

Changye Li

📍 Peking, China ✉ antoine031106@gmail.com 📞 +86 13860472996 🔗 antoinegg1.github.io
📄 antoinegg1

Summary

I am a junior undergraduate at Peking University (Class of '26), currently focused on Large Language Model, with a particular interest in learning algorithms and data. My research interests also cover Reinforcement Learning and mechanistic interpretability. My research is driven by the following questions:

- How can the gap between artificial intelligence and human-level intelligence be quantified through advanced AI system, and how can it be bridged using learning algorithms?
- How can the essence of intelligence be revealed through the modeling of AI system?

Education

Peking University Sept 2022 – May 2026
B.S. Student in Artificial Intelligence

Fellowships & Awards

Peking University Freshman Scholarship (¥25000 RMB)(2022)

Research Experience

Visiting Student Researcher at PAIR Lab: PKU Alignment and Interaction Research Lab Peking, China
2023 –

Currently working on Alignment and Interpretability of Language Models under the guidance from Dr. Yaodong Yang.

Publications

Language Models Resist Alignment, *Accepted at Neurips 2024 SoLar Program* Oct 2024
Jiaming Ji, Kaile Wang, Tianyi Qiu, Boyuan Chen, **Changye Li**, Hantao Lou, Jiayi Zhou, Josef Dai, Yaodong Yang,

Towards efficient collaboration via graph modeling in reinforcement learning, *In Submission* August 2024
Wenzhe Fa, Zishun Yu, Chengdong Ma, **Changye Li**, Yaodong Yang, Xinhua Zhang,

Projects

Sparse Autoencoder in Vision [antoinegg1/SAELens-V](#) 🔗

- Developed an sparse autoencoder (SAE) training repository based on [jbloomAus/SAELens](#) 🔗. Realizing the training and utilizing process of the SAE on vision language model. (eg.Llava, Chameleon)
- Tools Used: Python
- This project is expected to be used in an ongoing paper.

SUMON [antoinegg1/SUMON](#) 🔗

- A simplified traffic simulation program refers to [eclipse-sumo/sumo](#) 🔗, supports reinforcement learning and multi-agent training.
- Tools Used: Python, XML

- The project was applied in paper, **Towards efficient collaboration via graph modeling in reinforcement.**