empty then
      leftover ← leftover + (PQ_Tokens[q] - consumed);
  for each q in Queues do // Phase 2
    if leftover == 0 then
      break;
    consumed ← put_batch(q, leftover, batch);
    leftover ← leftover - consumed;
  return batch;

def put_batch:
  Inputs: Queue; Tokens; Batch
  Result: Tokens consumed by added requests from the queue.
  consumed ← 0;
  for each req in Queue do
    needed ← need_resources(req);
    if resources < needed then
      break;
    resources ← resources - needed;
    consumed ← consumed + needed;
    batch.append(req);
  queue ← [req for each req in queue if req not in batch]
  return consumed;

```

clustering. Given the distribution of request sizes, the scheduler computes K-Means clustering for values of K ranging from 1 to K_{max} . The intuition is that requests similar in size will be grouped within the same cluster and requests from different clusters are different enough to require separate resources. The scheduler calculates the Within-Cluster Sum of Squares (WCSS) for each configuration and picks K that yields minimal WCSS as the optimal number of queues. For practicality, we set the maximum number of queues, K_{max} , to four to keep queue management overheads tolerable.

Each request size is assigned to one of the K clusters based on its proximity to the cluster centroids. To determine the per-queue cutoffs, we exploit the results of K-Means clustering. We obtain the K centroids from the clustering result, and define the cluster boundaries as the midpoint between the centroids of two consecutive clusters. For example, the boundary between $Cluster_i$ and $Cluster_{i+1}$ is $\frac{Centroid_i + Centroid_{i+1}}{2}$. The boundaries represent the maximum request size for each queue: $Queue_1$ handles requests smaller than $Boundary_1$, $Queue_2$ handles requests larger than $Boundary_1$ but smaller than $Boundary_2$, and this continues to all K queues.

However, the distribution of request sizes changes over time due to the fluctuating load behavior. Hence, static queue configurations can lead to inefficiencies. Instead, Chameleon dynamically adjusts the number of queues based on the observed load patterns. Specifically, the system periodically gathers recent request data to analyze the distribution of request sizes, and, every $T_{refresh}$, it re-computes the optimal number of queues using the aforementioned method.

Note that changes in load patterns are not sharp [50], hence, changing the multi-queue organization can happen relatively infrequently, e.g., $T_{refresh}$ can be 5 minutes, inducing negligible overheads.

Assigning Quotas per Queue. After determining the number of queues in the system and the per-queue cut-offs, Chameleon Scheduler needs to assign the resource quotas for each queue. These quotas govern how much resources can be assigned for requests from a given queue. For that, we use queuing theory, modeling the $K * M/M/1$ queues [32]. Recall that, to adjust to the dynamic nature of the workload, Chameleon recomputes the number of queues, per-queue cut-offs, and per-queue resource quotas, every $T_{refresh}$.

We take the maximum allowed size of a request for a queue (S), the assigned resource quota to a queue (Tok), the expected duration of a request from a queue (D), the arrival rate for these requests (λ), and the requests' SLO . The requests' processing rate is then $\mu = \frac{Tok}{S * D}$, while the total time that a request spends in the system is $T_{total} = \frac{1}{\mu - \lambda}$. To meet the SLO , the system needs to satisfy the following equation: $T_{total} \leq SLO$. Combining these constraints, we compute the minimum number of tokens, Tok_{min} , required to meet the SLO of requests within a given queue:

$$Tok_{min} \geq S * D * \left(\frac{1}{SLO} + \lambda \right)$$

The total number of available tokens in the system, Tok_{total} , is greater or equal to $\sum Tok_{min}^q$ (sum of minimum tokens needed by each queue q). Each queue is assigned its minimal required tokens (Tok_{min}^q), and the remaining tokens ($Tok_{total} - \sum Tok_{min}^q$) are split across queues proportionally to their initial weights.

5 Evaluation

5.1 Evaluation Methodology

Hardware Platforms and LLMs. We run our experiments on a server equipped with an A40 NVIDIA GPU [36] and with an AMD EPYC 9454 CPU. The GPU has 48GB memory, while the CPU has 48 cores and 377GB of main memory. For the majority of experiments, we use the Llama-7B [54] model. For our scalability experiments, we use a server equipped with A100 NVIDIA GPU [35] configured with 24GB, 48GB, and 80GB of GPU memory. When memory capacity allows, in addition to the Llama-7B model, we also run the Llama-13B and Llama-30B models. While we present the results with Llama series of models, we also experimented with other models, such as Falcon [53], OPT [29], and Mixtral [34]. We observed similar performance trends.

Workload Configuration. We set the input and output lengths of requests based on the open-source production trace from Azure [39]. To vary the load on the system, we use the Poisson distribution for the request inter-arrival time [4, 25, 26]. To each request we attach an adapter following the power-law distribution [47]. Specifically, we set the number of different adapters used by the requests to N_a . There are five adapter ranks: 8, 16, 32, 64, and 128. Each rank has equal number of different adapters, i.e., $N_a/5$. Then, for a request we pick the adapter rank following the power-law distribution, smaller adapters are more likely to be chosen. After choosing the rank, we pick one of the adapters of a given rank with a uniform distribution. If not specified otherwise, throughout all experiments in the evaluation we set N_a to 100.

