

# Jingyuan Huang

✉ 1155173905@link.cuhk.edu.hk | 🏠 harveyyellow.github.io | 🎓 Huang Jingyuan

## Education

### The Chinese University of Hong Kong

UNDERGRADUATE STUDENT IN ARTIFICIAL INTELLIGENCE

GPA: 3.6/4

HONG KONG

Sep '21 - Present.

## Research Experience

### Summer Research with Assistant Professor Jieyu Zhao (USC)

RESEARCH COLLABORATOR

- Researched bias in Vision-Language Models (VLMs), currently writing for submission to ACL in February.
- Conducted the majority of the literature review and all experimental work, contributing ideas that were incorporated into the research.
- Designated as one of two co-first authors on the paper.

REMOTE

May '24 - Present.

### ARISE Lab (CUHK)

RESEARCH ASSISTANT

- Evaluating the Safety and Reliability of Multi-modal models and Large language models
- Developed the core idea of the first and second papers: distributing information across different modalities to circumvent AI moderation systems, revealing that only 10% of processed toxic content was flagged by AI in 2021.
- Led the majority of the experimental design and execution, dataset collection and generation, and authored the experimental section of the paper.
- Collaborated on creating a benchmark to evaluate the medical capabilities of Vision-Language Models (VLMs).
- Currently leading two students in their final year projects within the lab.

HONG KONG

May '22 - Present.

### Tencent AI Lab

RESEARCH COLLABORATOR

- Collaborated closely with researchers to produce a top-tier conference paper.
- Proposed the core idea: investigating the cultural biases of large language models across different languages.
- Designed the experiment, collecting all relevant datasets and establishing evaluation metrics.
- Completed the majority of the experimental work and authored the experimental and result section of the paper.

HONG KONG

Oct '22 - Sep '23.

## Publication

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

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

The Annual Meeting of the Association for Computational Linguistics

ACL 2024

(CCFA)

### A Picture is Worth a Thousand Toxic Words: A Metamorphic Testing Framework for Content Moderation Software

WENXUAN WANG, **JINGYUAN HUANG**, CHANG CHEN, PINJIA HE, JIAZHEN GU, MICHAEL R. LYU

The IEEE/ACM International Conference on Automated Software Engineering

ASE 2023

(CCFA)

### Validating Multimedia Content Moderation Software via Semantic Fusion

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

(CCFA)

### Asclepius: A Spectrum Evaluation Benchmark for Medical Multi-Modal Large Language Models

WENXUAN WANG\*, YIHANG SU, **JINGYUAN HUANG**, JIE LIU, WENTING CHEN, YUDI ZHANG, CHENG-YI LI, KAO-JUNG CHANG,

XIAOHAN XING, LINLIN SHEN, MICHAEL R. LYU

Under Review

### LLMs as GeoGuessr Masters: Exceptional Performance, Hidden Biases, and Privacy Risks

**JINGYUAN HUANG\***, JEN-TSE HUANG\*, ZIYI LIU, XIAOYUAN LIU, WENXUAN WANG, JIEYU ZHAO

Ongoing

## Skills

**Computer Skills**, Machine Learning, Deep Learning, Python, C, SQL

**Research Skills**, Self-motivated; skilled at summarizing papers and understanding key concepts.