

# HPC-GENIE Project: Towards Code Generative AI for HPC Programming

**Takahiro Katagiri** (Information Technology Center, Nagoya University, E-mail: katagiri@cc.nagoya-u.ac.jp)

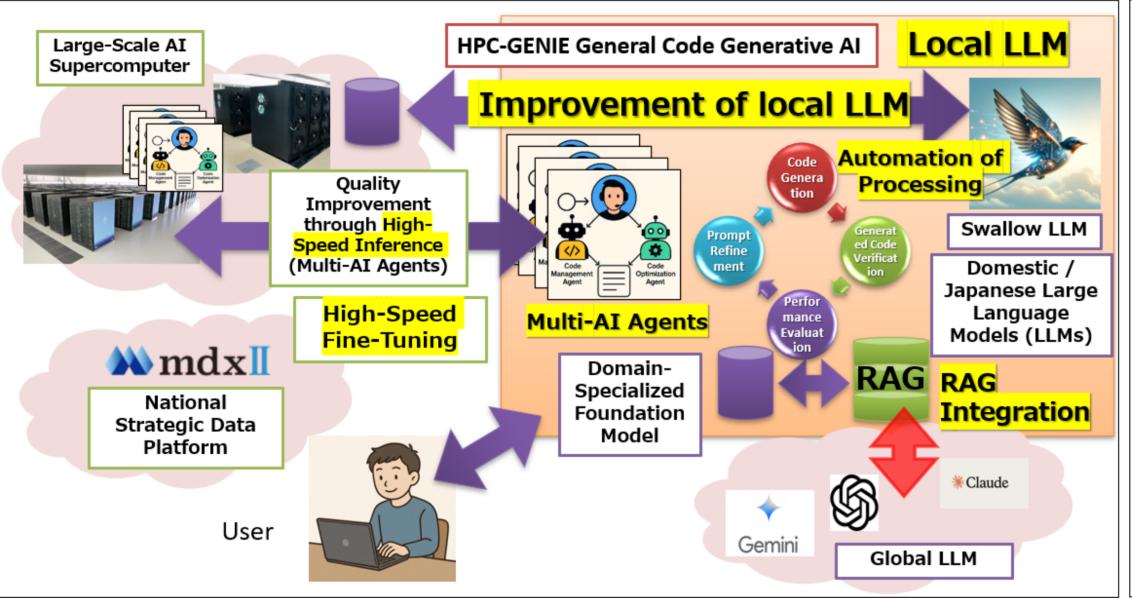
### **Aim of the Project**

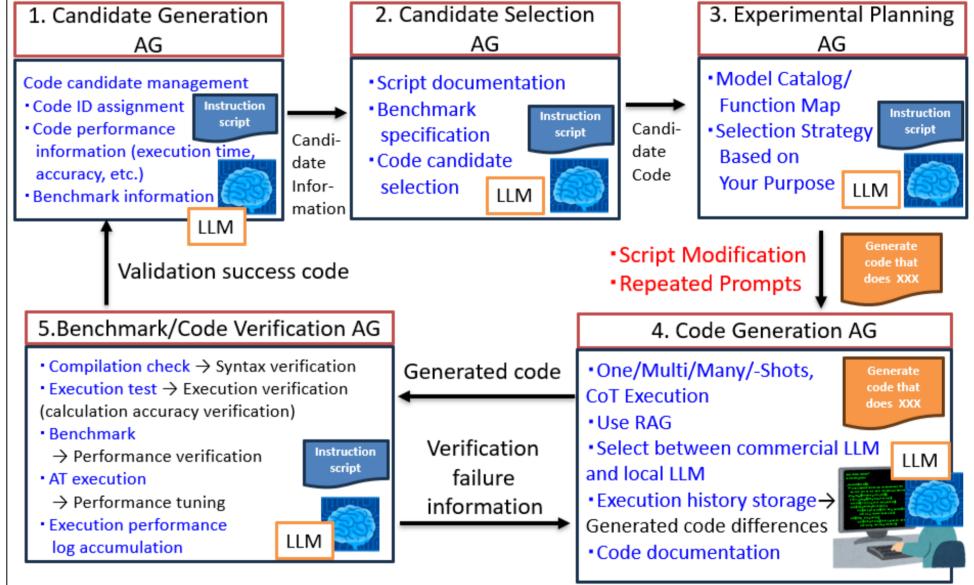
- The HPC-GENIE (<u>High-Performance Computing with GE</u>nerative <u>Neural Intelligence for Execution</u>) project is an initiative launched by members of Information Technology Center and Graduate School of Informatics at Nagoya University.
- It is a research project that leverages code-generating AI for the automatic generation of HPC programs. By integrating prompt engineering based on large language models (LLMs) with software auto-tuning (AT) technologies, the project aims to achieve automation that dramatically enhances the productivity of HPC software development.



