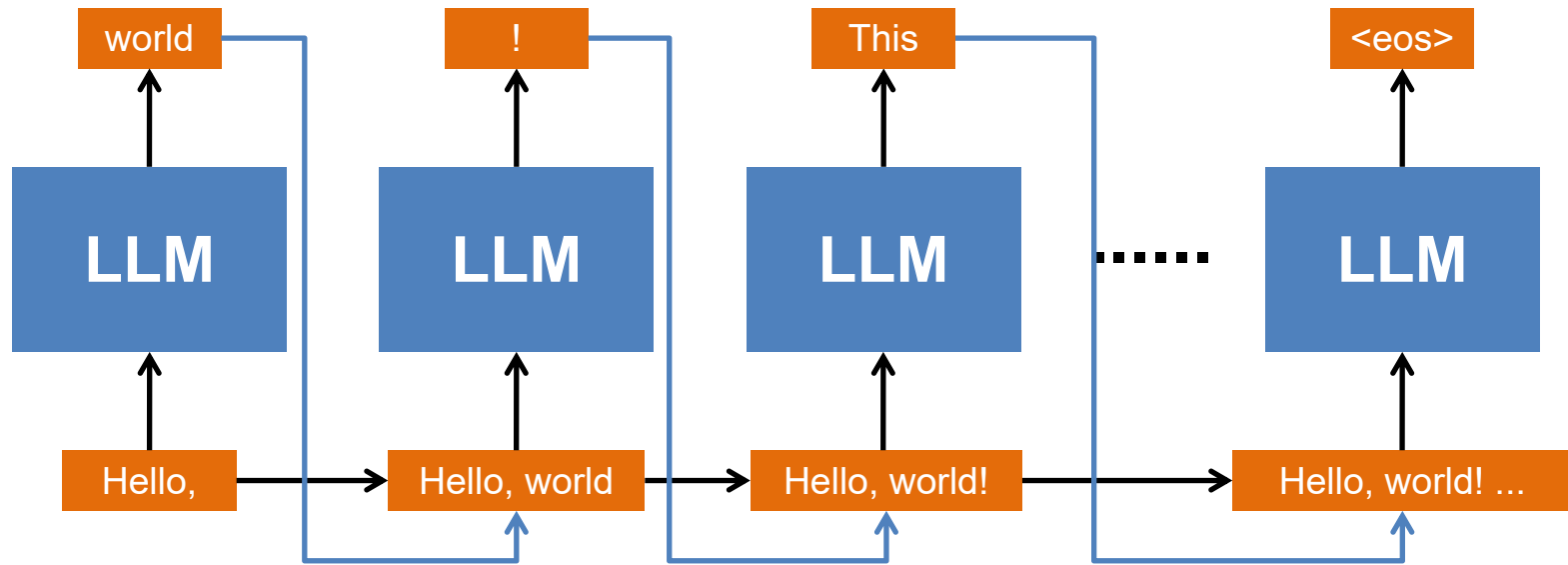


# Efficient LLM Scheduling By Learning To Rank

Yichao Fu, Siqi Zhu, Runlong Su, Aurick Qiao, Ion  
Stoica, Hao Zhang

# LLM Decoding Process

---



The output length of the LLM is unpredictable due to the autoregressive decoding process.

# Head-of-Line Blocking in LLM Serving

---

Requests:  $R_0$  10 tokens |  $R_1$  2 tokens |  $R_2$  1 token

System throughput: 1 token/s

FCFS: 