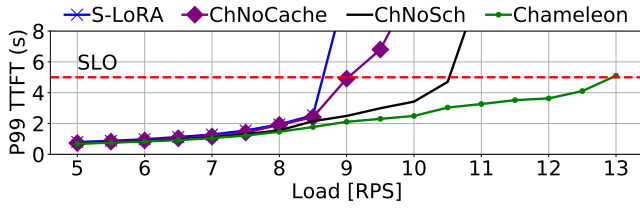


Figure 10: P99 TTFT tail latency for *S-LoRA*, *ChameleonNoCache*, *ChameleonNoSched* and *Chameleon* under different loads. Red dashed line indicates SLO.

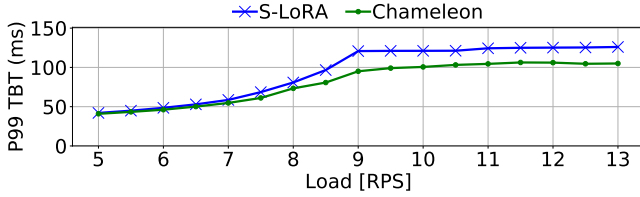


Figure 11: P99 TBT tail latency for *S-LoRA* and *Chameleon*.

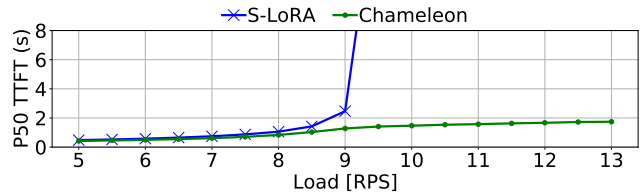


Figure 12: P50 TTFT latency for *S-LoRA* and *Chameleon*.

Baseline Systems. We run the experiments on open-source state-of-the-art LLM inference serving platform for adapter environments, *S-LoRA* [47], and compare *Chameleon* to *S-LoRA*. Note that *S-LoRA* performs iteration-level scheduling and asynchronous adapter prefetching. We also compare *Chameleon*'s scheduling component to the recently proposed SJF scheduling in μ Serve [44]. We measure TTFT, TBT, and end-to-end latency. As in prior work [26, 39, 50], we set the SLO to be 5 \times the average request execution time in a low-load system. In the following, we show our main results.

5.2 Performance Gains

Tail Latency. Figure 10 shows the P99 TTFT tail latency for *S-LoRA* and *Chameleon* under different loads. *Chameleon* constantly outperforms the baseline, and the benefits become more pronounced as the load increases. At low (6RPS), medium (8RPS), and high (9RPS) loads, *Chameleon* reduces the tail latency over the baseline by 14.7%, 24.6%, and 80.7%, respectively. There are two reason for *Chameleon*'s performance benefits. First, its caching mechanism reduces the adapter loading time and alleviates the bottlenecks on the PCIe link bandwidth. Second, its scheduling policy removes head-of-line blocking and prevents starvation, especially helping the requests at the tail.

Chameleon reduces both TTFT and TBT tail latencies. Figure 11 shows the P99 TBT latency for *S-LoRA* and *Chameleon* under different loads. *Chameleon* consistently outperforms the baseline. However, both systems keep their TBT latency under the SLO (150ms). The reason is that TBT latencies are less affected by the queuing effects, and they do not wait on adapter loading.

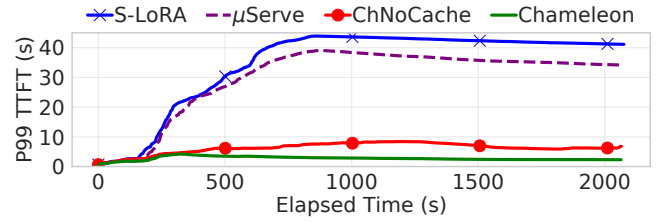


Figure 13: P99 TTFT latency over time with different scheduling policies: FIFO in *S-LoRA*, SJF in μ Serve, and our proposed policy in *ChameleonNoCache* and *Chameleon*.

Throughput. Figure 10 also shows the expected TTFT SLO as a red dashed line. We define the throughput as the load that a system can sustain without violating the SLO. From Figure 10, we can see that *S-LoRA*'s starts violating the SLO around 8.7 RPS, while *Chameleon*'s starts violating the SLO around 12.9 RPS. This results in 1.5 \times higher throughput for *Chameleon*.

Median Latency. Figure 12 shows the P50 TTFT latency for *S-LoRA* and *Chameleon* under different loads. At low (6RPS), medium (8RPS), and high (9RPS) loads, *Chameleon* reduces the median latency over the baseline by 13.9%, 20.9%, and 48.1%, respectively. The benefits of *Chameleon* are still significant, although not as pronounced as in the tail latency. The reason for this is that, on average, even with the baseline the requests do not experience head-of-line blocking and adapter loading does not happen in bursts.

Performance Breakdown. To understand the performance benefits of the two main *Chameleon*'s techniques, we run them in isolation. Figure 10 shows the P99 TTFT latency of *Chameleon* when running only with our proposed caching or scheduling techniques. We call these systems *ChameleonNoSched* and *ChameleonNoCache*, respectively. Both systems improve the throughput over the *S-LoRA* baseline: *ChameleonNoSched* and *ChameleonNoCache* have 1.2 \times and 1.1 \times higher throughput than the baseline, respectively. However, their performance is substantially lower than the performance of *Chameleon*. Hence, we need to synergistically perform both adapter caching and adapter-aware scheduling.

5.3 Different Scheduling/Caching Policies

In earlier experiments, we were comparing *Chameleon* with *S-LoRA* baseline that performs FIFO request scheduling and does not cache unused adapters. Here, we compare to a SJF (shortest-job-first) scheduling policy proposed by μ Serve [44]. Also, we augment the baseline with *Chameleon* Cache using a simple LRU eviction policy.

