# BurstGPT: A Real-World Workload Dataset to Optimize LLM Serving Systems

### Yuxin Wang[*]
HKRC
Hong Kong, CHINA

### Yuhan Chen[*]
HKUST (Guangzhou)
Guangzhou, Guangdong, CHINA

### Zeyu Li
HKUST (Guangzhou)
Guangzhou, Guangdong, CHINA

### Xueze Kang
HKUST (Guangzhou)
Guangzhou, Guangdong, CHINA

### Yuchu Fang
Huawei Technologies Co.
Shenzhen, Guangdong, CHINA

### Yeju Zhou
Huawei Technologies Co.
Shenzhen, Guangdong, CHINA

### Yang Zheng
Huawei Technologies Co.
Shenzhen, Guangdong, CHINA

### Zhenheng Tang
HKUST
Hong Kong, CHINA

### Xin He
Hong Kong Baptist University
Hong Kong, CHINA

### Rui Guo
Tsinghua University
Beijing, CHINA

### Xin Wang
Tsinghua University
Beijing, CHINA

### Qiang Wang
HIT (Shenzhen)
Shenzhen, Guangdong, CHINA

### Amelie Chi Zhou
Hong Kong Baptist University
Hong Kong, CHINA

### Xiaowen Chu[†]
HKUST (Guangzhou)
Guangzhou, Guangdong, CHINA

## Abstract

Serving systems for Large Language Models (LLMs) are often optimized to improve quality of service (QoS) and throughput. However, due to the lack of open-source LLM serving workloads, these systems are frequently evaluated under unrealistic workload assumptions. Consequently, performance may degrade when systems are deployed in real-world scenarios.

This work presents BurstGPT, an LLM serving workload with 10.31 million traces from regional Azure OpenAI GPT services over 213 days. BurstGPT captures LLM serving characteristics from user, model and system perspectives: (1) User request concurrency: burstiness variations of requests in Azure OpenAI GPT services, revealing diversified concurrency patterns in different services and model types. (2) User conversation patterns: counts and intervals within conversations for service optimizations. (3) Model response lengths: auto-regressive serving processes of GPT models, showing statistical relations between requests and their responses. (4) System response failures: failures of conversation and API services, showing intensive resource needs and limited availability of LLM services in Azure. The details of the characteristics can serve multiple purposes in LLM serving optimizations, such as system evaluation and trace provisioning. In our demo evaluation with BurstGPT, frequent variations in BurstGPT reveal declines in efficiency, stability, or reliability in realistic LLM serving. We identify that the generalization of KV cache management, scheduling and disaggregation optimizations can be improved under realistic workload evaluations. BurstGPT is publicly available now at https://github.com/HPMLL/BurstGPT and is widely used to develop prototypes of LLM serving frameworks in the industry.

## Keywords

LLM Serving, Workload Trace, Workload Management, Workload Provisioning, System Scheduling

[*]Both authors contributed equally to this research.

[†]Corresponding author. Email: xwchu@hkust-gz.edu.cn.

## 1 Introduction

Large language models (LLMs), particularly Generative Pre-trained Transformer (GPT) models, have significantly impacted the AI industry. Notable examples include OpenAI's ChatGPT [18][19], Google's Gemini [25], and Meta's Llama [26]. However, deploying and operating LLM services is costly due to the substantial use of AI accelerators [14]. The high volume of user requests and limited hardware resources further constrain the quality of service (QoS) [14]. Therefore, optimizing serving systems to reduce costs and improve user experience is urgent.
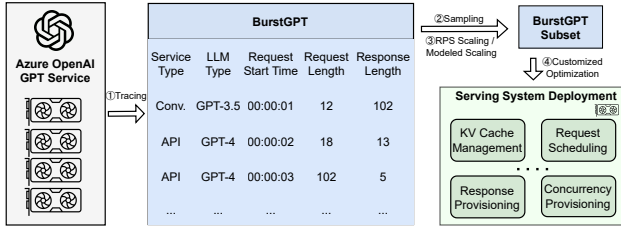
**Figure 1: Data collection and use method of `BurstGPT`. `BurstGPT` is a real-world workload trace from the Azure OpenAI GPT service. A scaled sample from a period of `BurstGPT` can be used to optimize serving systems using specific methods, considering realistic concurrency and response patterns. Note that we open-sourced two versions of `BurstGPT`: a cleaned trace and a raw trace, with failure logs excluded from the cleaned version.**

Existing systems primarily optimize system throughput and achieve better QoS. However, these studies often use non-LLM serving workloads in evaluations, derived from synthetic models [28] or non-LLM services such as Microsoft Azure Functions (MAF) [22][31]. This work identifies that these unrealistic synthetic workloads do not reflect the workload patterns of auto-regressive LLM serving and limit the effectiveness of evaluating these systems.

In addressing this gap, we introduce a real-world LLM serving workload dataset named `BurstGPT`. As shown in Figure 1, the traces are collected from the Azure OpenAI GPT Service [1], which serves GPT models on Azure and provides APIs for regional GPT service providers, such as enterprises and campuses, to build customized GPT services. We collected data from one of these regional GPT service providers with over 3,000 users. By deploying a logging engine on this regional service provider, we monitored privacy-independent metrics of each GPT request. In Figure 1, we collected the *request start time*, *request and response token lengths*, *service type* (API or conversational), and *LLM type* (ChatGPT, GPT-4 or GPT-4v) for 213 days. The dataset comprises 8.69 million traces of ChatGPT and 0.95 million GPT-4 traces using API services, and 0.30 million ChatGPT traces and 0.16 million traces of GPT-4 using conversational services. It also includes failure requests while the response length is zero.