Latency:  $R_0$  1s/tok |  $R_1$  6s/tok |  $R_2$  13s/tok



SRTF:  *Time* →

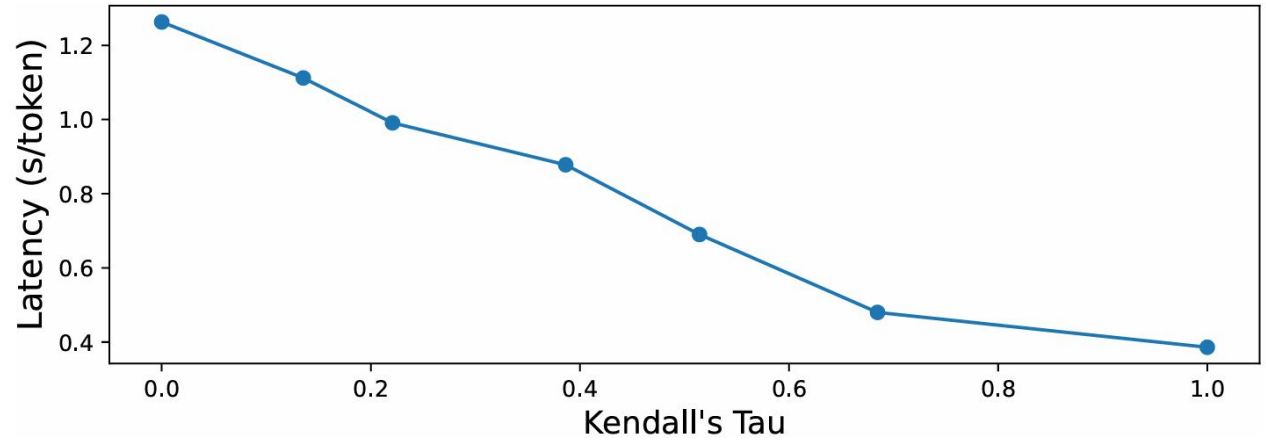
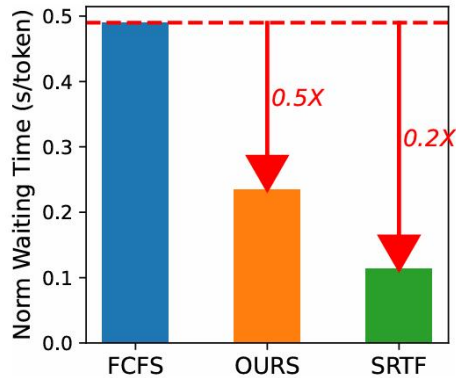
Latency:  $R_0$  1.3s /tok |  $R_1$  1.5s/tok |  $R_2$  1s/tok



LLM services that implement a first-come-first-serve (FCFS) scheduling inevitably face significant Head-Of-Line (HOL) blocking.

# Ranking is All You Need (to appx. SJF)

---



The precise generation length is not needed.

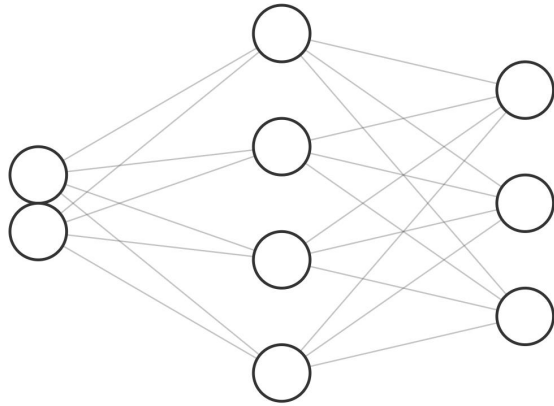
An accurate generation length ranking prediction is enough.