Scheduling Policies. Figure 13 shows the P99 TTFT latency over time with different scheduling policies at high system load (9 RPS). Specifically, we run *S-LoRA*'s FIFO scheduling policy [47] and μ Serve's SJF scheduling policy [44], as two state-of-the-art baselines. Additionally, we run our proposed adapter-aware multi-queue scheduling policy (*ChameleonNoCache*). This system does not include our caching mechanism. Finally, we run the full *Chameleon* design with both caching and scheduling proposals.

Both baselines (*S-LoRA* and μ Serve) have huge tail latency that increases over time due to the queuing bottlenecks. These designs have non-operational TTFT latencies. With FIFO scheduling (*S-LoRA*) the requests at the tail are the short ones blocked by the

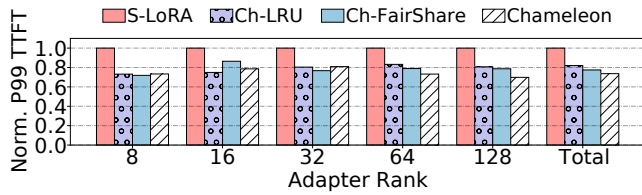


Figure 14: Normalized P99 TTFT latency for *S-LoRA*, *Chameleon-LRU*, *Chameleon-FairShare*, and *Chameleon*.

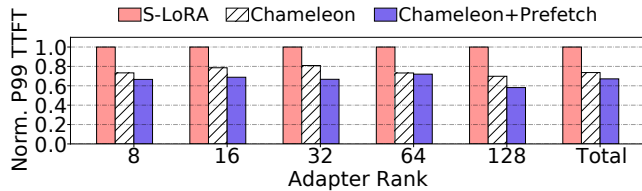


Figure 15: P99 TTFT latency for requests of different adapter ranks with *S-LoRA*, *Chameleon*, and *Chameleon+Prefetch*.

earlier longer requests, while with SJF scheduling (μ Serve) the requests at the tail are the long ones starved by the prioritization of shortest requests. Our proposed scheduling policy (*Chameleon-NoCache*) is very effective: it removes both head-of-line blocking effects and starvation, leading to much lower tail latencies.

The results in the figure are also aligned with observations from earlier experiments: to achieve the full performance potential the system needs to integrate both scheduling and caching approaches, as *Chameleon* has the lowest tail latency.

Caching Policies. Figure 14 shows the normalized P99 TTFT latency for requests of different adapter ranks with different caching policies at medium system load (8 RPS). We run the baseline system without an adapter cache (*S-LoRA*) and our proposed adapter cache mechanism with different replacement policies. *LRU* is a simple replacement policy found in conventional hardware caches: the least recently used adapter is evicted from the cache. *FairShare* uses our mechanism to consider adapter’s recency, frequency, and size while assigning each of the three knobs the same weight. Finally, *ChameleonCache* tunes the weights for the three knobs based on our extensive profiling (Section 4.1).

We observe that *Chameleon*’s proposed caching mechanism is very effective. All caching schemes reduce the P99 TTFT latency over the baseline by considerable amount for all adapter ranks. Additionally, our proposed replacement policy further improves the performance, especially for larger adapters. For example, for requests with adapter rank 128, *Chameleon* reduces the P99 TTFT latency over *Ch-FairShare* by 12%. In total, *Ch-LRU*, *Ch-FairShare*, and *Chameleon* reduce the P99 TTFT latency over *S-LoRA* by 18%, 22%, and 26%, respectively.

Prefetching Mechanism. To reduce the latency on cache misses, *Chameleon* uses a prefetching mechanism. Specifically, the system predicts which adapters are going to be used in the near future and prefetches them to the adapter cache ahead of time. The system uses a histogram-based technique to predict the future load of user-facing services [46]. As this mechanism is highly dependent on the prediction accuracy of the future loads, we do not include it in any of our main experiments. Here, we show the potential benefits of

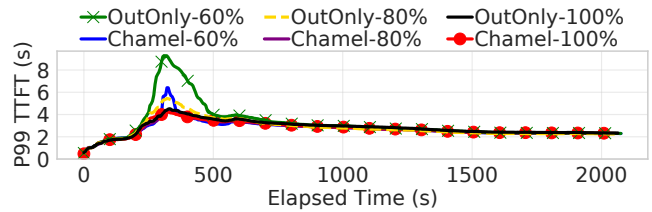


Figure 16: P99 TTFT latency over time for two scheduling policies (*OutputOnly* and *Chameleon*) under different output predictor accuracies.

such a prefetching mechanism. Figure 15 shows the normalized P99 TTFT latency for requests of different adapter ranks under medium load in three systems: *S-LoRA*, *Chameleon*, and *Chameleon+Prefetch*. Prefetching can further improve the performance of *Chameleon*. In total, across all requests, prefetching further reduces the P99 TTFT latency by 8.8%. As adapters are set to follow the power-law distribution across ranks and uniform distribution across adapters within the same rank, their predictability is very high. Hence, under these conditions, the accuracy of our prefetcher is above 95%.

5.4 Sensitivity Analysis

To evaluate the robustness of our design, we perform a sensitivity study. Specifically, we evaluate the efficiency of our scheduling policy under different accuracies of our output length predictor. Figure 16 shows the P99 TTFT latency for two scheduling policies: *OutputOnly* considers only the request’s output length (similar to prior work [44]), while *Chameleon* considers all request’s properties (input length, output length, and adapter size). We run with 100%, 80%, and 60% accuracy of the output length predictor. In all other experiments, the accuracy of our predictor is 80%.

