

# Changsheng SUN

Email: cssun@u.nus.edu

---

I am a Ph.D. candidate at the **PLSE Lab, NUS School of Computing**, working at the intersection of **Trustworthy AI, LLMs, and Graph Learning**. My research aims to make learning systems reliable and auditable by developing principled methods for robustness against adversarial threats and distribution shifts. Currently, I focus on **LLM safety and explainable graph modeling**, with specific interests in directionality- and risk-aware explanations. My work spans NLP, geometric learning, and large-scale time-series data analysis.

**Research Interests: Trustworthy AI, Graph Representation Learning, LLM Safety, Explainability (XAI).**

## LINKS

LinkedIn: [linkedin.com/in/changshengsun/](https://linkedin.com/in/changshengsun/)

GitHub: [github.com/cs-sun](https://github.com/cs-sun)

HomePage: [sunchangsheng.com](http://sunchangsheng.com) ( <http://comp.nus.edu.sg/~sun97/> )

Google Scholar: [scholar.google.com/citations?user=1M5\\_gpIAAAAJ](https://scholar.google.com/citations?user=1M5_gpIAAAAJ)

## EDUCATION

**National University of Singapore** Jan 2022 - Present (Expected Graduation: Early 2026)  
Doctor of Philosophy, Computer Science  
- Thesis: Trustworthy Geometric Learning: From Structural Biases to Risk-Aware Robustness  
- Advisor: [Prof. Dong Jin Song](#)

**National University of Singapore** Aug 2020 - Dec 2021  
Master of Computing, Artificial Intelligence

**Xidian University, Xi'an, China** Sep 2015 - Aug 2019  
Bachelor of Engineering, Computer Science

## WORK EXPERIENCES

**NUS School of Computing** Research Assistant Jan 2026 - Present  
- Working on AISG project: Trustworthy Geometric Learning;  
- Host: Prof. Dong Jin Song and Dr. Zhang Yedi.

Research Assistant July 2021 - Dec 2021  
- Built and maintained a reproducible evaluation pipeline for trustworthiness assessment of deep learning systems, emphasizing uncertainty-aware reliability analysis;  
- Contributed to work published at ASE 2022 (InputReflector).  
- Host: Prof. Dong Jin Song and Prof. Xiao Yan

**NUS-Singtel Cyber Security R&D Lab** Research Intern Jan 2020 - June 2020  
- Developed a white-box testing approach for deep neural networks by leveraging intermediate-layer signals and density estimation for reliability assessment.  
- Contributed to work published at ICSE 2021 (Self-Checking Deep Neural Networks).  
- Host: [Prof. David S. Rosenblum](#)

**JD Intelligent Cities Research, JD.com** Research Intern & Algorithm Engineer Jan 2018 - June 2018  
- Prototyped a distributed DQN research system on spatiotemporal data using the Ray execution framework; focused on training scalability and system performance profiling.  
- Host: [Prof. Yu Zheng](#) and Dr. Junbo Zhang

## SKILLS

Programming Languages: Python; Frameworks: PyTorch, TensorFlow, JAX, Ray, TorchGeometric.

## SERVICES & AWARDS

**Program Committee / Reviewer:** NeurIPS, ICML, AAAI, WWW, FSE, ASE (2022-2026).

**Teaching Assistant:** CS5232 Formal Specification and Design Techniques (2024); CS4218 Software Testing (2023); CS4211 Formal Methods for Software Engineering (2022)

**Supported and awarded by:** Graduate Student Travel Grants (IEEE S&P 2025, AAAI 2024, ASE 2022); NUS Research Scholarship (2022–2026).

## PUBLICATIONS

1. Ignoring Directionality Leads to Compromised Graph Neural Network Explanations. **Changsheng Sun**, Xinkle Li, Jin Song Dong. *2025 IEEE Security and Privacy: Workshops (SPW)*. April 2025.
2. PointCVaR: Risk-optimized Outlier Removal for Robust 3D Point Cloud Classification. Xinkle Li, Junchi Lu, Henghui Ding, **Changsheng Sun**, Joey Tianyi Zhou, Chee Yeow Meng. *Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI)*, 2024. Jan 2024.
3. Repairing Failure-inducing Inputs with Input Reflection. Yan Xiao, Yun Lin, Ivan Beschastnikh, **Changsheng Sun**, David S. Rosenblum, Jin Song Dong. *The 37th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2022. Oct 2022
4. PIANO: Influence Maximization Meets Deep Reinforcement Learning. Hui Li, Mengting Xu, Sourav S Bhowmick, Joty Shafiq Rayhan, **Changsheng Sun**, Jiangtao Cui. *IEEE Transactions on Computational Social Systems*. 2022.
5. Self-Checking Deep Neural Networks in Deployment. Yan Xiao, Ivan Beschastnikh, David S. Rosenblum, **Changsheng Sun**, Sebastian Elbaum, Y. Lin, Jin Song Dong. *The 43rd IEEE/ACM International Conference on Software Engineering (ICSE)*. Aug 2021.
6. Digraph Inception Convolutional Networks. Zekun Tong, Yuxuan Liang, **Changsheng Sun**, Xinkle Li, David S. Rosenblum, Andrew Lim. *Thirty-fourth Conference on Neural Information Processing Systems (NeurIPS)*, May 2020.

### Preprints:

7. Generalizing Neural Networks by Reflecting Deviating Data in Production. Yan Xiao, Yun Lin, Ivan Beschastnikh, **Changsheng Sun**, David S. Rosenblum, Jin Song Dong. *arXiv:2110.02718*.
8. Directed Graph Convolutional Network. Zekun Tong, Yuxuan Liang, **Changsheng Sun**, David S. Rosenblum, Andrew Lim. *arXiv:2004.13970*.
9. Disco: Influence Maximization Meets Network Embedding and Deep Learning. Hui Li, Mengting Xu, Sourav S Bhowmick, **Changsheng Sun**, Zhongyuan Jiang, Jiangtao Cui. *arXiv:1906.07378*.

### Under Review:

10. From Attack Surfaces to Actual Operations: A Survey of Modern LLM Jailbreaks. Ruikang Zhou, **Changsheng Sun**, Mark Huasong Meng. *Under Review*.
11. Risk-Aware Robust Graph Network Explanation. **Changsheng Sun**, Xinkle Li, Jin Song Dong. *Under Review*.