

# Zizhao Wang

zizhao.wang@utexas.edu · <https://wangzizhao.github.io/> · google scholar: <https://tinyurl.com/zizhaowangscholar>

## Research focus

- LLM post-training
- Reinforcement Learning (RL)
- Causal Reasoning

## Skills

- LLM
  - LLM Agents
  - RL post-training
  - Reasoning
  - Safety
- decision making
  - model-based RL
  - RLHF
  - offline RL
  - hierarchical RL
  - imitation learning
  - planning
- representation learning
- generalization and robustness
- interpretability and explainability
- ML frameworks (PyTorch, TensorFlow, vLLM, Transformers)
- distributed training (deepspeed, PEFT)
- Python
- data structure, algorithms

## EDUCATION

2020 - 25	<b>PhD</b> , Electrical and Computer Engineering expected graduation: Dec, advisor: <b>Peter Stone</b>	<b>University of Texas at Austin</b>
2018 - 19	<b>MS</b> , Computer Science	<b>Columbia University</b>
2016 - 18	<b>BS</b> , Computer Engineering (dual degree program)	<b>University of Michigan</b>
2014 - 18	<b>BS</b> , Electrical and Computer Engineering	<b>Shanghai Jiao Tong University</b>

## WORK EXPERIENCE

2025/03	<b>Student Researcher</b>	<b>Google</b>
	<ul style="list-style-type: none"><li>• Designed an adversarial <b>RL post-training</b> framework that enhance <b>LLM agent</b> security again prompt injections, by co-training two LLMs: an attacker that learns to create diverse prompt injections and an agent that learns to defend against them.</li><li>• Implemented the data collection pipeline with vLLM and parallel simulation environments for fast LLM agent rollout inference.</li><li>• Fine-tuned the LLM model with the GRPO algorithm, implemented with Transformer, deepspeed, and LoRA for fast and memory-efficient training.</li><li>• Reduced the attack success rate by 21% and improve task success rate by 18% compared to the untrained model.</li></ul>	
2024/06	<b>Research Intern</b>	<b>Microsoft Research</b>
	<ul style="list-style-type: none"><li>• Designed an <b>generative world model</b> that can synthesize images of novel scenarios, by using object-centric representations and disentangled representations.</li><li>• Enhanced the generalization of <b>reinforcement learning</b> policies by 30%, when learning with generated out-of-distribution data.</li></ul>	
2024/01	<b>Research Intern</b>	<b>Honda Research Institute</b>
	<ul style="list-style-type: none"><li>• Developed a <b>motion prediction</b> algorithm for <b>autonomous driving</b> that, reduced prediction error by 48%, by applying <b>causal reasoning</b> to vehicle interactions.</li><li>• Sped up model training with <b>distributed training</b> and <b>efficient CUDA implementations</b> for sparse attention.</li></ul>	

## RESEARCH EXPERIENCE

2021-22	<b>Causal World Model (ICML oral, AAAI oral)</b>	<b>University of Texas at Austin</b>
	<ul style="list-style-type: none"><li>• Developed a <b>world model</b> that can analyzes <b>causal relationships</b> between state factors (e.g., whether an object moves because of itself or other objects).</li><li>• Increased the model's <b>out-of-distribution generalization</b> by 46%, by leveraging the identified relationships and conditioning predictions only on relevant inputs.</li><li>• Derived a theoretically-grounded state abstraction for model-based RL, which improved sample efficiency and generalization in <b>planning</b> for <b>robotics</b> tasks.</li></ul>	
2022-23	<b>Unsupervised Skill Learning (NeurIPS)</b>	<b>University of Texas at Austin</b>
	<ul style="list-style-type: none"><li>• Proposed a skill discovery method for <b>structured decision-making</b> tasks, where reusable skills are learned to induce interactions between state factors.</li><li>• Implemented a novel <b>hierarchical RL</b> algorithm for skill learning in PyTorch – the high-level policy selects the interaction to induce and the low-level policy learns to induce it using primitive actions.</li><li>• Enhanced skill diversity and downstream task performance on long-horizon robotics tasks and structured decision-making tasks by 40%.</li></ul>	

## SELECTED PUBLICATIONS

See google scholar (<https://tinyurl.com/zizhaowangscholar>) for a complete list of publications.

1. Adversarial Reinforcement Learning for LLM Agent Safety, *In submission*  
**Z Wang**, D Li, V Keshava, P Wallis, A Balashankar, P Stone, L Rutishauser.
2. SkILD: Unsupervised Skill Discovery Guided by Local Dependencies, *NeurIPS 2024*  
**Z Wang\***, J Hu\*, C Chuck\*, S Chen, R Martin-Martin, A Zhang, S Niekum, P Stone.
3. Building Minimal and Reusable Causal State Abstractions for RL, *AAAI 2024 (oral)*  
**Z Wang\***, C Wang, X Xiao, Y Zhu, and P Stone.
4. ELDEN: Exploration via Local Dependencies, *NeurIPS 2023*  
**Z Wang\***, J Hu\*, R Martin-Martin, and P Stone.
5. Causal Dynamics Learning for Task-Independent State Abstraction, *ICML 2022 (oral)*  
**Z Wang**, X Xiao, Z Xu, Y Zhu, and P Stone.