To use `BurstGPT`, the user can sample a period of trace and scale it to fit the target system size. Users can apply vanilla `BurstGPT` using request-per-second (RPS) scaling or model the trace and adjust parameters, as shown in Figure 1. The subset can be used for various optimizations, such as identifying concurrency patterns for provisioning and predicting response lengths based on request-response length distributions. Additionally, this paper provides a detailed demo of `BurstGPT` for LLM serving evaluation. Key findings from our analysis and evaluations are listed below:

- User Request Concurrency. We identify unique patterns in request concurrency and lengths, varying by service types. This insight suggests that service-specific user request pattern evaluations are important for optimizing LLM serving. In our evaluation, (1) concurrency from non-LLM workloads does not accurately reflect system performance as

`BurstGPT` does. (2) Even for LLM workloads, generalization of system optimization across serving types is not guaranteed. For example, we verify that while optimized first-come-first-served scheduling is effective for conversation services, its efficiency can be reduced for API services.
- User Conversation Patterns. We share analyze traces for the first 2 months in 3.1, 3.3, 3.4 and conversation traces from the last 3 months in 3.2. The detailed analysis of conversation traces can help schedule LLM serving, e.g., KV cache offload strategies.
- Model Response Lengths. Longer response lengths lead to higher workloads during LLM serving. The unpredictable nature of response lengths per request introduces uncertainty to system pressure. The statistical study of the relationship between request and response lengths in `BurstGPT` sheds light on future response provisioning, which can help improve system performance.
- System Response Failures. We observe a relatively high failure rate in `BurstGPT`. In our evaluation, we identify that this is primarily due to inefficiencies in KV cache management. Variations in burstiness in `BurstGPT` lead to memory bottlenecks, causing spikes in failure rates and performance degradation.

With the dataset, we open-sourced a benchmark suite, BurstGPT-Perf. In our evaluation using BurstGPT-Perf, we first present our insights on using `BurstGPT` to assess system performance differences between LLM and non-LLM workloads. We also demonstrate use cases of `BurstGPT` in scheduling policy selection, workload provisioning. Then, we share our insights in prefill-decode(PD) disaggregation in industrial system prototypes. The evaluation results show that `BurstGPT` can effectively suggest optimizations for LLM serving systems in real-world scenarios.

## 2 Preliminary and Motivation

**Table 1: Comparison of BurstGPT with Other Traces**

| Serving Workloads | Real-world | | | Synthetic |
|---|---|---|---|---|
| | BurstGPT | MAF1 [22] | MAF2 [31] | FastServe [28] |
| LLM-User Req. Pat. | ✓ | ✗ | ✗ | ✗ |
| LLM-User Conv. Pat. | ✓ | ✗ | ✗ | ✗ |
| LLM Resp. Len. | ✓ | ✗ | ✗ | ✓ |
| Sys. Resp. Fail. | ✓ | ✗ | ✗ | ✗ |

### 2.1 Limitations of LLM Serving

Each request in LLM serving requires uncertain but large computational due to the model size and auto-regressive nature [12]. To reduce deployment and operational expenses, several specialized frameworks have been developed, including TensorRT-LLM [17], vLLM [12], DeepSpeed [21], and lightLLM [2]. In addition to these frameworks, various memory and compute optimizations have been proposed for efficient decoding [27][9][7][29][6][23][4][34][24], significantly improving LLM inference performance in deployments. Also, efficient request scheduling is crucial for optimizing LLM service. For example, Orca [30] introduces iteration-level scheduling, that dynamically adjusts batch size during iterations; FastServe [28]

proposes optimized Multilevel-Feedback-Queue (MLFQ) scheduling in the request queue to improve system efficiency.

However, these optimizations have yet to be evaluated on real-world LLM serving workloads. We identify that generalizing serving optimizations from non-LLM or synthetic workloads to LLM workloads can be challenging [16][15] for ensuring serving efficiency, stability, and reliability.

## 2.2 Towards Workload-aware LLM Serving

Stable workloads lead to consistent system performance across various metrics. However, in the real world, the workload of LLM is diversified due to user, system, and model behaviors. Neglecting any aspect of workload traces during system optimization may yield an incomplete understanding of a framework's behavior in real-world deployment.

*User Concurrency Patterns.* Currently, the assessment of performance in LLM serving frameworks relies on synthetic concurrency [28] or non-LLM concurrency, such as Microsoft Azure Functions (MAF) [22][31]. As concluded in Table 1, these approaches do not reflect real-world LLM workloads due to a lack of empirical data on system and model patterns, as well as user behaviors. For example, the request-per-second (RPS) metric in MAF is much higher (1.64 RPS on average) compared to LLM services (0.019 RPS on average in conversation service and 0.21 RPS on average in API service when using ChatGPT) due to the non-autoregressive and lightweight nature of function services in Azure. However, burstiness in concurrency is much less severe in MAF, likely leading to performance divergence from LLM workloads.

*Model Patterns.* The computational complexities of prefilling and decoding in LLM systems are $O(s^2)$ and $O(s)$ relative to the request length $s$, respectively [11]. This relationship between request length and system load significantly differs from that observed in other lightweight cloud functions and synthetic workloads [22][31], highlighting a unique resource occupation of LLM serving. These complexities are also reflected in the request concurrency, failure rates, and latency of each request in LLM serving.

## 3 Introduction to BurstGPT

### 3.1 User Request Concurrency

*Long-term Patterns: Periodicity and Aperiodicity.* The long-term usage patterns of conversation services, spanning hours and days, exhibit periodic characteristics. As illustrated in Figure 2, conversation volumes for both ChatGPT and GPT-4 peak during weekdays and diminish during weekends, implying heightened user interaction primarily on workdays. Notably, GPT-4's overall usage is lower compared to ChatGPT. Figure 4 shows that conversation volumes for both ChatGPT and GPT-4 exhibit periodic highs during working hours and lows during night hours.

