

# SHILIN ZHANG

•+86-189-8883-9683•torbjornclancy@gmail.com•Guangzhou, China

[Github](#) – [Google Scholar](#) – Last updated: July 10, 2024

---

## EDUCATION

### South China University of Technology

Bachelor of Engineering, Intelligent Manufacturing

Guangzhou

Sep 2021 - Jun 2025

- **Core Courses**

Introduction to Engineering, Data structure & Algorithms, Design & Manufacturing, Artificial Intelligence Technology and Applications, Fluid Mechanics, Embedded system design

---

## PUBLICATION

### MLLM-as-a-Judge: Assessing Multimodal LLM-as-a-Judge with Vision-Language Benchmark

GuangZhou

#### First Co-author

2024.1-2024.6

- Project Background: This work focuses on exploring the processing and logical capabilities of various state-of-the-art MLLMs on text, images, and videos, as well as evaluating the processing capabilities of existing MLLMs. At the same time, this work also emphasizes the differences between the judgments of MLLMs and humans, as well as the feasibility of MLLMs themselves.
- Accepted as **ICML 2024 Oral**

### A Blurred Border: The Dilemma in Identifying LLM-Human Mixcase

GuangZhou

#### Co-author

2023.9-2023.12

- Project Background: This study introduces the MixSet (mixcase dataset) to unravel the challenges in distinguishing mixcases of human-written and machine-generated texts, revealing significant gaps in current detection methods for mixed-case scenarios.
- Accepted as **NAACL 2024 Findings**

### PyramidInfer: Pyramid KV Cache Compression for High-throughput LLM Inference

GuangZhou

#### Co-author

2023.12-2024.5

- Project Background: Large Language Models (LLMs) face challenges in GPU memory usage, limiting their scalability for real-time applications. To address this, our project, PyramidInfer, introduces a method that compresses the Key-Value (KV) cache by selectively retaining crucial context layer-wise, significantly reducing memory needs and improving inference speed without compromising performance. Our tests demonstrate that PyramidInfer achieves a 2.2x increase in throughput and reduces GPU memory usage by over 54% compared to existing methods.
  - Accepted as **ACL 2024 Findings**
- 

## PROJECT EXPERIENCES

### Applied Optimization and Learning Research

Mentored by Prof. Peter Zhang, CMU

Shanghai Jiao Tong University, Shanghai

2023.7-2023.11

- Project Background: Focused on operations research and AI, particularly Q-learning and DQN, for hospital outpatient data coordination and doctor scheduling.
- Project Outcomes: Under Prof. Zhang's guidance, simulated problem-solving in a real hospital setting with both operational and AI techniques; evaluated the impact of different algorithms on coordination; visualized results via scripting and composed an experimental report and academic paper, all received high praise from the professor.

**Porsche**

New Energy Vehicle Charging Process Technology Research

**Guangzhou**

Apr 2023 – Jul 2023

**Research Assistant**

- **Research Background:** Tasked with research on EV market and charging technologies based on Porsche's strategic goals, in line with the burgeoning trend of electric vehicles.
- **Research Outcomes:** Selected EV charging processes as research focus; conducted thorough study using literature, market surveys, and presented findings related to automotive charging processes and market research, earning mentor's recognition.

**Stanley Electric Co., Ltd.**

A company producing high-tech automotive components for GAC Group

**Guangzhou**

Jan 2023

**Researcher**

- **Analysis & Research:** Investigated production line issues within Stanley Electric, especially cache and efficiency problems in warehouse management; proposed solutions integrating existing practices with Industry 4.0 technologies and implemented them.
- **Reporting:** Analyzed industrial process inefficiencies, comparing with Industry 4.0 insights and personal project experience, and authored a comprehensive report.

---

**ADDITIONAL INFORMATION**

- **Language Skills:** English (TOEFL: 96), Mandarin, Cantonese
- **Hardware skills:** Embedded system design experiences on stm32 and Arduino
- **Interests:** Rugby (Third place in the 2019 NFL Flag Football China High School League, Third place in the 2019 Guangzhou High School Rugby League), Swimming, Guitar, Violin