

XIAOMENG HU

PH.D. STUDENT

(+852) 53931720

greghxm@link.cuhk.edu.hk

<https://gregxmhu.github.io/>

| | | |
|---------------------------------------|--|---|
| EDUCATION | The Chinese University of Hong Kong <i>Ph.D. in Computer Science and Engineering</i> • Advisor: Prof. Tsung-Yi Ho and Dr. Pin-Yu Chen (IBM Research) • Research area: LLMs & Trustworthy Machine Learning | Hong Kong, China 2023 - 2027 (<i>expected</i>) |
| | Northeastern University <i>B.E. in Artificial Intelligence</i> • GPA: 4.22/5.00, Rank: 1/123. | Shenyang, China 2019 - 2023 |
| PUBLICATIONS | <ol style="list-style-type: none">Xiaomeng Hu, Pin-Yu Chen, Tsung-Yi Ho. Token Highlighter: Inspecting and Mitigating Jailbreak Prompts for Large Language Models. <i>Under Review</i>, 2024.Xiaomeng Hu, Pin-Yu Chen, Tsung-Yi Ho. Gradient Cuff: Detecting Jailbreak Attacks on Large Language Models by Exploring Refusal Loss Landscapes. <i>NeurIPS 2024</i>, 2024.Xiaomeng Hu, Pin-Yu Chen, Tsung-Yi Ho. RADAR: Robust AI-Text Detection via Adversarial Learning. <i>NeurIPS 2023</i>, 2023.Xiaomeng Hu, Shi Yu, Chenyan Xiong, Zhenghao Liu, Zhiyuan Liu, Ge Yu. P³ Ranker: Mitigating the Gaps between Pre-training and Ranking Fine-tuning with Prompt-based Learning and Pre-finetuning. <i>SIGIR 2022</i>, 2022. | |
| PROJECTS | Jailbreak Prompt Detection <i>Detecting Jailbreak prompts via checking the refusal loss landscape.</i> | 2023.08 - 2023.12 |
| | AI-Text Detection <i>Training a robust ai-text detector via adversarial learning</i> | 2023.03 - 2023.05 |
| | Prompt-based document reranking <i>Using prompt-based learning to achieve few-shot document reranking.</i> | 2022.08 - 2022.12 |
| AWARDS AND HONORS (SELECTED) | <ul style="list-style-type: none">• NeurIPS Scholar Award, NeurIPS 2024 2024.10• Full Postgraduate Scholarship, The Chinese University of Hong Kong, 2023.08• First Class Scholarship, Northeastern University 2022.09• Outstanding Student Leader, Northeastern University 2021.09• First Prize Award, The Chinese Mathematics Competitions 2020.12• Second Prize Award, The Chinese Physics Olympiad 2018.09 | |
| SKILLS | Languages: Chinese, English, Cantonese. Programming: Python, C++, MATLAB, Latex. | |
| ACADEMIC SERVICES | Reviewers for: <i>ICLR 2025, NeurIPS Adv-ML 2023</i> Invited speaker for: <i>AIS 2022</i> | |