The system is robust to predictor accuracy for most of the time. However, during load spikes (at around 300s), the systems with lower accuracy have high tail latency. Importantly, the scheduler that uses only the predicted output length as knob to prioritize requests is more sensitive to the predictor accuracy than our proposed scheme, as shown in the difference between *OutputOnly-60%* and *Chameleon-60%*. With the 80% predictor accuracy (as observed in our evaluation), the system has negligible performance loss compared to an ideal predictor (100% accuracy).

5.5 Scalability Analysis

To assess the scalability of *Chameleon*, we run experiments with larger models (Llama-7B, Llama-13B, and Llama-30B) and with different memory capacities (24GB, 48GB, and 80GB). In this section, we run all the experiments on an A100 NVIDIA GPU that by default has 80GB of memory. Given the available memory space, we use 500, 100, and 10 different adapters with 7B, 13B, and 30B models.

Scalability to LLM Size. First, we consider performance scaling with larger LLMs. Figure 17 on the left shows the normalized P99 TTFT latency of *Chameleon* over *S-LoRA* with different LLMs (Llama-7B, Llama-13B, and Llama-30B). *Chameleon*’s latency for a given model and load is normalized to the latency of *S-LoRA* for the same model and at the same load. *Chameleon* outperforms the baseline across all model sizes and system loads.

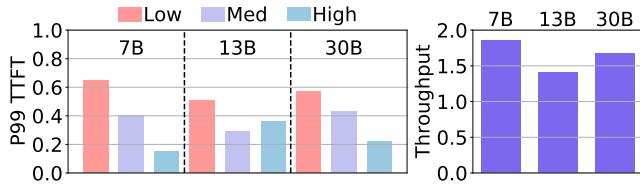


Figure 17: Normalized P99 TTFT latency (left) and throughput (right) for *Chameleon* over *S-LoRA* with different LLMs (Llama-7B, 13B, and 30B). Latency is normalized to that of *S-LoRA* for a given model and load level (Low, Medium, High).

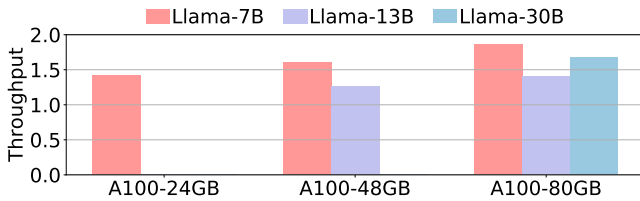


Figure 18: Normalized throughput for *Chameleon* over *S-LoRA* with different LLMs (Llama-7B, 13B, and 30B), when running with different memory configurations (24GB, 48GB, 80GB) on an A100 GPU.

As the model size increases, the heterogeneity across requests’ execution times becomes more pronounced, i.e., the difference in requests’ sizes creates even larger gap in their execution times. Thus, efficient scheduling provides more benefits. On the other hand, as the load increases, the system with larger models uses more GPU memory, leaving less space available for caching. Therefore, a sweet-spot for *Chameleon* is medium load with the medium size LLM (Llama-13B). Overall, averaged across all loads, *Chameleon* reduces the P99 TTFT latency over the baseline by 60.0%, 61.3%, and 59.3% for Llama-7B, Llama-13B, and Llama-30B models, respectively.

Figure 17(right) shows the normalized throughput of *Chameleon* over *S-LoRA* with different LLMs (Llama-7B, Llama-13B, and Llama-30B). *Chameleon* improves the throughput over the baseline by 1.86 \times , 1.41 \times , and 1.67 \times for Llama-7B, Llama-13B, and Llama-30B models, respectively. For larger models, *Chameleon* has less space to improve the throughput as such models are bounded by the memory capacity. For example, for both Llama-13B and Llama-30B, *Chameleon* increases the throughput of *S-LoRA* from 2.5K and 1.5K tokens per second to 3.5K and 2.5K tokens per second, respectively.

Scalability to GPU Memory Capacity. We consider performance scaling with different memory capacity. Figure 18 shows the normalized throughput of *Chameleon* over *S-LoRA* with different memory configurations of an A100 GPU (24GB, 48GB, 80GB) while running different LLMs (Llama with 7B, 13B, and 30B). Llama-30B fits only in 80GB memory, Llama-13B fits in 48GB and 80GB memory, while Llama-7B fits in all memory configurations. As memory capacity increases, *Chameleon* improves its relative performance over the baseline to a larger extent. *Chameleon* improves the throughput of Llama-7B model over *S-LoRA* by 1.4 \times , 1.6 \times , and 1.9 \times with 24GB, 48GB, and 80GB of GPU memory, respectively. The reason for the larger throughput improvements of *Chameleon* is that GPUs with larger memory capacity create more space for adapter caching.

Scalability to GPU Compute Capability. In addition, we analyze the performance scaling of *Chameleon* when running on different hardware platforms with the same memory capacity. Figure 10 shows that *Chameleon* improves the throughput over the baseline by 1.5 \times on an A40 GPU with 48GB of memory, while Figure 17 shows that *Chameleon* improves the throughput over the baseline by 1.9 \times on an A100 GPU with 48GB of memory. The reason for larger performance savings on A100 over A40, is that A100 improves the request’s execution time. Thus, adapter loading overheads substantially affect the request’s critical path.

