

Yu Wu

☎ 551-229-5305 ✉ yw828@rutgers.edu <https://yuwuofrutgers.github.io>

PhD candidate at Rutgers University working on LLM agents, multi-LLM collaboration, and resource-efficient AI systems. My research focuses on agent orchestration, model routing and cascades, retrieval-augmented decision making, memory-augmented systems, and reliable inference under compute constraints.

EDUCATION

- Rutgers, The State University of New Jersey, New Jersey, USA** *Aug. 2020 - Present*
PhD candidate in Electrical and Computer Engineering (Advisor: Prof. [Anand D. Sarwate](#))
- University of Science and Technology of China (USTC), Anhui, China** *Jun. 2017 - Jun. 2020*
Master's in Electronic Engineering and Information Science (Advisor: Prof. [Bin Liu](#))
- University of Science and Technology of China (USTC), Anhui, China** *Jul. 2013 - Jun. 2017*
Bachelor's in Information Security

WORK EXPERIENCE

- **Agentic Data Curation for LLM Training Data** | LLM agents, tool calling, multi-agent systems *May 2026 – Present*
- *Research Scientist Intern at **AnalogyAI***
- Built a **Claude Code-like, ReAct-style curation agent** in a multi-agent data pipeline: perceives unexplored data and turn raw sources into training-ready LLM datasets; shipped as the team default.
- **Modality- and schema-agnostic**: dynamically handles **multimodal inputs** and composes operators to deliver **novel, user-specified data kinds**.
- Built its **execution runtime: budget-aware, large-batch concurrent, with safe operator execution**.
- **Online Learning Algorithm for Audio Echo Cancellation** | least mean square, adaptive filtering, acoustic signal processing *Jun. 2022 - Aug. 2022*
- *Researcher (Intern) at **Nokia Bell Lab***
- Toward the misalignment error of online estimation system, we analyze the accuracy and effectiveness of “delay and extrapolate” algorithm, which is widely used in adaptive filter.
- Our method proves the **optimality** and **corrects misuse** in previous works. The experimental results are **consistent** with the theoretical derivation. [[paper](#)]

SELECTED RESEARCH PROJECTS

- **Interactive LLM Cascade for Agentic Model Routing** | RAG, In-context learning, orchestration, distillation, memory mechanisms *Sep. 2024 - Dec. 2025*
- ***Project Leader, Research Assistant at Rutgers***
- Designed an interactive **multi-LLM collaboration** framework that combines **retrieval-augmented routing** with adaptive deferral from weaker LLMs to stronger LLMs.
- Formulated model selection as a long-term **collaboration** process, enabling weaker models to query stronger models only when needed and improve through feedback over time.
- Built a LLM cascade for agentic decision making under resource constraints, achieving up to **33.06%** accuracy improvement while reducing stronger-model calls by up to **48.05%**. Studied the **trade-off** between quality, latency, and API cost in multi-model **orchestration** for real-world deployment. [[paper](#)]
- **Learning to Help: Expert Routing for Resource-Constrained AI Systems** | real-time inference, distributed system, ViT *Sep. 2020 - Sep. 2024*
- ***Project Leader, Research Assistant at Rutgers***
- Developed a framework to jointly train external machine/human experts and a query **router** that efficiently guides data samples to suitable models, enhancing legacy ML system performance.
- Proved **Bayes optimality** theoretically, improved overall system accuracy by **4%-12%**, and demonstrated effectiveness across different **distributed** settings.

- Demonstrated the benefits of structured defer-and-route mechanisms for collaborative AI under limited compute and communication budgets. [[paper1](#), [paper2](#), [code](#)]

