

# BOYI WEI

312 Sherrerd Hall, Princeton, NJ 08540

+1 (949) 678-3985 [wby@princeton.edu](mailto:wby@princeton.edu) [boyiwei.com](https://boyiwei.com) [github.com/boyiwei](https://github.com/boyiwei)

## Education

### Princeton University

*Doctor of Philosophy in Electrical and Computer Engineering*

Advisor: Peter Henderson

**August 2023 – Present**

*Princeton, NJ*

### Princeton University

*Master of Arts in Electrical and Computer Engineering*

**August 2023 – September 2025**

*Princeton, NJ*

### University of Science and Technology of China

*Bachelor of Science in Applied Physics (Summa Cum Laude)*

GPA: 4.00/4.30 (Top 3%)

**September 2019 – July 2023**

*Hefei, Anhui, China*

## Research Interest

My research focuses on building capable, reliable, and trustworthy language systems, including:

- Self-improving agent systems.
- Alignment in AI systems (large language models, AI agents), especially safety alignment.
- Interpretation, understanding, and mitigation of law/policy issues for language models.

## Publications and Preprints

1. **Boyi Wei\***, Benedikt Stroebel\*, Jiacen Xu, Joie Zhang, Zhou Li, Peter Henderson. Dynamic Risk Assessments for Offensive Cybersecurity Agents. *NeurIPS 2025 Datasets and Benchmarks*.
2. **Boyi Wei\***, Kaixuan Huang\*, Yangsibo Huang\*, Tinghao Xie, Xiangyu Qi, Mengzhou Xia, Prateek Mittal, Mengdi Wang, and Peter Henderson. Assessing the brittleness of safety alignment via pruning and low-rank modifications. *ICLR 2024 Set-LLM (Best Paper) / ICML 2024*.
3. **Boyi Wei\***, Weijia Shi\*, Yangsibo Huang\*, Noah A.Smith, Chiyuan Zhang, Luke Zettlemoyer, Kai Li, Peter Henderson. Evaluating Copyright Takedown Methods for Language Models. *NeurIPS 2024 Datasets and Benchmarks*.
4. **Boyi Wei\***, Zora Che\*, Nathaniel Li, Udari Madhushani Sehwag, Jasper Götting, Samira Nedungadi, Julian Michael, Summer Yue, Dan Hendrycks, Peter Henderson, Zifan Wang, Seth Donoughe, Mantas Mazeika. Best Practices for Biorisk Evaluations on Open-Weight Bio-Foundation Models. *NeurIPS 2025 Biosecurity Safeguards for Generative AI. arXiv preprint arXiv:2510.27629*
5. Xiangyu Qi\*, **Boyi Wei\***, Nicolas Carnili, Yangsibo Huang, Tinghao Xie, Luxi He, Matthew Jagielski, Milad Nasr, Prateek Mittal, Peter Henderson. On Evaluating the Durability of Safeguards for Open-Weight LLMs. *ICLR 2025*.
6. Jakub Lucki, **Boyi Wei**, Yangsibo Huang, Peter Henderson, Florian Tramèr, Javier Rando. An Adversarial Perspective on Machine Unlearning for AI Safety. *NeurIPS 2024 SoLaR Workshop (Best Paper) / TMLR*.
7. Tinghao Xie\*, Xiangyu Qi\*, Yi Zeng\*, Yangsibo Huang\*, Udari Madhushani Sehwag, Kaixuan Huang, Luxi He, **Boyi Wei**, Dacheng Li, Ying Sheng, Ruoxi Jia, Bo Li, Kai Li, Danqi Chen, Peter Henderson, Prateek Mittal. SORRY-Bench: Systematically Evaluating Large Language Model Safety Refusal Behaviors. *ICLR 2025*.

