

# Jiawei ZHANG

UserID: jiaweiz7@illinois.edu  $\diamond$  [javyduck.github.io](https://javyduck.github.io)  
(+1) 217-200-3511  $\diamond$  Unit 1129, Thomas M. Siebel Center for Computer Science, Urbana, IL 61801

## EDUCATION

---

### University of Illinois Urbana-Champaign (UIUC)

Ph.D. in Computer Science, GPA: 4.0/4.0.

Advisor: *Prof. Bo Li*

*August. 2023 – May. 2027 (Expected)*

### University of Illinois Urbana-Champaign (UIUC)

Master of Science in Computer Science (Research-Oriented), GPA: 4.0/4.0.

*August. 2021 – May. 2023*

### Zhejiang University (ZJU)

Bachelor of Engineering (Excellent Class) [Rank: 6/130].

Hangzhou, China

*Sept. 2017 – Jun. 2021*

## RESEARCH INTEREST:

---

My current research predominantly centers on **trustworthy large language models (LLMs)**. I'm particularly interested in enhancing their trustworthiness by mitigating issues like hallucination, using external knowledge sources as leverage. While my foundation in **robustness, privacy, fairness, and explainability** remains intact, my renewed focus aims at the integration of these principles into the development and understanding of LLMs, thereby ensuring they align more closely with human values and expectations.

## PUBLICATION

---

- **Jiawei Zhang**, Zhongzhu Chen, Huan Zhang, Chaowei Xiao, Bo Li. DiffSmooth: Certifiably Robust Learning via Diffusion Models and Local Smoothing. *32th USENIX Security Symposium 2023*. [[paper](#)]
- **Jiawei Zhang**, Linyi Li, Ce Zhang, Bo Li. CARE: Certifiably Robust Learning with Reasoning via Variational Inference. *IEEE Conference on Secure and Trustworthy Machine Learning (SatML) 2023*. [[arxiv](#)]
- Zhuolin Yang\*, Zhikuan Zhao\*, Boxin Wang, **Jiawei Zhang**, Linyi Li, Hengzhi Pei, Bojan Karlas, Ji Liu, Heng Guo, Ce Zhang, Bo Li. Improving Certified Robustness via Statistical Learning with Logical Reasoning. *Advances in Neural Information Processing Systems (NIPS) 2022*. [[arxiv](#)]
- Linyi Li, **Jiawei Zhang**, Tao Xie, Bo Li. Double Sampling Randomized Smoothing. *International Conference on Machine Learning (ICML) 2022*. [[arxiv](#)]
- **Jiawei Zhang\***, Linyi Li\*, Huichen Li, Xiaolu Zhang, Shuang Yang, Bo Li. Progressive-Scale Boundary Blackbox Attack via Projective Gradient Estimation. *International Conference on Machine Learning (ICML) 2021*. [[arxiv](#)]

## RESEARCH EXPERIENCE

---

### Detecting the Hallucination for LLM

Research Assistant | Secure Learning Lab, UIUC

*Sept. 2023 – Present*

Advised by Prof. Bo Li

- Developed methods to evaluate the self-consistency, knowledge consistency, and logical consistency of text generated by LLMs.
- Leverage a retrieval system to externally cross-check claims made in LLM responses against trusted knowledge bases, ensuring the validity of generated content.

### Safety-Critical Driving Scenario Generation Based on LLM

Research Assistant | Secure Learning Lab, UIUC

*Oct. 2022 – March. 2023*

Advised by Prof. Bo Li

- Aim to enrich the safety-critical testing scenarios in SafeBench [[link](#)] for Autonomous Vehicles.
- Train an adversarial agent (vehicle/pedestrian/bicyclist) via specifically designed multi-agent reinforcement learning to cause the unexpected collision of the ego vehicle.

### Enhance Robustness via Diffusion Models and Local Smoothing

Research Assistant | Secure Learning Lab, UIUC

*Jun. 2022 – Oct. 2022*

Advised by Prof. Bo Li & Postdoc. Huan Zhang (CMU)

- Prove that the “one-shot” denoising of DDPM can approximate the mean of the generated posterior distribution by continuous-time diffusion models, which is an approximation of the original instance under mild conditions.
- Propose a local smoothing technique based on the diffusion models, achieve the **SOTA 43.6%** certified accuracy on CIFAR-10 under  $\ell_2$  radius 1.0 and the **SOTA 53.0%** certified accuracy on ImageNet under the  $\ell_2$  radius 1.5.