In contrast, the long-term pattern of API services, over the same hourly and daily intervals, follows an aperiodic pattern characterized by burstiness. Figures 3 and 5 reveal irregular, dense request submissions to API services, significantly diverging from the more predictable patterns observed in conversation services. The higher variance in daily request volumes, indicated by a greater standard
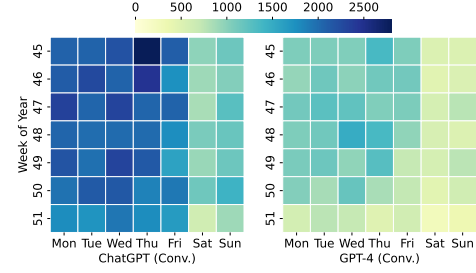


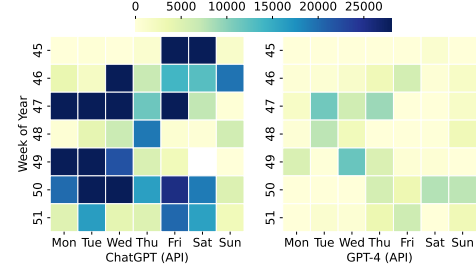**Figure 2: Weekly Periodicity of Conversation Services in BurstGPT.**



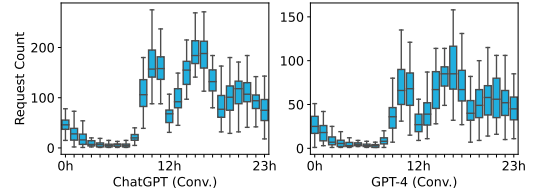**Figure 3: Weekly Aperiodicity API Services in BurstGPT.**



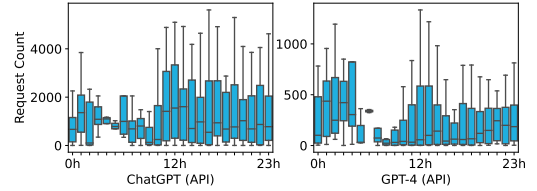**Figure 4: Daily Periodicity Conversation Services in BurstGPT.**



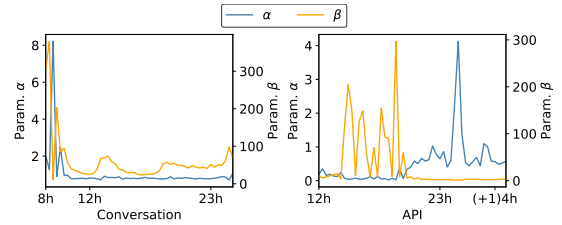**Figure 5: Daily Aperiodicity API Services in BurstGPT.**



**Figure 6: Variation of Burstiness in Gamma Distribution in BurstGPT.**

deviation compared to Figure 4, suggests these requests may be automated.

Yuxin Wang, Yuhan Chen, Zeyu Li, Xueze Kang, Yuchu Fang, Yeju Zhou, Yang Zheng, Zhenheng Tang, Xin He, Rui Guo, Xin Wang, Qiang Wang, Amelie Chi Zhou, and Xiaowen
Conference'17, July 2017, Washington, DC, USA                                                                                                                                   Chu

*Short-term Patterns: Variant Burstiness.* Figure 6 presents the modeled burstiness every twenty minutes in ChatGPT's conversation and API services during working hours on weekdays. The burstiness is modeled using the Gamma distribution, as supported in [14][28]. We observe that the shape parameter $\alpha$ and the rate parameter $\beta$ of the Gamma distribution vary between conversation and API services. A smaller $\alpha$ indicates a higher coefficient of variation (CV) of the Gamma distribution, meaning the workload is more bursty [28]. In conversation services, the parameter $\alpha$ varies sharply during working hours. In API services, both $\alpha$ and $\beta$ vary sharply during the day. These observations underscore the instability of burstiness in the system.

## 3.2 User Conversation Patterns

*Distribution.* Over 35% of conversations end with only one request. The left of Figure 9 shows that the distribution of conversation requests exhibits an exponential decline, with a median of 2 and 75% of conversations with four or fewer requests.

*Interval Time.* The left of Figure 9 shows the average interval time is not directly related to the number of requests. The P90 interval time increases when the number of requests is less than 5 and then fluctuates after that.

*The More Requests, the Longer Intervals.* The right of Figure 9 illustrates the number of intervals exceeding a certain duration $c$ in each conversation, which we define as long intervals. We make $c$ equal to 3 different time deltas and the interquartile range (IQR) upper bound. These intervals increase as the number of conversation requests increases.

## 3.3 Model Response Patterns

To identify differences in response length distributions across LLMs, this subsection employs conversational traces in BurstGPT and uses Llama-2-13b-chat[26] for comparison. For Llama-2-13b-chat,
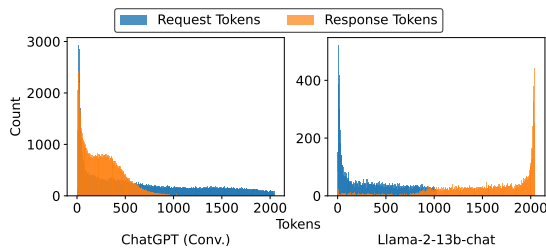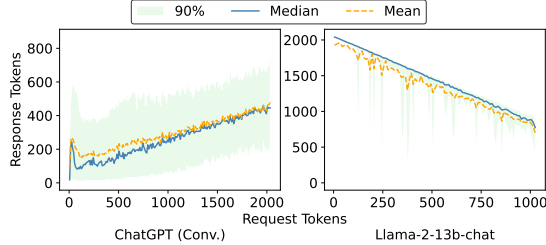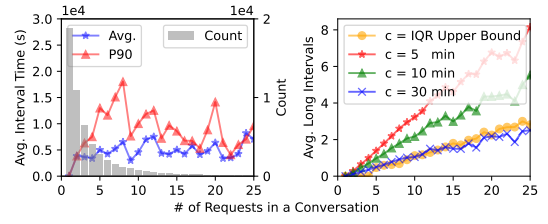


Figure 9: Count, Average Interval times, and Average Long Intervals as the Number of Requests in a Conversation in BurstGPT.

we feed 10k real-world prompts to the model for inference and gather the response lengths. The prompts are randomly truncated from the ShareGPT dialogues [8] with aligninled distribution to BurstGPT's conversational traces [32]. This alignment allows for comparing the response length distributions between the two LLMs, facilitating a deeper understanding of their behaviors in responding to conversational prompts.

*Distribution of Request and Response Lengths.* In Figure 7, we observe similarities in request distributions and differences in response distributions. Both ChatGPT and Llama-2-13b-chat request distributions adhere to a Zipf distribution [28], characterized by a peak in the frequency of shorter requests. The histogram of ChatGPT's response tokens reveals a bimodal distribution of length variability. In contrast, Llama-2-13b-chat exhibits a Zipf distribution with a higher frequency of longer requests. These distributions provide insights into the conversational dynamics of request and response lengths for the models. ChatGPT handles longer contexts compared to Llama, reflecting differences in model capacities or interactions.