8. Shuai Shao, Qihan Ren, Chen Qian, **Boyi Wei**, Dadi Guo, Jingyi Yang, Xinhao Song, Linfeng Zhang, Weinan Zhang, Dongrui Liu, Jing Shao. Your Agent May Misevolve: Emergent Risks in Self-evolving LLM Agents. *arXiv preprint arXiv:2509.26354* (2025).
9. Sayash Kapoor, Benedikt Stroebel, Peter Kirgis, Nitya Nadgir, Zachary S Siegel, **Boyi Wei**, Tianci Xue, Ziru Chen, Felix Chen, Saiteja Utpala, Franck Ndzhomga, Dheeraj Oruganty, Sophie Luskin, Kangheng Liu, Botao Yu, Amit Arora, Dongyoong Hahn, Harsh Trivedi, Huan Sun, Juyong Lee, Tengjun Jin, Yifan Mai, Yifei Zhou, Yuxuan Zhu, Rishi Bommasani, Daniel Kang, Dawn Song, Peter Henderson, Yu Su, Percy Liang, Arvind Narayanan. Holistic Agent Leaderboard: The Missing Infrastructure for AI Agent Evaluation. *arXiv preprint arXiv:2510.11977* (2025).
10. Rui-Jie Zhu\*, Zixuan Wang\*, Kai Hua\*, Tianyu Zhang\*, Ziniu Li\*, Haoran Que\*, **Boyi Wei\***, Zixin Wen\*, Fan Yin\*, He Xing\*, Lu Li, Jiajun Shi, Kaijing Ma, Shanda Li, Taylor Kergan, Andrew Smith, Xingwei Qu, Mude Hui, Bohong Wu, Qiyang Min, Hongzhi Huang, Xun Zhou, Wei Ye, Jiaheng Liu, Jian Yang, Yunfeng Shi, Chenghua Lin, Enduo Zhao, Tianle Cai, Ge Zhang, Wenhao Huang, Yoshua Bengio, Jason Eshraghian. Scaling Latent Reasoning via Looped Language Models. *arXiv preprint:2510.25741* (2025)
11. Xiangyu Qi, Yangsibo Huang, Yi Zeng, Edoardo Debenedetti, Jonas Geiping, Luxi He, Kaixuan Huang, Udari Madhushani Sehwag, Vikash Sehwag, Weijia Shi, **Boyi Wei**, Tinghao Xie, Danqi Chen, Pin-Yu Chen, Jeffrey Ding, Ruoxi Jia, Jiaqi Ma, Arvind Narayanan, Weijie J. Su, Mengdi Wang, Chaowei Xiao, Bo Li, Dawn Song, Peter Henderson, Prateek Mittal. AI Risk Management Should Incorporate Both Safety And Security. *arXiv preprint:2405.19524* (2024).

## Experience

<b>Scale AI</b>	<b>May 2025 – September 2025</b>
<i>Research Scientist Intern</i> (Advisor: Nathaniel Li, Zifan Wang, Julian Michael)	<i>San Francisco, CA</i>
<b>University of California, Irvine</b>	<b>March 2023 – June 2023</b>
<i>Research Intern</i> (Advisor: Sitao Huang)	<i>Irvine, CA</i>

  

<b>Georgia Institute of Technology</b>	<b>January 2022 – November 2022</b>
<i>Research Intern</i> (Advisor: Cong Hao)	<i>Atlanta, GA</i>

## Teaching

<b>COS 568: Systems and Machine Learning</b>	<b>Spring 2025</b>
--	--------------------

## Honors and Awards

<b>Francis Robbins Upton Fellowship</b>	<b>September 2023</b>
<b>China National Scholarship</b>	<b>September 2021, September 2022</b>
<b>USTC Outstanding Student Scholarship (Gold Prize)</b>	<b>October 2021</b>
<b>Yan Jici Outstanding Student Scholarship</b>	<b>November 2021</b>

## Services

### Reviews:

- Conference: ICLR (2025, 2026), NeurIPS (2025)
- Workshop: SeT-LLM (ICLR 2024), SoLaR (NeurIPS 2024), L2M2 (ACL 2025), DIG-BUG (ICML 2025)

**Tutorials:** LLMs and Copyright Risks: Benchmarks and Mitigation Approaches (AAAI 2025 / NAACL 2025)