

# Wanli Yang

## Curriculum Vitae

☎ (+86) 18759861072

✉ yangyywl@gmail.com

👤 WanliYoung

🏠 yangwl.site

🎓 Google Scholar

## Education

2024–Present **Institute of Computing Technology, Chinese Academy of Sciences**

**Ph.D. Candidate in Cyberspace Security**

- Research Advisor: Prof. Fei Sun & Prof. Xinran Liu

2020–2024 **Nankai University, College of Software**

**B.E. in Software Engineering**

- Rank: 5/140, CET4: 616, CET6: 614

## Research Experience

12/2023 **Reliable & Practical Evaluation of Model Editing**

–Present

- Advisor:** Prof. Fei Sun

- Institution:** Institute of Computing Technology, CAS

- Uncovering the potential of model editing to trigger LLMs collapse and proposing employing perplexity as a surrogate metric to monitor model collapse.

Accepted to **Findings of ACL 2024**.

- Revealing existing model editing evaluation adopts inappropriate strategies, such as teacher forcing during testing, which substantially overestimate the effectiveness of existing techniques.

Accepted to **Main Conference of ACL 2025**.

12/2022 **Graph Contrastive Learning for Recommendation**

–03/2023

- Advisor:** Prof. Wayne Xin Zhao

- Institution:** Gaoling School of Artificial Intelligence, RUC

- Contributing implementations of the latest recommendation models for **RecBole** (GitHub 3.3k Stars).

## Internship Experience

11/2023 **Reasons Behind LLMs Collapse Caused by Editing**

–04/2025

- Mentor:** Dr. Xinyu Ma

- Company:** Baidu Inc.

- Revealing the root causes behind the model collapse triggered by editing and proposes a straightforward solution to prevent collapse and achieve remarkable editing performance.

- Accepted to **Findings of EMNLP 2024**.

04/2025 **Effective Model Editing through Fine-tuning**

–Present

- Mentor:** Dr. Hongyu Zang

- Company:** Meituan Inc.

- Investigating the underlying reasons for the failure of fine-tuning in model editing and devising effective tuning strategies.

- Working on progress.

## Publications and Preprints

- ACL 2024 **Wanli Yang**, Fei Sun, Xinyu Ma, Xun Liu, Dawei Yin, Xueqi Cheng. The Butterfly Findings Effect of Model Editing: Few Edits Can Trigger Large Language Models Collapse.
- ACL 2024 Hexiang Tan, Fei Sun, **Wanli Yang**, Yuanzhuo Wang, Qi Cao, Xueqi Cheng. Blinded Main by Generated Contexts: How Language Models Merge Generated and Retrieved Contexts When Knowledge Conflicts?
- EMNLP 2024 **Wanli Yang**, Fei Sun, Jiajun Tan, Xinyu Ma, Du Su, Dawei Yin, Huawei Shen. The Findings Fall of ROME: Understanding the Collapse of LLMs in Model Editing.
- ACL 2025 **Wanli Yang**, Fei Sun, Jiajun Tan, Xinyu Ma, Qi Cao, Dawei Yin, Huawei Shen, Main Xueqi Cheng. The Mirage of Model Editing: Revisiting Evaluation in the Wild.

## Skills

- Programming Python, C/C++, Java, SQL
- Frameworks PyTorch, TensorFlow, JAX
- Others Git,  $\text{\LaTeX}$ , Linux

## Honors and Awards

- |             |   |   |
|-------------|---|---|
| 2024        | <b>Outstanding Graduates (Top 5%)</b>           | Nankai University                         |
| 2023        | <b>College Scholarship</b>                      | University of Chinese Academy of Sciences |
| 2022 & 2023 | <b>Academic Excellence Scholarship (Top 5%)</b> | Nankai University                         |
| 2021 & 2022 | <b>Arts and Sports Scholarship (Top 2%)</b>     | Nankai University                         |
| 2021        | <b>Innovation Scholarship (Top 3%)</b>          | Nankai University                         |

## Hobbies

- Sports Fitness, Swimming, Cycling, Dragon Boat Racing (former paddler for the Dragon Boat Team of Nankai University), etc.
- Reading Works of Haruki Murakami.