

# 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 in general-purpose applications.

## Publication

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

## Education

<b>University of California, Davis</b> , Ph.D. in Computer Science	Sept 2024 - Present
• <b>GPA:</b> 3.95/4.00	
<b>Huazhong University of Science and Technology</b> , M.S. in Computer Science	Sept 2021 – June 2024
• <b>GPA:</b> 3.90/4.00	
• <b>Awards:</b> First-Class Graduate Scholarship	
<b>Dalian University of Technology</b> , B.Eng in Software Engineering	Sept 2017 – June 2021
• <b>GPA:</b> 3.89/4.00	
• <b>Awards:</b> Mitsubishi Scholarship, Outstanding Scholarship, Meritorious Winner of 2019 ICM	