

BOYI WEI

312 Sherrerd Hall, Princeton, NJ 08540

☎ +1 (949) 678-3985 ✉ wby@princeton.edu 🌐 boyiwei.com 🐙 [GitHub](#) 🌐 [LinkedIn](#)

Education

Princeton University

August 2023 – Present

Doctor of Philosophy in Electrical and Computer Engineering

Princeton, NJ

Advisor: Peter Henderson

Princeton University

August 2023 – September 2025

Master of Arts in Electrical and Computer Engineering

Princeton, NJ

University of Science and Technology of China

September 2019 – July 2023

Bachelor of Science in Applied Physics (Summa Cum Laude)

Hefei, Anhui, China

GPA: 4.00/4.30 (Top 3%)

Research Interest

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

- Self-evolving 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*, Jiachen 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 Sehwan, 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*.
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. Hadas Orgad, **Boyi Wei**, Kaden Zheng, Martin Wattenberg, Peter Henderson, Seraphina Goldfarb-Tarrant, Yonatan Belinkov. Large Language Models Generate Harmful Content Using a Unified Mechanism. *ICLR 2026 Workshop on Logical Reasoning of Large Language Models*.
8. Tinghao Xie*, Xiangyu Qi*, Yi Zeng*, Yangsibo Huang*, Udari Madhushani Sehwan, 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*.

9. 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 Misedevolve: Emergent Risks in Self-evolving LLM Agents. *ICLR 2026*.
10. Sayash Kapoor, Benedikt Stroebel, Peter Kirgis, Nitya Nadgir, Zachary S Siegel, **Boyi Wei**, Tianci Xue, Ziru Chen, Felix Chen, Saiteja Utpala, Franck Ndzomga, Dheeraj Oruganty, Sophie Luskin, Kangheng Liu, Botao Yu, Amit Arora, Dongyoon Hahm, 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. *ICLR 2026*.
11. 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)*
12. 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

Scale AI	May 2025 – September 2025
<i>Research Scientist Intern</i> (Advisor: Nathaniel Li, Zifan Wang, Julian Michael)	<i>San Francisco, CA</i>
University of California, Irvine	March 2023 – June 2023
<i>Research Intern</i> (Advisor: Sitao Huang)	<i>Irvine, CA</i>
Georgia Institute of Technology	January 2022 – November 2022
<i>Research Intern</i> (Advisor: Cong Hao)	<i>Atlanta, GA</i>

Teaching

COS 568: Systems and Machine Learning	Spring 2025
--	--------------------

Honors and Awards

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

Services

Reviews:

- Conference: ICML (2026), ICLR (2025, 2026), NeurIPS (2025, 2026)
- 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)