

DSA Research Experiences for Undergraduates

Research Project

Section1: Faculty Information

Full Name	Zeyi WEN	Tel	
Thrust/Hub	DSA Thrust, Info Hub	Office	E2 312
Email	wenzeyi@ust.hk		

Section2: Research Project Proposal

Project Title	Efficient Hyperparameter Optimization for LLM Inference
Project Description (max 800 words)	<p><i>Provide a brief summary of the project, objectives, and expected outcomes.</i></p> <p>Large language models (LLMs) have demonstrated remarkable capabilities in text generation, summarization, and dialogue systems. However, deploying these models effectively requires careful tuning of inference hyperparameters—such as temperature and top_p—which are crucial for balancing text quality, creativity, and computational efficiency. Hyperparameter optimization for LLMs presents significant challenges, as traditional methods like grid search, random search, and standard Bayesian optimization are often time-consuming, token-intensive, and struggle with the exploration-exploitation trade-off, particularly when evaluating high-cost configurations (e.g., full-scale inference with long prompts). To address these challenges, this project aims to develop an efficient hyperparameter optimization method for LLM deployment, focusing on minimizing token consumption while maintaining or improving text quality. Key research questions include designing advanced searching strategies and implementing multi-fidelity evaluation techniques. The project seeks to establish a principled framework for hyperparameter optimization in LLMs, emphasizing efficiency, scalability, and adaptability. Students will explore both theoretical foundations (e.g., convergence guarantees) and practical implementations, validated through extensive benchmarks on diverse NLP tasks.</p>
Proposed Research Duration	Start Date: <u>2025</u> / <u>03</u> / <u>15</u> End Date: <u>2025</u> / <u>09</u> / <u>01</u>
Student/Researcher Duties	<p><i>List the primary responsibilities and tasks expected from the student/researcher during the project.</i></p> <p>Conduct benchmarking on NLP tasks in terms of text quality and efficiency. Establish an LLM serving system using SOTA frameworks. Design and implement adaptive search strategies to balance exploration-exploitation trade-offs.</p>
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	<i>List preferred coursework, experience, or skills (e.g., statistics, programming, AI).</i> Math, programming
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	
Preferred Time	
Meeting Mode	<input type="checkbox"/> In-Person <input type="checkbox"/> Online
Venue (if in-person)	
Meeting Link (if online)	