6 Related Work

LLM Inference Optimizations. Large body of work proposed hardware [5, 15, 19, 22, 23, 39, 42, 43, 60–62, 65], algorithm [8, 13, 18] and system-level [1, 20, 30, 31, 59] optimizations for performance and energy-efficiency [38, 50] of LLM inference systems. All these works consider only a single monolithic LLM and do not optimize for a multi-adapter LLM inference environment. *Chameleon* is orthogonal to such techniques and can be combined with them for better performance or energy-efficiency.

LLM Inference with Parameter-Efficient Fine Tuning. Since the adoption of parameter-efficient fine tuning techniques [16, 24, 27, 56] researchers have been working on optimizing the system stack for efficient LLM inference in such multi-adapter environments [4, 25, 47, 58, 64]. *S-LoRA* [47] and *Punica* [4] decouple the base model from task-specific adapters and fetch the required adapters on the fly from the host to the GPU memory. *dLoRA* [58] dynamically merges and unmerges adapters with the base model based on the current system state. *Chameleon* builds on top of these works and proposes two new techniques: adapter caching and adapter-aware scheduling. Throughout the paper, we quantitatively compare *Chameleon* to *S-LoRA* as state-of-the-art baseline.

LLM Inference Scheduling. [10, 37, 44, 48, 51, 57, 59]. *μ Serve* [44] reduces the head-of-line blocking effects via a shortest-job-first (SJF) scheduling policy. In the paper, we compare quantitatively to *μ Serve*. Based on the input and output request length distribution, *ExeGPT* [37] and *DynamoLLM* [50] allocate resources (batch size and model parallelism) and schedule the requests for the minimal cost and energy consumption, respectively. *Llumnix* [51] reschedules the requests across worker replicas to improve load balance. These works focus on request scheduling in a multi-node environment, while using iteration-level scheduling [59] within a node. In contrast, *Chameleon* redesigns the scheduling policy within a node, and can be combined with cluster-level schedulers.

General-Purpose Workload Scheduling. [6, 14, 33] *Size-Interval Task Assignment (SITA)* [6, 14] addresses head-of-line blocking by providing an “express-lane” for short tasks. *Q-Zilla* [33] capitalizes on this idea and proposes a *Server-Queue Decoupled Size-Interval Task Assignment* for highly diverse microservice invocations. *Chameleon* adopts the algorithm to a new domain, multi-adapter LLM inference serving. Moreover, *SITA* assumes perfect knowledge of task size while *Q-Zilla* relies on request preemption. On the other hand, *Chameleon* uses a predictor for a request’s output length (unknown ahead of time), and it does not use preemption due to its high cost in LLM inference environments [20, 44, 57].

7 Conclusion

This paper presents *Chameleon*, an efficient LLM inference serving system in many-adapter environments. Chameleon relies on two main principles: adapter caching and adapter-aware request scheduling. Caching minimizes the overhead of loading the adapter weights on the request's critical path, while scheduling alleviates the head-of-line blocking for requests with highly diverse execution times. Under high loads, Chameleon reduces P99 TTFT latency by 80.7% and P50 TTFT latency by 48.1%, while improving the throughput by 1.5× over state-of-the-art-baseline.

References

- [1] Amey Agrawal, Nitin Kedia, Ashish Panwar, Jayashree Mohan, Nipun Kwatra, Bhargav Gulavani, Alexey Tumanov, and Ramachandran Ramjee. 2024. Taming Throughput-Latency Tradeoff in LLM Inference with Sarathi-Serve. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI '24)*. USENIX Association, Santa Clara, CA, 117–134. <https://www.usenix.org/conference/osdi24/presentation/agrawal>
- [2] Keivan Alizadeh, Iman Mirzadeh, Dmitry Belenko, Karen Khatamifard, Minsik Cho, Carlo C Del Mundo, Mohammad Rastegari, and Mehrdad Farajtabar. 2024. LLM in a flash: Efficient Large Language Model Inference with Limited Memory. arXiv:2312.11514 [cs.CL]
- [3] Jairus Bowne. 2024. Using Large Language Models in Learning and Teaching. <https://biomedsciences.unimelb.edu.au/study/dlh/assets/documents/large-language-models-in-education/flms-in-education>
- [4] Lequn Chen, Zihao Ye, Yongji Wu, Danyang Zhuo, Luis Ceze, and Arvind Krishnamurthy. 2024. Punica: Multi-Tenant LoRA Serving. In *Proceedings of Machine Learning and Systems*, P. Gibbons, G. Pekhimenko, and C. De Sa (Eds.), Vol. 6. 1–13. https://proceedings.mlsys.org/paper_files/paper/2024/file/054de805fceb78a201f5e9d53c85908-Paper-Conference.pdf
- [5] Jaewan Choi, Jaehyun Park, Kwanhee Kyung, Nam Sung Kim, and Jung Ho Ahn. 2024. Unleashing the Potential of PIM: Accelerating Large Batched Inference of Transformer-Based Generative Models. In *Proceedings of the International Symposium on High-Performance Computer Architecture (HPCA '24)*.
- [6] Mark E. Crovella, Mor Harchol-Balter, and Cristina D. Murta. 1998. Task assignment in a distributed system (extended abstract): improving performance by unbalancing load. In *Proceedings of the 1998 ACM SIGMETRICS Joint International Conference on Measurement and Modeling of Computer Systems (Madison, Wisconsin, USA) (SIGMETRICS '98/PERFORMANCE '98)*. Association for Computing Machinery, New York, NY, USA, 268–269. <https://doi.org/10.1145/277851.277942>
- [7] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness. arXiv:2205.14135 [cs.LG]
- [8] Jyotikrishna Dass, Shang Wu, Huihong Shi, Chaojian Li, Zhifan Ye, Zhongfeng Wang, and Yingyan Lin. 2023. ViTALiTy: Unifying Low-rank and Sparse Approximation for Vision Transformer Acceleration with a Linear Taylor Attention. In *Proceedings of the International Symposium on High-Performance Computer Architecture (HPCA '23)*.