# Learning To Rank

---



Data



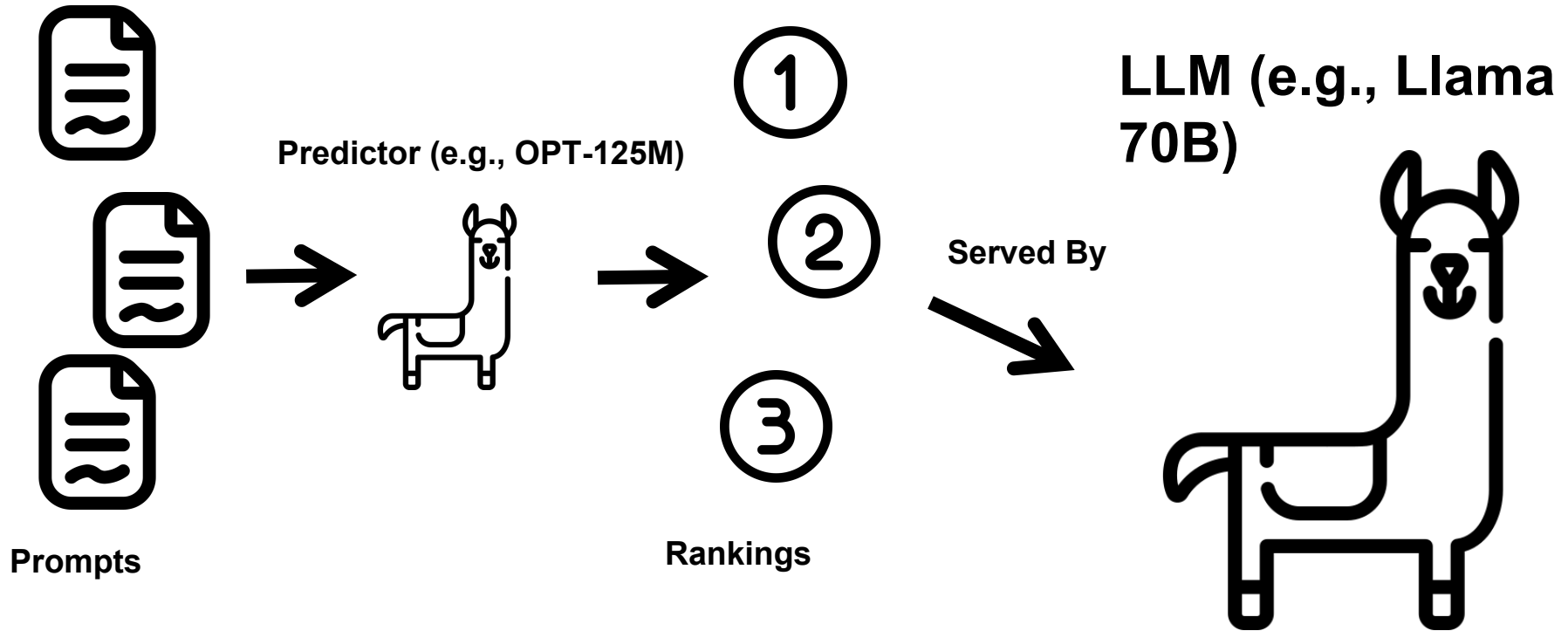
Neural Network



Rankings

*Learning to rank* (LTR) applies machine learning methods to ranking supervised data.

