

Yunfan Zhang

530 W 120th St, New York, NY 10024 | (919) 564-9552 | yunfan.z@columbia.edu | www.yunfanzhang.me

Education

Columbia University – Ph.D. Student in Computer Science
Duke University – B.S. in Computer Science, Minor in History

Expected May 2027
December 2020

Research Experience

Research Associate, Columbia University – New York, NY January 2024 – Present

Advisors: Prof. Kathleen McKeown and Prof. Smaranda Muresan

- Conduct research on LLM reasoning, alignment, and agent evaluation.
- Improve LLM alignment with CoT-focused training, including RLVR with GRPO and iterative SFT (STaR).
- Propose enhancements to RL training for LLMs, including new credit-assignment strategies and self-distillation.
- Develop evaluation benchmarks for LLM web-search agents: design ReAct-style agent scaffolding; implement data-crawling pipelines to create real-world, time-sensitive, contamination-free evaluation sets.
- Measure and optimize distributed training and inference efficiency with ver1, PyTorch FSDP, and vLLM.

Research Associate, Columbia University – New York, NY Aug 2021 – May 2023

Advisor: Prof. Ethan Katz-Bassett

- Conducted Internet measurement research on user behavior, performance, reliability, and security.
- Developed LLM agent methods for open Internet measurement problems, including Autonomous System (AS) classification, rDNS inference, and device identification.
- Inferred user activity patterns from public Internet services and datasets (e.g., Google Public DNS, M-Lab, Censys, root DNS logs); trained Gradient Boosting models to predict user activity from DNS-derived signals.
- Analyzed TB-scale datasets in Google BigQuery; implemented high-performance networking code in Golang; deployed measurements across 20+ geographically distributed vantage points.

Research Associate, Duke University – Durham, NC May 2019 – July 2021

Advisors: Prof. Maria Gorlatova and Dr. Guohao Lan

- Conducted research on augmented reality (AR), computer vision, and applied deep learning.
- Proposed a real-time, deep-learning-based depth inpainting and denoising method for consumer RGB-Depth cameras; demonstrated end-to-end user experience gains in mobile AR applications due to refined depth inputs.
- Implemented computer vision algorithms in PyTorch, OpenCV, and SciPy.
- Built Unity-based AR application prototypes for AR headsets (HoloLens, Magic Leap) and mobile (ARCore).

Publications

- **LiveNewsBench: Evaluating LLM Web Search Capabilities with Freshly Curated News**

Yunfan Zhang, Kathleen McKeown, Smaranda Muresan.

Under review at the International Conference on Learning Representations, 2026 (Under review, ICLR 2026).

- **Exploring Chain-of-Thought Reasoning for Steerable Pluralistic Alignment**

Yunfan Zhang, Kathleen McKeown, Smaranda Muresan.

Conference on Empirical Methods in Natural Language Processing, 2025 (EMNLP 2025)

- **Forecasting Communication Derailments Through Conversation Generation**

Yunfan Zhang, Kathleen McKeown, Smaranda Muresan.

International Natural Language Generation Conference, 2025 (INLG 2025)

- **Who Squats IPv4 Addresses?**

Loqman Salamatian, Todd Arnold, Italo Cunha, Jiangchen Zhu, **Yunfan Zhang**, Ethan Katz-Bassett, Matt Calder.

ACM SIGCOMM Computer Communication Review, 2023 (CCR 2023)

- **InDepth: Real-time Depth Inpainting for Mobile Augmented Reality**
Yunfan Zhang, Tim Scargill, Ashutosh Vaishnav, Gopika Premsankar, Mario Di Francesco, Maria Gorlatova.
Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies, 2022 (IMWUT/UbiComp 2022)
- **Optimal Network Protocol Selection for Competing Flows via Online Learning**
Xiaoxi Zhang, Siqi Chen, **Yunfan Zhang**, Youngbin Im, Maria Gorlatova, Sangtae Ha, Carlee Joe-Wong.
IEEE Transactions on Mobile Computing, 2022 (TMC 2022)
- **Edge-assisted Collaborative Image Recognition for Mobile Augmented Reality**
Guohao Lan, Zida Liu, **Yunfan Zhang**, Tim Scargill, Jovan Stojkovic, Carlee Joe-Wong, Maria Gorlatova.
ACM Transactions on Sensor Networks, 2022 (TOSN 2022)
- **Towards Identifying Networks with Internet Clients Using Public Data**
Weifan Jiang, Tao Luo, Thomas Koch, **Yunfan Zhang**, Ethan Katz-Bassett, Matt Calder.
ACM Internet Measurement Conference, 2021 (IMC 2021)
- **Towards a Traffic Map of the Internet: Connecting the Dots between Popular Services and Users**
Thomas Koch, Weifan Jiang, Tao Luo, Petros Gigis, **Yunfan Zhang**, Kevin Vermeulen, Emile Aben, Matt Calder, Ethan Katz-Bassett, Lefteris Manassakis, Georgios Smaragdakis, Narseo Vallina-Rodriguez.
ACM Workshop on Hot Topics in Networks, 2021 (HotNets 2021)
- **CollabAR: Edge-assisted Collaborative Image Recognition for Mobile Augmented Reality**
Zida Liu, Guohao Lan, Jovan Stojkovic, **Yunfan Zhang**, Carlee Joe-Wong, Maria Gorlatova.
ACM/IEEE International Conference on Information Processing in Sensor Networks, 2020 (IPSN 2020)
- **Edge-based Provisioning of Holographic Content for Contextual and Personalized Augmented Reality**
Michael Glushakov, **Yunfan Zhang**, Yuqi Han, Tim Scargill, Guohao Lan, Maria Gorlatova.
IEEE International Conference on Pervasive Computing and Communications Workshops, 2020 (PerCom Workshops 2020)

