

Yao Tong

Ph.D. Candidate • Trustworthy ML, Copyright, Privacy Auditing for LLMs

tongyao@u.nus.edu • [Google Scholar](#) • Singapore

Research Interests

Trustworthy machine learning; copyright for data and models; privacy leakage auditing; memorization and generalization; ai security.

Education

Ph.D. in Computer Science

2022 – Present

National University of Singapore • Singapore

- Advisor: Prof. Reza Shokri.

B.S. in Computer Science

2018 – 2022

The Chinese University of Hong Kong • China

- Graduated with First-class Honors.

Preprints

2. Decomposing Extrapolative Problem Solving: Spatial Transfer and Length Scaling with Map Worlds

Yao Tong, Jiayuan Ye, Anastasia Borovykh, Reza Shokri

Under review at ICLR, 2025

1. Identifying Optimal Output Sets for Differential Privacy Auditing

Yao Tong, Jiayuan Ye, Reza Shokri

Publications & Workshops

6. When Transformers Can (or Can't) Generalize Compositionally? A Data-Distribution Perspective

Yao Tong, Jiayuan Ye, Anastasia Borovykh, Reza Shokri

NeurIPS Workshop on What Can('t) Transformers Do?, 2025

5. SeedPrints: Fingerprints Can Even Tell Which Seed Your Large Language Model Was Trained From

Yao Tong*, Haonan Wang*, Siquan Li, Kenji Kawaguchi, Tianyang Hu

NeurIPS Workshop on Prevent Unauthorized Knowledge Use from Large Language Models, 2025

Under review at ICLR, 2025

4. Cut the Deadwood Out: Training-Free Backdoor Purification via Guided Module Substitution

Yao Tong*, Weijun Li*, Xuanli He, Haolan Zhan, Qionghai Xu

In Findings of Association for Computational Linguistics EMNLP, 2025

3. **How much of my dataset did you use? Quantitative Data Usage Inference in Machine Learning**

Yao Tong*, Jiayuan Ye*, Sajjad Zarifzadeh, Reza Shokri

In International Conference of Learning Representations (ICLR), 2025

Oral Presentation (Top ~1.5% among submissions)

2. **The Stronger the Diffusion Model, the Easier the Backdoor: Data Poisoning to Induce Copyright Breaches Without Adjusting Finetuning Pipeline**

Haonan Wang, Qianli Shen, **Yao Tong**, Yang Zhang, Kenji Kawaguchi

NeurIPS Workshop on Backdoors in Deep Learning, 2023

Oral Presentation

In International Conference on Machine Learning (ICML), 2024

Oral Presentation (Top ~2% among submissions)

1. **Towards Regulatable AI Systems: Technical Gaps and Policy Opportunities**

Xudong Shen, Hannah Brown, Jiashu Tao, Martin Strobel, **Yao Tong**, Akshay Narayan, Harold Soh, Finale Doshi-Velez

Communications of the ACM (CACM), 2024

Work conducted during Finale's visit to NUS as an outcome of the RRAI Workshop; middle authors order is alphabetical.

* denotes equal contribution.

Privacy Meter

An open-source library to audit data privacy in statistical and machine learning algorithm via membership inference. 2025

Open-source • https://github.com/privacytrustlab/ml_privacy_meter

- Implemented privacy auditing tools such as DUCI and RMIA.
- Contributed to the development and long-term maintenance of the library as one of the organizers.

Teaching

Teaching Assistant, CS5562 Trustworthy Machine Learning

National University of Singapore

2023 Fall

Teaching Assistant, CS3244 Machine Learning

National University of Singapore

2024 Spring

Teaching Assistant, CS6208 Advanced Topics in Artificial Intelligence

National University of Singapore

2024 Fall

Teaching Assistant, Data and Knowledge Management, Software Engineering

The Chinese University of Hong Kong

2021-2022

Honors & Awards

President Graduate Fellowship - NUS

2024

Dean's List - CUHK

2019-2022

University Research Award - CUHK

2021, 2022

School Academic Scholarship (for Top 2% students) - CUHK

2020-2022

Bowen Scholarship - *CUHK*

2018-2022

Service

Reviewer: NeurIPS 2025, ICLR 2025, ICML Workshop 2025

Sub-reviewer: CCS 2024, USENIX Security 2024