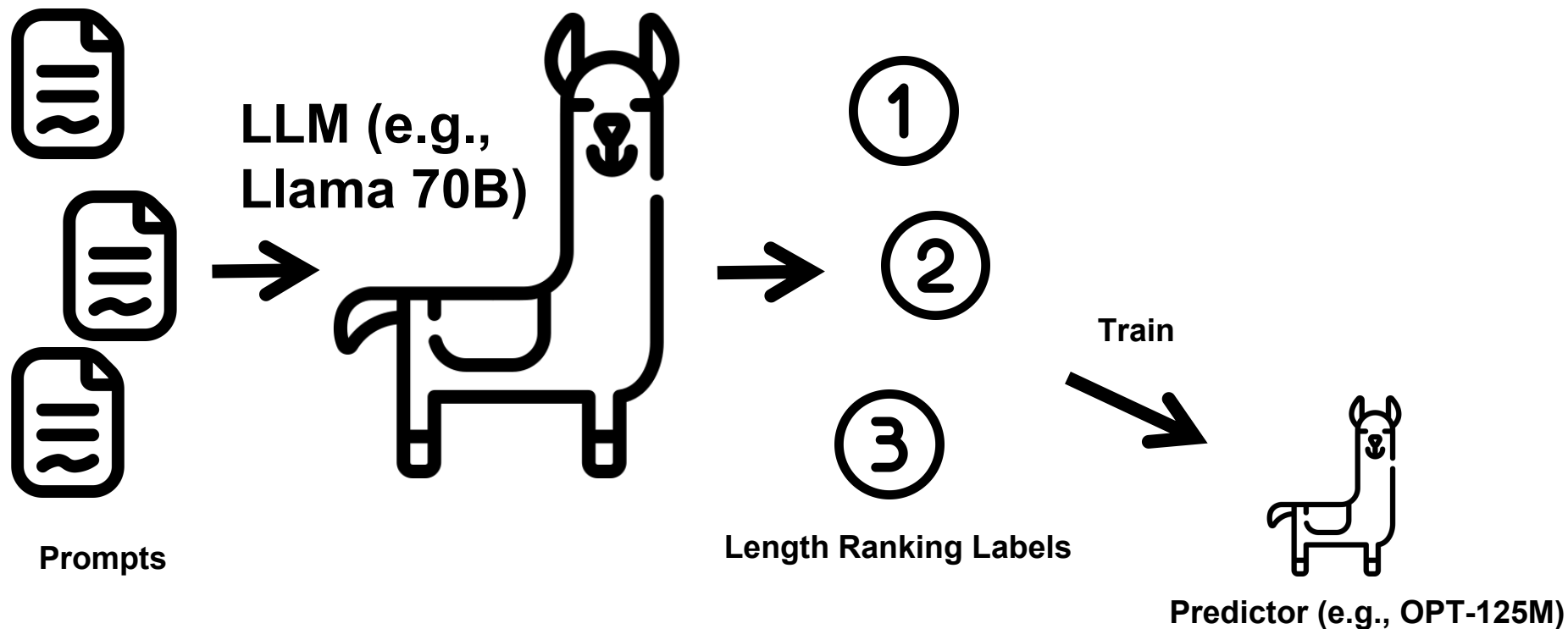
# Generation Length Ranking Predictor



The predictor estimates the relative length of responses for incoming prompts, allowing them to be efficiently ordered before being processed by the target language model.

# Train Length Ranking Predictor

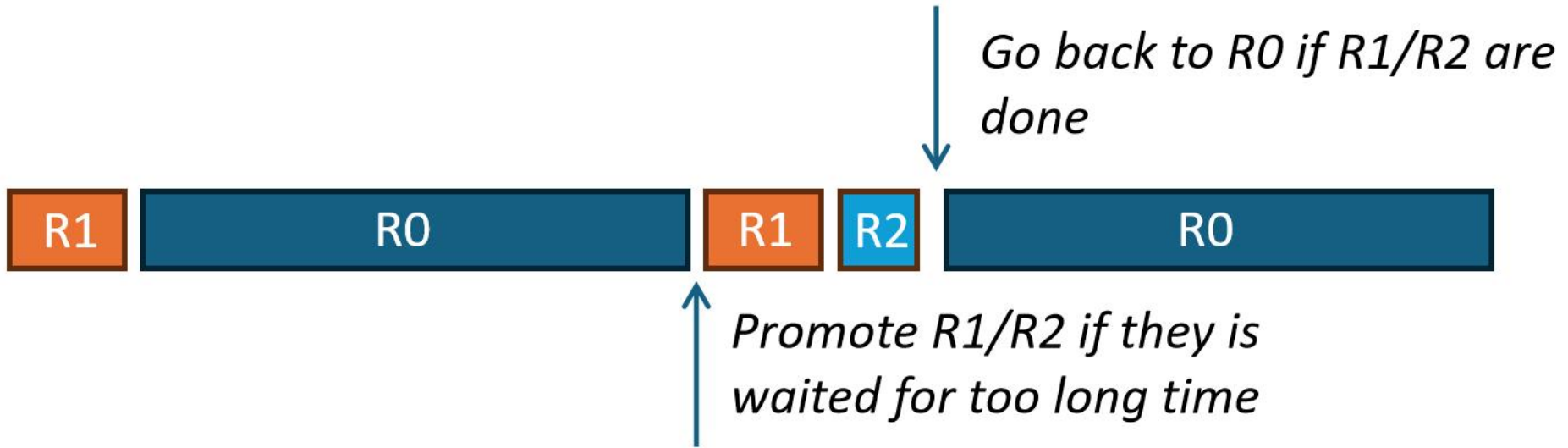
---



The predictor can be rapidly trained on live production data within minutes (e.g., 5 minutes) during actual LLM deployment.

