

Wenxiao Wang

✉ www@relai.ai | Head of AI at RELAI, Inc

Education

University of Maryland, College Park

PH.D. IN COMPUTER SCIENCE (DISSERTATION: TOWARDS RELIABLE AGENTIC LLMs [\[url\]](#))

Advisor: Prof. [Soheil Feizi](#)

College Park, US

Sept. 2021 - 2025

Institute for Interdisciplinary Information Sciences (a.k.a. Yao class), Tsinghua University

B.ENG. IN COMPUTER SCIENCE AND TECHNOLOGY.

Peking, China

Sept. 2016 - June. 2020

Experience

RELA, Inc

HEAD OF AI

Maryland, US

July 2025 - Present

Sony AI

RESEARCH INTERN MENTORED BY **WEIMING ZHUANG** AND **LINGJUAN LYU**

Remote, US

May. 2023 - Aug. 2023

Bytedance

RESEARCH INTERN MENTORED BY **LINJIE YANG**, **HENG WANG** AND **YU TIAN**

Remote, US

June 2022 - Nov. 2022

Institute for Interdisciplinary Information Sciences, Tsinghua University

RESEARCH ASSISTANT MENTORED BY PROF. **HANG ZHAO**

Peking, China

Sep. 2020 - Aug. 2021

University of California, Berkeley

VISITING STUDENT RESEARCHER ADVISED BY **XINYUN CHEN**, **RUOXI JIA** AND PROF. **DAWN SONG**

Berkeley, US

Apr. 2019 - Aug. 2019

Bytedance AI Lab

INTERN IN VISUAL SEARCH GROUP MENTORED BY **YI HE** AND **LEI LI**

Peking, China

May. 2018 - Nov. 2018

Preprints

Maestro: Joint Graph & Config Optimization for Reliable AI Agents [\[url\]](#)

WENXIAO WANG, PRIYATHAM KATTAKINDA, SOHEIL FEIZI

2025

Chain-of-Defensive-Thought: Structured Reasoning Elicits Robustness in Large Language Models against Reference Corruption [\[url\]](#)

WENXIAO WANG, PARSA HOSSEINI, SOHEIL FEIZI

2025

Publications

Tool Preferences in Agentic LLMs are Unreliable [\[url\]](#)

KAZEM FAGHIH*, **WENXIAO WANG***, YIZE CHENG*, SIDDHANT BHARTI, GAURANG SRIRAMANAN, SRIRAM
BALASUBRAMANIAN, PARSA HOSSEINI, SOHEIL FEIZI (*EQUAL CONTRIBUTION)
Conference on Empirical Methods in Natural Language Processing (EMNLP)

2025

DyePack: Provably Flagging Test Set Contamination in LLMs Using Backdoors [\[url\]](#)

YIZE CHENG*, **WENXIAO WANG***, MAZDA MOAYERI, SOHEIL FEIZI (*EQUAL CONTRIBUTION)
Conference on Empirical Methods in Natural Language Processing (EMNLP)

2025

Can AI-Generated Text be Reliably Detected? [\[url\]](#)

VINU SANKAR SADASIVAN, AOUNON KUMAR, SRIRAM BALASUBRAMANIAN, **WENXIAO WANG**, SOHEIL FEIZI
Transactions on Machine Learning Research (TMLR)
Media Coverage: [\[Washington Post\]](#) [\[Wired\]](#) [\[New Scientist\]](#) [\[The Register\]](#) [\[TechSpot\]](#) [\[UMD Science\]](#)

2025

Robustness of AI-Image Detectors: Fundamental Limits and Practical Attacks [\[url\]](#)

MEHRDAD SABERI, VINU SANKAR SADASIVAN, KEIVAN REZAEI, AOUNON KUMAR, ATOOSA CHEGINI, **WENXIAO WANG**,
SOHEIL FEIZI
International Conference on Learning Representations (ICLR)
Media Coverage: [\[Wired\]](#) [\[MIT Tech Review\]](#) [\[Bloomberg News\]](#) [\[The Register\]](#)

2024

DRSM: De-Randomized Smoothing on Malware Classifier Providing Certified Robustness [\[url\]](#)

SHOUMIK SAHA, **WENXIAO WANG**, YIGITCAN KAYA, SOHEIL FEIZI, TUDOR DUMITRAS
International Conference on Learning Representations (ICLR)