- [9] Yao Fu, Leyang Xue, Yeqi Huang, Andrei-Octavian Brabete, Dmitrii Ustiugov, Yuvraj Patel, and Luo Mai. 2024. ServerlessLLM: Low-Latency Serverless Inference for Large Language Models. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI '24)*.
- [10] Yichao Fu, Siqi Zhu, Runlong Su, Aurick Qiao, Ion Stoica, and Hao Zhang. 2024. Efficient LLM Scheduling by Learning to Rank. arXiv:2408.15792 [cs.LG] <https://arxiv.org/abs/2408.15792>
- [11] Alexander Fuerst and Prateek Sharma. 2021. FaasCache: keeping serverless computing alive with greedy-dual caching. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '21)*. New York, NY, USA. <https://doi.org/10.1145/3445814.3446757>
- [12] GitHub. 2024. The world's most widely adopted AI developer tool. <https://github.com/features/copilot>.
- [13] Cong Guo, Jiaming Tang, Weiming Hu, Jingwen Leng, Chen Zhang, Fan Yang, Yunxin Liu, Minyi Guo, and Yuhao Zhu. 2023. Olive: Accelerating Large Language Models via Hardware-friendly Outlier-Victim Pair Quantization. In *Proceedings of the 50th Annual International Symposium on Computer Architecture (ISCA '23)*.
- [14] Mor Harchol-Balter, Mark E. Crovella, and Cristina D. Murta. 1999. On Choosing a Task Assignment Policy for a Distributed Server System. *J. Parallel and Distrib. Comput.* 59, 2 (1999), 204–228. <https://doi.org/10.1006/jpdc.1999.1577>
- [15] Guseul Heo, Sangyeop Lee, Jaehong Cho, Hyunmin Choi, Sanghyeon Lee, Hyungkyu Ham, Gwangsun Kim, Divya Mahajan, and Jongse Park. 2024. Neupims: NPU-PIM Heterogeneous Acceleration for Batched LLM Inference. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*.
- [16] Edward J Hu, Weizhu Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations (ICLR)*. <https://arxiv.org/abs/2106.09685>
- [17] Huggingface. 2023. PEFT: Parameter-Efficient Fine-Tuning of Billion-Scale Models on Low-Resource Hardware. <https://github.com/huggingface/peft>
- [18] Ranggi Hwang, Jianyu Wei, Shijie Cao, Changho Hwang, Xiaohu Tang, Ting Cao, and Mao Yang. 2024. Pre-gated MoE: An Algorithm-System Co-Design for Fast and Scalable Mixture-of-Expert Inference. In *Proceedings of the 51st Annual International Symposium on Computer Architecture (ISCA '24)*.
- [19] Hongsun Jang, Jaeyong Song, Jaewon Jung, Jaeyoung Park, Youngsok Kim, and Jinho Lee. 2024. Smart-Infinity: Fast Large Language Model Training using Near-Storage Processing on a Real System. In *Proceedings of the International Symposium on High-Performance Computer Architecture (HPCA '24)*.
- [20] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of the 29th Symposium on Operating Systems Principles (Koblenz, Germany) (SOSP '23)*. Association for Computing Machinery, New York, NY, USA, 611–626. <https://doi.org/10.1145/3600006.3613165>
- [21] Marina Lammertyn. 2024. 60+ ChatGPT Statistics And Facts You Need to Know in 2024. <https://blog.invgate.com/chatgpt-statistics>.
- [22] Jungi Lee, Wonbeom Lee, and Jaewoong Sim. 2024. Tender: Accelerating Large Language Models via Tensor Decomposition and Runtime Requantization. In *Proceedings of the 51st Annual International Symposium on Computer Architecture (ISCA '24)*.
- [23] SeungYul Lee, Hyunseung Lee, Jihoon Hong, SangLyu Cho, and Jae W. Lee. 2024. VGA: Hardware Accelerator for Scalable Long Sequence Model Inference. In *Proceedings of the International Symposium on Microarchitecture (MICRO '24)*.
- [24] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. arXiv:2104.08691 [cs.CL] <https://arxiv.org/abs/2104.08691>
- [25] Suyi Li, Hanfeng Lu, Tianyuan Wu, Minchen Yu, Qizhen Weng, Xusheng Chen, Yizhou Shan, Binhang Yuan, and Wei Wang. 2024. CaraServe: CPU-Assisted and Rank-Aware LoRA Serving for Generative LLM Inference. arXiv:2401.11240 [cs.DC] <https://arxiv.org/abs/2401.11240>
- [26] Zhuohan Li, Lianmin Zheng, Yinmin Zhong, Vincent Liu, Ying Sheng, Xin Jin, Yanping Huang, Zhifeng Chen, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. AlpaServe: Statistical Multiplexing with Model Parallelism for Deep Learning Serving. In *Proceedings of the 17th USENIX Symposium on Operating Systems Design and Implementation (OSDI '23)*.
- [27] Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks. arXiv:2110.07602 [cs.CL] <https://arxiv.org/abs/2110.07602>
- [28] Meta. 2024. Llama3-70B. <https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct>.
- [29] Meta AI. 2024. Open Pre-trained Transformer Language Models. https://huggingface.co/docs/transformers/model_doc/opt.
