香港科技大学（广州）
THE HONG KONG
UNIVERSITY OF SCIENCE AND
TECHNOLOGY (GUANGZHOU)

数据科学与分析学域
DATA SCIENCE AND ANALYTICS THRUST
信息枢纽
INFORMATION HUB

DSA Research Experiences for Undergraduates

## Research Project

### Section1: Faculty Information

| Full Name | Jia Li | Tel | |
|---|---|---|---|
| Thrust/Hub | Thrust of Data Science and Analytics | Office | W2 L6 605 |
| Email | jialee@hkust-gz.edu.cn | | |

### Section2: Research Project Proposal

| Project Title | Enhancing Reasoning in Large Language Models: Data, Algorithms, and Applications |
|---|---|
| Project Description (max 800 words) | Our project focuses on the exploration and optimization of the reasoning capabilities of large language models (LLMs), a domain that has become a core direction in natural language processing research in recent years. From early methods based on prompting to the current gradual shift toward post-training approaches, the potential of LLMs in System 2 thinking—characterized by slower, more deliberate logical reasoning—has been significantly demonstrated. Whether represented by open-source models like Deepseek R1 and Qwen QwQ or closed-source benchmarks such as OpenAI o1/o3, these models have exhibited robust reasoning capabilities in high-precision data-driven post-training tasks across mathematics, coding, and scientific domains. However, current research is predominantly confined to structured data within specific fields, leaving ample room for new exploration: enhancing the logical reasoning, search, and planning abilities of models through broader and more diverse data sources.<br><br>Our work primarily revolves around several key aspects. First, we aim to transcend the reliance on domain-specific data in existing research by exploring the application of graph data or other underutilized datasets in reinforcement learning training. By introducing diversified datasets and training paradigms, we seek to enable models not only to solve specific tasks but also to exhibit stronger generalization and flexibility. Second, we are committed to optimizing algorithms and reinforcement learning methodologies to mitigate common issues such as "semantic repetition" and "hallucination" during the reasoning process, thereby improving the efficiency of model thinking and the quality of outputs. For instance, we aim to guide models toward forming more efficient reasoning pathways, avoiding verbose and semantically redundant responses. |

香港科技大学（广州）
THE HONG KONG
UNIVERSITY OF SCIENCE AND
TECHNOLOGY (GUANGZHOU)

数据科学与分析学域
DATA SCIENCE AND ANALYTICS THRUST
信息枢纽
INFORMATION HUB

|  | Furthermore, building upon our previous achievements, we strive to push the boundaries of model capabilities. For example, our prior work developed Graph Arena—a benchmarking tool for graph reasoning that supports natural language queries. This tool not only evaluates models' performance in graph-structured reasoning but also provides clear guidance for subsequent training efforts. Additionally, we plan to deepen the capabilities of LLMs in code generation and comprehension, enabling them to demonstrate higher levels of logical reasoning and problem-solving in programming tasks.

In terms of technical methodology, we train models ranging from 1.5B to 7B parameters using specialized synthetic datasets while drawing inspiration from the concept of Chain-of-Thought Prompting and integrating reinforcement learning techniques. These approaches aim to guide models in constructing clear reasoning chains and demonstrating generalization capabilities on unseen datasets. Building on Graph Arena, we will continue to refine evaluation frameworks and apply them to real-world scenarios to validate the practical effectiveness of the models.

Through in-depth exploration of data, algorithms, and application scenarios, we are dedicated to advancing the development of next-generation LLMs with powerful reasoning capabilities. By transcending data boundaries, we aim to equip models with a deeper understanding of the underlying patterns governing the world's complexities. |
|---|---|
| Proposed Research Duration | Start Date: __March__ / __01__ / __2025__<br>End Date: __Dec__ / __31__ / __2025__ |
| Student/Researcher Duties | Maintaining code repositories, designing and conducting experiments, optimizing algorithms, managing datasets, and developing evaluation tools, while actively contributing to research documentation, team collaboration, and result dissemination. |
| Technical Skills Required | ☑ Python     ☑ Machine Learning     ☐ Big Data<br>☐ R     ☑ Deep Learning     ☐ SQL<br>☐ C/C++     ☐ Other: _____ |
| Preferred Student/Researcher Background | Machine learning, programming (Python, PyTorch, Triton), and experience with large language models, graph data, or reinforcement learning, |
| Maximum Number of Students/Researchers | ☑ 1          ☐ 2 |

**Section3: Pre-Application Research Exposure Meeting**

Faculty members are encouraged to schedule a Research Exposure Meeting to introduce students to their projects.

香港科技大学（广州）
THE HONG KONG
UNIVERSITY OF SCIENCE AND
TECHNOLOGY (GUANGZHOU)

数据科学与分析学域
DATA SCIENCE AND ANALYTICS THRUST
信息枢纽
INFORMATION HUB

| | |
|---|---|
| Preferred Date | |
| Preferred Time | |
| Meeting Mode | ☐ In-Person        ☐ Online |
| Venue (if in-person) | |
| Meeting Link (if online) | |

香港科技大学（广州）
THE HONG KONG
UNIVERSITY OF SCIENCE AND
TECHNOLOGY (GUANGZHOU)

数据科学与分析学域
DATA SCIENCE AND ANALYTICS THRUST
信息枢纽
INFORMATION HUB