- **Reasoning-Guided Teacher-Student Learning for Autonomous Driving** | Reasoning VLM, Distillation, GRU
- *Project collaborator* Apr. 2025 - Nov. 2025
- Built a reasoning-guided teacher-student framework for end-to-end autonomous driving, where the teacher model generates refined reasoning and the student model distills this knowledge to predict numerical trajectories in the absence of explicit reason-action annotations.
- Achieved a **24%** performance improvement over non-reasoning baselines, while providing enhanced **interpretability** and maintaining suitability for deployment on **resource-constrained devices**.
- **Enhancing Model-Based Reinforcement Learning with Data Filter** | Out-of-distribution, RL, MuJoCo, MBPO, Actor-Critic
- *Project collaborator* Jun. 2024 - Sep. 2024
- To bridge model-free and model-based reinforcement learning more efficiently, we introduced an out-of-distribution (**OOD**) data filter into **Dyna-style model-based RL** to remove unreliable samples generated by the learned dynamics model.
- Our theorems provide a **tighter bound** on estimation error, and experiments reduce the number of training epochs required for convergence by up to **25%**. [[paper](#)]
- **Human-computer Interactive Sensing** | wearable device, human-centered computing, human-computer interaction
- *Researcher (Intern) at X-discovery Lab, Dartmouth College* Dec. 2018 - Apr. 2019
- To achieve human-computer interaction through soft material, we propose an inductive sensing based prototype called Tessutivo. We yielded **93.9% real-time accuracy** for object recognition. [[paper](#), [demo](#)]
- **Anti-interference for WiFi-based Human Activity Recognition (HAR)** | CSI, non-intrusive sensing, PyTorch
- *Research Assistant at EEIS department of USTC* Sep. 2017 - May 2020
- To mitigate the interference components in WiFi signals, we propose PhaseAnti system. Our method **improves up to 16%** on accuracy and **9× faster** recognition speed. [[paper1](#), [paper2](#)]

PUBLICATIONS

- **Y Wu**, S Wu, et al., "Online In-Context Knowledge Distillation for LLM Cascades." [Under review](#)
- **Y Wu**, Y Li, et al., "Learning to Help in Multi-Class Settings." [ICLR 2025](#)
- **Y Wu**, and Anand Sarwate, "Learning to Help: Training Models to Assist Legacy Devices." [ISIT 2024 Workshop IT-TML](#)
- Y Tao, B Baker, **Y Wu**, et al., "Batch Effects In Brain Foundation Model Embeddings." [ICML 2026 Workshop SD4H](#)
- Z Dong, Y Zhu, **Y Wu**, et al., "FROST-Drive: Scalable and Efficient End-to-End Driving with a Frozen Vision Encoder." [WACV 2026 Workshop LLVM-AD](#)
- J Huang, B Liu, C Miao, Y Lu, **Y Wu**, et al., "PhaseAnti: An anti-interference WiFi-based activity recognition system using interference-independent phase component." [IEEE Transactions on Mobile Computing 2021](#)
- J Huang, B Liu, P Liu, C Chen, N Xiao, **Y Wu**, et al., "Towards anti-interference WiFi-based activity recognition system using interference-independent phase component." [INFOCOM 2020](#)
- J Gong, **Y Wu**, et al., "Tessutivo: Contextual interactions on interactive fabrics with inductive sensing." [UIST 2019](#)
- T Liu, Z Li, X Liu, Y Li, **Y Wu**, et al., "LaTeX Compilation: Challenges in the Era of LLMs" [arXiv \(2026\)](#)
- W Zhang, Y Li, Z Dong, **Y Wu**, et al., "Renaissance of Literate Programming in the Era of LLMs: Enhancing LLM-Based Code Generation in Large-Scale Projects." [arXiv \(2024\)](#)
- Z Yu, L An, Y Li, **Y Wu**, et al., "EAPCR: A Universal Feature Extractor for Scientific Data without Explicit Feature Relation Patterns." [arXiv \(2024\)](#)
- Y Li, Z Dong, E Luo, **Y Wu**, et al., "When to trust your data: enhancing dyna-style model-based reinforcement learning with data filter." [arXiv \(2024\)](#)

SKILLS

Programming: Python, C, C++, MATLAB, Processing, SQL, Java, R

ML / Systems: PyTorch, NumPy, SciPy, Pandas, scikit-learn, Slurm, Git, LoRA, HuggingFace

Course: Qishi Deep Learning in Recommender Systems [Bootcamp](#)

SERVICE

Reviewer: ICLR, ICML, TMLR, IJCNLP-AAACL, Journal of Knowledge-Based Systems

Teaching Assistant: Probability and Random Processes, Linear Signal and System, Discrete Mathematics, etc.