# LLM-based physics analysis assistant at BESIII - 'Dr.Sai'

Yipu Liao (廖一朴)

Inistitute of High Energy Physics, CAS, Beijing

On behalf of **Dr.Sai** working group

2025.7.17

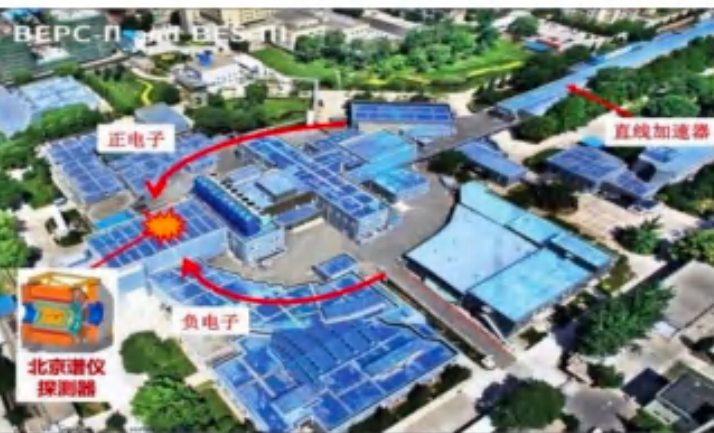第六届 粒子物理前沿研讨会

@ Jilin University, Changchun

# Outline

1. Motivation

2. Dr.Sai project

3. Dr.Sai-BESIII

4. Prospect

5. Summary

# Challenges in HEP

- Data magnitude: **~10 PB**

- Growth rate: accelerated with new technology

  - BEPCII → BEPCII-U (May 17, 2025): **Luminosity increase ~300% at 2.35 GeV**
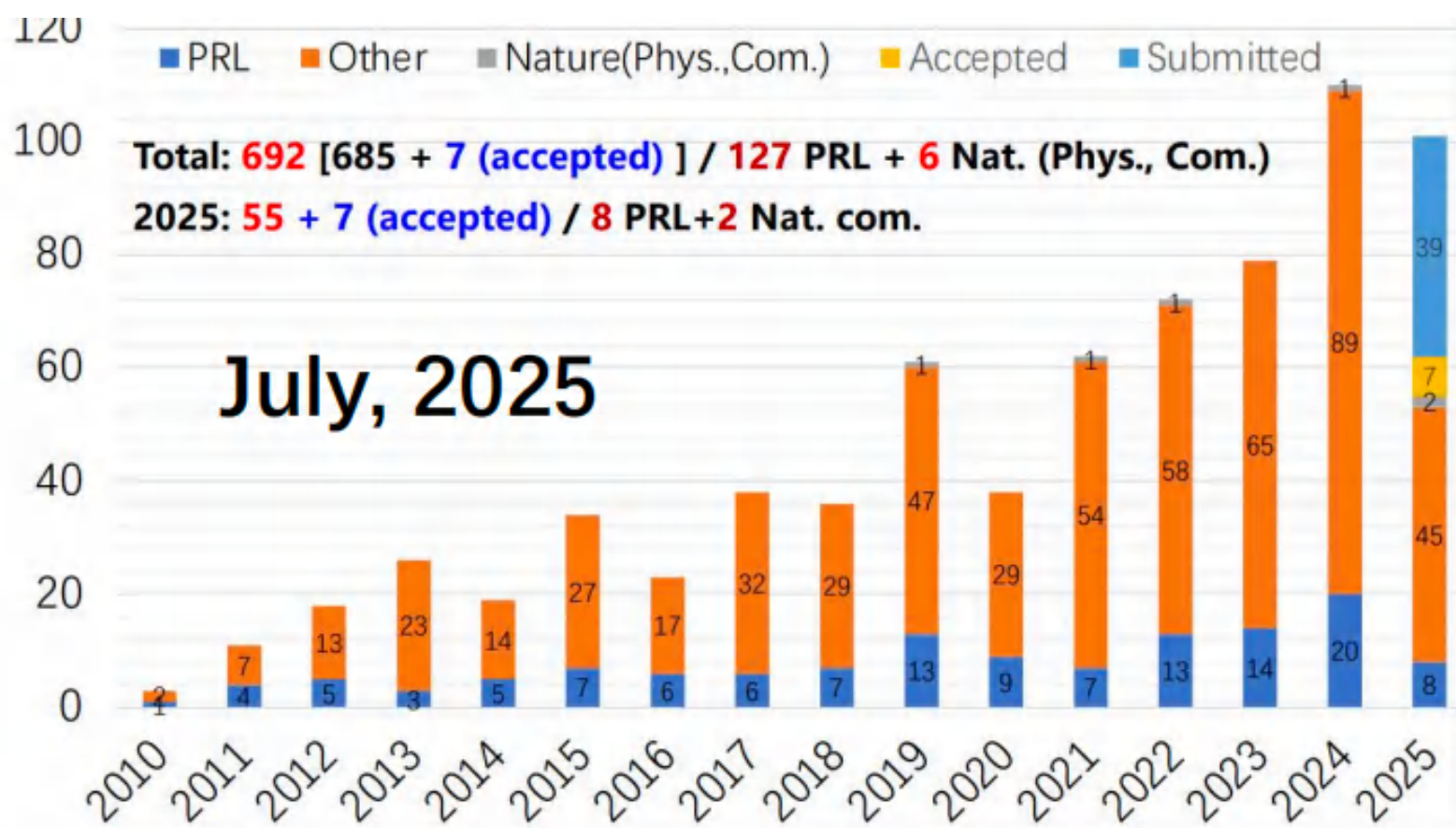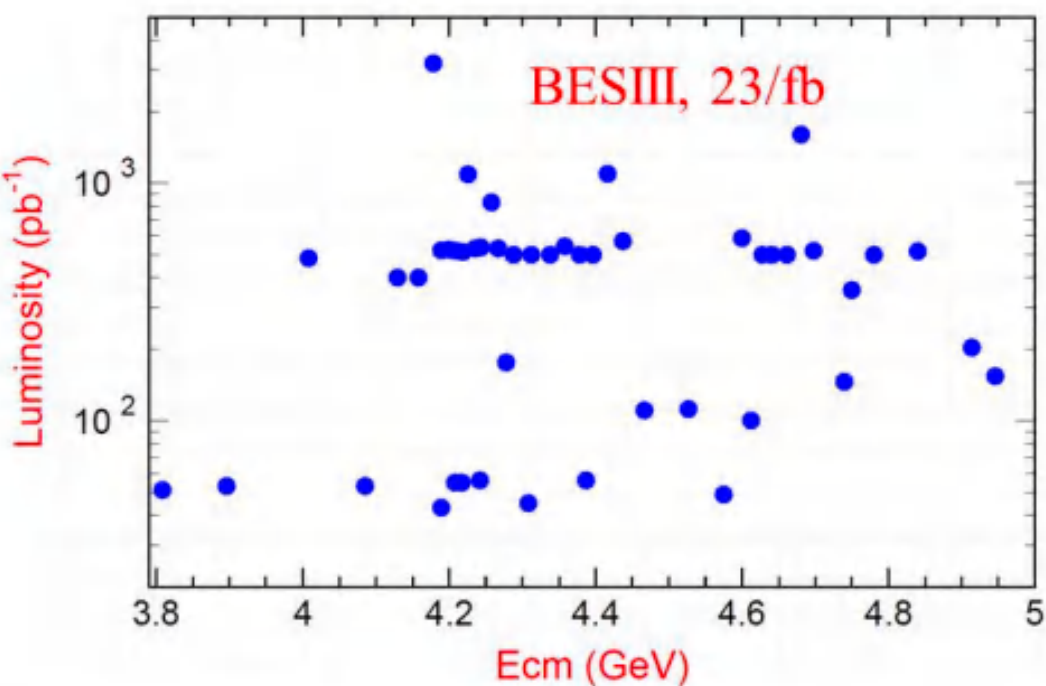
  - HEPS (2025): **~800 TB / day**



- ~700 physics results with ~700 people in 16 years
  - One result normally took ~3 years
- HEP need a more efficient tool to investigate data

# Challenges in HEP

**Data process workflow at HEP experiment** (take Hadron collider as an example)

| Collision, MC generator | Trigger, Simulation | Reconstruction | Statistical analysis | Extract physics variables |
|---|---|---|---|---|

| | | | | |
|---|---|---|---|---|
| • Accelerator control | • Data acquisition | • Track & vertex Finding and fitting | • Event selection | • Systematic uncertainty |
| • Initial-state-radiation | • Fast reconstruction | • Clusterization & reconstruction of jet | • Optimizations | • Fitting |
| • Parton showering | • Data input/output | • Jet tagging | • Background analysis | • Uncertainty propagation |
| • Hadronization | • Online monitoring | • Kinematic fit | • Injection test | • Radiation and VP corrections |
| • NP-correction | • Detector geometry | • Detector calibration | • Reweighting | • ... |
| • Pileup | • Detector noise | • ... | • Correlation corrections | |
| • ... | • Calibration | | • ... | |
| | • Multi-scattering | | | |
| | • ... | | | |

- **Too complicated, similar lines of code as windows/macOS**
- **One small task need: several people + several years!**

