

# Prefix and output-length aware scheduling for efficient online LLM inference



Iñaki Arango,<sup>1,2,3</sup> Ayush Noori,<sup>1,2,3</sup> Yepeng Huang,<sup>3</sup> Rana Shahout,<sup>2</sup> Minlan Yu <sup>3</sup> <sup>1</sup> Harvard College; <sup>2</sup> Harvard John A. Paulson School of Engineering and Applied Sciences; <sup>3</sup> Harvard Medical School

# Motivation

LLM inference in large data centers can benefit from data parallelization, where models are replicated across GPU devices that can serve requests in parallel. How should we assign requests to GPU workers? Instead of dividing requests evenly and randomly, in real-world applications, requests exhibit patterns that can be exploited to improve performance. These include:



# **Benchmarking PREBLE**

The first approach to benefit from prompt sharing under a distributed LLM serving system with data parallelism across multiple GPUs was PREBLE (Srivatsa *et al.*, 2024). However, PREBLE sets the expected decode output length equal to the average output length of requests during scheduling. Therefore, PREBLE may induce imbalanced workloads on different GPUs. Here, we extend PREBLE **by integrating prefix-aware scheduling with output length-aware scheduling** of S<sup>3</sup> (Jin *et al.*, 2024). We extend PREBLE's simulator of LLM inference to benchmark PREBLE vs. baseline prefix-unaware schedulers and identify opportunities for improvement. We find that:

- At high request rates, PREBLE is outperformed by several prefix-unaware schedulers.
- 2 This effect is exacerbated by an increase in the number of GPUs available for inference.
  - PREBLE suffers from overhead introduced by its E2 scheduler, which scales with both # of GPUs and request rate.
  - ) Decode length heterogeneity worsens PREBLE performance.

### We evaluated and benchmarked PREBLE across...

1

3



## Adding output-length aware scheduling to PREBLE



Having carefully characterized the scalability challenges associated with PREBLE, we sought to improve its performance by leveraging both prefix-aware and output length-aware scheduling. We build on the E2 scheduler of PREBLE by considering prefix sharing, fairness, and output length. As a proof-of-concept, we use a perfect oracle of true output length. However, decode length **can also be predicted**.

To incorporate output length, we modify the global prefix tree of the E2 scheduler. We evaluate this modified version of PREBLE on a dataset with high variance in token lengths, created using our **ReAct-based variance-customizable dataset generator**.

When using output length for per-GPU load calculation, we improve the performance of PREBLE in high-demand settings, with 14.31% reduced latency at 64 RPS (0.1223 vs. 0.1427) and 28.89% reduced latency at 128 RPS (0.1820 vs. 0.2559).

### **Predicting decode length**

We trained a lightweight 6-layer BERT-based all-MiniLM-L6-v2 language model to predict decode length on the Alpaca-52K dataset, achieving 4.8× performance over random.

Metric	Score	Metric	Score
Accuracy	0.24	Recall	0.24
F <sub>1</sub> score	0.22	AUROC	0.85
Precision	0.23	MCC	0.20



#### Generating a variance-maximizing benchmark

To understand how decode length heterogeneity impacts PREBLE, we generated artificial benchmarking datasets from ReAct whose output lengths can follow any discrete distribution, then used this to create a benchmark with high variance in token lengths (1 vs. 300 tokens).

🔭 <u>inaki.io</u> <u>www.ayushnoori.com</u>

☑ {inakiarango, anoori}@college.harvard.edu

R.S. and M.Y. are partially supported by NSF CNS NeTS 2107078. This work was supported in part by ACE, one of the seven centers in JUMP 2.0, a Semiconductor Research Corporation (SRC) program sponsored by DARPA.