## Automation of Local LLM-Specialized Processing by HPC-GENIE

- Automation of Code Tuning Process: iterative prompting of code generation  $\rightarrow$  generated code verification  $\rightarrow$  performance evaluation  $\rightarrow$  prompt refinement
- Enhancement of code quality through multi-Al-agent (A2A) collaboration
- RAG-based reinforcement of local LLMs
- Improvement of generated code quality via supercomputer-accelerated fine-tuning and inference (multi-AI-agent execution) for local LLMs





**HPC-GENIE**: A General Code Generative AI System

**Iterative Promompt** for the HPC-GENIE System

#### **Overview of VibeCodeHPC**

• VibeCodeHPC: Multi-Agent System for Auto-Tuning of HPC Code Optimization



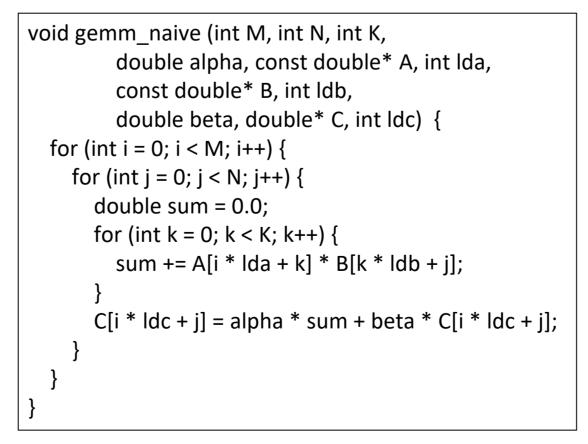
Launch screen

# **Roles of Agents**

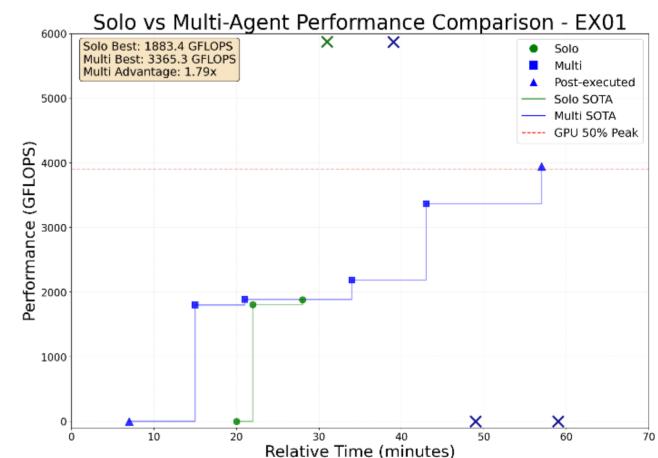
<mark>Agents</mark>	Roles	Scope of Responsibility
<mark>PM</mark>	Project Management	Requirements Definition, Resource Allocation, and Budget Management
SE	System Design	Agent Monitoring, Statistical Analysis, and Report Generation
PG	Code Generation and Execution	Parallel Implementation, SSH/SFTP Connection, Job Execution, Performance Measurement, and SOTA Evaluation
CD	Deployment Management	Publication and Anonymization of SOTA-Achieving Code

A Case Study: Multi-AI agents code opination for matrix-matrix multiplication

# Original Code (Before Optimization)



# Optimizaztion History



A multi-agent system generates high-performance code faster than a solo-agent system and additionally achieves even higher performance.

#### VibeCodeHPC:

Source Code: <a href="https://github.com/Katagiri-Hoshino-Lab/VibeCodeHPC-jp">https://github.com/Katagiri-Hoshino-Lab/VibeCodeHPC-jp</a> arXiv: <a href="https://www.arxiv.org/abs/2510.00031">https://www.arxiv.org/abs/2510.00031</a>

Information Technology Center, Nagoya University <a href="http://www.icts.nagoya-u.ac.jp/en/center/">http://www.icts.nagoya-u.ac.jp/en/center/</a>

