# GABRIELE OLIARO

## EDUCATION

Carnegie Mellon University

Ph.D. in Computer Science. Advisor: Zhihao JiaTsinghua UniversityM.S in Advanced Computing. Advisor: Jidong Zhai

Harvard University B.S. in Electrical Engineering. Advisor: Minlan Yu

## EXPERIENCE

### **CMU** Catalyst Group

PhD Researcher with Prof. Zhihao Jia

- (SpecInfer project): Co-design ML and systems optimizations to reduce LLM inference latency by  $1.5 \times -3.5 \times$ . Build multi-node, multi-GPU distributed runtime with efficient CUDA kernels. Lead development of open-source FlexFlow Serve (1.8K stars on GitHub), deploying the SpecInfer algorithm. **Prints author ASPLOS** '24 paper, cited 250+ times.
- (FlexLLM project): Designed FlexLLM, the first system to co-serve LLM inference and parameter-efficient finetuning, leveraging token-level finetuning, dependent parallelization, and graph pruning to reduce GPU memory overhead by up to 8× and preserve over 80% finetuning throughput under heavy inference loads.
- (AdaServe project): Co-developed AdaServe, a novel LLM serving system that enables fine-grained SLO-customized inference through speculative decoding, achieving up to 73% higher SLO attainment and 74% higher goodput compared to state-of-the-art systems.
- (QST project): Mentor the design of a quantization-based LLM fine-tuning technique requiring up to 2.3× less GPU memory and up to 3× lower fine-tuning time compared to PEFT. ♥ Outstanding Paper Award at ACL '24

### Snowflake AI Research

Research Intern

• Designed and implemented SuffixDecoding, a model-free speculative decoding method that accelerates LLM inference using suffix trees to predict token sequences, achieving up to  $2.9 \times$  higher throughput and  $3 \times$  lower latency compared to state-of-the-art methods.

#### Whist Technologies, Inc

Software Engineer

• Implemented client-server profiling suite for a streaming application. Built scalable back-end components in Go to efficiently and securely handle user data in the cloud. Profiled and improved performance of low-latency network streaming protocol in C/C++.

## Systems + Theory Lab @ Harvard University

Undergraduate Researcher with Prof. Minlan Yu

• Designed scalable RDMA protocol for real-time in-band network telemetry in large datacenters. Implemented protocol in P4, deployed on Tofino, and ran simulations in Mininet and ns-3. Co-authored papers at SIGCOMM '23 and HotNets '21.

#### **RISE Lab @ UC Berkeley**

Undergraduate Researcher with Prof. Ion Stoica

• Designed a work-stealing mechanism, together with task pipelining technique to improve the throughput of distributed Ray framework. Committed to the GitHub repo (35.2K stars): ray-project/ray

## RECENT PUBLICATIONS

## **First-author** papers

- (ArXiv '24) SuffixDecoding: A Model-Free Approach to Speeding Up Large Language Model Inference <u>Gabriele Oliaro</u>, Zhihao Jia, Daniel Campos, Aurick Qiao. [pdf]
- (ArXiv '24) FlexLLM: A System for Co-Serving Large Language Model Inference and Parameter-Efficient Finetuning. Xupeng Miao\*, <u>Gabriele Oliaro</u>\*, Xinhao Cheng, Mengdi Wu, Colin Unger, Zhihao Jia. [pdf]
- (ASPLOS '24) SpecInfer: Accelerating Generative Large Language Model Serving with Speculative Inference and Token Tree Verification. Xupeng Miao\*, <u>Gabriele Oliaro</u>\*, Zhihao Zhang\*, Xinhao Cheng\*, Zeyu Wang, Rae Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, Chunan Shi, Zhuoming Chen, Daiyaan Arfeen, Reyna Abhyankar, Zhihao Jia. [pdf] Tited 250+ times

goliaro@cs.cmu.edu

Aug 2022 - Present

Aug 2021 - June 2023 Awarded full-ride merit scholarship Aug 2017 - May 2021 Cum Laude with High Honors (GPA 3.83)

> Aug 2022 - present Carnegie Mellon University

> > May 2024 - Present San Mateo, CA

Sept 2021 - Aug 2022

New York, NY

Cambridge, MA

Jun 2020 - Aug 2020

Sept 2020 - Feb 2022

Berkeley, CA

Collaborations

- (ArXiv '25) SpecReason: Fast and Accurate Inference-Time Compute via Speculative Reasoning Rui Pan, Yinwei Dai, Zhihao Zhang, <u>Gabriele Oliaro</u>, Zhihao Jia, Ravi Netravali. [pdf]
- (ArXiv '25) AdaServe: SLO-Customized LLM Serving with Fine-Grained Speculative Decoding Zikun Li<sup>\*</sup>, Zhuofu Chen<sup>\*</sup>, Remi Delacourt, <u>Gabriele Oliaro</u>, Zeyu Wang, Qinghan Chen, Shuhuai Lin, April Yang, Zhihao Zhang, Zhuoming Chen, Sean Lai, Xupeng Miao, Zhihao Jia. [pdf]
- (ACL '24 Oral) Quantized Side Tuning: Fast and Memory-Efficient Tuning of Quantized Large Language Models . Zhengxin Zhang, Dan Zhao, Xupeng Miao, <u>Gabriele Oliaro</u>, Qing Li, Yong Jiang, Zhihao Jia. [pdf] 🏆 Outstanding Paper Award
- (ASPLOS '24) Optimal Kernel Orchestration for Tensor Programs with Korch. Muyan Hu, Ashwin Venkatram, Shreyashri Biswas, Balamurugan Marimuthu, Bohan Hou, <u>Gabriele Oliaro</u>, Haojie Wang, Liyan Zheng, Xupeng Miao, Jidong Zhai, Zhihao Jia
- (ACM Comput. Surv.) Towards Efficient Generative Large Language Model Serving: A Survey from Algorithms to Systems. Xupeng Miao, <u>Gabriele Oliaro</u>, Zhihao Zhang, Xinhao Cheng, Hongyi Jin, Tianqi Chen, Zhihao Jia. [pdf]
- (SIGCOMM '23) Direct Telemetry Access. Jonatan Langlet, Ran Ben Basat, <u>Gabriele Oliaro</u>, Michael Mitzenmacher, Minlan Yu, Gianni Antichi. [pdf]
- (HotNets '21) Zero-CPU Collection with Direct Telemetry Access. Jonatan Langlet, Ran Ben Basat, Sivaram Ramanathan, <u>Gabriele Oliaro</u>, Michael Mitzenmacher, Minlan Yu, Gianni Antichi. [pdf]
- \* Equal contribution

### SELECT PROJECTS

 $\label{eq:chickadee OS: Designed a multi-core OS kernel in C/C++. Implemented virtual memory, buddy allocator, processes, threads, wait queues, file system, disk support, buffer cache, signals and system calls.$ 

**LSM-Tree based Key-Value Store**: Designed and implemented in C++ a NoSQL key-value store using a Log-Structured Merge Tree and an in-memory skip list.

**TMALL Repeat Buyers Prediction**: Built an ensemble model based on gradient boosting to predict repeat buyers. Achieved Top 0.6% ranking on Tianchi challenge portal on aliyun.com

XRDict: Implemented a basic reverse dictionary leveraging XLM-R, a cross-lingual understanding language model.