*Response Lengths per Request and Request Bin.* ChatGPT shows a linear correlation between request and response lengths, with a symmetric distribution of response lengths in Figure 8. In contrast, Llama-2-13b-chat decreases response length as request length increases, with higher variability for longer requests due to context length limits.

In Figure 10, ChatGPT tends to produce shorter responses than Llama-2-13b-chat, indicated by a shifted Gaussian distribution with a peak at shorter lengths and a longer spread as request length increases. In contrast, Llama-2-13b-chat shows less variability in response length, peaking at 2048 characters.

## 3.4 System Response Failures

We observe a high failure rate in GPT services. Figure 11 shows that ChatGPT has a higher average service failure rate but better stability compared to GPT-4 in both conversational and API interfaces. Overall, GPT-4 has better service reliability over ChatGPT, demonstrating greater consistency and stability with fewer outliers. Additionally, the average failure rate of the conversation service is consistently high, exceeding 5% for ChatGPT, which is significantly higher than that of regular cloud services [3]. We will discuss the reliability problem in detail in Section 4.
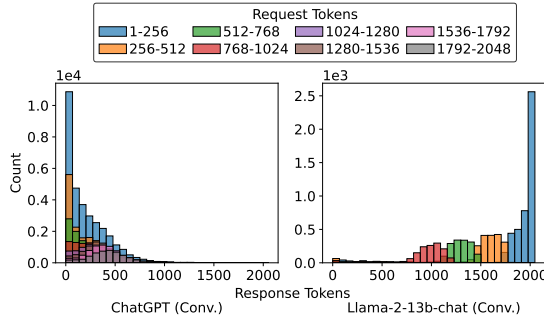


Figure 7: Distribution of Request and Response Tokens.



Figure 8: Response Lengths per Request

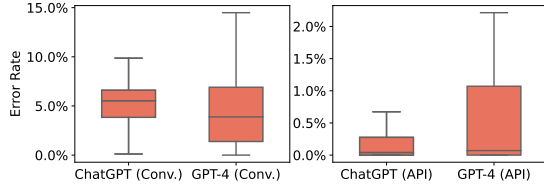**Figure 10: Response Lengths per Request Bin**



**Figure 11: Comparison of Statistical Failure Rates in BurstGPT.**

## 4 BurstGPT-Perf: A Benchmark Suite

This section presents our open-sourced benchmark suite: BurstGPT-Perf to evaluate LLM serving systems under BurstGPT with streaming, stochastic, and bursty workloads. The suite is lightweight and modular, requiring minimal code integration for effective deployment. Additionally, it serves as a user-friendly example of how to implement BurstGPT for evaluation and will be integrated and open-sourced within BurstGPT. With our use case, we advocate for using BurstGPT on more serving systems to explore system limitations and optimization opportunities further.
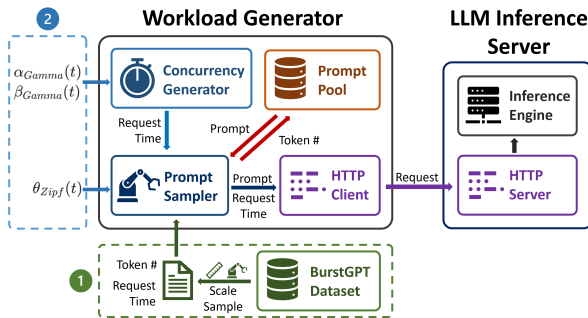


**Figure 12: Overview of BurstGPT-Perf. It generates simulations of BurstGPT in a burst manner with two scaling methods: 1. *RPS Scaling* scales the original BurstGPT data; 2. *Modeled Scaling* uses Gamma distribution parameters to generate request times and Zipf distribution parameter to generate prompt token lengths.**

## 4.1 Implementation

*Workload Generator.* Figure 12 presents an overview of the workload generator. The system comprises four main components: the Prompt Sampler, the Prompt Pool, the Concurrency Generator, and the HTTP client. The HTTP client is included using the efficient and asynchronous aiohttp framework. An index with token numbers as keys and prompt indices as values will be established in the Prompt Pool to mitigate the overhead of prompt sampler queries and a prompt is randomly selected for return. Once the prompt sampler retrieves a prompt from the Prompt Pool to enable workload concurrency, the concurrency generator times requests.

*Evaluation Workflow.* Requests of specified lengths are sampled from the Prompt Pool at stochastic intervals that are generated by Concurrency Generator in the Workload Generator. This concurrency pattern is treated as a time series, which inputs it into the Inference Engine via HTTP for performance evaluation purposes. The benchmark suite systematically adjusts RPS if we use the vanilla trace, or the parameters $\lambda$, $\alpha$, and $\beta$ according to a predetermined sequence, if we use modeled traces, thereby generating average and instantaneous performance metrics for users.

## 4.2 Scaling BurstGPT to Any Scale

In Section 3, we identify specific burstiness patterns from workload traces. To assess LLM serving systems on BurstGPT at any scale, we scale BurstGPT to match the system size, enabling scalable evaluations. We use *RPS Scaling* or *Modeled Scaling* in the evaluations. (1) For RPS scaling, we use a scaling parameter $c$ to multiply the timestamp, making the average RPS in the scaled trace equal to the given RPS. (2) Modeled scaling offers more flexibility in adjusting the scaling parameters. Our analysis in Section 3 reveals that within BurstGPT, concurrency follows a varied Gamma distribution, characterized by the shape parameter $\alpha$ and the scale parameter $\beta$. The parameter $\alpha$ predominantly influences the coefficient of variation (*CV*) of the distribution, mathematically expressed as $CV = 1/\sqrt{\alpha}$.