Selected Industry Experience and Projects

- External Research Collaborator, Fireworks AI – New York, NY** December 2024 – Present
- Design benchmarks to evaluate LLM instruction-following capabilities under realistic API user workflows.
 - Develop evaluations to quantify modality gaps between textual and visual inputs in VLMs.
 - Coordinate paper submission, benchmark leaderboard launches, and technical blog releases with the team.
- Research Intern, Cloudflare – New York, NY** June 2023 – September 2023
- Developed a prototype bot traffic detection system using Language Models and natural language features; among the first LM-based network traffic analysis efforts at Cloudflare.
 - Improved bot detection accuracy from ~60% to ~90% in offline evaluation on global production traffic.
 - Built high-performance feature engineering and data pipelines in ClickHouse SQL, processing O(100B) database rows per run.
- Software Engineer, Duke Electric Vehicles Team – Durham, NC** October 2016 – August 2019
- **Guinness World Record:** Most efficient electric vehicle: 27,482 MPGe (battery-electric).
 - **Guinness World Record:** Most fuel-efficient vehicle: 14,573 MPGe (hydrogen fuel cell).
 - Built a real-time telemetry and data analytics pipeline (Bluetooth Communication, Android dashboard, Python backend / strategy planners, JS visualization) to inform driving strategies in races and record attempts..
 - Contributed to the vehicle's trapezoidal motor control algorithms and vehicle sensor drivers. Reconstructed vehicle dynamics parameters (speed/accel/elevation) from noisy sensor readings using LOWESS regression.
- Software Engineering Intern, Red Hat – Durham, NC** May 2018 – August 2018
- Contributed to Ansible and Ansible Tower, Red Hat's distributed server management and orchestration tool.
 - Implemented features across the stack, including DSL parsing, OAuth authentication, user management, auditing, and CLI tooling.

- Diagnosed and fixed concurrency bugs (race conditions, deadlocks) in distributed job execution.
- Triaged and reviewed community issues/PRs for the Ansible GitHub repository (40K+ stars, 4K+ contributors) and Ansible Tower (8K+ stars, 200+ contributors).

Teaching Experience

Teaching Assistant, CSEE 4119: Computer Networks – Columbia University

Fall 2022

Teaching Assistant, CS 356: Computer Network Architecture – Duke University

Fall 2019, Spring 2020

Honors and Awards

- Reviews: ACL Rolling Review (ARR) 2025, COLING 2025, IEEE Transactions on Computational Imaging 2024

As a member of Duke Electric Vehicles Team:

- Guinness World Record: Most efficient electric vehicle: 27,482 MPGe (battery-electric).
- Guinness World Record: Most fuel-efficient vehicle: 14,573 MPGe (hydrogen fuel cell).
- Shell Eco-Marathon Americas 2018: First place in battery-electric prototype. Best of 25 teams.
- Shell Eco-Marathon Americas 2018: First place in hydrogen prototype. Best of 7 teams.
- Shell Eco-Marathon Americas 2018: Technical Innovation Award
- Shell Eco-Marathon Americas 2017: First place in battery-electric prototype. Best of 30 teams.

Technical Skills

- **LLM Research:** RL with LLMs, verl, RLVR, RLHF, GRPO, PPO, STaR, SFT, LoRA, Distributed Model Training (with PyTorch FSDP, DeepSpeed, JAX), Model Serving Optimizations (with vLLM), Model Evaluations & Benchmarks, Agent Scaffolding, Hugging Face Transformers, PEFT, Accelerate.
- **Computer Vision:** OpenCV, Torch Vision, Hugging Face Diffusers, SciPy, NumPy
- **Software Engineering:** Flask, Django, Puppet, JavaScript, Vue.js, Bootstrap, Golang, Java, Android