

# ZHAORUN CHEN

(+1) 773-952-0790 Email: [zhaorun@uchicago.edu](mailto:zhaorun@uchicago.edu) [◇ Homepage](#) [◇ Google Scholar](#)

## EDUCATION

---

<b>University of Chicago</b> Ph.D. in Computer Science – Advisor: Prof. <a href="#">Bo Li</a>	2024.09 - now
<b>Purdue University</b> M.S in Computer Engineering	2022.08 - 2023.08
<b>Shanghai Jiao Tong University</b> B.E in Computer Engineering	2018.09 - 2022.06

## WORK EXPERIENCES

---

<b>Virtue AI</b> <b>Tech Lead of Agent Red-Teaming, Founding Member of Technical Staff</b> - Led the development of Virtue AgentSuite, Virtue Agent ForgingGround, Virtue Action Guard, Virtue AgentGateway	2024.01 - now
<b>Meta</b> <b>Research Scientist Intern working with <a href="#">Jason Weston</a> at Meta Superintelligence Labs</b> - Led the research of the first synthetic training framework for general AI agents ( <a href="#">DreamGym</a> , <a href="#">Early Experience</a> )	2025.06 - 2025.10

## PUBLICATIONS & PREPRINTS

---

Full publication list is in [Google Scholar](#).

- [1] Kai Zhang, Xiangchao Chen, Bo Liu, Tianci Xue, Zeyi Liao, Zhihan Liu, Xiyao Wang, Yuting Ning, **Zhaorun Chen**, Xiaohan Fu, Jian Xie, Yuxuan Sun, Boyu Gou, Qi Qi, Zihang Meng, Jianwei Yang, Ning Zhang, Xian Li, Ashish Shah, Dat Huynh, Hengduo Li, Zi Yang, Sara Cao, Lawrence Jang, Shuyan Zhou, Jiacheng Zhu, Huan Sun, Jason Weston, Yu Su, Yifan Wu, [Agent Learning via Early Experience Authors](#), in Proceedings of the Forty-Third International Conference on Machine Learning (**ICML 2026**), Seoul, South Korea, July 2026. [[Paper](#)] [[Agent RL](#)]
- [2] Nuoya Xiong, Yuhang Zhou, Hanqing Zeng, **Zhaorun Chen**, Furong Huang, Shuchao Bi, Lizhu Zhang, Zhuokai Zhao, [Token-Level LLM Collaboration via FusionRoute](#), in Proceedings of the Forty-Third International Conference on Machine Learning (**ICML 2026**), Seoul, South Korea, July 2026. [[Paper](#)] [[LLM Infra](#)]
- [3] **Zhaorun Chen**, Zhuokai Zhao, Kai Zhang, Bo Liu, Qi Qi, Yifan Wu, Tarun Kalluri, Sara Cao, Yuanhao Xiong, Haibo Tong, Huaxiu Yao, Hengduo Li, Jiacheng Zhu, Xian Li, Dawn Song, Bo Li, Jason Weston, Dat Huynh, [Scaling Agent Learning via Experience Synthesis](#), in Proceedings of the 14th International Conference on Learning Representations (**ICLR 2026**), Rio de Janeiro, Brazil, April 2026. [[Paper](#)] [[Agent RL](#)]
- [4] **Zhaorun Chen**, Xun Liu, Mintong Kang, Jiawei Zhang, Minzhou Pan, Shuang Yang, Bo Li, [ARMs: Adaptive Red-Teaming Agent against Multimodal Models with Plug-and-Play Attacks](#), in Proceedings of the 14th International Conference on Learning Representations (**ICLR 2026**), Rio de Janeiro, Brazil, April 2026. [[Paper](#)] [[Red-Teaming](#)][[Multi-Modal](#)]
- [5] **Zhaorun Chen**, Zichen Wen, Yichao Du, Yiyang Zhou, Chenhang Cui, Zhenzhen Weng, Haoqin Tu, Chaoqi Wang, Zhengwei Tong, Qinglan Huang, Canyu Chen, Qinghao Ye, Zhihong Zhu, Yuqing Zhang, Jiawei Zhou, Zhuokai Zhao, Rafael Rafailov, Chelsea Finn, Huaxiu Yao, [MJ-Bench: Is Your](#)