2024

Temporal Robustness against Data Poisoning [\[url\]](#)

WENXIAO WANG, SOHEIL FEIZI
Conference on Neural Information Processing Systems (NeurIPS)

2023

Spuriousity Rankings: Sorting Data for Spurious Correlation Robustness [\[url\]](#)

MAZDA MOAYERI, **WENXIAO WANG**, SAHIL SINGLA, SOHEIL FEIZI
Conference on Neural Information Processing Systems (NeurIPS)[\[spotlight\]](#)

2023

Lethal Dose Conjecture on Data Poisoning [\[url\]](#)

WENXIAO WANG, ALEXANDER LEVINE, SOHEIL FEIZI
Conference on Neural Information Processing Systems (NeurIPS)

2022

Improved Certified Defenses against Data Poisoning with (Deterministic) Finite Aggregation [\[url\]](#)

WENXIAO WANG, ALEXANDER LEVINE, SOHEIL FEIZI
International Conference on Machine Learning (ICML)

2022

On Feature Decorrelation in Self-Supervised Learning [\[url\]](#)

TIANYU HUA*, **WENXIAO WANG***, ZIHUI XUE, SUCHENG REN, YUE WANG, HANG ZHAO
(*EQUAL CONTRIBUTION)
International Conference on Computer Vision (ICCV)[\[oral\]](#)

2021

DPLis: Boosting Utility of Differentially Private Deep Learning via Randomized Smoothing [\[url\]](#)

WENXIAO WANG, TIANHAO WANG, LUN WANG, NANQING LUO, PAN ZHOU, DAWN SONG, RUOXI JIA
Privacy Enhancing Technologies Symposium (PETS)

2021

REFIT: A Unified Watermark Removal Framework For Deep Learning Systems With Limited Data [\[url\]](#)

XINYUN CHEN*, WENXIAO WANG*, YIMING DING, CHRIS BENDER, RUOXI JIA, BO LI, DAWN SONG
(*EQUAL CONTRIBUTION)
ACM Asia Conference on Computer and Communications Security (AsiaCCS)

2021

The Secret Revealer: Generative Model Inversion Attacks Against Deep Neural Networks [\[url\]](#)

YUHENG ZHANG*, RUOXI JIA*, HENGZHI PEI, WENXIAO WANG, BO LI, DAWN SONG
(*EQUAL CONTRIBUTION)
Conference on Computer Vision and Pattern Recognition (CVPR)[\[oral\]](#)

2020

Leveraging Unlabeled Data for Watermark Removal of Deep Neural Networks [\[url\]](#)

XINYUN CHEN*, WENXIAO WANG*, YIMING DING, CHRIS BENDER, RUOXI JIA, BO LI, DAWN SONG
(*EQUAL CONTRIBUTION)
ICML2019 Workshop on Security and Privacy of Machine Learning

2019

Talks

- **Temporal Robustness against Data Poisoning**, AI TIME Youth PhD Talk, November 2023.
- **Lethal Dose Conjecture: From Few-shot Learning to Potentially Nearly Optimal Defenses against Data Poisoning**, TMLR Group, Hong Kong Baptist University, December 2022.
- **Lethal Dose Conjecture on Data Poisoning**, AI TIME Youth PhD Talk, November 2022.
- **Improved Certified Defenses against Data Poisoning with (Deterministic) Finite Aggregation**, AI TIME Youth PhD Talk, August 2022.

Services

Program Committee / Reviewer of: NeurIPS, ICML, ICLR, ICCV, CVPR, TPAMI, TMLR, ...

Awards

Gold Medal(4th place)	National Olympiad in Informatics	2015
Gold Medal(1st place)	Asia and Pacific Informatics Olympiad in China District	2015
Gold Medal(10th place)	China Team Selection Competition	2015
Gold Medal	National Olympiad in Informatics	2014
Bronze Medal	Asia and Pacific Informatics Olympiad in China District	2014
Silver Medal	China Team Selection Competition	2014

Teaching

- Teaching Assistant of *CMSC828W: Foundations of Deep Learning*, Fall 2022, University of Maryland, College Park.

- Teaching Assistant of *CMSC422: Introduction to Machine Learning*, Spring 2022, University of Maryland, College Park.
- Teaching Assistant of *CMSC351: Algorithms*, Fall 2021, University of Maryland, College Park.