- [30] Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Zhengxin Zhang, Rae Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, Chunan Shi, Zhuoming Chen, Daiyaan Arfeen, Reyna Abhyankar, and Zhihao Jia. 2024. SpecInfer: Accelerating Large Language Model Serving with Tree-based Speculative Inference and Verification. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*.
- [31] Xupeng Miao, Chunan Shi, Jiangfei Duan, Xiaoli Xi, Dahua Lin, Bin Cui, and Zhihao Jia. 2024. SpotServe: Serving Generative Large Language Models on Preemptible Instances. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*.
- [32] Amirhossein Mirhosseini and Thomas Wenisch. 2021. μ Steal: a theory-backed framework for preemptive work and resource stealing in mixed-criticality microservices. In *Proceedings of the 35th ACM International Conference on Supercomputing (Virtual Event, USA) (ICS '21)*. Association for Computing Machinery, New York, NY, USA, 102–114. <https://doi.org/10.1145/3447818.3463529>
- [33] Amirhossein Mirhosseini, Brendan L. West, Geoffrey W. Blake, and Thomas F. Wenisch. 2020. Q-Zilla: A Scheduling Framework and Core Microarchitecture for Tail-Tolerant Microservices. In *2020 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. 207–219. <https://doi.org/10.1109/HPCA47549.2020.00026>
- [34] Mistral AI. 2024. The Mixtral-8x22B Large Language Model. <https://huggingface.co/mistralai/Mixtral-8x22B-Instruct-v0.1>.
- [35] NVIDIA. 2024. Introduction to the NVIDIA DGX A100 System. <https://docs.nvidia.com/dgx/dgxa100-user-guide/introduction-to-dgxa100.html>.
- [36] NVIDIA. 2024. NVIDIA A40 Data Center GPU. <https://www.nvidia.com/en-us/data-center/a40/>.

- [37] Hyungjun Oh, Kihong Kim, Jaemin Kim, Sungkyun Kim, Junyeol Lee, Du-seong Chang, and Jiwon Seo. 2024. ExeGPT: Constraint-Aware Resource Scheduling for LLM Inference. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*.
- [38] Pratyush Patel, Esha Choukse, Chaojie Zhang, Ínigo Goiri, Brijesh Warrier, Nithish Mahalingam, and Ricardo Bianchini. 2024. Characterizing Power Management Opportunities for LLMs in the Cloud. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*.
- [39] Pratyush Patel, Esha Choukse, Chaojie Zhang, Aashaka Shah, Ínigo Goiri, Saeed Maleki, and Ricardo Bianchini. 2024. Splitwise: Efficient generative LLM inference using phase splitting. In *Proceedings of the 51st Annual International Symposium on Computer Architecture (ISCA '24)*.
- [40] Cheng Peng, Xi Yang, Aokun Chen, Kaleb Smith, Nima PourNejatian, Anthony Costa, Cheryl Martin, Mona Flores, Ying Zhang, Tanja Magoc, Gloria Lipori, Mitchell Duane, Naykky Ospina, Mustafa Ahmed, William Hogan, Elizabeth Shenkman, Yi Guo, Jiang Bian, and Yonghui Wu. 2023. A study of generative large language model for medical research and healthcare. *npj Digital Medicine* (2023).
- [41] Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2020. AdapterHub: A Framework for Adapting Transformers. *arXiv preprint arXiv:2007.07779* (2020). <https://arxiv.org/abs/2007.07779>
- [42] Yubin Qin, Yang Wang, Dazheng Deng, Zhiren Zhao, Xiaolong Yang, Leibo Liu, Shaojun Wei, Yang Hu, and Shouyi Yin. 2023. FACT: FFN-Attention Co-optimized Transformer Architecture with Eager Correlation Prediction. In *Proceedings of the 50th Annual International Symposium on Computer Architecture (ISCA '23)*.
- [43] Yubin Qin, Yang Wang, Zhiren Zhao, Xiaolong Yang, Yang Zhou, Shaojun Wei, Yang Hu, and Shouyi Yin. 2024. MECLA: Memory-Compute-Efficient LLM Accelerator with Scaling Sub-matrix Partition. In *Proceedings of the 51st Annual International Symposium on Computer Architecture (ISCA '24)*.
- [44] Haoran Qiu, Weichao Mao, Archit Patke, Shengkun Cui, Saurabh Jha, Chen Wang, Hubertus Franke, Zbigniew Kalbarczyk, Tamer Başar, and Ravishankar K. Iyer. 2024. Power-aware Deep Learning Model Serving with μ -Serve. In *2024 USENIX Annual Technical Conference (USENIX ATC '24)*. USENIX Association, Santa Clara, CA, 75–93. <https://www.usenix.org/conference/atc24/presentation/qiu>
- [45] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- [46] Mohammad Shahrad, Rodrigo Fonseca, Inigo Goiri, Gohar Chaudhry, Paul Batum, Jason Cooke, Eduardo Laureano, Colby Tresness, Mark Russinovich, and Ricardo Bianchini. 2020. Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider. In *Proceedings of the 2020 USENIX Annual Technical Conference (USENIX ATC '20)*.
