

EDUCATION	<p><b>Zhejiang University</b> Hangzhou, China  <i>Undergraduate Student</i> 2022 - 2026 (<i>expected</i>)</p> <ul style="list-style-type: none"> <li>• <b>Major:</b> Computer Science and Technology.</li> <li>• <b>Minor Program:</b> Advanced Honor Class for Engineering Education(only 50 science&amp;engineering students selected in each Grade), Chu Kochen Honors College, Zhejiang University.</li> <li>• <b>GPA:</b> 4.03/4.3</li> </ul>
RESEARCH INTERESTS	<p>My research interests lie broadly in deep learning, spanning both theory and applications. I am particularly interested in understanding the <b>mechanisms</b> behind deep learning models. Additionally, I am interested in <b>empirical methodology</b>, believing that intelligence can be understood through construction.</p> <p>Regarding the working principles of deep learning models, I am currently focusing on <b>mechanistic interpretability</b> of LLMs—a rapidly growing research area that aims to reverse-engineer the internal operations of neural networks to uncover how they function. On the empirical side, I plan to explore aspects such as <b>RL</b> and <b>foundation models</b> in the future.</p>
PUBLICATIONS AND PREPRINTS	<ol style="list-style-type: none"> <li>1. Xu Cheng*, Lei Cheng*, Zhaoran Peng, <b>Yang Xu</b>, 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.</li> <li>2. Qihan Ren*, Junpeng Zhang*, <b>Yang Xu</b>, Yue Xin, Dongrui Liu, Quanshi Zhang. <i>Towards the Dynamics of a DNN Learning Symbolic Interactions</i>. Neural Information Processing Systems (NeurIPS), 2024. <ul style="list-style-type: none"> <li>• Originally second author for theoretical contributions; authorship adjusted after merging experimental paper's first author.</li> </ul> </li> <li>3. <b>Yang Xu</b>, Yi Wang, Hao Wang. <i>Tracking the Feature Dynamics in LLM Training: A Mechanistic Study</i>. arXiv preprint arXiv:2412.17626, 2024.</li> </ol>
INTERNSHIPS	<p><b>Shanghai Jiao Tong University</b>   Shanghai, China 2023.05-Present(remote)</p> <ul style="list-style-type: none"> <li>• <b>Remote Research Intern</b> in the John Hopcroft Center for Computer Science, School of electronic information and electrical engineering.</li> <li>• <b>Advisor:</b> Prof. Quanshi Zhang.</li> <li>• <b>Study:</b> Interpretability of Neural Networks and Deep Learning Theory. <b>2 conference papers</b> were published.</li> </ul> <p><b>Rutgers University</b>   New Jersey, USA 2024.07-Present(remote since 2024.9)</p> <ul style="list-style-type: none"> <li>• <b>Visiting Student and Research Intern</b> in the Department of Computer Science.</li> <li>• <b>Advisor:</b> Prof. Hao Wang.</li> <li>• <b>Study:</b> Mechanistic Interpretability Study of LLMs.</li> </ul>
AWARDS AND HONORS	<ul style="list-style-type: none"> <li>• <b>2<sup>st</sup> Scholarship</b> in Zhejiang University.</li> <li>• <b>Win 1<sup>st</sup> Prize for twice(2022, 2023)</b> in Zhejiang Division of National Mathematics Competition for College Students.</li> </ul>
ENGLISH PROFICIENCY	<p><b>TOEFL iBT:</b> 102 (Listening: 30, Reading: 28, Speaking: 20, Writing: 24) <i>March 2024</i></p> <p><b>Activities:</b> Member of ZJUFLA (Zhejiang University Foreign Language Association), English Corner Organizer for 2 semesters.</p> <p><i>Sept. 2023 – June 2024</i></p>