To make BurstGPT applicable for evaluating LLM serving systems of different sizes, we adjust these parameters in BurstGPT's concurrency distribution model at intervals of 20 minutes. For assessments of BurstGPT on a serving system of a particular scale, users can modify the $\beta$ parameter to determine a warm-up arrival rate compatible with their system's capacity, such as 30% to 40% GPU utilization of KV cache in our setting. Following this adjustment, users can evaluate their system within any chosen time range in BurstGPT.

## 4.3 Matching BurstGPT with Dialogues

For request lengths, in Section 3 and related simulations [14], it is observed that the request length follows a Zipf distribution. The sampler samples request lengths using a variable parameter $\lambda$, influencing the frequency of shorter requests in input streams.

For response lengths in BurstGPT, they are not predefined, allowing the model to determine them dynamically. This strategy avoids restrictions that could degrade system runtime, especially with specific LLMs. Over many requests, the distribution of output lengths will naturally align with the model's behavior, as detailed in Section 3.

Yuxin Wang, Yuhan Chen, Zeyu Li, Xueze Kang, Yuchu Fang, Yeju Zhou, Yang Zheng, Zhenheng Tang, Xin He, Rui Guo, Xin Wang, Qiang Wang, Amelie Chi Zhou, and Xiaowen Chu
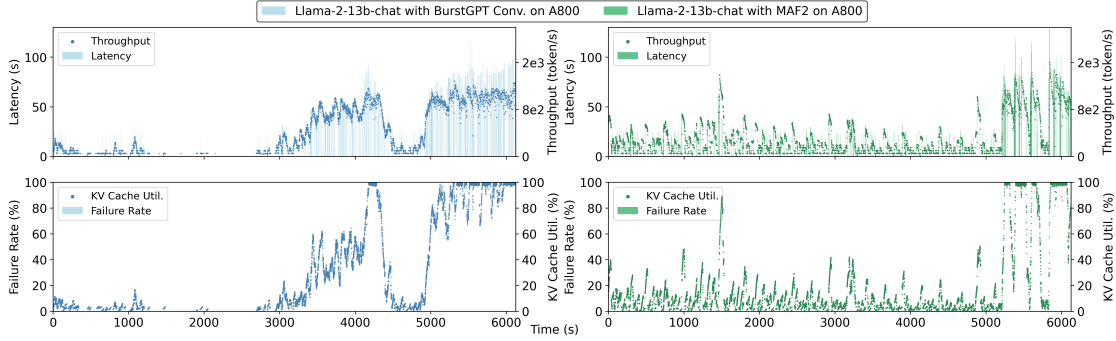


Figure 13: The impact of concurrency of BurstGPT Conversation and BurstGPT API traces on the latency, throughput, KV cache usage, and failure rates of vLLM serving on a single A800 GPU. Using the first 2000 queries for both traces, we scale BurstGPT by a factor of 10 to reach an RPS of 0.326. To reach the same RPS, MAF2 is scaled with a parameter of 1.234, ensuring both are complete in 6126 seconds.

After determining request lengths, the suite initializes a CPU buffer for loading a truncated prompt dataset pool. This dataset can be sampled from existing GPT conversation dialogues, such as Alpaca [13]. The dataset truncation and request length sampling adhere to the predefined Zipf distribution. The buffer size is user-configurable to balance with the CPU's memory capacity. Requests of specific lengths are then selected from the prompt dataset and sent to the server at intervals, where they queue for processing.

## 4.4 Metrics and Setups

In serving systems, Quality of Service (*QoS*) is often measured by *latency* and *throughput* [16]. Latency refers to the completion time of each token, while throughput is the number of tokens processed within a given period. In addition to efficiency, the stability and reliability of LLM serving are particularly important to users, yet these aspects are often overlooked. In this evaluation, we also focus on latency jitters to assess stability and failure rates to assess the reliability of LLM serving.

We categorize metrics into two types: average values (denoted with a superscript $^{avg}$) and instantaneous values (marked with a superscript $^{ins}$). This categorization provides a comprehensive understanding of both long-term trends and immediate system behavior. Based on the trace study, the following key metrics are employed for a thorough evaluation of a serving system's performance:

- $R^{avg}$ / $R^{ins}$: Request Failure Rate;
- $L^{avg}$ / $L^{ins}$: Token Latency;
- $\sigma_L^{avg}$ / $\sigma_L^{ins}$: Token Latency Jitters (Standard Deviation);
- $P^{avg}$ / $P^{ins}$: System Throughput.

## 5 Demo Evaluations

Our evaluation is conducted on Llama-2-13b-chat on an A800 and A6000 GPU server. In experiments that use modeled scaling, we set the initial parameters $\alpha_{Gamma}$ and $\beta_{Gamma}$ of the Gamma distribution in BurstGPT to 0.5 and 2, respectively, for 30% GPU KV Cache utilization. The request length parameter $\theta_{zipf}$ is set to 1.1. We use requests in ShareGPT as the request pool.
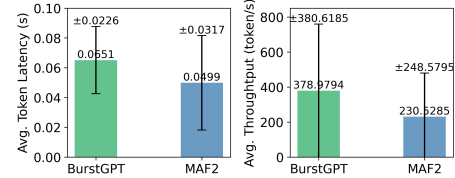


Figure 14: Average token latency and throughtput of vLLM serving Llama-2-13b-chat on an A800 GPU server with BurstGPT and MAF2 workloads.

## 5.1 Evaluation: BurstGPT v.s. MAF2

In this evaluation, we assess the efficiency and stability of vLLM [11] using BurstGPT and MAF2. We set the same RPS to 0.326 to scale both MAF2 and BurstGPT. Both experiments were completed in 6126 seconds. Figure13 shows the serving system's instantaneous metrics under the two workloads. The results indicate that the request burstiness of concurrency in BurstGPT is higher, whereas the request distribution of the MAF2 workload is more uniform. As we can see from Figure 14, the average latency and throughput are higher for BurstGPT compared to MAF2, indicating a larger system resource occupation by BurstGPT. In conclusion, although both methods share the same RPS, the more bursty concurrency in BurstGPT introduces more challenges to the serving system, particularly in terms of KV cache management, compared to MAF2.

