

# HAOTIAN SUN

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

**Research Interests:** Large-scale multi-modal foundation models;  
LLM adaptation and personalization; Agentic planning

## EDUCATION

---

<b>Georgia Institute of Technology</b>   GPA – 4.0/4.0 <i>Ph.D. Student in Machine Learning</i>	Fall 2022 - Present <i>Atlanta, GA</i>
<b>École CentraleSupélec</b>   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>
<b>Xi'an Jiaotong University</b>   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

---

<b>Apple</b>   AI/ML Intern, Foundation Model Team Manager: Nan Du <b>Research Focus:</b> Building unified diffusion models for multi-modal generation	May 2025 - August 2025 <i>Cupertino, CA</i>
<b>Apple</b>   AI/ML Intern, Foundation Model Team Manager: Nan Du <b>Research Focus:</b> Scaling up multi-modal generative models with efficient adaptive computation	May 2024 - August 2024 <i>Cupertino, CA</i>
<b>Georgia Institute of Technology</b>   Graduate Research Assistant Advisors: Bo Dai, Chao Zhang <b>Research Focus:</b> Multi-modal representation learning; Adaptation and personalization for black-box LLMs; Agentic planning with LLMs.	August 2022 - Present <i>Atlanta, GA</i>

## FEATURED PROJECTS

---

<b>AmorLIP: Efficient Language-Image Pretraining via Amortization</b> – Proposed an amortized CLIP framework that reduces large-batch contrastive costs via lightweight networks; – Delivered up to 12.24% relative improvement in zero-shot classification and retrieval across 38 tasks, outperforming standard CLIP baselines with higher efficiency.	Spring 2025
<b>EC-DIT: Scaling Diffusion Transformers with Adaptive Expert-Choice Routing</b> – Scaled up 97B-parameter diffusion transformer with expert-choice routing, enabling adaptive heterogeneous compute for efficient text-to-image generation; – Achieved state-of-the-art GenEval score of 71.68% with faster convergence, stronger text-image alignment, and superior scalability over conventional sparse models.	Fall 2024
<b>HYDRA: Model Factorization Framework for Black-Box LLM Personalization</b> – 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; – Delivered 9.01% average improvement over SoTA prompt-based methods across diverse personalization tasks.	Spring 2024
<b>BBox-Adapter: Lightweight Adapting for Black-Box Large Language Models</b> – Proposed an effective domain adaptation approach that is transparent, privacy-conscious, and cost-effective for customizing commercial black-box LLMs with only APIs; – 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
<b>AdaPlanner: Adaptive Planning from Feedback with Language Models</b> – Proposed AdaPlanner, a closed-loop planning approach allowing the LLM agent to refine its self-generated plan adaptively in response to environmental feedback; – Designed a skill discovery mechanism that leverages successful plans as few-shot exemplars, boosting sample efficiency by up to 600x.	Spring 2023
<b>ToolQA: A Dataset for LLM Question Answering with External Tools</b> – Proposed a new dataset to faithfully evaluate LLMs' ability to use external tools for question answering; – Conducted an in-depth diagnosis of existing tool-use LLMs to highlight their strengths, weaknesses, and potential improvements.	Spring 2023

- [1] **Haotian Sun**, Yitong Li, Yuchen Zhuang, Niao He, Hanjun Dai, Bo Dai, "AMORLIP: Efficient Language-Image Pretraining via Amortization," *Thirty-ninth Conference on Neural Information Processing Systems (NeurIPS 2025)*.
- [2] **Haotian Sun**, Tao Lei, Bowen Zhang, Yanghao Li, Haoshuo Huang, Ruoming Pang, Bo Dai, Nan Du. EC-DIT: Scaling Diffusion Transformers with Adaptive Expert-Choice Routing, *Thirteenth International Conference on Learning Representations (ICLR 2025)*.
- [3] **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%)**.
- [4] **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)*.
- [5] **Haotian Sun**, Antoine Moulin, Tongzheng Ren, Arthur Gretton, Bo Dai. Spectral Representation for Causal Estimation with Hidden Confounders, *The 28th International Conference on Artificial Intelligence and Statistics (AISTATS 2025)*.
- [6] Yuchen Zhuang, **Haotian Sun**, Yue Yu, Rushi Qiang, Qifan Wang, Chao Zhang, Bo Dai. HYDRA: Model Factorization Framework for Black-Box LLM Personalization, *Thirty-eighth Conference on Neural Information Processing Systems (NeurIPS 2024)*.
- [7] Lingkai Kong<sup>\*</sup>, **Haotian Sun**<sup>\*</sup>, 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)*.
- [8] Changhao Li, Yuchen Zhuang, Rushi Qiang, **Haotian Sun**, Hanjun Dai, Chao Zhang, and Bo Dai. "Matryoshka Pilot: Learning to Drive Black-Box LLMs with LLMs." *Thirty-ninth Annual Conference on Neural Information Processing Systems (NeurIPS 2025)*.
- [9] 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)*.
- [10] 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)*.
- [11] Ran Xu, Yuchen Zhuang, Yue Yu, **Haotian Sun**, Hang Wu, Carl Yang, May D Wang. MedAdapter: Efficient Test-Time Adaptation of Large Language Models Towards Medical Reasoning. *Conference on Empirical Methods in Natural Language Processing (EMNLP 2024)*.