

Wenxin Ding

Research Interest

My research interest lies in machine learning security and privacy. Specifically, I focus on bridging the gap between theoretical understanding and empirical practice. My research studies the safety behavior of machine learning models under strategically optimized training data. Recently, I have been working on problems regarding vulnerabilities of text-to-image diffusion models and developing tools for content creators against copyright infringement.

Education

University of Chicago	Chicago, IL
Ph.D. in Computer Science (4.0 / 4.0)	2021.9 – 2026.6
Carnegie Mellon University	Pittsburgh, PA
M.S. in Computer Science – Research Thesis (4.11 / 4.33)	2020.8 – 2021.8
B.S. in Computer Science and B.S. in Mathematical Sciences (3.92 / 4.0)	2016.8 – 2020.5
Minor in Computational Finance	

Work Experience

Research Assistant, University of Chicago (Chicago, IL)	2021.9 – Present
Advisors: Prof. Heather (Haitao) Zheng and Prof. Ben Y. Zhao	
<ul style="list-style-type: none">• Design effective data poisoning attacks against text-to-image generative models as a tool to protect content creators from copyright infringement• Develop analytical framework to quantify the change in diffusion model performance under data poisoning attacks• Propose and implement methodology to train robust model versions by selecting augmented training data• Theoretically formulate the optimal loss of robust multi-class classifiers and derive theorems to approximate the loss	
Research Engineer, Qualcomm (San Diego, CA)	2025.6 – 2025.9
Mentors: Dr. Jonathan Petit and Dr. Cong Chen	
<ul style="list-style-type: none">• Evaluate black-box and white-box jailbreak attacks against vision language models (VLMs)• Develop effective jailbreak attacks against VLMs• Build new benchmark datasets for evaluating VLM safety commissioned by MLCommons	
Research Assistant, Carnegie Mellon University (Pittsburgh, PA)	2019.2 – 2021.8
Advisors: Prof. Nihar B. Shah and Prof. Weina Wang	
<ul style="list-style-type: none">• Design mathematical modeling for score calibration of peer-review data• Apply differential privacy to protect anonymity of reviewer identity• Derive Pareto-optimal calibration method to trade-off data utility and user privacy	

Peer-Reviewed Publications

Conferences

- **Wenxin Ding**, Cathy Li, Shawn Shan, Ben Y. Zhao, Haitao Zheng. “Understanding Implosion in Text-to-Image Generative Models.” in *Proceedings of ACM SIGSAC Conference on Computer and Communications Security (CCS)*, 2024.
- Shawn Shan, **Wenxin Ding**, Josephine Passananti, Haitao Zheng, Ben Y. Zhao. “Prompt-Specific Poisoning Attacks on Text-to-Image Generative Models.” in *Proceedings of IEEE Symposium on Security and Privacy (S&P)*, 2024.

- **Wenxin Ding**, Arjun Nitin Bhagoji, Ben Y. Zhao, and Haitao Zheng. “Towards Scalable and Robust Model Versioning.” in *Proceedings of IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, 2024.
- Sihui Dai*, **Wenxin Ding***, Arjun Nitin Bhagoji, Daniel Cullina, Ben Y. Zhao, Haitao Zheng, and Prateek Mittal. “Characterizing the Optimal 0-1 Loss for Multi-class Classification with a Test-time Attacker.” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2023. **Spotlight paper.** (* for equal contribution)
- Shawn Shan, **Wenxin Ding**, Emily Wenger, Haitao Zheng, and Ben Y. Zhao. “Post-breach recovery: Protection against white-box adversarial examples for leaked DNN models.” in *Proceedings of ACM SIGSAC Conference on Computer and Communications Security (CCS)*, 2022.
- **Wenxin Ding**, Gautam Kamath, Weina Wang, and Nihar B. Shah. “Calibration with privacy in peer review.” in *Proceedings of IEEE International Symposium on Information Theory (ISIT)*, 2022.

Workshops

- Wenxin Ding, Nihar B. Shah, and Weina Wang. “On the privacy-utility tradeoff in peer-review data analysis.” *AAAI Privacy-Preserving Artificial Intelligence (PPAI) workshop*, 2021. **Spotlight paper.**

Teaching Experience

Teaching Assistant

University of Chicago

- CMSC 25800 Adversarial Machine Learning
- CMSC 25300/35300 Mathematical Foundations of Machine Learning

Carnegie Mellon University

- 15110 Principles of Computing (Head Teaching Assistant)
- 15213 Introduction to Computer Systems
- 15440 Distributed Systems

Mentor

- Strong Women Strong Girls, Pittsburgh, PA

Service

Technical Program Committee

- 2025 ACM Conference on Computer and Communications Security (CCS)
- 2025, 2026 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)
- 2024 ACM Workshop on Artificial Intelligence and Security (AISec)

Reviewer

- Nature
- 2025 ACM SIGGRAPH Asia
- 2024, 2025 The Conference on Uncertainty in Artificial Intelligence (UAI)
- SIAM Journal on Mathematics of Data Science (SIMODS)

Awards

- 2025 University of Chicago William Rainey Harper Dissertation Fellowship
- 2024 University of Chicago UU Fellowship
- 2021 University of Chicago Eckhardt Scholar
- 2020 Carnegie Mellon University Senior Leadership Recognition
- 2019 Mark Stehlik SCS Alumni Undergraduate Impact Scholarship
- 2017 William Lowell Putnam Mathematical Competition (Rank: 255 / 4638)