



MLSys 2025

The Eighth Annual Conference on Machine Learning and Systems

SOLA: Optimizing SLO Attainment for Large Language Model Serving with State-Aware Scheduling

Ke Hong¹², Xiuhong Li^{2*}, Lufang Chen², Qiuli Mao², Xuefei Ning¹, Guohao Dai^{23*}, Shengen Yan¹², Yun Liang⁴, Yu Wang^{1*} ¹Tsinghua University, ²Infinigence-Al ³Shanghai Jiao Tong University, ⁴Peking University *Corresponding Authors





Background



Model as a Service (MaaS) and Service-Level Objective (SLO)



[1] C. Zhang, et al. MArk: Exploiting Cloud Services for Cost-Effective, SLO-Aware Machine Learning Inference Serving, ATC'19.
 [2] M. Haque, et al. Few-to-Many: Incremental Parallelism for Reducing Tail Latency in Interactive Services, ASPLOS'15.
 2025/5/11 NICS-efc Lab





The Goal of Service Providers: Maximize RPS Adhering to SLO







LLM Serving: Two Distinct Phases - Prefill and Decode



*Here LLM refers to the mainstream LLMs with Decoder-only Transformer architecture







LLM Serving: Two Distinct Phases - Prefill and Decode



[1]K. Hong, et al. FlashDecoding++: Faster Large Language Model Inference with Asynchronization, Flat GEMM Optimization, and Heuristics, MLSys'24.



Overview



What We Do: Optimize SLO Attainment in Large Language Model (LLM) Serving





Motivation



The bias and the variance of the latency distribution



BIAS

one of TTFT and TPOT has much lower attainment against its SLO

VARIANCE

many requests fail the SLOs while others still having redundant budgets to spare



Motivation



Two-fold trade-offs exist in request-level latency

BIAS

one of TTFT and TPOT has much lower attainment against its SLO



Between TTFT and TPOT

VARIANCE

many requests fail the SLOs while others still having redundant budgets to spare





Motivation



Two-fold trade-offs exist in request-level latency









SOLA: State-Aware Scheduling







SOLA: State-Aware Scheduling



State-aware State update Strategy optimization scheduler ((i-1)-th iteration) (i-th iteration) i-th iteration (i-1)-th iteration Iterative-level **Formulation and** Q_i^{wait} Q_{i+1}^{wait} Q_i^{run} Q_{i-1}^{run} scheduling Schedule Schedule Execution Execution **Design Space** optimization Iteration-level state update Iteration-level strategy optimization ④Optimize 1 Update Request Constrained **SLOs** State optimization **Real-time** Aware of both 2 Statistics Scheduling LLM/ Cost System system and **Optimization** Parallelism State model 3 Tune request states based on States State Monitor Strategy Generator Deployment (Sec. 4.4) (Sec. 4.3) Setting

2025/5/11







SOLA: State-Aware Scheduling State-aware State update Strategy optimization (i-th iteration) scheduler ((i-1)-th iteration) (i-1)-th iteration *i*-th iteration Iterative-level **Formulation and** Q_i^{wait} Q_{i+1}^{wait} Q_i^{run} $Q_{i-1}^{\operatorname{run}}$ scheduling Schedule Schedule Execution Execution **Design Space** optimization Iteration-level state update Iteration-level strategy optimization 1 Update ④Optimize Request Constrained SLOs State optimization 2 Statistics Cost LLM/ System Parallelism State model 3 Tune State Monitor Strategy Generator Deployment (Sec. 4.3) (Sec. 4.4) Setting

Modeling and Design Space



Algorithm 1 Scheduling at the *i*-th iteration.

