

Jaechul Roh

jrohsc.github.io · [Github](#) · [Google Scholar](#)

+1 (470) 915 - 1137 · jroh@umass.edu

EDUCATION

University of Massachusetts Amherst

Ph.D. in Computer Science

Advisor: Prof. [Amir Houmansadr](#)

GPA: 4.0/4.0

September 2023 – Present

Amherst, Massachusetts, USA

Hong Kong University of Science and Technology

B.Eng. in Computer Engineering, School of Engineering

Final Year Thesis Advisor: Prof. [Jun Zhang](#)

*2 years of Compulsory Korean Military Duty

September 2017 – May 2023

Clear Water Bay, Hong Kong

RESEARCH INTERESTS

My research is centered on **Privacy & Security** of AI models and agents. I have been conducting research on the trustworthiness of multi-modal generative models across various domains, including text, vision, and audio modalities, under the supervision of Prof. Amir Houmansadr. I am currently doing my Summer Research Internship at Brave Software working on AI web-agent privacy under the guidance of Dr. Ali Shahin Shamsabadi.

PUBLICATIONS

Preprints

1. **Bob's Confetti: Phonetic Memorization Attacks in Music and Video Generation**

Jaechul Roh, Zachary Novack, Yuefeng Peng, Niloofar Miresghallah, Taylor Berg-Kirkpatrick, Amir Houmansadr

Preprint at arXiv

[\[paper\]](#) [\[demo page\]](#)

2. **Chain-of-Code Collapse: Reasoning Failures in LLMs via Adversarial Prompting in Code Generation**

Jaechul Roh, Varun Gandhi, Shivani Anilkumar, Arin Garg

Preprint at arXiv

[\[paper\]](#) [\[code\]](#)

3. **R1dacted: Investigating Local Censorship in DeepSeek's R1 Language Model**

Ali Naseh, Harsh Chaudhari, **Jaechul Roh**, Mingshi Wu, Alina Oprea, Amir Houmansadr

Preprint at arXiv

[\[paper\]](#)

4. **OverThink: Slowdown Attacks on Reasoning LLMs**

Abhinav Kumar, **Jaechul Roh**, Ali Naseh, Marezna Karpinska, Mohit Iyyer, Amir Houmansadr, and Eugene Bagdasarian

Preprint at arXiv

[\[paper\]](#) [\[code\]](#)

5. **FameBias: Embedding Manipulation Bias Attack in Text-to-Image Models**

Jaechul Roh^{*}, Andrew Yuan^{*}, Jinsong Mao^{*}

*Equal Contribution**

Preprint at arXiv

[\[paper\]](#)

6. **Understanding (Un)Intended Memorization in Text-to-Image Generative Models**

Ali Naseh, **Jaechul Roh**, Amir Houmansadr

Preprint at arXiv

[\[paper\]](#)

Conference

1. **Multilingual and Multi-Accent Jailbreaking of Audio LLMs**

Jaechul Roh, Virat Shejwalkar, Amir Houmansadr

COLM 2025

[\[paper\]](#)

2. **Backdooring Bias (B^2) into Stable Diffusion Models**
Ali Naseh, **Jaechul Roh**, Eugene Bagdasarian, Amir Houmansadr
USENIX Security '25
[\[paper\]](#) [\[code\]](#)
3. **OSLO: One-Shot Label-Only Membership Inference Attacks**
Yuefeng Peng, **Jaechul Roh**, Subhransu Maji, Amir Houmansadr
NeurIPS 2024
[\[paper\]](#)
4. **Memory Triggers: Unveiling Memorization in Text-To-Image Generative Models through Word-Level Duplication**
Ali Naseh, **Jaechul Roh**, Amir Houmansadr
The 5th AAAI Workshop on Privacy-Preserving Artificial Intelligence
[\[paper\]](#)
5. **Robust Smart Home Face Recognition under Starving Federated Data**
Jaechul Roh, Yajun Fang
IEEE International Conference on Universal Village (IEEE UV2022)
Oral Presentation
[\[paper\]](#)[\[code\]](#)[\[slides\]](#)[\[video\]](#)
6. **MSDT: Masked Language Model Scoring Defense in Text Domain**
Jaechul Roh, Minhao Cheng, Yajun Fang
IEEE International Conference on Universal Village (IEEE UV2022)
Oral Presentation
[\[paper\]](#)[\[code\]](#)[\[slides\]](#)[\[video\]](#)
7. **Impact of Adversarial Training on the Robustness of Deep Neural Networks**
Jaechul Roh
2022 IEEE 5th International Conference on Information Systems and Computer Aided Education (ICISCAE)
[\[paper\]](#)[\[code\]](#)

INVITED TALKS

Google Speech Technologies Group

Paper presentation

July 2025

Google DeepMind

- Presented our "*Multilingual and Multi-Accent Jailbreaking of Audio LLMs*" paper to the NFM Reading Group led by the Speech Technologies Group at Google DeepMind.
[\[slides\]](#)

WORK EXPERIENCE

Brave Software

Research Intern, Supervisor: Ali Shahin Shamsabadi

June 2025 – September 2025

London, United Kingdom (Remote)

- Working on privacy & security of AI agents.

Super Chain AI (Conard International)

NLP Engineer Intern, Supervisor: Cat Yung

June 2021 – August 2021

Kowloon Bay, Hong Kong

- In charge of topic modeling and semantic analysis based on customer reviews and assigning specific semantics to the topics extracted.
- Competitors' analysis through web-scraping customer reviews from other drop-shipping websites.

Military Service at Head Quarter of 12th Infantry Division

Sergeant of Republic of Korea Army

July 2018 – March 2020

Injae, Kang Won Do, Republic of Korea

- Officer Administrative Clerk Specialist
- Squad Leader of the Head Quarter

PROFESSIONAL SERVICES

Program Committee Member (Reviewer)

- **Main conferences:** ICLR (2025)

SKILLS / LANGUAGES

Programming Language: Python, C++

Languages: Korean (Native), English (Native), Chinese (Fluent)