

# ZHAORUN CHEN

University of Chicago, 5801 S Ellis Ave, Chicago, IL 60637, USA

Homepage: [billchan226.github.io](https://billchan226.github.io) | [zhaorun@uchicago.edu](mailto:zhaorun@uchicago.edu)

## EDUCATION

<b>Shanghai Jiao Tong University (SJTU)</b> School of Electronic Information and Electrical Engineering (SEIEE)	<b>Shanghai, China</b>
<b>Bachelor of Engineering in Automation</b>	<i>Sept. 2018 – June 2022</i>
<b>Purdue University</b> Elmore Family School of Electrical and Computer Engineering	<b>West Lafayette, IN, US</b>
<b>Master in Electrical and Computer Engineering (2 year program, full-funded)</b>	<i>Aug. 2022 – Aug. 2023</i>
<b>University of Chicago</b> Department of Computer Science.	<b>Chicago, IL, US</b>
<b>PhD in Computer Science</b>	<i>Sept. 2024 – June 2030</i>

## PUBLICATIONS

- [1] **Chen, Z.**, Zhao, Z., Luo, H., Yao, H., Li, B., & Zhou, J. (2024). HALC: Object Hallucination Reduction via Adaptive Focal-Contrast Decoding, in Proceeding of the Forty-first International Conference on Machine Learning (ICML 2024), Vienna, Austria, July 2024.
- [2] **Chen, Z.**, Zhao, Z., Zhu, Z., Zhang, R., Li, X., Raj, B., & Yao, H. (2024). AutoPRM: Automating Procedural Supervision for Multi-Step Reasoning via Controllable Question Decomposition, in Proceeding of 2024 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2024), Mexico City, Mexico, Jun 2024.
- [3] **Chen, Z.**, Zhao, Z., He, T., Chen, B., Zhao, X., Gong, L., & Liu C. (2024). Safe Reinforcement Learning via Hierarchical Adaptive Chance-Constraint Safeguards, in Proceeding of 2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2024), Abu Dhabi, UAE, October 2024.

## PREPRINTS

- [4] **Chen, Z.**, Xiang, Z., Xiao, C., Song, D., & Li, B. (2024). AgentPoison: Red-teaming LLM Agents via Memory or Knowledge Base Backdoor Poisoning.
- [5] **Chen, Z.**, Du, Y., ..., Rafailov, R., Finn, C., & Yao, H. (2024). MJ-Bench: Is Your Multimodal Reward Model Really a Good Judge?

## RESEARCH EXPERIENCE

**AgentPoison: Red-teaming LLM Agents via Memory or Knowledge Base Backdoor Poisoning** Feb. 2024 – June 2024

**Advisor:** Prof. Bo Li, CS, UIUC/University of Chicago; Prof. Chaowei Xiao, University of Wisconsin, Madison

- Propose the first backdoor attack framework targeting generic and RAG-based LLM agents by poisoning their long-term memory or knowledge base; propose an integral trigger optimization algorithm to map triggered instances to unique embedding space to achieve high backdoor attack success rate while preserving agent utility in non-triggered cases;
- On three types of real-world LLM agents (e.g. autonomous driving agent, QA agent, healthcare EHRAgent), AgentPoison achieves an average ASR  $\geq 80\%$  with minimal impact on benign performance ( $\leq 1\%$ ) with a minor poison ratio  $< 0.1\%$ .

**HALC: Object Hallucination Reduction via Adaptive Focal-Contrast Decoding**

Aug. 2023 – Jan. 2024

**Advisor:** Prof. Bo Li, CS, UIUC/University of Chicago; Prof. Jiawei Zhou, CS, TTIC

- Propose the key pattern to mitigating object hallucination is identifying the optimal visual context to safely predict each token; propose an adaptive decoding algorithm HALC which employ a sampling based interpolated JSD-filter method to select proxies to approximate the optimal visual context and then perform contrastive decoding to reveal their context-specific information.
- HALC outperforms all SOTA (e.g. OPERA, VCD, woodpecker, LURE) in terms of CHAIR, POPE and MME benchmarks.

**Automating Procedural Supervision for LLMs Reasoning via Controllable Problem Decomposition** May 2023–Aug. 2023

**Advisor:** Prof. Huaxiu Yao, CS, UNC Chapel Hill; Prof. Bhiksha Raj, CS, Carnegie Mellon University

- Propose AutoPRM, a process-supervision pipeline that automatically decomposes long-form reasoning into sub-questions and assign fine-grained RL rewards to intermediate steps to enhance internal logic consistency and transfer of meta-knowledge.

- To allow fine-grained control during question solving, we propose a user-controlled switch to manage decomposition granularity and introduce FIM-guided-decoding to guide each subproblem solver towards solution of the primary problem. AutoPRM improves the accuracy on GSM8K (+4.4%), MATH (+4.2%) and StrategyQA (+6.2%) tasks over SOTA without external data.

INTERNSHIP EXPERIENCE

---

Guardrail Models for Video Content Moderation   Research Intern	June 2024 – Present
Advisor: Prof. Bo Li, VirtueAI	
Self-supervised Fine-grained Feedback for LLMs RLHF   Research Intern	June 2023 – Feb. 2024
Advisor: Prof. Huaxiu Yao, UNC-Chapel Hill	
Motion Planning for Decentralized Multi-manipulator Assembly System   AI Intern	Feb. 2022 - Aug. 2022
Advisor: Prof. Cewu Lu, Flexiv Robotics Inc.	

SKILLS

---

- **English Standard Test:** TOEFL: 111 (R: 28, L: 27, S: 28, W: 28)
- **Programming Languages:** solid expertise in Python, Matlab, C++, and various algorithms and data structures
- **Tools:** PyTorch, Matlab, LaTeX, TensorFlow, OpenCV, ROS, Linux, QT5, MySQL

AWARDS & HONOR

---

Full-funded two-year Guaranteed Fellowship at Purdue ECE (Very Rare)	2022
Best paper award in IEEE ISRIMT 2021	2021
Third Prize in the Contemporary Undergraduate Mathematical Contest in Modeling	2020
First Prize in RoboCup China Open	2020
Academic Excellence Scholarship (Top 5%)	2019-2022

SERVICE

---

Graduate Teaching Assistant for Data Structure (ECE368)	2022-2023
Conference Reviewer: NeurIPS'24, ICLR'24, COLM'24, ARR'24, IROS'24	