

# HONGXIANG ZHANG

Lafayette, IN, USA ◊ (+1) 530-564-2580 ◊ [hxxzhang@gmail.com](mailto:hxxzhang@gmail.com) ◊ [LinkedIn](#) ◊ [Homepage](#) ◊ [Google Scholar](#)

## EDUCATION BACKGROUND

Purdue University Visiting Scholar	Augest 2025 - Present
University of California, Davis (GPA: 3.9/4.0) M.Eng of Science, Computer Science	Sep 2022 - June 2025
Australian National University Bachelor of Engineering (Honours), Computer Science	Jul 2020 - Sep 2022
Shandong University Bachelor of Engineering, Computer Science	Sep 2018 - Jun 2022

## ACADEMIC EXPERIENCE

<b>Multi-agent System alignment</b> <i>Hongxiang Zhang, Yuan Tian, Tianyi Zhang</i>	Ongoing
<ul style="list-style-type: none"><li>Designing an alignment framework for <b>multi-agent systems</b> to improve inter-agent consistency and task coordination.</li><li>Focusing on <b>inter-agent communication</b>, we make attention steering, allowing the agent to maintain focus throughout the generation.</li></ul>	

<b>Self-Anchor: Large Language Model Reasoning via Step-by-step Attention Alignment</b> <i>Hongxiang Zhang, Yuan Tian, Tianyi Zhang</i>	[arXiv]
<ul style="list-style-type: none"><li>Proposed Self-Anchor, a reasoning-time alignment method that leverages the inherent structure of reasoning to steer LLM attention, allowing the model to maintain focus throughout generation.</li><li>Self-Anchor significantly reduces the performance gap between “<b>non-reasoning</b>” models and specialized reasoning models, with the potential to enable most LLMs to tackle complex reasoning tasks without retraining.</li></ul>	

<b>Active Layer-Contrastive Decoding Reduces Hallucination in Large Language Model Generation</b> <i>Hongxiang Zhang, Hao Chen, Muhan Chen, Tianyi Zhang</i>	EMNLP 2025 Main [arXiv] [Project][Code]
<ul style="list-style-type: none"><li>Proposed Active Layer-Contrastive Decoding (ActLCD), a decoding algorithm that <b>actively contrasts model layers</b> to suppress hallucination.</li><li>By casting decoding as a sequential decision-making problem, ActLCD employs a <b>reinforcement learning</b> policy guided by a reward-aware classifier to optimize factuality <b>beyond the token level</b>.</li><li>Achieved SOTA hallucination reduction across five benchmarks.</li></ul>	

<b>SteerDiff: Steering towards Safe Text-To-Image Diffusion Model</b> <i>Hongxiang Zhang, Yifeng He, Hao Chen</i>	Under review [arXiv]
<ul style="list-style-type: none"><li>Developed <b>SteerDiff</b>, a lightweight plug-in module that filters unsafe or inappropriate concepts in diffusion model prompts via latent-space steering.</li><li>Proposed a semantic-preserving projection mechanism between text encoder and UNet, enabling safety control without retraining.</li><li>Demonstrated scalability to concept removal and fairness alignment tasks with minimal fine-tuning.</li></ul>	

## LLAMAFUZZ: Large Language Model Enhanced Greybox Fuzzing

*Hongxiang Zhang, Yuyang Rong, Yifeng He, Hao Chen*

*Under review [arXiv]*

- Introduced **LLAMAFUZZ**, a hybrid fuzzing framework that leverages LLM-driven mutation to learn structured input formats for both binary and text-based targets.
- Integrated an adaptive seed-selection policy that balances exploration and structure preservation, achieving significant coverage improvements on real-world benchmarks.
- Provided interpretability analyses showing how language-based mutation heuristics enhance fuzzing efficiency and semantic validity.

## INDUSTRY EXPERIENCE

---

**Software and Operation Engineer Intern** Volkswagen, Beijing, China

Feb 2022 - Jul 2022

- Achieved data masking and auto-populated to **JIRA** log in Python.
- Intuitively analyzed and provided services and finance data to CIO with **Tableau**, helping the management level adjust the company service strategy agilely.

**Quality Assurance Intern** Didi Global, Beijing, China

Dec 2020 - Feb 2021

- Developed distributed Java invoice services to enable the interaction between server and end-user devices that served more than 1 million users per day.
- Achieved unit testing automation by using **JUnit** and the test case in **Redis**, reduced over 15% of the testing engineer's workload.

## TECHNICAL SKILLS

---

**Programming Languages**

Python, Java, C/C++, JavaScript, HTML, CSS

**Software and skills**

MetaGPT, Autogen, MySQL, Numpy, Pytorch, Numpy, Pandas, **Azure**

## TEACHING

---

ECS 153 Computer Security - UC Davis

2024 Spring

ECS 036A Programming & Problem Solving - UC Davis

2024 Winter

ECS 036C Data Structures, Algorithms, & Programming - UC Davis

2023 Fall

ECS 140A Programming Language - UC Davis

2022 Winter

ENGN 4528 Computer Vision - Australian National University

2021

## ACADEMIC SERVICES

---

**Conference reviewer**

2025 IEEE 43rd International Conference on Consumer Electronics (ICCE 2025)

2025 ACL Rolling Review (ARR)