Input:
$$M$$
, Q_i^{wait} , Q_i^{run} , m_i^{ratio} , n_i , k_i , \mathcal{F}_i
1: Initialize $Q_i^{\text{run}} \leftarrow \emptyset$
2: if get_req_num $(Q_i^{\text{wait}}) = 0$ then
3: Return $(Q_i^{\text{wait}}, Q_i^{\text{run}})$
4: end if
5: $Q_i^{\text{wait}} \leftarrow \mathcal{F}_i(Q_i^{\text{wait}})$
6: for $r \in Q_i^{\text{wait}}$ do
7: $m^{\text{peak}} \leftarrow \text{cal_peak_mem}(Q_i^{\text{run}} \cup r)$
8: if $m^{\text{peak}} \ge M \times (1 - m_i^{\text{ratio}})$ then
9: continue
10: end if
11: $Q_i^{\text{run}} \leftarrow Q_i^{\text{run}}$.push_req_with_token_num $(r, k_{i,r}^{\text{new}})$
12: $Q_i^{\text{wait}} \leftarrow Q_i^{\text{wait}}$.pop_req_with_token_num $(r, k_{i,r}^{\text{new}})$
13: if get_req_num $(Q_i^{\text{run}}) \ge n_i$ then
14: break
15: end if
16: if get_token_num $(Q_i^{\text{run}}) \ge k_i$ then
17: break
18: end if
19: end for
20: Return $(Q_i^{\text{wait}}, Q_i^{\text{run}})$

Design space of scheduling

Execution Order Which requests are prioritized

X

Workload Size

How many requests (or tokens) to run in this iteration







SOLA: State-Aware Scheduling



Constrained Optimization



Problem: Constrained Optimization

Solution: Order and Workload Control

Condition (or)	Constrained Optimization
$(1) 1 > p_i^{\text{TPOT}} > p_i^{\text{TTFT}}$ $(2) p_i^{\text{TPOT}} > 1 > p_i^{\text{TTFT}}$	$\frac{\min\max_{r}(t_{i,r}^{\text{TPOT}})}{s.t.\max_{r}(t_{i,r}^{\text{TTFT}}) \leq T^{\text{TTFT}}}$
(1) $1 > p_i^{\text{TTFT}} > p_i^{\text{TPOT}}$ (2) $p_i^{\text{TTFT}} > 1 > p_i^{\text{TPOT}}$	$\begin{vmatrix} \min \max_{r} (t_{i,r}^{\text{TTFT}}) \\ s.t. \max_{r} (t_{i,r}^{\text{TPOT}}) \leq T^{\text{TPOT}} \end{vmatrix}$
$p_i^{\text{TTFT/TPOT}}$: the ratio between maximum	

real-time TTFT/TPOT and its SLO

- Objective: the less fulfilled one of TTFT and TPOT SLOs, "doing the best"
- Constraint: the other one of TTFT and TPOT SLOs, "ensure the constraint is met"









Order Control

Workload Control



Optimize TPOT subject to TTFT







Results1: Higher SLO Attainment

- **Baseline:** vLLM-S (vLLM with SplitFuse strategy), vLLM-D (vLLM with default strategy), SJF (Shortest-Job-First strategy)
- Testbed: NVIDIA A100 GPUs
- Benchmark: ShareGPT and LongBench datasets



(a) Llama3-8B, ShareGPT. (b) Qwen1.5-14B, ShareGPT. (c) Llama3-70B, ShareGPT. (d) Vicuna-7B, LongBench. (e) Qwen1.5-72B, LongBench.



Other Results





- SLO setting: 10×/15× of the single request latency to be the tight/loose SLOs for 7B-14B models, and 5×/10× of the single request latency to be the tight/loose SLOs for 70B+ models
- More compact distribution with the center close to linearly scaled SLO line

TTFT [s]

4

(b) Qwen1.5-72B, LongBench, 0.35 request/s.

345

TTFT [s]

0

2345

TTFT [s]

2345

TTFT [s]

01







The goal of this work:

Higher SLO attainment, higher RPS, lower costs

The motivation of this work:

We can exploit two-fold trade-offs (between TTFT and TPOT, among different requests) to optimize SLO attainment in LLM serving

What we do in this work:

1. The formulation of iterative-level scheduling and design space

- 2. State-aware scheduling that is aware of both system and request states
 - **3.** Constrained optimization on the scheduling strategy





SOLA: Optimizing SLO Attainment for Large Language Model Serving



Thank you, Q&A

Ke Hong, hk24@mails.tsinghua.edu.cn