- [47] Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper, Nicholas Lee, Shuo Yang, Christopher Chou, Banghua Zhu, Lianmin Zheng, Kurt Keutzer, Joseph Gonzalez, and Ion Stoica. 2024. SLoRA: Scalable Serving of Thousands of LoRA Adapters. In *Proceedings of Machine Learning and Systems*, P. Gibbons, G. Pekhimenko, and C. De Sa (Eds.), Vol. 6. 296–311. https://proceedings.mlsys.org/paper_files/paper/2024/file/906419cd502575b617cc489a1a696a67-Paper-Conference.pdf
- [48] Ying Sheng, Shiyi Cao, Dacheng Li, Banghua Zhu, Zhuohan Li, Danyang Zhuo, Joseph E. Gonzalez, and Ion Stoica. 2024. Fairness in Serving Large Language Models. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI '24)*. USENIX Association, Santa Clara, CA, 965–988. <https://www.usenix.org/conference/osdi24/presentation/sheng>
- [49] Jovan Stojkovic, Esha Choukse, Chaojie Zhang, Inigo Goiri, and Josep Torrellas. 2024. Towards Greener LLMs: Bringing Energy-Efficiency to the Forefront of LLM Inference. *arXiv e-prints*, Article arXiv:2403.20306 (March 2024), arXiv:2403.20306 pages. <https://doi.org/10.48550/arXiv.2403.20306> arXiv:2403.20306 [cs.AI]
- [50] Jovan Stojkovic, Chaojie Zhang, Inigo Goiri, Josep Torrellas, and Esha Choukse. 2024. DynamoLLM: Designing LLM Inference Clusters for Performance and Energy Efficiency. *CoRR abs/2408.00741* (2024). arXiv:2408.00741 <http://arxiv.org/abs/2408.00741>
- [51] Biao Sun, Ziming Huang, Hanyu Zhao, Wencong Xiao, Xinyi Zhang, Yong Li, and Wei Lin. 2024. Llumnix: Dynamic Scheduling for Large Language Model Serving. In *Proceedings of the 18th USENIX Symposium on Operating Systems Design and Implementation (OSDI '24)*.
- [52] Kartik Talamadupula. 2024. A Guide to LLM Inference Performance Monitoring. <https://syml.ai/developers/blog/a-guide-to-llm-inference-performance-monitoring/>.
- [53] Technology Innovation Institute (TII). 2024. Falcon-180B. <https://huggingface.co/tiiuae/falcon-180B>. (2024).
- [54] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).
- [55] Tanay Varshney. 2024. Build an LLM-Powered Data Agent for Data Analysis. <https://developer.nvidia.com/blog/build-an-llm-powered-data-agent-for-data-analysis/>.
- [56] Zhengbo Wang, Jian Liang, Ran He, Zilei Wang, and Tieniu Tan. 2024. LoRA-Pro: Are Low-Rank Adapters Properly Optimized? <https://arxiv.org/abs/2407.18242>
- [57] Bingyang Wu, Yinmin Zhong, Zili Zhang, Gang Huang, Xuanzhe Liu, and Xin Jin. 2023. Fast Distributed Inference Serving for Large Language Models. arXiv:2305.05920 [cs.LG]
- [58] Bingyang Wu, Ruidong Zhu, Zili Zhang, Peng Sun, Xuanzhe Liu, and Xin Jin. 2024. dLoRA: Dynamically Orchestrating Requests and Adapters for LoRA LLM Serving. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI '24)*. USENIX Association, Santa Clara, CA, 911–927. <https://www.usenix.org/conference/osdi24/presentation/wu-bingyang>
- [59] Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. 2022. Orca: A Distributed Serving System for Transformer-Based Generative Models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI '22)*. USENIX Association, Carlsbad, CA, 521–538. <https://www.usenix.org/conference/osdi22/presentation/yy>
- [60] Zhongkai Yu, Shengwen Liang, Tianyun Ma, Yunke Cai, Ziyuan Nan, Di Huang, Xinkai Song, Yifan Hao, Jie Zhang, Tian Zhi, Yongwei Zhao, Zidong Du, Xing Hu, Qi Guo, and Tianshi Chen. 2024. FlashLLM: A Chiplet-Based In-Flash Computing Architecture to Enable On-Device Inference of 70B LLM. In *Proceedings of the International Symposium on Microarchitecture (MICRO '24)*.
- [61] Sungmin Yun, Kwanhee Kyung, Juhwan Cho, Jaewan Choi, Jongmin Kim, Byeongho Kim, Sukhan Lee, Kyomin Sohn, and Jung Ho Ahn. 2024. Duplex: A Device for Large Language Models with Mixture of Experts, Grouped Query Attention, and Continuous Batching. In *Proceedings of the International Symposium on Microarchitecture (MICRO '24)*.
- [62] Hengrui Zhang, August Ning, Rohan Baskar Prabhakar, and David Wentzloff. 2024. LLMCompass: Enabling Efficient Hardware Design for Large Language Model Inference. In *Proceedings of the 51st Annual International Symposium on Computer Architecture (ISCA '24)*.
- [63] Youpeng Zhao, Di Wu Wu, and Jun Wang. 2024. ALISA: Accelerating Large Language Model Inference via Sparsity-Aware KV Caching. In *Proceedings of the 51st Annual International Symposium on Computer Architecture (ISCA '24)*.
- [64] Zhe Zhou, Xuechao Wei, Jiejing Zhang, and Guangyu Sun. 2022. PetS: A Unified Framework for Parameter-Efficient Transformers Serving. In *Proceedings of the USENIX Annual Technical Conference (USENIX ATC '22)*.
- [65] Hanqing Zhu, Jiaqi Gu, Hanrui Wang, Zixuan Jiang, Zhekai Zhang, Rongxing Tang, Chenghao Feng, Song Han, Ray T. Chen, and David Z. Pan. 2024. Lightening-Transformer: A Dynamically-Operated Optically-Interconnected Photonic Transformer Accelerator. In *Proceedings of the International Symposium on High-Performance Computer Architecture (HPCA '24)*.