- Multimodal Reward Model Really a Good Judge for Text-to-Image Generation?, in Proceeding of the Thirty-Ninth Conference on Neural Information Processing Systems (**NeurIPS 2025**), San Diego, the United States, Dec 2025. [Paper] [Code] [RL Post-Training]
- [6] **Zhaorun Chen**, Mintong Kang, Bo Li, **ShieldAgent: Shielding Agents via Verifiable Safety Policy Reasoning**, in Proceedings of the Forty-Second International Conference on Machine Learning (**ICML 2025**), Vancouver, Canada, July 2025. [Paper] [Code] [Agent Safety] [Guardrail]
- [7] **Zhaorun Chen**, Francesco Pinto, Minzhou Pan, Bo Li, **SafeWatch: An Efficient Safety-Policy Following Video Guardrail Model with Transparent Explanations**, in Proceedings of the 13th International Conference on Learning Representations (**ICLR 2025**), Singapore, Apr 2025. [Paper] [Code] [Video Safety Reasoning] [RL Post-Training]
- [8] **Zhaorun Chen**, Zhen Xiang, Chaowei Xiao, Dawn Song, Bo Li, **AgentPoison: Red-teaming LLM Agents via Poisoning Memory or Knowledge Bases**, in Proceeding of the Thirty-Eighth Conference on Neural Information Processing Systems (**NeurIPS 2024**), Vancouver, Canada, Dec 2024. [Paper] [Code] [LLM Agent Safety]
- [9] **Zhaorun Chen**, Zhuokai Zhao, Hongyin Luo, Huaxiu Yao, Bo Li, Jiawei Zhou, **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. [Paper] [Code] [Multimodal Hallucination]
- [10] Chengquan Guo, Chulin Xie, Yu Yang, **Zhaorun Chen**, Zinan Lin, Xander Davies, Yarin Gal, Dawn Song, Bo Li, **RedCodeAgent: Automatic Red-teaming Agent against Diverse Code Agents**, in Proceedings of the 14th International Conference on Learning Representations (**ICLR 2026**), Rio de Janeiro, Brazil, April 2026. [Paper] [Red-Teaming][Coding Agent]
- [11] Zijian Zhang, Kaiyuan Zheng, **Zhaorun Chen**, Joel Jang, Yi Li, Chaoqi Wang, Mingyu Ding, Dieter Fox, and Huaxiu Yao, **GRAPE: Generalizing Robot Policy via Preference Alignment**, in Proceeding of the IEEE International Conference on Robotics and Automation (**ICRA 2026**), Vienna, Austria, June 2026. [Paper] [Code] [RL Post-Training] [Robotics]
- [12] Siwei Han, Haonian Ji, Siyang Xin, Juanquan Shi, Shi Qiu, Xinyu Ye, Peng Xia, Jiaqi Liu, **Zhaorun Chen**, Yiyang Zhou, Linjie Li, Lijuan Wang, Huaxiu Yao, **Paper2Figure: A Multi-Agent Collaborative System for Figure Generation Towards Academic Research Paper**, in Proceeding of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (**CVPR 2026**), Denver, CO, June 2026. [Paper] [Code] [Image Generation]
- [13] Zichen Wen, Jiashu Qu, **Zhaorun Chen**, Dongrui Liu, Zhiyuan Liu, Ruixi Wu, Yicun Yang, Xiangqi Jin, Haoyun Xu, Xuyang Liu, Weijia Li, Chaochao Lu, Jing Shao, Conghui He, Linfeng Zhang, **The Devil behind the mask: An emergent safety vulnerability of Diffusion LLMs (ICLR 2026)**, Rio de Janeiro, Brazil, April 2026. [Paper] [Red-Teaming][Diffusion LLMs]
- [14] Mintong Kang, **Zhaorun Chen**, Bo Li, **C-SafeGen: Certified Safe LLM Generation with Claim-Based Streaming Guardrails**, in Proceeding of the Thirty-Ninth Conference on Neural Information Processing Systems (**NeurIPS 2025**), San Diego, the United States, Dec 2025. [Paper] [Code] [Guardrail]
- [15] Mintong Kang, **Zhaorun Chen**, Chejian Xu, Jiawei Zhang, Chengquan Guo, Minzhou Pan, Ivan Revilla, Yu Sun, Bo Li, **PolyGuard: Massive Multi-Domain Safety Policy-Grounded Guardrail Dataset**, in Proceeding of the Thirty-Ninth Conference on Neural Information Processing Systems (**NeurIPS 2025**), San Diego, the United States, Dec 2025. [Paper] [Code] [Guardrail] [Dataset]
- [16] Haibo Tong, Zhaoyang Wang, **Zhaorun Chen**, Haonian Ji, Shi Qiu, Siwei Han, Zhongkai Xue, Yiyang Zhou, Peng Xia, Kexin Geng, Mingyu Ding, Rafael Rafailov, Chelsea Finn, Huaxiu Yao, **MJ-Video: Fine-Grained Benchmarking and Rewarding Video Preferences in Video Generation**, in Proceeding of the Thirty-Eighth Conference on Neural Information Processing Systems (**NeurIPS 2025**), San Diego, the United States, Dec 2025. [Paper] [Code] [RL Post-Training]
- [17] Andy Zhou, Kevin Wu, Francesco Pinto, **Zhaorun Chen**, Yi Zeng, Yu Yang, Shuang Yang, Sanmi Koyejo, James Zou, Bo Li, **AutoRedTeamer: Autonomous Red Teaming with Lifelong Attack In-**

