

Xiaofei Wen

Davis, CA | xfwe@ucdavis.edu | (279)766-9275 | Homepage | LinkedIn | GitHub

Bio

Xiaofei Wen is a Computer Science Ph.D. student at University of California, Davis, advised by Prof. Muhaao Chen. His research interest lies in improving the reliability and behavior of language and video models.

Publication

- [a] **OmniGuard: Unified Omni-Modal Guardrails with Deliberate Reasoning** Dec 2025
Boyu Zhu, Xiaofei Wen, Wenjie Jacky Mo, Tinghui Zhu, Yanan Xie, Peng Qi, Muhaao Chen
[Preprint](#) [paper] [webpage]
- [b] **Towards Policy-Compliant Agents: Learning Efficient Guardrails For Policy Violation Detection** Oct 2025
Xiaofei Wen, Wenjie Jacky Mo, Yanan Xie, Peng Qi, Muhaao Chen
[Preprint](#) [paper] [webpage]
- [c] **REDCODER: Automated Multi-Turn Red Teaming for Code LLMs** July 2025
Wenjie Jacky Mo, Qin Liu, Xiaofei Wen, Dongwon Jung, Hadi Askari, Wenxuan Zhou, Zhe Zhao, Muhaao Chen
[Preprint](#) [paper] [code]
- [d] **Diagnosing and Mitigating Modality Interference in Multimodal Large Language Models** May 2025
Rui Cai, Bangzheng Li, Xiaofei Wen, Muhaao Chen, Zhe Zhao
[Preprint](#) [paper] [code]
- [e] **ThinkGuard: Deliberative Slow Thinking Leads to Cautious Guardrails** Feb 2025
Xiaofei Wen, Wenxuan Zhou, Wenjie Jacky Mo, Muhaao Chen
[ACL 2025 Findings](#) [paper] [code]
- [f] **Red Teaming Language Models for Processing Contradictory Dialogues** Nov 2024
Xiaofei Wen, Bangzheng Li, Tenghao Huang, Muhaao Chen
[EMNLP 2024](#) [paper] [code]
- [g] **Personalized Topic Selection Model for Topic-Grounded Dialogue** Aug 2024
Shixuan Fan, Wei Wei, Xiaofei Wen, Xianling Mao, Jixiong Chen, Danyang Chen
[ACL 2024 Findings](#) [paper]
- [h] **Sequential Topic Selection Model with Latent Variable for Topic-Grounded Dialogue** Dec 2022
Xiaofei Wen, Wei Wei, Xianling Mao
[EMNLP 2022 Findings](#) [paper]

Education

- University of California, Davis**, Ph.D. in Computer Science Sept 2024 - Present
• **GPA:** 3.95/4.00
- Huazhong University of Science and Technology**, M.S. in Computer Science Sept 2021 – June 2024
• **GPA:** 3.90/4.00
• **Awards:** First-Class Graduate Scholarship
- Dalian University of Technology**, B.Eng in Software Engineering Sept 2017 – June 2021
• **GPA:** 3.89/4.00

- **Awards:** Mitsubishi Scholarship, Outstanding Scholarship, Meritorious Winner of 2019 ICM

Experience

Research Intern, NetEase – Hangzhou, China

June 2023 – Sept 2023

- Collected and standardised the data used to address the problem of dialogue hallucinations, and validated the validity through fine-tuning methods.
- Participated in the construction of Chinese persona dialogue model based on LLaMA and LLaMA2. The finetuned model was used for character dialogue fragments generation for the game Treacherous Waters.

Research Intern, Ant Group – Hangzhou, China

Sept 2022 – Dec 2022

- Implemented a RL-based conversational(interactive) recommendation system [code].
- Improved the performance of the interactive recommendation system for the real insurance business in Alipay.

Service

Reviewer: ACL, EMNLP, NAACL, since 2023.

Vice President, Model United Nations of DUT, 2018-2019