

Jianliang He

Kline Tower, 219 Prospect Street, New Haven, CT 06511

Phone: (+1) 203-410-5714 | Email: jianliang.he@yale.edu | Website: jlianghe.github.io

RESEARCH INTERESTS

Machine Learning Theory; Mechanistic Interpretability; Large Language Model.

EDUCATION

Yale University

Department of Statistics and Data Science
Ph.D. in Statistics. Advisor: Zhuoran Yang.

2024.9 - Present

New Haven, CT

Fudan University

Department of Statistics and Data Science, School of Management
B.S. in Statistics

2020.9 - 2024.6

Shanghai, China

RESEARCH PAPERS

* stands for equal contribution or alphabetical ordering.

1. **He, J.**, Wang, L., Chen, S., Yang, Z. “On the Mechanism and Dynamics of Modular Addition: Fourier Features, Lottery Ticket, and Grokking”. Submitted, 2025.
2. **He, J.**, Pan, X., Chen, S., Yang, Z. “In-Context Linear Regression Demystified: Training Dynamics and Mechanistic Interpretability of Multi-Head Softmax Attention”. [arXiv.2503.12734](https://arxiv.org/abs/2503.12734). *International Conference on Machine Learning (ICML)*, 2025.
3. Qin, S.*, **He, J.***, Kuang, Q*, Gang, B, Xia, Y. “Data-light Uncertainty Set Merging with Admissibility”. [arXiv.2410.12201](https://arxiv.org/abs/2410.12201). Submitted to *Journal of Machine Learning Research (JMLR)*, 2024.
4. **He, J.***, Chen, S.*., Zhang, F., Yang, Z. “From Words to Actions: Unveiling the Theoretical Underpinnings of LLM-Driven Autonomous Systems”. [arXiv.2405.19883](https://arxiv.org/abs/2405.19883). *International Conference on Machine Learning (ICML)*, 2024.
5. **He, J.**, Zhong, H., Yang, Z. “Sample-Efficient Learning of Infinite-Horizon Average-Reward MDPs with General Function Approximation”. [arXiv.2404.12648](https://arxiv.org/abs/2404.12648). *International Conference on Learning Representations (ICLR)*, 2024.
6. Banerjeea, T*, Gang, B*, **He, J.***. “Harnessing the Collective Wisdom: Fusion Learning using Decision Sequences from Diverse Sources”. [arXiv.2308.11026](https://arxiv.org/abs/2308.11026). *Biometrika*, 2025+.

INDUSTRIAL EXPERIENCE

Machine Learning Engineer, Cisco Foundation AI Team, San Francisco

2025.6 - Present

- Applied post-training pipelines to develop a reasoning Large Language Model (LLM) for cybersecurity domain.

TEACHING

Teaching Assistant

- MANA130083.01 Nonparametric Statistics, Fudan University Spring, 2023
- S&DS 2410 Probability Theory, Yale University Fall, 2025

SERVICE

Conference Reviewer: NeurIPS (2024, 2025), ICLR (2025,2026), ICML 2025.

Journal Reviewer: Management Science.