

EDUCATION	Zhejiang University <i>Undergraduate Student</i>	Hangzhou, China 2022 - 2026 (expected)
RESEARCH INTERESTS	<p>Recently, my research philosophy has undergone a significant evolution from Deconstruction to Construction. Previously, I focused on Interpretability and Trustworthy ML, driven by a desire to decipher the “physics” of neural networks and ensure their controllability.</p> <p>Now, my research interests have shifted toward new model architectures (e.g., Sparse/Linear Attention, DeltaNet) and new learning paradigms (e.g., Continual Learning, Test-time Learning), grounded in two core beliefs:</p> <ul style="list-style-type: none"> (1) Understanding through construction, not just deconstruction. As Richard Feynman famously said: ”What I cannot create, I do not understand.” (2) The next paradigm won’t emerge from reverse-engineering Transformers alone. Relying solely on that path risks trapping our understanding in the local minima of existing frameworks. <p>While this represents a substantial shift in direction and I am still in an exploratory phase, I believe my previous research experience provides valuable priors for this new terrain. My long-term goal is to build models with true agency and efficiency.</p>	
PUBLICATIONS AND PREPRINTS	<ol style="list-style-type: none"> 1. Xu Cheng*, Lei Cheng*, Zhaoran Peng, Yang Xu, Tian Han, Quanshi Zhang. <i>Layerwise Change of Knowledge in Neural Networks</i>. Proceedings of the 41st International Conference on Machine Learning (ICML), PMLR 235:8038-8059, 2024. 2. Qihan Ren*, Junpeng Zhang*, Yang Xu, Yue Xin, Dongrui Liu, Quanshi Zhang. <i>Towards the Dynamics of a DNN Learning Symbolic Interactions</i>. Neural Information Processing Systems (NeurIPS), 2024. 3. Yang Xu, Yi Wang, Hao Wang. <i>Tracking the Feature Dynamics in LLM Training: A Mechanistic Study</i>. arXiv preprint arXiv:2412.17626, 2024. 4. Yang Xu*, Xuanming Zhang*, Samuel Yeh, Jwala Dhamala, Ousmane Dia, Rahul Gupta, Sharon Li. <i>Simulating and Understanding Deceptive Behaviors in Long-Horizon Interactions</i>. arXiv preprint arXiv:2510.03999, 2024. NeurIPS 2025 Workshop @ResponsibleFM Accepted. ICLR 2026 Under Review. 	
INTERNSHIPS	<p>Shanghai Jiao Tong University Shanghai, China</p> <ul style="list-style-type: none"> • Remote Research Intern in the John Hopcroft Center for Computer Science, School of electronic information and electrical engineering. • Advisor: Prof. Quanshi Zhang. • Study: Interpretability of Neural Networks and Deep Learning Theory. 2 conference papers were published. <p>Rutgers University New Jersey, USA</p> <ul style="list-style-type: none"> • Visiting Student and Research Intern in the Department of Computer Science. • Advisor: Prof. Hao Wang. • Study: Mechanistic Interpretability Study of LLMs. <p>UW-Madison Madison, USA</p> <ul style="list-style-type: none"> • Remote Research Intern in the School of Computer, Data & Information Sciences • Advisor: Prof. Sharon Li. 	2023.05-2024.09(remote) 2024.07-2025.02(remote since 2024.9) 2025.03-2025.11(remote)

AWARDS
AND
HONORS

- 2st Scholarship in Zhejiang University.
- Win 1st Prize for twice(2022, 2023) in Zhejiang Division of National Mathematics Competition for College Students.

ENGLISH
PROFICIENCY

TOEFL iBT: 102 (Listening: 30, Reading: 28, Speaking: 20, Writing: 24) *March 2024*
Activities: Member of ZJUFLA (Zhejiang University Foreign Language Association), English Corner Organizer for 2 semesters.

Sept. 2023 – June 2024

ACADEMIC
SERVICES

Reviewer for: *International Conference on Learning Representations (ICLR) 2025, North American Chapter of the Association for Computational Linguistics (NAACL) 2025, ICLR 2026.*