

EDUCATION

•Bachelor of Engineering in AI: Systems & Technology

Sept 2021 - July 2025

The Chinese University of Hong Kong

CGPA: 3.637

- Dean's List 2021-2022
- Dean's List 2023-2024
- ELITE Stream Scholarship 2022-2023
- Relevant Coursework: Fundamentals of Machine Learning, Introduction to Computer Systems, Introduction to NLP, Foundations of Optimization, Foundation of Applied Deep Learning, Data Comm & Computer Networks

PERSONAL PROJECTS

•LLMs for Software

Apr. 2024 - present

A comprehensive evaluation framework on current LLMs' capabilities to help with software development.

- Investigated multimodal-related papers within seven years of those that can benefit from LLMs and accordingly formulated a task tree that integrates all potential tasks.
- Building an evaluation framework that constructs several LLMs, typical datasets, and corresponding tasks' evaluation metrics.

•LLM Safety

Jul. 2022 - Aug. 2023

Consist of three projects focusing on safety evaluation frameworks: DUO & OASIS & Multilingual Safety.

- DUO: Developed a general, effective methodology and tool for testing multimedia content moderation software. Successfully attacked 7 popular software through 1,200 pieces of data, measurably indicating their weaknesses in multimedia content moderation.
- OASIS: Conducted a preliminary study on over 5,000 real-world image messages and summarised over 20 images' transform rules to hide toxic content. Designed a comprehensive testing framework for textual toxic content spread via images, whose results can effortlessly evade software moderation.
- Multilingual Safety: Built a multilingual safety benchmark consisting of 10 different languages and 11 tasks. Demonstrated that LLMs potentially exist in different biases between languages with different security audit strengths.

EXPERIENCE

•Research Collaboration

Jun - Aug 2023

Tencent AI Lab, Shenzhen, China

Online

- Supervisor: Dr. Zhaopeng Tu
- Project Title: Multilingual Safety of Large Language Model
- Built a multilingual safety benchmark containing more than 26,000 pieces of data and 10 different languages for testing large language models' multilingual safety ability.

•UG Summer Research Internship

Jun - Aug 2024

The Chinese University of Hong Kong, Hong Kong

Offline

- Supervisor: Prof. Michael Rung Tsong Lyu
- Project Title: Automatically Detecting Audio-Visual Inconsistencies in 3D Software Systems
- Proposed an automatic tool to detect audio and video inconsistencies in software, which utilized LLMs and LLMs to extract information and analysis.

PUBLICATIONS

•All Languages Matter: On the Multilingual Safety of LLMs

Wang, Wenxuan, Zhaopeng Tu, Chang Chen, Youliang Yuan, Jen-tse Hung, Wenxiang Jiao, and Michael Lyu

Findings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL), 2024

•Validating Multimedia Content Moderation Software via Semantic Fusion

Wang, Wenxuan, Jingyuan Huang, Chang Chen, Jiazhen Gu, Jianping Zhang, Weibin Wu, Pinjia He, and Michael Lyu

Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA), 2023

•An Image is Worth a Thousand Toxic Words: A Metamorphic Testing Framework for Content Moderation Software

Wang, Wenxuan, Jingyuan Huang, Jen-tse Huang, Chang Chen, Jiazhen Gu, Pinjia He, and Michael R Lyu

2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE), 2023