

Xirui Li

4244679380 | xiruili@g.ucla.edu | LinkedIn | GitHub

EDUCATION

| | |
|---|--|
| University of California, Los Angeles (UCLA) <i>Master of Electrical and Computer Engineering</i> <ul style="list-style-type: none">GPA: 3.85 | Sep.2022-Jul.2024 <i>Los Angeles, USA</i> |
| Technical University of Munich (TUM) <i>Bachelor of Electrical Engineering and Information Technology</i> <ul style="list-style-type: none">GPA: 3.84 | Oct.2018-Jul.2022 <i>Munich, Germany</i> |

TECHNICAL SKILLS

Programming Languages: Python, C/C++, Matlab, Java, VBA, RestAPI, MySQL, CUDA, NoSQL
Frameworks: Django, Pytorch, HTML5, CSS **System/Tools:** Linux, Docker, Git, Confluence, Jira

RESEARCH

| | |
|---|--|
| Research Assistant <i>University of California, Los Angeles</i> <ul style="list-style-type: none">Create an oversensitive benchmark to evaluate MLLMs' safety alignment balance.Propose an novel adversarial attack on LLMs by prompt decomposition and reconstruction. | Oct.2023-Present <i>Los Angeles, California</i> |
| Research Assistant <i>Technical University of Munich</i> <ul style="list-style-type: none">Investigated visual interpretations for DETection TRansformers and human-in-the-loop workflows of DETection TRansformers. [GitHub]Interpreted attention mechanism in Deformable DETection TRansformer in Local Interpretable Model-Agnostic Explanations (LIME) architecture.Implemented the Deformable DETection TRansformer in Caltech Pedestrian dataset.[GitHub] | Jul.2021-Dec.2021 <i>Munich, Germany</i> |

PREPRINTS

| | |
|--|----------|
| DrAttack: Prompt Decomposition and Reconstruction Makes Powerful LLMs Jailbreakers <ul style="list-style-type: none">A novel jailbreaking method that leverages prompt decomposition and reconstruction to make LLMs generate harmful contents. DrAttack obtains a substantial gain of success rate on powerful LLMs over attackers (15% to 80% on GPT-4).Project details could be found in this [website]. | Feb.2024 |
| MOSSBench: Is Your MLLM Oversensitive to Safe Queries? <ul style="list-style-type: none">The first comprehensive MLLM benchmark to evaluate the over-sensitivity issue in MLLMs. Prevalent oversensitivity behaviors have been evaluated in SOTA MLLMs with Claude-3 of 76.33% oversensitivity rate.Project details could be found in this [website]. | Jun.2024 |

WORK EXPERIENCE

| | |
|---|---|
| Software Engineer Intern <i>Mathworks</i> <ul style="list-style-type: none">Developed HTML Verifier for HDL code generation reports for pattern inspection to improved 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 0 and improve robustness for kernel HDL coder 5 mostly-used test constraints calculation by refactoring for both Simulink and MATLAB HDL code generation workflow. | Jul.2023-Sep.2023 <i>Natick, MA</i> |
| Software Engineer Intern <i>BMW Group</i> <ul style="list-style-type: none">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. | Feb.2021-Jul.2021 <i>Munich, Germany</i> |