

Yuzhen Huang

Email: yhuanghj@cse.ust.hk

Homepage: hyz17.github.io

Research Interests: I am primarily focused on large language models, particularly in advancing their reasoning capabilities and multimodal understanding. To achieve this, my research interests lie in: (1) enhancing reasoning and planning abilities through self-improvement and RL techniques, (2) improving the architecture and training methods of multimodal models to strengthen their understanding across multiple modalities and (3) developing reliable evaluation methods for language models.

EDUCATION

The Hong Kong University of Science and Technology, Hong Kong SAR, China.

Feb. 2024 - present

– **Ph.D in Computer Science**

– Advisor: Prof. Junxian He

Shanghai Jiao Tong University, Shanghai, China

Sep. 2019 – Jul. 2023

– **B.Eng. in Computer Science**

– *GPA: 3.89/4.3, Score: 90.27/100*

RESEARCH PROJECTS

- [1] B-STaR: Monitoring and Balancing Exploration and Exploitation in Self-Taught Reasoner
 - **Quantitatively analyze the dynamics of exploration and exploitation during self-improvement.**
 - Introduce B-STaR, a Self-Taught Reasoning framework that autonomously adjusts its configurations.
 - Balance exploration and exploitation, leading to superior performance.
- [2] Compression Represents Intelligence Linearly
 - **Investigate the linear correlation between compression and intelligence in LLMs.**
 - Provide evidence for the belief that superior compression is indicative of greater intelligence.
 - Propose compression efficiency serves as an unsupervised and reliable metric to assess LLMs' abilities.
 - Published in COLM 2024 as the first author.
- [3] C-Eval: A Multi-Level Multi-Discipline Chinese Evaluation Suite for Foundation Models
 - **The first comprehensive Chinese evaluation suite for LLMs.**
 - Conduct a thorough evaluation of the most advanced LLMs.
 - Over 9.8M downloads on Hugging Face and more than 100 models on leaderboard.
 - Published in NeurIPS 2023 as the first author.

PUBLICATIONS

* denotes equal contribution

- [1] W Zang*, **Y Huang***, L zhao, Y Wang, Z Shan, J He. B-STaR: Monitoring and Balancing Exploration and Exploitation in Self-Taught Reasoner. ICLR 2025.
- [2] **Y Huang***, J Zhang*, Z Shan, J He. Compression Represents Intelligence Linearly. Conference on Language Modeling (COLM), 2024.
- [3] **Y Huang***, Y Bai*, Z Zhu, J Zhang, J Zhang, T Su, J Liu, C Lv, Y Zhang, Y Fu, M Sun, J He. C-Eval: A Multi-Level Multi-Discipline Chinese Evaluation Suite for Foundation Models. NeurIPS (Datasets and Benchmarks track), 2023

PAST EMPLOYMENT

Research Intern, Tencent
Mentor: Zifei Shan

Nov. 2023 – Jan. 2024

PROFESSIONAL ACTIVITIES

Reviewer: NeurIPS 2024, ICLR 2025

TEACHING

Teaching Assistant, The Hong Kong University of Science and Technology
COMP 5212 Machine Learning

Fall 2024

STANDARD TESTS

[1] **TOEFL** – 102