- tegration, in Proceeding of the Thirty-Ninth Conference on Neural Information Processing Systems (**NeurIPS 2025**), San Diego, the United States, Dec 2025. [Paper] [Code] [Red-Teaming]
- [18] Shaobo Wang, Yantai Yang, Guo Chen, Peiru Li, Kaixin Li, Yufa Zhou, **Zhaorun Chen**, Linfeng Zhang, [Grounding and Enhancing Informativeness and Utility in Dataset Distillation \(ICLR 2026\)](#), Rio de Janeiro, Brazil, April 2026. [Paper] [Multi-Modal][Data Distillation]
- [19] Haonian Ji, Shi Qiu, Siyang Xin, Siwei Han, **Zhaorun Chen**, Dake Zhang, Hongyi Wang, Huaxiu Yao, [From EduVisBench to EduVisAgent: A Benchmark and Multi-Agent Framework for Reasoning-Driven Pedagogical Visualization \(ICLR 2026\)](#), Rio de Janeiro, Brazil, April 2026. [Paper] [Multi-Modal][Benchmark]
- [20] Yixiong Fang, Ziran Yang, **Zhaorun Chen**, Zhuokai Zhao, Jiawei Zhou, [Enhancing Vision-Language Model Reliability with Uncertainty-Guided Dropout Decoding](#), in Proceeding of the Thirty-Ninth Conference on Neural Information Processing Systems (**NeurIPS 2025**), San Diego, the United States, Dec 2025. [Paper] [Code] [Hallucination]
- [21] Zichen Wen, Shaobo Wang, Yufa Zhou, Junyuan Zhang, Qintong Zhang, Yifeng Gao, **Zhaorun Chen**, Bin Wang, Weijia Li, Conghui He, Linfeng Zhang, [Efficient Multi-modal Large Language Models via Progressive Consistency Distillation](#), in Proceeding of the Thirty-Ninth Conference on Neural Information Processing Systems (**NeurIPS 2025**), San Diego, the United States, Dec 2025. [Paper] [Code] [Multimodal Understanding]
- [22] Yiming Zhang, Zhuokai Zhao, **Zhaorun Chen**, Zhili Feng, Zenghui Ding, Yining Sun, [RankClip: Ranking-Consistent Language-Image Pretraining](#), in Proceeding of the International Conference on Computer Vision (**ICCV 2025**), Honolulu, the United States, Dec 2025. [Paper] [Code] [Multimodal Pre-training]
- [23] Yiming Zhang, Zhuokai Zhao, **Zhaorun Chen**, Zenghui Ding, Xianjun Yang, Yining Sun, [Beyond Training: Dynamic Token Merging for Zero-Shot Video Understanding](#), in Proceeding of the International Conference on Computer Vision (**ICCV 2025**), Honolulu, the United States, Dec 2025. [Paper] [Code] [Multimodal Understanding]
- [24] Chejian Xu, Jiawei Zhang, **Zhaorun Chen**, Chulin Xie, Mintong Kang, Zhuowen Yuan, Zidi Xiong, Chenhui Zhang, Lingzhi Yuan, Yi Zeng, Peiyang Xu, Chengquan Guo, Andy Zhou, Jeffrey Ziwei Tan, Zhun Wang, Alexander Xiong, Xuandong Zhao, Yu Gai, Francesco Pinto, Yujin Potter, Zhen Xiang, Zinan Lin, Dan Hendrycks, Dawn Song, Bo Li, [MMDT: Decoding the Trustworthiness and Safety of Multimodal Foundation Models](#), in Proceedings of the 13th International Conference on Learning Representations (**ICLR 2025**), Singapore, Apr 2025. [Paper] [Code] [Multi-modal Safety]
- [25] Peng Xia, Siwei Han, Shi Qiu, Yiyang Zhou, Zhaoyang Wang, Wenhao Zheng, **Zhaorun Chen**, Chenhang Cui, Mingyu Ding, Linjie Li, Lijuan Wang, Huaxiu Yao, [MMIE: Massive Multimodal Interleaved Comprehension Benchmark for Large Vision-Language Models](#), in Proceedings of the 13th International Conference on Learning Representations (**ICLR 2025**) (**Oral Presentation**), Singapore, Apr 2025. [Paper] [Code] [Multi-modal Reasoning]
- [26] Yiyang Zhou, Zhaoyang Wang, Tianle Wang, Shangyu Xing, Peng Xia, Bo Li, Kaiyuan Zheng, Zijian Zhang, **Zhaorun Chen**, Wenhao Zheng, Xuchao Zhang, Chetan Bansal, Weitong Zhang, Ying Wei, Mohit Bansal, Huaxiu Yao, [AnyPrefer: An Automatic Framework for Preference Data Synthesis](#), in Proceedings of the 13th International Conference on Learning Representations (**ICLR 2025**), Singapore, Apr 2025. [Paper] [RL Post-Training]
- [27] Chenhang Cui, An Zhang, Yiyang Zhou, **Zhaorun Chen**, Gelei Deng, Huaxiu Yao, Tat-Seng Chua, [Fine-Grained Verifiers: Preference Modeling as Next-token in Vision-Language Alignment](#), in Proceedings of the 13th International Conference on Learning Representations (**ICLR 2025**), Singapore, Apr 2025. [Paper] [RL Post-Training]
- [28] Yiyang Zhou, Zhiyuan Fan, Dongjie Cheng, Sihan Yang, **Zhaorun Chen**, Chenhang Cui, Xiyao Wang, Yun Li, Linjun Zhang, Huaxiu Yao, [Calibrated Self-Rewarding Vision Language Models](#), in Proceeding of the Thirty-Eighth Conference on Neural Information Processing Systems (**NeurIPS**