# Beijing Electron-Positron Collider (BEPC II)

Ground breaking: 1984
CM energy : 2 - 5 GeV
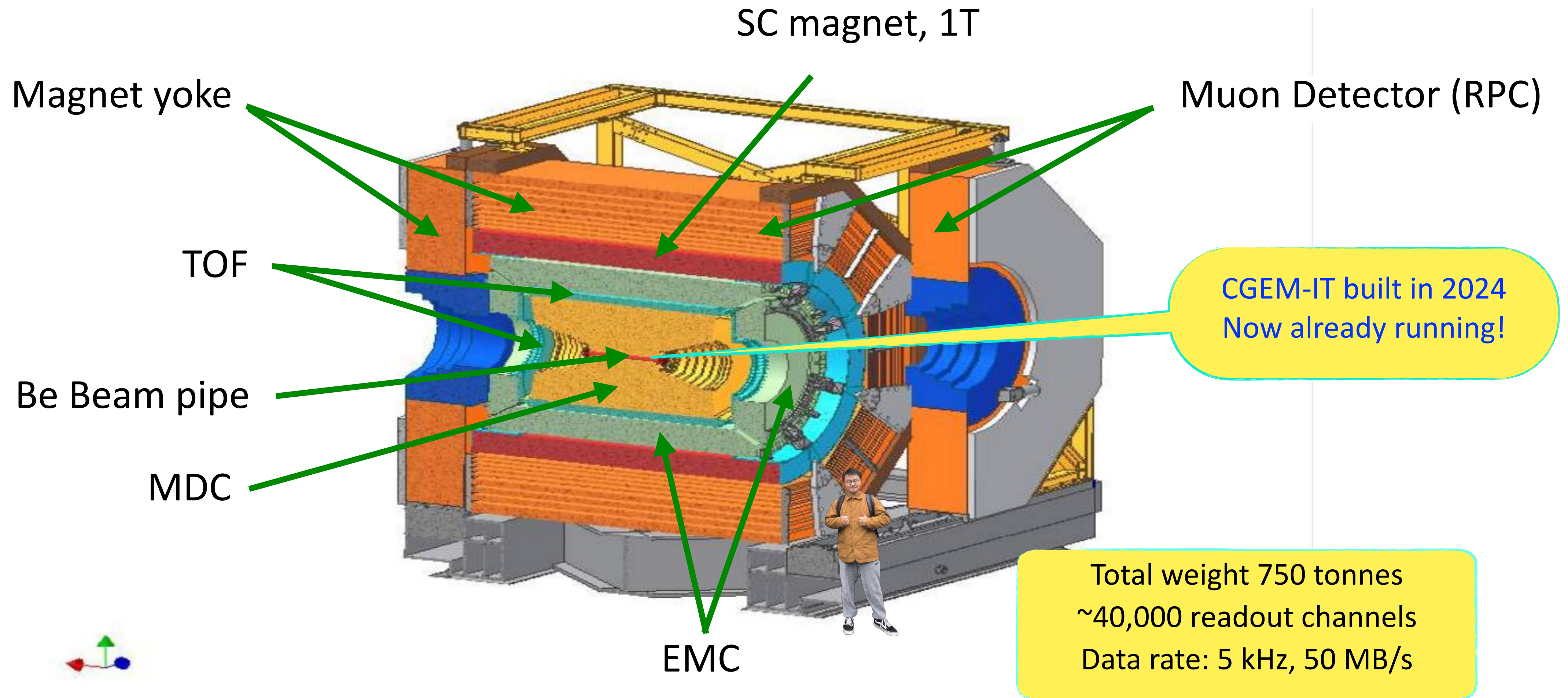Major upgrade: 2004
Energy upgrade: 2024

1989-2005 (BEPC): $L_{peak}=1.0 \times 10^{31}$ /cm$^2$s
2008-now (BEPCII): $L_{peak}=1.0 \times 10^{33}$/cm$^2$s (Apr. 5, 2016)

World unique e$^+$e$^-$ accelerator in $\tau$-charm energy region



北京谱仪
BES

IHEP, Beijing

# The BESIII experiment



SC magnet, 1T

Magnet yoke

Muon Detector (RPC)

TOF

CGEM-IT built in 2024
Now already running!

