

# Wenxuan Wang

✉ jwxwang@gmail.com | 🌐 <https://jarviswang94.github.io/>

## Summary

I am a **postdoc** researcher at **the Hong Kong University of Science and Technology**, advised by Prof. Shuai Wang. I received my Ph.D. degree at The Chinese University of Hong Kong, advised by Prof. Michael R. Lyu. I was a visiting scholar at UCLA NLP Group, worked with Prof. Nanyun Peng and Harvard Medical School, worked with Prof. Xiang Li.

My research interests are the **trustworthy AI models**, such as Large Language Models, Multi-Modal AI Models, Code Intelligent Models, and Medical Intelligent Models. My works have over 2000 citations and an h-index of 21. My GitHub projects have received over 3000 stars.

## Education

### The Chinese University of Hong Kong

PHD'S STUDENT IN COMPUTER SCIENCE AND ENGINEERING

Hong Kong

2020 - 2024

### Huazhong University of Science and Technology

B.E. IN COMPUTER SCIENCE AND TECHNOLOGY

Wuhan, Hubei, China

2013 - 2017

## Experience

### Hong Kong University of Science and Technology

POSTDOC RESEARCHER

- Security and Safety of LLMs

Hong Kong, China

2024.9 - Now

### UCLA NLP Group

VISITING RESEARCHER

- Evaluation of Large Language Models and Multi-modal AI Models

Los Angeles, CA, USA

2023 - 2024

### Tencent AI Lab

RESEARCH INTERN

- Safety and Reliability of Language Generation Models

Shenzhen, Guangdong, China

2019 - 2023

### Turing Robot Inc.

RESEARCH AND DEVELOPMENT ENGINEER

- Multi-modal AI

Beijing, China

2017 - 2019

## Selected Publication

### Apathetic or Empathetic? Evaluating LLMs' Emotional Alignment with Humans

JEN-TSE HUANG, MAN HO LAM, ERIC JOHN LI, SHUJIE REN, **WENXUAN WANG (CORRESPONDING)**, WENXIANG JIAO, ZHAOPENG TU, MICHAEL R. LYU

the Thirty-Eighth Annual Conference on Neural Information Processing Systems

NeurIPS 2024

(CCF A)

### On the Reliability of Psychological Scales on Large Language Models

JEN-TSE HUANG, WENXIANG JIAO, MAN HO LAM, ERIC JOHN LI, **WENXUAN WANG (CORRESPONDING)**, MICHAEL R. LYU

The 2024 Conference on Empirical Methods in Natural Language Processing

EMNLP 2024

(CAAI A, CCF B)

### LogicAsker: Triggering the Logical Reasoning Failures in Large Language Models

YUXUAN WAN\*, **WENXUAN WANG (CO-FIRST)\***, YILIU YANG, PINJIA HE, WENXIANG JIAO, MICHAEL R. LYU

The 2024 Conference on Empirical Methods in Natural Language Processing

EMNLP 2024

(CAAI A, CCF B)

### New Job, New Gender? Measuring the Social Bias in Image Generation Models

**WENXUAN WANG**, HAONAN BAI, JEN-TSE HUANG, JINGYUAN HUANG, HAOWEI QIU, NANYUN PENG, MICHAEL R. LYU

ACM Multimedia 2024 (Oral)

ACM MM 2024 (Oral)

(CCF A)

### Not All Countries Celebrate Thanksgiving: On the Cultural Dominance in Large Language Models

**WENXUAN WANG**, WENXIANG JIAO, RUYI DAI, JINGYUAN HUANG, JEN-TSE HUANG, ZHAOPENG TU, MICHAEL R. LYU

The Annual Meeting of the Association for Computational Linguistics

ACL 2024

(CCF A)