## Certifiably Robust Learning with Reasoning via Variational Inference May. 2022 – Sep. 2022

Research Assistant | Secure Learning Lab, UIUC

Advised by Prof. Bo Li & Prof. Ce Zhang (ETH Zürich)

- Propose a scalable and certifiably robust learning with reasoning pipeline CARE, which is able to integrate knowledge rules to enable reasoning ability for reliable prediction
- Propose an efficient Expectation Maximization (EM) algorithm to approximate the reasoning based on Markov Logic Network (MLN) via variational inference using Graph Convolutional Network (GCN).
- Extensive experiments on different datasets show that the proposed method achieves significantly higher certified robustness than SOTA baselines, for example, the certified accuracy could be improved from 36.0% (SOTA) to 61.8% under  $\ell_2$  radius 2.0 on AWA2.

## Boundary Blackbox Attack via Projection Based Gradient Estimation Jun. 2020 – May. 2021

Research Intern | Cooperate with Ant Financial

Advised by Prof. Bo Li

- Propose the first theoretical framework to analyze boundary blackbox attacks with general projection functions.
- Characterize the key characteristics and trade-offs for a good projective gradient estimator.
- Propose Progressive-Scale based projective Boundary Attack via progressively searching for the optimal scale in a self-adaptive way under spatial, frequency, and spectrum scales.
- The extensive experiments show that our method outperforms the state-of-the-art boundary attacks on MNIST, CIFAR-10, CelebA, and ImageNet against different blackbox models and an online API (MEGVII Face++).

## 6DoF Pose Estimation Sep. 2018 – Apr. 2019

Research Intern | State Key Laboratory of CAD&CG, ZJU

Advised by Prof. Xiaowei Zhou

- Use Unity to build an aircraft carrier deck and estimate the 6DoF pose of the moving planes on it with PVNet.
- Employ a coarse-to-fine prediction scheme to predict per voxel likelihoods for each human joint by ConvNet.

## SELECTED COURSE PROJECTS

---

### Verification of Neural Networks Based on DeepPoly Sep. 2021 – Nov. 2021

Team Leader | Course - Logic and Artificial Intelligence, UIUC

Advised by Prof. Gagandeep Singh

- Propose a much more efficient certification method based on the widely used bound-propagation algorithm DeepPoly (CROWN-IBP) via a recursive refinement of the linear constraints.
- Achieve the **Top-1** certification performance in class, and achieve 100% verification accuracy on all testing cases.

## INTERNSHIP

---

### Sea AI Lab (Singapore) May 2023 – Present

Research Intern

Advised by Dr. Tianyu Pang & Dr. Chao Du

- Conducted evaluations on cross-modal models (e.g., Stable Diffusion, Whisper) to assess their consistency under a range of common data corruptions.
- Developed a rigorous benchmark for assessing the self-consistency of these models. The benchmark was designed to provide a comprehensive understanding of model behavior, incorporating a wide range of scenarios and inputs to measure their resilience and accuracy.

## TEACHING

---

### CS 307 - Modeling and Learning in Data Science (Spring 2022)

- Teaching Assistant with Prof. Bo Li and Prof. David Forsyth

## SELECTED HONOR & AWARDS

---

Meritorious Winner, MCM COMAP's Mathematical Contest in Modeling Apr. 2020

First Prize, China Harbour Scholarship (1/130) Jan. 2020

First Prize, The Chinese Mathematics Competitions (non-math major), Zhejiang Province Nov. 2018

First Prize, National High School Mathematics League, Zhejiang Province Nov. 2016

## SKILLS & OTHERS

---

**Computer Languages** Python, C/C++, Java, Matlab, Shell script, Markdown

**Frameworks & Packages** PyTorch, TensorFlow, OpenCV, MySQL, Hadoop, Unity, SPSS