## 5.2 Evaluation: Micro Workload Variations

In this experiment, we model the request lengths to a Zipf distribution, as in [14]. The sampler samples the request lengths in Zipf distribution with a variable parameter $\lambda$. We investigate the impact of varying specific variables: $\alpha$, $\beta$, and $\lambda$. We utilized the Llama-2-7b-chat and Llama-2-13b-chat models in our testbed. The results demonstrate a correlation between the adjustments in these variables and the observed system behavior.

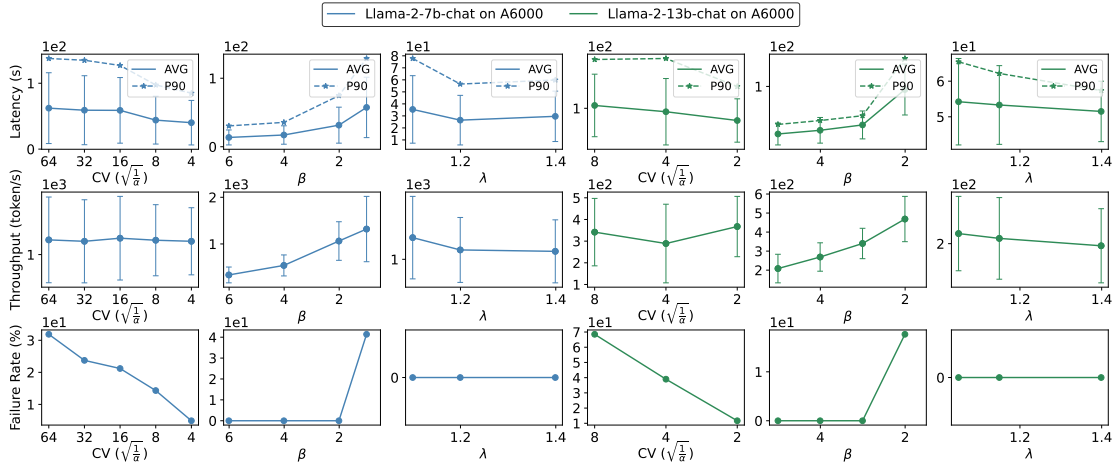We observe that minor changes in the parameter $\alpha$ or the *CV* can lead to sudden request failures, significantly undermining the

**Figure 15: Impact of parameters $\alpha$, $\beta$, and $\lambda$ on average and p90 latency, average throughput, and average failure rates in vLLM. The evaluation is conducted with Llama-2-7b-chat and Llama-2-13b-chat models using A6000 and A800 GPU Servers. Experiment Settings: Larger $\alpha$ or smaller $CV$ and $\beta$ increase system burstiness; higher $\lambda$ correlates with shorter request likelihood. Note: Due to space limits, only results from the A6000 server are shown.**

system's reliability. However, these changes have less effect on system performance metrics such as average latency and throughput. This is because such changes in $\alpha$ are temporary and behave as quadratic functions. In contrast, alterations in $\alpha$ in a linear decline lead to ongoing performance degradation, particularly noticeable during request failures. These findings underscore the value of our trace analysis and provide essential insights for improving system efficiency in similar computational environments.

## 5.3 Demo Use: Request Scheduler Selection

We randomly sample a time range on `BurstGPT` and build two subsets with conversation workloads and API workloads and then used Modeled Scaling to derive the corresponding parameter sequences for generating workloads. We use these two generated workloads separately to evaluate three scheduling strategies i.e. fisrt-come-first-serve (FCFS), short-request-first (SRF), and long-request-fisrt (LRF), on vLLM.

In Figure 16, the efficiency of the various scheduling methods may shift across different types of real-world workloads. Comparing FCFS with SRF and LRF on average values in terms of efficiency and stability, we can find that the optimization on conversation service when LRF is superior to FCFS in latency and and latency jitters does not stand in API service. This interesting finding inspires us that even using real-world workloads, an optimization on one workload does not necessarily generalize to other workloads.

## 5.4 Demo Use: Workload Provisioning

To adjust the resources allocated to services in a timely manner, cloud systems can predict system load to proactively scale resources up or down. The workload provisioning task is a time series forecasting problem. Specifically, we select the request count per time bin and the average number of tokens per request as the prediction targets. To construct features for the prediction model, we use historical load data from both short-term and long-term perspectives.

For each time point, we incorporate load values from the past 3 time points (lag window) to capture short-term trends, along with statistical features (e.g., mean, variance) computed over a rolling window of up to 60 minutes to reflect long-term patterns. As shown in 3.1, the day of the week and the hour of the day also impact the workload, so these two features are included. The prediction target is the actual load value for that time point. We employ XGBoost [5] as the prediction model due to its simplicity and effectiveness. We compute the normalized mean absolute error (NMAE) and normalized mean square error (NMSE) between the prediction and the ground truth. The results are shown in Figure 17.

As shown, even without careful tuning, the models can predict future system load easily and accurately. The predicted load trends generally align well with the actual load, indicating significant opportunities for fine-grained load prediction modeling in real-world inference services, which lead to more precise service delivery.

We also investigate the performance of the prediction algorithm under different time intervals. As shown in Figure 18, predicting long-term patterns (10 minutes) is generally easier than predicting short-term patterns (1 minute), with the former achieving a 58% lower NMSE compared to the latter. This suggests that when using load prediction for serving, it is essential to balance the granularity of scheduling with the accuracy of predictions.

