

## DSA Research Experiences for Undergraduates

### Research Project

#### Section1: Faculty Information

Full Name	Guoming Tang	Tel	13618480331
Thrust/Hub	DSA/INFO	Office	W3-306
Email	guomingtang@hkust-gz.edu.cn		

#### Section2: Research Project Proposal

Project Title	Green GPU Computing for Large-Scale LLM Inferences
Project Description (max 800 words)	<p>As the scale of large language models (LLMs) continues to grow, so does the demand for <b>more efficient and adaptive GPU-based inference</b>. Traditional GPU resource allocation often relies on static or heuristic-based methods, which can lead to <b>suboptimal energy consumption</b> and increased latency. These challenges become even more pronounced in large-scale online inference scenarios, where diverse tasks—each with distinct <b>computational complexities</b>—must be processed simultaneously.</p> <p>This project will explore the following questions:</p> <ol style="list-style-type: none"> <li>1. How can we define <b>new metrics regarding the GPU computing potentials</b> in evaluating LLM performance and overhead?</li> <li>2. How do we seamlessly <b>match new tasks to the most suitable hardware</b> config class in an online environment with multiple GPU nodes?</li> <li>3. Can LLMs <b>facilitate dynamic adjustments</b> to maintain efficiency amid variable workloads?</li> <li>4. Can LLMs <b>generalize the hardware-tuning methodology</b> to new architectures, tasks, or even different hardware platforms?</li> </ol> <p>You will gain the following opportunities:</p> <ol style="list-style-type: none"> <li>1. <b>Hands-on HPC &amp; GPU Experience</b> Work with real GPU environments and gather practical skills in hardware tuning, optimization algorithm and cluster-based analysis. Gain a deep understanding of balancing performance metrics such as energy consumption and inference latency.</li> <li>2. <b>LLM Integration &amp; Agent Development</b> Explore how large language models can be harnessed for scheduling, optimization, and interpretability. Experiment with prompt engineering, fine-tuning (e.g., LoRA), and LLM-based reasoning to push the boundaries of HPC scheduling methods.</li> </ol>

	<p><b>3. Research &amp; Publications</b> Identify real-world problems within GPU hardware resource management and transform these insights into publishable research. There is a strong opportunity to work on cutting-edge topics and produce high-quality academic papers (if possible).</p> <p><b>4. Professional Growth &amp; Internships</b> Students in this project may access internship opportunities with industry partners, opening doors to broader collaborations and future career development.</p>
Proposed Research Duration	Start Date: Now End Date: Aug. 31, 2025
Student/Researcher Duties	Students should read papers in related fields, actively engage in team discussions, and contribute ideas to the project. Additionally, they will be responsible for part of coding tasks. Exploring different cutting-edge directions to solve problems is encouraged.
Technical Skills Required	<input checked="" type="checkbox"/> Python <input checked="" type="checkbox"/> Machine Learning <input type="checkbox"/> Big Data <input type="checkbox"/> R <input checked="" type="checkbox"/> Deep Learning <input type="checkbox"/> SQL <input type="checkbox"/> C/C++ <input type="checkbox"/> Other: _____
Preferred Student/Researcher Background	Programming, Academic Reading & Writing Skills, Teamwork Ability
Maximum Number of Students/Researchers	<input type="checkbox"/> 1 <input checked="" type="checkbox"/> 2

### Section3: Pre-Application Research Exposure Meeting

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

Preferred Date	Mar. 14-16, 2025
Preferred Time	Morning or afternoon time, 1 hour
Meeting Mode	<input checked="" type="checkbox"/> In-Person <input type="checkbox"/> Online
Venue (if in-person)	W3-306
Meeting Link (if online)	