

Zicong Zhang

Shanghai, China | zhangzico@sjtu.edu.cn | +(86) 188 9181 2036
<https://zhangzicong518.github.io>

Education

Shanghai Jiao Tong University , B.S. in Computer Science	Sept 2023 – Present
• John Hopcroft Class, Zhiyuan College (Honor Program)	
• GPA: 89.5/100	
• Coursework:	
Math: Mathematical Logic(97), Information Theory(97), Linear Algebra(94)	
CS: Foundations of Data Science(100), Machine Learning Theory(98), Computer System Design(98), Program Design and Data Structure(96), Machine Learning(95), Algorithm(94)	

Research Experience

Multimodal Reasoning and Generic Agent , Microsoft Research Asia, intern	July 2025 – Jan 2026
• Proposed and implemented a tool-integrated RL paradigm that unifies long-chain-of-thought reasoning and tool-calling into a single hybrid thinking mode, enabling vision–language models to “think with images” more effectively	
• Built and modified a large-scale, high-fidelity simulation environment that generates diverse, realistic agent–tool interaction trajectories, providing the massive training data required for robust tool-calling agents.	
Continual Learning with Model Fusion , SJTU MIFA Lab, intern	July 2024 – April 2025
• Co-led development of a novel continual learning framework for VLMs by introducing model fusion	
• Proposed aggregating the results of multiple decoupled task-specific models for prediction in zero-shot scenarios.	
• Conducted extensive experiments on multiple benchmarks, demonstrate that outperforming the original pre-trained VLM and other state-of-the-art continual learning methods.	

Awards

- Zhiyuan Honor Awards (Top 10% in SJTU) 2023, 2024, 2025
- The Third Prize of Academic Scholarship (Top 30% in major) 2024, 2025

Publications

Enhanced Continual Learning of Vision-Language Models with Model Fusion

Haoyuan Gao*, Zicong Zhang*, Yuqi Wei, Linglan Zhao, Guilin Li, Yexin Li, Linghe Kong, Weiran Huang
Fisrt Workshop SCOPE The Thirteenth International Conference on Learning Representations

Tracing the Dance of Embeddings: Visualizing High-Dimensional Sample-wise Trajectory for Training Analysis

Yiming Liu*, Zicong Zhang*, Yun Lin, Ruofan Liu, Yuhuan Huang, Weiyu Kong, Jinsong Dong
[Under Review]

Projects

Environment Scaling for General Agentic Model – Microsoft Research Asia	Jul 2025 – Jan 2026
• Collected 150k real-world APIs from RapidAPI, MCP-Server and ToolBench; embedded api info to build a similarity-weighted dependency tool graph.	
• Applied Louvain community detection to partition the graph into hundreds functional domains, and automatically derived a unified database schema per domain following read-write paradigm.	
• Materialised each API as deterministic read/write operators on the domain-specific schema, enabling fully verifiable environment states	
• Sampled long-horizon, logically coherent tool sequences on dependency graph, and generate diverse, high-fidelity agent–environment interaction trajectories for training	

Scaling Active Perception for Multimodal Reasoning – Microsoft Research Asia

Jul 2025 – Jan 2026

- Unified long-chain reasoning and tool-calling into a single hybrid thinking mode through SFT and RL, enabling vision-language models to autonomously switch thinking modes while “thinking with images”
- Curated high-quality long-CoT trajectories via teacher distillation and iterative rejection sampling and relieve tool-use hallucinations and inefficiencies
- Designed mode-aware process rewards that encourage the model to leverage both internal reasoning and external tools in a complementary manner
- Evaluated on MathVision, MathVista, HRBench, V* and other math or perception-oriented benchmarks, surpassing the original VLM baseline across various datasets

Continual Learning Framework for VLMs

July 2024 - Jan 2025

Results in papers Enhanced Continual Learning of Vision-Language Models with Model Fusion

- Designed a continual learning framework for VLMs by introducing model fusion
- Deployed a pipeline and conducted extensive experiments on multiple benchmarks, achieving up to 2% improvement over other state-of-the-art continual learning methods.

Technologies

Language: Mandarin (native), English (CET-6 600, TOEFL under preparation)

Programming: C/C++, Python, Rust

Technologies: Git, Pytorch, Latex, vLLM, VeRL, LLaMA-Factory, Deepspeed