# Starvation Prevention

---



$$\text{max\_waiting\_time} = \max(TTFT, \max(TPOT)).$$



# Evaluation

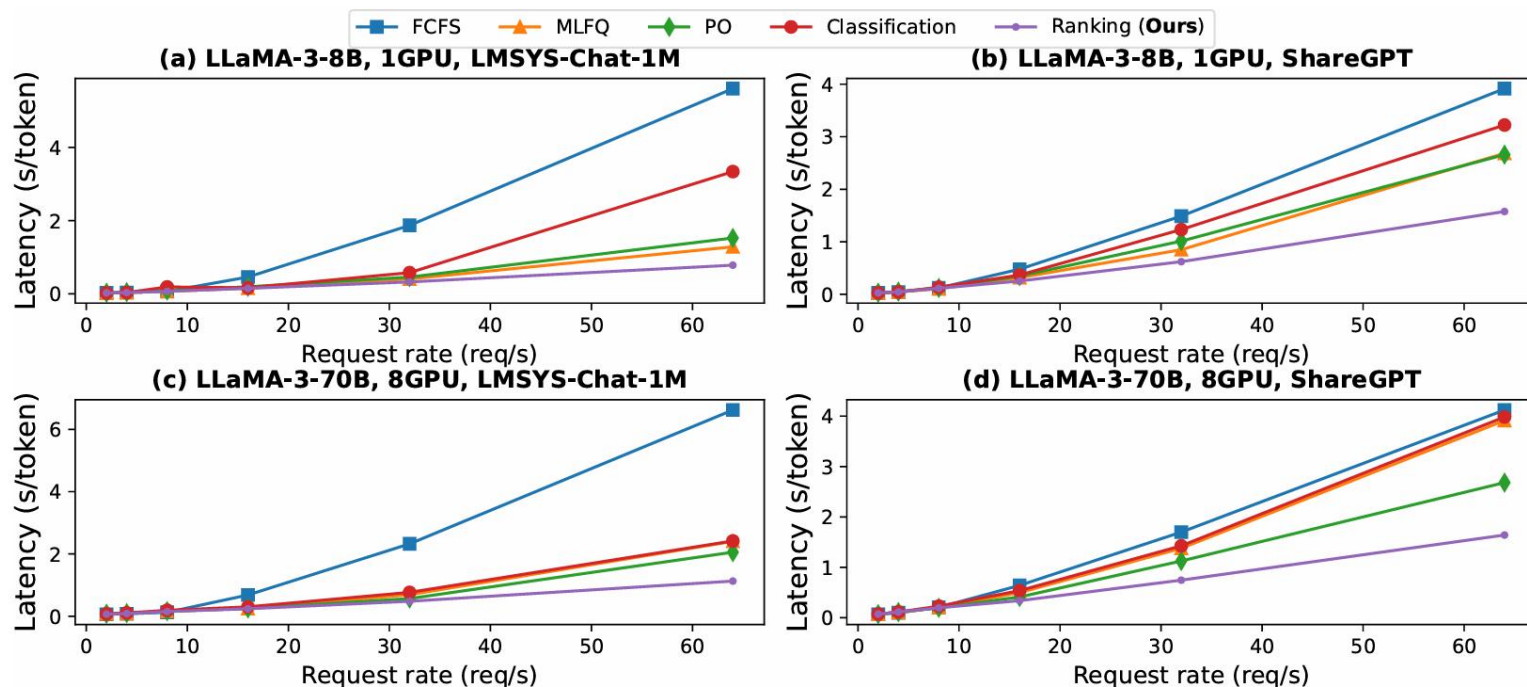


Figure 3: Mean latency of different schedulers with Llama-3 models on real workloads.

Our proposed method improve the mean latency by up to 6.9 $\times$  compared with FCFS and from 1.5 $\times$ –1.9 $\times$  compared with PO in Chatbot Serving.

---

**Thank You!**