

HAOTIAN SUN

E-mail: haotian.sun@gatech.edu | **Website:** haotiansun.tech

Research Interests: Adaptive Computation for Foundation Models;
Black-box Adaptation and Personalization; Agentic Planning

EDUCATION

Georgia Institute of Technology GPA – 4.0/4.0 <i>Ph.D. Student in Machine Learning</i>	Fall 2022 - Present <i>Atlanta, GA</i>
École CentraleSupélec GPA – 3.9/4.0 <i>Dual Degree Program in collaboration with Xi'an Jiaotong University</i>	Fall 2015 – Summer 2017 <i>Paris, France</i>
Xi'an Jiaotong University Average Grade – 91/100 <i>Honors Youth Program</i> <i>Highly selective nationwide program accepting under 120 students each year</i>	Fall 2013 – Summer 2020 <i>Xi'an, China</i>

EXPERIENCE

Apple <i>AI/ML Intern, Foundation Model Team</i> Manager: Nan Du Research Focus: Scaling up multi-modal generative models with adaptive computation	May 2024 - August 2024 <i>Cupertino, CA</i>
Georgia Institute of Technology <i>Graduate Research Assistant</i> Advisors: Bo Dai, Chao Zhang Research Focus: Adaptation and personalization for black-box LLMs; Agentic planning with LLMs.	August 2022 - Present <i>Atlanta, GA</i>

FEATURED PROJECTS

HYDRA: Model Factorization Framework for Black-Box LLM Personalization – Developed a model factorization framework that personalizes black-box LLMs by capturing user-specific patterns and shared knowledge without relying on access to the model's inherent parameters; – Designed a reranker and adapter that prioritize relevant historical records and align with user preferences; – Delivered 9.01% average improvement over SoTA prompt-based methods across diverse personalization tasks.	Spring 2024
BBox-Adapter: Lightweight Adapting for Black-Box Large Language Models – Proposed an effective domain adaptation approach that is transparent, privacy-conscious, and cost-effective for customizing commercial black-box LLMs with only APIs; – Designed an online adaptation framework iteratively sampling from previous inferences and optimizing the backend lightweight adapter (up to 0.3B); – Achieved up to 6.77% improvement over the base model with 31.30x less training cost and 1.84x less inference cost than the official SFT service.	Fall 2023
AdaPlanner: Adaptive Planning from Feedback with Language Models – Proposed AdaPlanner, a closed-loop planning approach allowing the LLM agent to refine its self-generated plan adaptively in response to environmental feedback; – Developed a code-style LLM prompt structure that facilitates plan generation across a variety of tasks; – Designed a skill discovery mechanism that leverages successful plans as few-shot exemplars, boosting sample efficiency by up to 600x.	Spring 2023
ToolQA: A Dataset for LLM Question Answering with External Tools – Proposed a new dataset to faithfully evaluate LLMs' ability to use external tools for question answering; – Minimized the overlap between our benchmark data and LLMs' pre-training data, enabling a more precise evaluation of LLMs' tool-use reasoning abilities; – Conducted an in-depth diagnosis of existing tool-use LLMs to highlight their strengths, weaknesses, and potential improvements.	Spring 2023
Autoregressive Diffusion Model for Graph Generation – Designed a diffusion network that learns an optimal node absorbing ordering from graph topology and a denoising network that uses the reverse node order to reconstruct the graph efficiently; – Achieved better generation performance than previous SoTA and guaranteed fast generation speed.	Fall 2022

- [1] **Haotian Sun**^{*}, Yuchen Zhuang^{*}, Wei Wei, Chao Zhang, Bo Dai. BBox-Adapter: Lightweight Adapting for Black-Box Large Language Models, *Forty-first International Conference on Machine Learning (ICML 2024)*, (**Spotlight, top 3.5%**).
- [2] **Haotian Sun**^{*}, Yuchen Zhuang^{*}, Lingkai Kong, Bo Dai, Chao Zhang. AdaPlanner: Adaptive Planning from Feedback with Language Models, *Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS 2023)*.
- [3] Yuchen Zhuang, **Haotian Sun**, Yue Yu, Rushi Qiang, Qifan Wang, Chao Zhang, Bo Dai. HYDRA: Model Factorization Framework for Black-Box LLM Personalization, *ArXiv abs/2406.02888 (2024)*.
- [4] Yuchen Zhuang, Yue Yu, Kuan Wang, **Haotian Sun**, Chao Zhang. ToolQA: A Dataset for LLM Question Answering with External Tools, *Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS 2023)*.
- [5] Lingkai Kong^{*}, **Haotian Sun**^{*}, Yuchen Zhuang, Haorui Wang, Wenhao Mu, Chao Zhang. Two Birds with One Stone: Enhancing Calibration and Interpretability with Graph Functional Neural Process, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics (AISTATS 2024)*.
- [6] Tongzheng Ren, **Haotian Sun**, Antoine Moulin, Arthur Gretton, Bo Dai. Spectral Representation for Causal Estimation with Hidden Confounders, *ArXiv abs/2407.10448 (2024)*.
- [7] Lingkai Kong, Jiaming Cui, **Haotian Sun**, Yuchen Zhuang, B Aditya Prakash, Chao Zhang. Autoregressive Diffusion Model for Graph Generation, *Fortieth International Conference on Machine Learning (ICML 2023)*.