

YIBIN WANG

✉ wyb896409234@gmail.com · <https://yibinwang.netlify.app>

Forth-year undergraduate seeking a **Research Internship in Trustworthy AI**.

🎓 EDUCATION

Huazhong University of Science and Technology (HUST), Hubei, China 2019 – 2024

B.E. in Computer Science (CS), GPA: 3.81/4.00 - Transcripts

I got injured and took a one-year leave of absence from school in 2019, and now I am about to graduate.

👥 EXPERIENCE

Generalization and Robustness of LLMs (Remote Research Intern) Sep. 2023 – Present

Advisor: Hao Wang - Rutgers University

- Conducted extensive research on adversarial robustness of large language models.
- Conducted in-depth research on Bayesian algorithms and their applications in large language models.

Adversarial Robustness (Research Intern) Sep. 2021 – Aug. 2023

Advisor: Kun He - Huazhong University of Science and Technology

- Conducted extensive research on adversarial attack and defense in machine learning
- Conducted in-depth research on certified robustness based on convex relaxation

👥 PUBLICATIONS

* indicates equal contribution

Robustness-Aware Word Embedding Improves Certified Robustness to Adversarial Word Substitutions In Findings of ACL 2023

- Authors: **Yibin Wang***, Yichen Yang*, Di He, Kun He
- In my research on adversarial word substitution in NLP, I have discovered the crucial impact of word embeddings on certified adversarial robustness in convex relaxation-based methods. Given the findings, I transformed the optimization problem of the model's certified robustness into an optimization problem of word embeddings through theoretical proofs. I independently completed all coding, experiments, and the main part of the paper writing.

Continual Learning of Large Language Models: A Comprehensive Survey Preprint

- Authors: Haizhou Shi, Zihao Xu, Hengyi Wang, Weiyi Qin, Wenyan Wang, **Yibin Wang**, Hao Wang
- Responsible for writing the parts related to large language models in the Preliminaries section.

BLoB: Bayesian Low-Rank Adaptation by Backpropagation for Large Language Models Under review of NeurIPS 2024

- Authors: **Yibin Wang***, Haizhou Shi*, Ligong Han, Dimitris Metaxas, Hao Wang
- We introduce a principled Bayesian framework for improving large language models' generalization and uncertainty estimation. I contributed to the design of the algorithm and the writing of the paper, independently optimized the algorithm, implemented the code, and conducted the primary experiments.

📖 MISCELLANEOUS

- Languages: English - IELTS overall score 7.0, Mandarin - Native speaker