

# Jiaxu Zhao


PhD Candidate,  
TU Eindhoven


 j.zhao@tue.nl


 +31 687 286 408

 zhao1402072392.github.io

## Languages

 Chinese

 English

 Dutch

## Skills

**Coding:** Python, C  
**Frameworks:** Keras, Pytorch, Mat-lab

## Research Interests

- Large Language Models
- Fairness and Bias in NLP
- AI Safety
- Causal Reasoning
- Sparse Neural Network
- Graph Neural Network
- Robustness
- Adversarial Attack

## About Me

I am passionate about fairness and robustness in NLP, with a strong interest in interdisciplinary col-laboration. I value diverse per-spectives and believe that respon-sible AI starts from inclusive de-sign and transparent evaluation.

Beyond the academic adventure, I also dive into photography, horse riding, squash, and badminton.

## Looking For

Actively seeking research scien-tist or applied scientist roles fo-cusing on NLP, LLMs, and AI Safety, in either industry research labs or academia.

## Education

- 2021.9 **PhD - Eindhoven University of Technology**
- Research Group: Mathematics & Computer Science, Data Mining
- 2025.9 Supervisors: Mykola Pechenizkiy, Meng Fang, Yulong Pei
- Research Focus: Fairness and Bias in Generative Language Models
- 2018.9 **Master - University of Electronic Science and Technology of China**
- Research Group: Computer Technology, CI Lab
- 2021.7 Supervisors: Xiaobin Wang, Hong Qu
- Thesis: The Research and Solution of Exposure Bias in Neural Machine Translation.

## Selected Publications

- Jiaxu Zhao, et al. Understanding Large Language Model Vulnerabilities to Social Bias At-tacks **ACL 2025**
- Jiaxu Zhao, et al. Unmasking Style Sensitivity: A Causal Analysis of Bias Evaluation Insta-bility in Large Language Models **ACL 2025**
- Jiaxu Zhao, et al. FS-GNN: Improving Fairness in Graph Neural Networks via Joint Sparsi-fication **Neurocomputing 2025**
- Qin Zhang, Sihan Cai, **Jiaxu Zhao**, Mykola Pechenizkiy, Meng Fang. CHAmbi: A New Bench-mark on Chinese Ambiguity Challenges for Large Language Models. **EMNLP findings 2024**
- Turbal, Bohdan, Anastasiia Mazur, **Jiaxu Zhao**, and Mykola Pechenizkiy. On Adversarial Robustness of Language Models in Transfer Learning. **NeurIPS 2024 SoLaR Workshop**
- Shenghui Li, Fanghua Ye, Meng Fang, **Jiaxu Zhao**, Yun-Hin Chan, Edith C-H Ngai, Thiemo Voigt. Synergizing Foundation Models and Federated Learning: A Survey **arXiv preprint arXiv:2406.12844**
- Jiaxu Zhao**, et al. More than Minorities and Majorities: Understanding Multilateral Bias in Language Generation. **ACL findings 2024**
- Jiaxu Zhao**, et al. CHBias: Bias Evaluation and Mitigation of Chinese Conversational Lan-guage Models. **ACL 2023**
- Jiaxu Zhao**, et al. Gptbias: A comprehensive framework for evaluating bias in large lan-guage models. **arXiv preprint arXiv:2312.06315**
- Jiaxu Zhao**, et al. Rest: Enhancing group robustness in DNNs through reweighted sparse training. **ECML PKDD 2023**

Tianjin Huang, Tianlong Chen, Meng Fang, Vlado Menkovski, **Jiaxu Zhao**, Lu Yin, Yulong Pei, Decebal Constantin Mocanu, Zhangyang Wang, Mykola Pechenizkiy, Shiwei Liu. You Can Have Better Graph Neural Networks by Not Training Weights at All: Finding Untrained GNNs Tickets. **LoG 2022 (Oral and Best Paper Award)**

## Internship Experience

- 2023.3 **Digital Brain Laboratory, Shanghai** Student Researcher
- **Project:** Chinese Large Language Models.
- 2023.6 **Responsibilities:** 1) Evaluated fairness and bias in Chinese LLMs using Statistical Parity, KL divergence, and group-wise accuracy across demographic groups. 2) Applied Proximal Policy Optimization (PPO) for Reinforcement Learning from Hu-man Feedback (RLHF) to guide behavior alignment.

## Selected Projects

- 2023.9 **Responsible AI Fall 2023 Research Program**
- Led by Prof. Julia Stoyanovich, New York University
- 2024.8 Co-supervise with Prof. Mykola Pechenizkiy
- Research Project:** Investigated robustness of LLMs under adversarial attacks post transfer learning. Found that larger models are more resilient to adversarial at-tacks, revealing trade-offs between adaptation and security.
- 2024.7 **OpenML AI Search Project**
- Led by OpenML and Prof. Joaquin Vanschoren
- 2024.8 **Research Project:** Developed a retrieval-augmented generation (RAG) framework for dataset recommendation. Integrated semantic similarity and query reformula-tion modules to enhance recommendation quality.
- 2022.11 **Supervised Project**
- Co-supervise with Prof. Meng Fang
- 2023.7 **Research Project:** Exploring the synergy within a context-aware and domain-flexible pipeline for neural machine translation.