- 2024), Vancouver, Canada, Dec 2024. [Paper] [Code] [RL Post-Training]
- [29] **Zhaorun Chen**, Zhuokai Zhao, Zhihong Zhu, Ruiqi Zhang, Xiang Li, Bhiksha Raj, Huaxiu Yao, [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. [Paper] [Code] [RL Post-Training]
- [30] Zhihong Zhu, Kefan Shen, **Zhaorun Chen**, Yunyan Zhang, Yuyan Chen, Xiaoqi Jiao, Zhongwei Wan, Shaorong Xie, Wei Liu, Xian Wu, Yefeng Zheng, [DGLF: A Dual Graph-based Learning Framework for Multi-modal Sarcasm Detection](#), in Proceeding of 2024 Conference on Empirical Methods in Natural Language Processing (**EMNLP 2024**), Miami, Florida, Nov 2024. [Multimodal Safety]
- [31] Siyue Wang, **Zhaorun Chen**, Zhuokai Zhao, Chaoli Mao, Yiyang Zhou, Jiayu He, Albert Sibo Hu, [EscIRL: Evolving Self-Contrastive IRL for Trajectory Prediction in Autonomous Driving](#), in Proceeding of 8th Annual Conference on Robot Learning (**CoRL 2024**), Munich, Germany, Nov 2024. [Paper] [Code] [RL for Robotics]
- [32] **Zhaorun Chen**, Zhuokai Zhao, Tairan He, Binhao Chen, Xuhao Zhao, Liang Gong, Chengliang Liu, [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. [Paper] [Code] [Safe RL] [Robotics]

## ACADEMIC SERVICES

---

### Conference Area Chair

- EMNLP 2025

### Conference Reviewer

- NeurIPS, ICLR, COLM, ARR, IROS 2024
- ICLR, CVPR, ICML, NeurIPS 2025

### Organizer

- NeurIPS CLAS 2024: The Competition for LLM and Agent Safety [Link] 2024
- COLM 2025 Workshop on AI Agents: Capabilities and Safety (AIA 2025) [Link] 2025