

Hongpeng Lin

<https://hopelin99.github.io> | Google Scholar | hopelin@ruc.edu.cn | (+86) 187-597-90891

Research Interests

My research interests are in Natural Language Processing and Multi-modal Understanding, with the hope of making machines perceive, understand, and express like humans. Currently, my research focus is committed to adversarial robustness (jailbreaking LLMs and mitigation), long video understanding (natural language query localization), and creative multi-modal text generation (humor and association).

Education

Renmin University of China, Gaoling School of Artificial Intelligence <ul style="list-style-type: none">Master of Science in Artificial Intelligence. <i>Advisor:</i> Ruihua Song	Sep. 2021 – Jun. 2024
Xidian University <ul style="list-style-type: none">Bachelor of Engineering in Software Engineering.	Sep. 2017 – Jun. 2021

Experience

Research Intern, Stanford University <ul style="list-style-type: none">Social Influence Jailbreak Program. <i>Advisor:</i> Weiyan Shi, Yi Zeng	Sep. 2023 - Jan. 2024
Intern, Xianyuan Technology <ul style="list-style-type: none">Studies of LLMs and MLLMs.	Jun. 2023 - Aug. 2023
Intern, Kuaishou Technology <ul style="list-style-type: none">WenLan project of large-scale multi-modal pre-training.	Nov. 2021 - Jun. 2022
Research Intern, Beijing Academy of Artificial Intelligence (BAAI) <ul style="list-style-type: none">Multi-modal data crawling and pre-processing.	Oct. 2021 - Jan. 2022
Research Intern, Tencent AI Lab <ul style="list-style-type: none">Development of a multi-modal dialogue WeChat mini-program.	Jul. 2021 - Aug. 2021

Publications

[*: Equal Contribution]

- (1) How Johnny Can Persuade LLMs to Jailbreak Them: Rethinking Persuasion to Challenge AI Safety by Humanizing LLMs
Yi Zeng*, **Hongpeng Lin***, Jingwen Zhang, Diyi Yang, Ruoxi Jia, Weiyan Shi
arXiv Preprint 2024.
- (2) TikTok: A Video-Based Dialogue Dataset for Multi-Modal Chitchat in Real World
Hongpeng Lin*, Ludan Ruan*, Wenke Xia*, Peiyu Liu, Jingyuan Wen, Yixin Xu, Di Hu, Ruihua Song, Wayne Xin Zhao, Qin Jin, Zhiwu Lu.
ACM MM 2023 (Oral).
- (3) An Analysis of Collection Methods and Simulation Algorithms for Chinese Speech Errors
Hao Pu, **Hongpeng Lin**, Ruihua Song, Mei Yan.
CLSW 2022.

Awards

Graduate Student First Class Academic Scholarship, 2023

Skills

Technical Skills: Proficient in Python and Pytorch, and has practical experience in finetuning (M)LLMs.

Characteristic: Reliable and responsible, with a strong sense of teamwork and a passion for research.

Academic Services

Conference Reviewer: EMNLP, ACM MM, CVPR, ICML