<b>All Languages Matter! A Multilingual Safety Benchmark for Large Language Models</b> WENXUAN WANG, ZHAOPENG TU, CHANG CHEN, YOULIANG YUAN, JEN-TSE HUANG, WENXIANG JIAO, MICHAEL R. LYU Findings of The Annual Meeting of the Association for Computational Linguistics	ACL 2024 Findings (CCF A)
<b>Understanding and Mitigating the Uncertainty in Zero-Shot Translation</b> WENXUAN WANG, WENXIANG JIAO, SHUO WANG, ZHAOPENG TU, MICHAEL R. LYU The IEEE/ACM Transactions on Audio, Speech, and Language Processing	TASLP (JCR Q1, SCI Q1, CAAI A)
<b>On the Shortcut Learning in Multilingual Neural Machine Translation</b> WENXUAN WANG, WENXIANG JIAO, JEN-TSE HUANG, ZHAOPENG TU, MICHAEL R. LYU Neurocomputing	Neurocomputing (JCR Q2, SCI Q1)
<b>A Picture is Worth a Thousand Toxic Words: A Metamorphic Testing Framework for Content Moderation Software</b> WENXUAN WANG, JINGYUAN HUANG, CHANG CHEN, PINJIA HE, JIAZHEN GU, MICHAEL R. LYU The IEEE/ACM International Conference on Automated Software Engineering	ASE 2023 (CCF A)
<b>BiasAsker: Measuring the Bias in Conversational AI System</b> YUXUAN WAN*, WENXUAN WANG* (CO-FIRST), PINJIA HE, JIAZHEN GU, HAONAN BAI, MICHAEL R. LYU The ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering	FSE 2023 (CCF A)
<b>Validating Multimedia Content Moderation Software via Semantic Fusion</b> WENXUAN WANG, JINGYUAN HUANG, CHANG CHEN, JIAZHEN GU, JIANPING ZHANG, WEIBIN WU, PINJIA HE, MICHAEL R. LYU The ACM SIGSOFT International Symposium on Software Testing and Analysis	ISSTA 2023 (CCF A)
<b>MTTM: Metamorphic Testing for Textual Content Moderation Software</b> WENXUAN WANG, JEN-TSE HUANG, WEIBIN WU, JIANPING ZHANG, YIZHAN HUANG, SHUQING LI, PINJIA HE, MICHAEL R. LYU The IEEE/ACM International Conference on Software Engineering	ICSE 2023 (CCF A)
<b>Understanding and Improving Sequence-to-Sequence Pretraining for Neural Machine Translation</b> WENXUAN WANG, WENXIANG JIAO, YONGCHANG HAO, XING WANG, SHUMING SHI, ZHAOPENG TU, MICHAEL R. LYU The Annual Meeting of the Association for Computational Linguistics	ACL 2022 (CCF A)
<b>Rethinking the Value of Transformer Components</b> WENXUAN WANG, ZHAOPENG TU The International Conference on Computational Linguistics	COLING 2020 (CCF B)
<b>FPETS: Fully Parallel End-to-end Text-to-speech System</b> DABIAO MA*, ZHIBA SU*, WENXUAN WANG* (CO-FIRST), YUHAO LU AAAI Conference on Artificial Intelligence	AAAI 2020 Oral (CCF A)
<b>Who is ChatGPT? Benchmarking LLMs’ Psychological Portrayal Using PsychoBench</b> JEN-TSE HUANG, WENXUAN WANG, ERIC JOHN LI, MAN HO LAM, SHUJIE REN, YOULIANG YUAN, WENXIANG JIAO, ZHAOPENG TU, MICHAEL R. LYU The Twelfth International Conference on Learning Representations	ICLR 2024 Oral (CAAI A)
<b>Gpt-4 is Too Smart to be Safe: Stealthy Chat with LLMs via Cipher</b> YOU LIANG YUAN, WENXIANG JIAO, WENXUAN WANG, JEN-TSE HUANG, PINJIA HE, SHUMING SHI, ZHAOPENG TU The Twelfth International Conference on Learning Representations	ICLR 2024 (CAAI A)
<b>Boosting Adversarial Transferability by Block Shuffle and Rotation</b> KUNYU WANG, XUANRAN HE, WENXUAN WANG, XIAOSEN WANG The IEEE/CVF Conference on Computer Vision and Pattern Recognition	CVPR 2024 (CCF A)
<b>Does ChatGPT Know that It Does Not Know? Evaluating the Black-Box Calibration of ChatGPT.</b> YOU LIANG YUAN, WENXUAN WANG, QINGSHUO GUO, YIMING XIONG, CHIHAI SHEN, PINJIA HE International Conference on Computational Linguistics	COLING 2024 (CCF B)
<b>Generative Type Inference for Python</b> YUN PENG, CHAOZHENG WANG, WENXUAN WANG, CUIYUN GAO, MICHAEL R. LYU The IEEE/ACM International Conference on Automated Software Engineering (ACM SIGSOFT Distinguished Paper Award)	ASE 2023 (CCF A)
<b>Constructing Effective In-Context Demonstration for Code Intelligence Tasks: An Empirical Study</b> SHUZHENG GAO, XIN-CHENG WEN, CUIYUN GAO, WENXUAN WANG, MICHAEL R. LYU The IEEE/ACM International Conference on Automated Software Engineering	ASE 2023 (CCF A)

### **Improving the Transferability of Adversarial Samples by Path-Augmented Method**

JIANPING ZHANG, JEN-TSE HUANG, **WENXUAN WANG**, YICHEN LI, WEIBIN WU, XIAOSEN WANG, YUXIN SU, MICHAEL R. LYU

The IEEE/CVF Conference on Computer Vision and Pattern Recognition

CVPR 2023

(CCF A)

### **Parrot: Translating during chat using large language models**

WENXIANG JIAO, JEN-TSE HUANG, **WENXUAN WANG**, XING WANG, ZHAOPENG TU

Findings of The Conference on Empirical Methods in Natural Language Processing

Findings of EMNLP 2023

(CAAI A, CCF B)

### **LFAA: Crafting Transferable Targeted Adversarial Examples with Low-Frequency Perturbations**

KUNYU WANG, JULUAN SHI, **WENXUAN WANG**

26th European Conference on Artificial Intelligence

ECAI 2023

(CCF B)

### **Improving Adversarial Transferability via Neuron Attribution-Based Attacks**

JIANPING ZHANG, WEIBIN WU, JEN-TSE HUANG, YIZHAN HUANG, **WENXUAN WANG**, YUXIN SU, MICHAEL R. LYU

The IEEE/CVF Conference on Computer Vision and Pattern Recognition

CVPR 2022

(CCF A)

### **AEON: A Method for Automatic Evaluation of NLP Test Cases**

JEN-TSE HUANG, JIANPING ZHANG, **WENXUAN WANG**, PINJIA HE, YUXIN SU, MICHAEL R. LYU

The ACM SIGSOFT International Symposium on Software Testing and Analysis

ISSTA 2022

(CCF A)

### **Revisiting, Benchmarking and Exploring API Recommendation: How Far Are We?**

YUN PENG, SHUQING LI, WENWEI GU, YICHEN LI, **WENXUAN WANG**, CUIYUN GAO, MICHAEL R. LYU

IEEE Transactions on Software Engineering

TSE 2021

(CCF A)