## 6 BurstGPT in Industry

We provide an example of how we use `BurstGPT` for industrial use in Huawei during our development of *PD* disaggregation systems. As a common practice in the industry, given a fixed serving system scale, developers use goodput [33] to evaluate the system's capacity, i.e., the maximum request rate that can be served while meeting the SLO goal (e.g., 90%) for each GPU. Recently, disaggregating the prefilling and decoding processes during LLM inference has become the standard approach in LLM service deployment [10, 20, 33]. However, scheduling the disaggregation of prefilling and decoding can be challenging. The prefilling instances can be seen as the

Yuxin Wang, Yuhan Chen, Zeyu Li, Xueze Kang, Yuchu Fang, Yeju Zhou, Yang Zheng, Zhenheng Tang, Xin He, Rui Guo, Xin Wang, Qiang Wang, Amelie Chi Zhou, and Xiaowen Chu
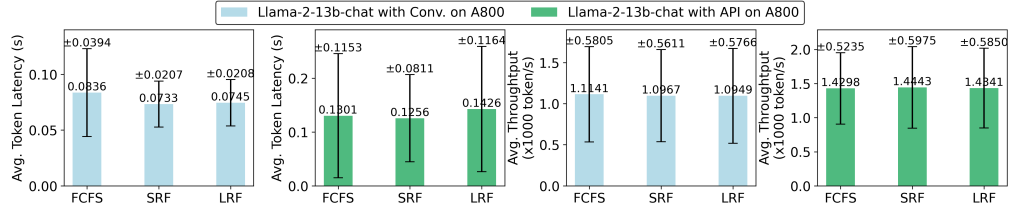
**Figure 16: Average token latency and throughput of FCFS, SRF, and LRF on vLLM serving Llama-2-13b-chat on an A800 GPU server with `BurstGPT` Conversations and API workloads.**
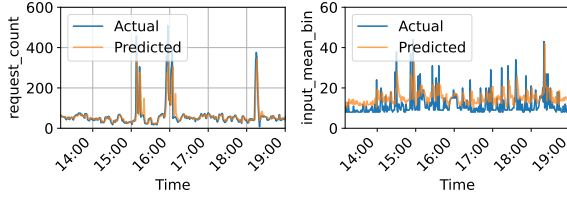


**Figure 17: Request count and mean request token number predictions. The figure only shows a random interval of 6 hours. The time bin interval is 1 minute. Left: Request count predictions, the NMAE is 0.73 and the NMSE is 0.48. Right: Mean request token number predictions, the NMAE is 0.72 and the NMSE is 0.67.**
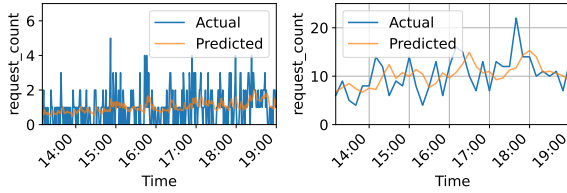


**Figure 18: Request count predictions. The figure only shows a random interval of 6 hours. Left: Request count predictions of Conv. with time bin interval 1 minute, the NMAE is 0.66 and the NMSE is 0.50. Right: Request count predictions of Conv. with time bin interval 10 minutes, the NMAE is 0.32 and the NMSE is 0.21.**

"Provider," which generates the first KV cache and token, while the decoding instances act as the "Consumer," accepting the KV cache and token to continue the auto-regressive generation. Thus, for a given system scale (i.e., the number of prefilling ($P$) and decoding ($D$) instances), the concurrency of requests, request lengths, and response lengths will determine the goodput of the system. As workloads vary rapidly in real-world serving, a fixed $PD$ ratio is not always ideal. Instead, a dynamic $PD$ ratio should be used.

Figure 20 and Figure 19 show examples of goodput at different PD ratios. Since we lack preliminary workload traces before the real-world deployment of our serving system, we use `BurstGPT` in our simulation platform to determine the optimal PD ratio under various user, model, and system workloads. The result is used to develop prototypes of the $PD$ disaggregation serving system at

Huawei. In the simulation, we select the $P$ and $D$ instance ratio, then run simulations, gradually increasing QPS while monitoring the delay SLO achievement rate until the SLO is no longer met. The QPS at this point is considered the system's capacity. We then fix the total number of $P$ and $D$ instances. First, considering a simple situation, we adjust the $PD$ ratio one time to fit the workload divergence in Figure 19, we observe that at 6:00, the $PD$ ratio switches from 2:6 to 6:2, improving the goodput of the serving system. This simulation serves as a solid prototype for real-world $PD$ disaggregation scheduling policies. Then, we develop a dynamic $PD$ ratio strategy. We use beam search to identify the best $PD$ ratio given the provisioned workload in BurstGPT. Instead of directly scaling the number of instances, we simulate scaling by adjusting the load size per instance, assuming a predefined instance launch time in seconds. As shown in Figure 20, with the provisioned workload, the system searches for the optimal $PD$ ratio to improve the goodput.
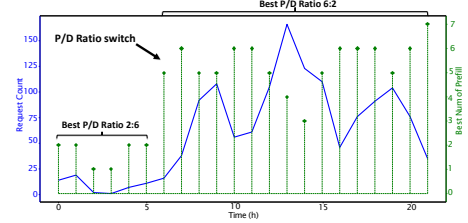


**Figure 19: At 6:00, the PD ratio switches from 2:6 to 6:2, improving the goodput of the serving system.**
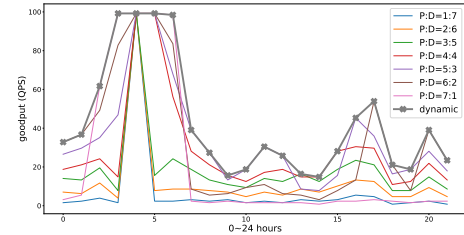


**Figure 20: Given 8 instances, comparison of different PD ratio and dynamic PD ratio based on workload provisioning and beam search during LLM serving.**

These findings highlight the significant improvements in LLM serving achieved through data-driven modeling of user and model

behavior. By carefully designing the scheduling algorithm, we can greatly enhance both system efficiency and QoS.

## 7 Broader Impacts and Conclusions

This research emphasizes incorporating real-world workload data to enhance LLM serving systems. It identifies a critical gap: the need for such data in optimizing these systems. In this work, we introduce the real-world LLM workload `BurstGPT` from Azure OpenAI GPT service. The traces and evaluations reveal possible performance degradation of serving systems under real-world workloads, highlighting challenges and opportunities in workload-aware LLM serving optimizations.

We encourage using `BurstGPT` in optimizing and evaluating serving systems to improve the efficiency, stability, and reliability of LLM services under realistic workloads. We also advocate for data-driven methodologies in developing LLM systems in the future.

