

Xirui Li

4244679380 | xirui@g.ucla.edu | xirui-li.github.io | [in](#) xirui-li | [G](#) xirui-li | [tw](#) xirui-li

EDUCATION

- | | |
|---|--|
| • University of Maryland, College Park
<i>PhD of Computer Science</i> | <i>Sep 2025 - Now</i>
Washington DC, USA |
| • University of California, Los Angeles
<i>Master of Electrical and Computer Engineering</i> | <i>Sep 2022 - Jul 2024</i>
Los Angeles, USA |
| • Technical University of Munich
<i>Bachelor of Electrical Engineering and Information Technology</i> | <i>Oct 2018 - Jul 2022</i>
Munich, Germany |

WORK EXPERIENCE

- | | |
|---|---|
| • Mathworks []
<i>Software Engineer Intern</i> | <i>Jul 2023 – Sep 2023</i>
Natick, USA |
| ◦ Developed HTML Verifier for HDL code generation reports for pattern inspection that improves use cases from 1 to 7 with optimized user experience. | |
| ◦ Performed unit test and system test on individual kernel HDL coder QE test constraints and achieved 100% code coverage for the constraints. | |
| ◦ Reduced coupling degree to zero and improve robustness for kernel HDL coder five mostly-used test constraints calculation by refactoring for both Simulink and MATLAB HDL code generation workflow. | |
| • BMW Group []
<i>Software Engineer Intern</i> | <i>Feb 2021 - Jul 2021</i>
Munich, Germany |
| ◦ Accelerated <i>Ticket Maker</i> script from 5 steps to 3 steps for automated Jira tickets generation by optimizing tickets generation logic and algorithm. | |
| ◦ Developed <i>Budget Viewer</i> script to generate ticket-related budgets visualization with customization filter based on VBA and Jira Rest-API, which reduce half-day work to 5 minutes. | |
| ◦ Optimized <i>Budget Viewer</i> , reducing reaction time by 96.67% (from 5 minutes to 10 seconds) and streamlined functional redundancy of <i>Ticket Maker</i> software. | |

FIRST AUTHOR PROJECTS AND PUBLICATIONS

- | | |
|--|---|
| • VisualThinker: R1-Zero's "Aha Moment" in Visual Reasoning [ArXiv] [] | <i>Sep 2025 - Present</i>
Submitted to ICLR 2026 |
| ◦ UMD, UCLA, TurningPoint AI Zhou, Hengguang, Li, Xirui, et al. | |
| ◦ Reproduced the first ever visual "aha moment" on a 2B non-SFT model by RLVR, demonstrating emergent reasoning capabilities in vision-language models. | |
| ◦ Developed multimodal agentic reasoning pipeline integrating visual feedback with RLVR to enable systematic reasoning on visual tasks. | |
| ◦ Created comprehensive gradient analysis toolkit to monitor memorization versus grokking behaviors during RLVR training, including metrics for effective rank, nuclear norms, and SVD-based statistics. | |
| • MOSSBench: Oversensitivity in Multimodal LLMs [ICLR 2025][] | <i>Oct 2023 - Jun 2024</i>
ICLR 2025 |
| ◦ UCLA, TurningPoint AI Li, Xirui, et al. | |
| ◦ Proposed the first benchmark to reveal and analyze the oversensitivity prevalence on vision-language models (VLMs) to safe queries. | |
| ◦ Identified three key types of visual stimuli that trigger oversensitivity in multimodal LLMs: Exaggerated Risk, Negated Harm, and Counterintuitive Interpretation. | |
| ◦ Revealed widespread oversensitivity across 20 SOTA MLLMs with refusal rates. | |
| • DrAttack: Decomposition-based Jailbreaking [EMNLP 2024][] | <i>Oct 2023 - Aug 2024</i>
EMNLP 2024 |
| ◦ UCLA, TurningPoint AI Li, Xirui, et al. | |
| ◦ Developed the first decomposition-based jailbreaking attacks on large language models (LLMs), achieving state-of-the-art attack success rate on GPT-4. | |
| ◦ Implemented prompt decomposition and reconstruction techniques to bypass safety mechanisms in SOTA LLMs | |
| ◦ Conducted evaluation across multiple commercial and open-source LLMs to demonstrate attack effectiveness. | |
| ◦ Analyzed defense mechanisms and proposed improved safety alignment strategies based on attack insights. | |

SKILLS

- | | |
|--|--|
| • Programming: Python (Package: PyTorch, Tensorflow, PySpark, PyTest), Java, SQL, Shell Script, MATLAB, R | |
| • Languages: English (Professional), German (Professional), Mandarin (Native) | |