## 8 Acknowledgement

## References

[1] [n. d.]. Azure OpenAI Service. https://azure.microsoft.com/en-us/products/ai-services/openai-service/

[2] 2023. Lightllm: A python-based large language model inference and serving framework. https://github.com/ModelTC/lightllm

[3] Muhammad Ahsan, Sacheendra Talluri, and Alexandru Iosup. 2023. Failure Analysis of Big Cloud Service Providers Prior to and During Covid-19 Period. arXiv:2210.08006 [cs.DC]

[4] Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. 2023. GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints. *arXiv preprint arXiv:2305.13245* (2023).

[5] Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*. 785–794.

[6] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. Llm. int8 (): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339* (2022).

[7] Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. 2023. Model Tells You What to Discard: Adaptive KV Cache Compression for LLMs. *arXiv preprint arXiv:2310.01801* (2023).

[8] Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. The False Promise of Imitating Proprietary LLMs. arXiv:2305.15717 [cs.CL]

[9] Ke Hong, Guohao Dai, Jiaming Xu, Qiuli Mao, Xiuhong Li, Jun Liu, Kangdi Chen, Hanyu Dong, and Yu Wang. 2023. FlashDecoding++: Faster Large Language Model Inference on GPUs. *arXiv preprint arXiv:2311.01282* (2023).

[10] Cunchen Hu, Heyang Huang, Junhao Hu, Jiang Xu, Xusheng Chen, Tao Xie, Chenxi Wang, Sa Wang, Yungang Bao, Ninghui Sun, and Yizhou Shan. 2024. MemServe: Context Caching for Disaggregated LLM Serving with Elastic Memory Pool. arXiv:2406.17565 [cs.DC] https://arxiv.org/abs/2406.17565

[11] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Yu, Joey Gonzalez, Hao Zhang, and Ion Stoica. 2023. vLLM: Easy, Fast, and Cheap LLM Serving with PagedAttention. https://vllm.ai/

[12] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient

[13] Zhihui Li, Max Gronke, and Charles Steidel. 2023. ALPACA: A New Semi-Analytic Model for Metal Absorption Lines Emerging from Clumpy Galactic Environments. arXiv:2306.11089 [astro-ph.GA]

[14] Zhuohan Li, Lianmin Zheng, Yinmin Zhong, Vincent Liu, Ying Sheng, Xin Jin, Yanping Huang, Zhifeng Chen, Hao Zhang, Joseph E Gonzalez, et al. 2023. AlpaServe: Statistical Multiplexing with Model Parallelism for Deep Learning Serving. *arXiv preprint arXiv:2302.11665* (2023).

[15] Xupeng Miao and Gabriele Oliaro et.al. 2023. SpecInfer: Accelerating Generative Large Language Model Serving with Speculative Inference and Token Tree Verification. arXiv:2305.09781 [cs.CL]

[16] Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Hongyi Jin, Tianqi Chen, and Zhihao Jia. 2023. Towards Efficient Generative Large Language Model Serving: A Survey from Algorithms to Systems. arXiv:2312.15234 [cs.LG]

[17] Oh Neal, Vaidya anf Fred and Comly Nick. 2023. Optimizing inference on large language models with nvidia tensorrt-llm. https://github.com/NVIDIA/TensorRT-LLM

[18] OpenAI. 2022. Introducing ChatGPT. https://openai.com/blog/chatgpt

[19] OpenAI. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774* (2023).

[20] Ruoyu Qin, Zheming Li, Weiran He, Mingxing Zhang, Yongwei Wu, Weimin Zheng, and Xinran Xu. 2024. Mooncake: A KVCache-centric Disaggregated Architecture for LLM Serving. arXiv:2407.00079 [cs.DC] https://arxiv.org/abs/2407.00079

[21] Jeff Rasley and Rajbhandari et.al. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 3505–3506.

[22] Mohammad Shahrad, Rodrigo Fonseca, Íñigo Goiri, Gohar Chaudhry, Paul Batum, Jason Cooke, and Laureano et.al. [n. d.]. Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider. ([n. d.]).

[23] Noam Shazeer. 2019. Fast transformer decoding: One write-head is all you need. *arXiv preprint arXiv:1911.02150* (2019).

[24] Zhenheng Tang, Yonggang Zhang, Peijie Dong, Yiu ming Cheung, Amelie Chi Zhou, Bo Han, and Xiaowen Chu. 2024. FuseFL: One-Shot Federated Learning through the Lens of Causality with Progressive Model Fusion. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. https://openreview.net/forum?id=E7fZOoiEKl

[25] Gemini Team, Machel Reid, Nikolay Savinov, and etc. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv:2403.05530 [cs.CL]

[26] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. arXiv:2302.13971 [cs.CL]

[27] Dao Tri, Haziza Daniel, Massa Francisco, and Sizov Grigory. 2023. Flash-decoding for long-context inference. https://crfm.stanford.edu/2023/10/12/flashdecoding.html

[28] 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]

[29] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*. PMLR, 38087–38099.

[30] 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/yu

[31] Yanqi Zhang, Íñigo Goiri, Gohar Irfan Chaudhry, Rodrigo Fonseca, Sameh Elnikety, Christina Delimitrou, and Ricardo Bianchini. 2021. Faster and Cheaper Serverless Computing on Harvested Resources. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*. doi:10.1145/3477132.3483580

[32] Zangwei Zheng, Xiaozhe Ren, Fuzhao Xue, Yang Luo, Xin Jiang, and Yang You. 2023. Response Length Perception and Sequence Scheduling: An LLM-Empowered LLM Inference Pipeline. arXiv:2305.13144 [cs.CL]

[33] Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu, Yibo Zhu, Xuanzhe Liu, Xin Jin, and Hao Zhang. 2024. DistServe: disaggregating prefill and decoding for goodput-optimized large language model serving. In *Proceedings of the 18th USENIX Conference on Operating Systems Design and Implementation* (Santa Clara, CA, USA) *(OSDI'24)*. USENIX Association, USA, Article 11, 18 pages.

[34] Jiahang Zhou, Yanyu Chen, Zicong Hong, Wuhui Chen, Yue Yu, Tao Zhang, Hui Wang, Chuanfu Zhang, and Zibin Zheng. 2024. Training and Serving System of Foundation Models: A Comprehensive Survey. arXiv:2401.02643 [cs.AI]

memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*. 611–626.