Be Beam pipe

MDC

EMC
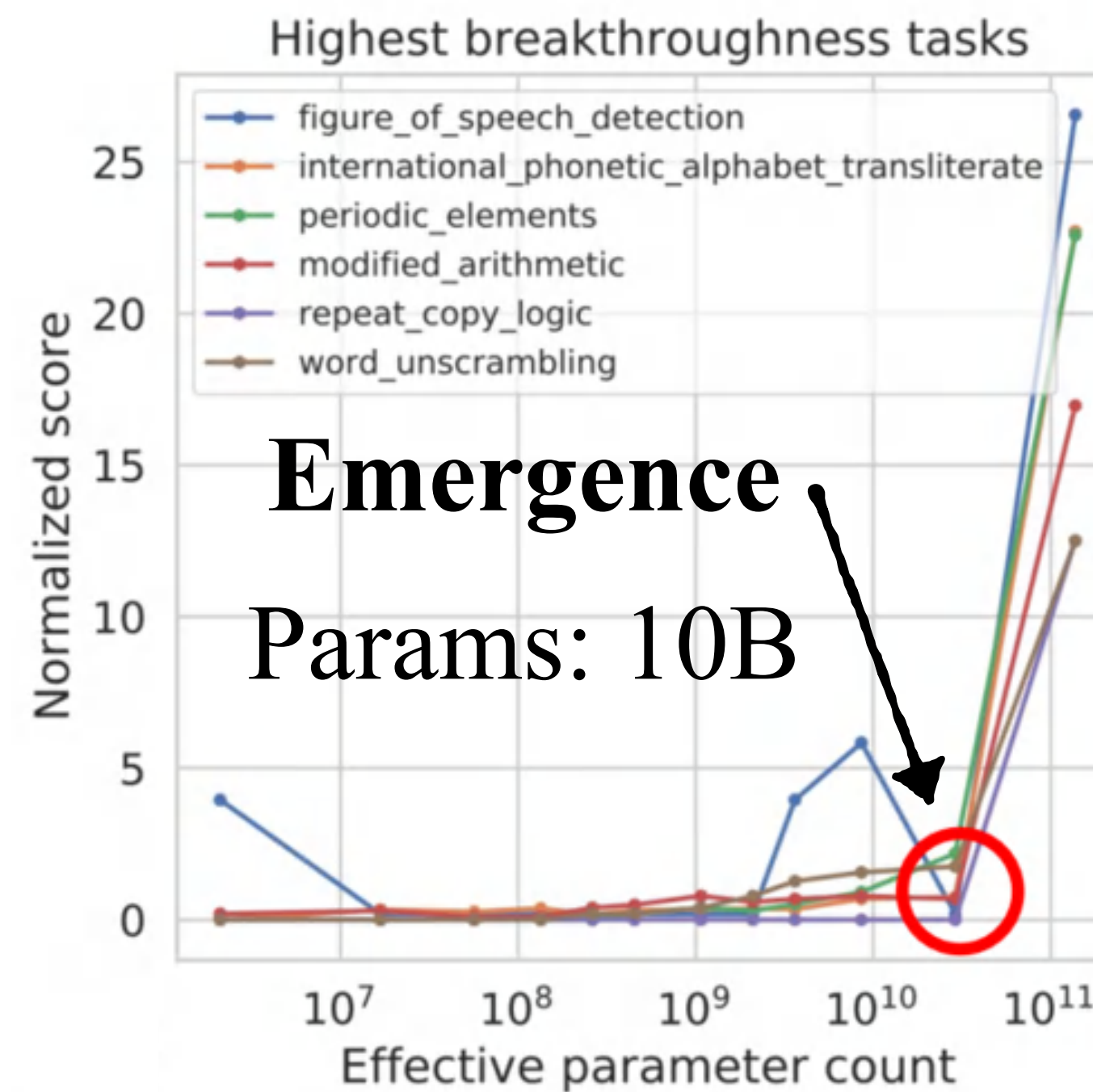
Total weight 750 tonnes
~40,000 readout channels
Data rate: 5 kHz, 50 MB/s

Has been in full operation since 2008, all sub-detectors are in good status!
**A lot of data to be analyzed → Use the more intelligent automatic workflow**

# What is Large Language Model (LLM)

- **Large language models (LLMs), normally build on Transformer architecture (Deep Learning) and based on the next-token prediction**

- **It has demonstrated impressive performance in <span style="color:red">text / code generation</span>**

  - GPT4o, Gemini, LLaMa3 …
  - Could be used for HEP studies
  - Game changer

- A foundation model (large, computing intensive) + fine tuning for each task individually (smaller data set)



Highest breakthroughness tasks

- figure_of_speech_detection
- international_phonetic_alphabet_transliterate
- periodic_elements
- modified_arithmetic
- repeat_copy_logic
- word_unscrambling

**Emergence**

Params: 10B

Normalized score / Effective parameter count

Difficulty Groups

Questions answered correctly (%)

Easy — Difficult

GPT-4.0  GPT-3.5  Humans

# Retrieval Augment Generation (RAG)

- RAG is the most cheap & **promising** solution to reduce hallucinations (幻觉)

- Store private data so no need for retraining the model

# Some applications based on RAG

*\* Focus on HEP*



- An AI assistant of the ATLAS Collaboration with decades of institutional documentation (wiki, docs, indico, ...)
- Focus on the search system and questions about ATLAS

# Agent(ic): LLM with tools

A LLM agent framework consists of:

- User request: a user question or request

- **Agent: a LLM-based agent core**

- **Tools: manage tools except for text/ image reply; more professional tools to finish specific tasks**

- Sensor: text/video/image interface

- Planning: choose method to reply

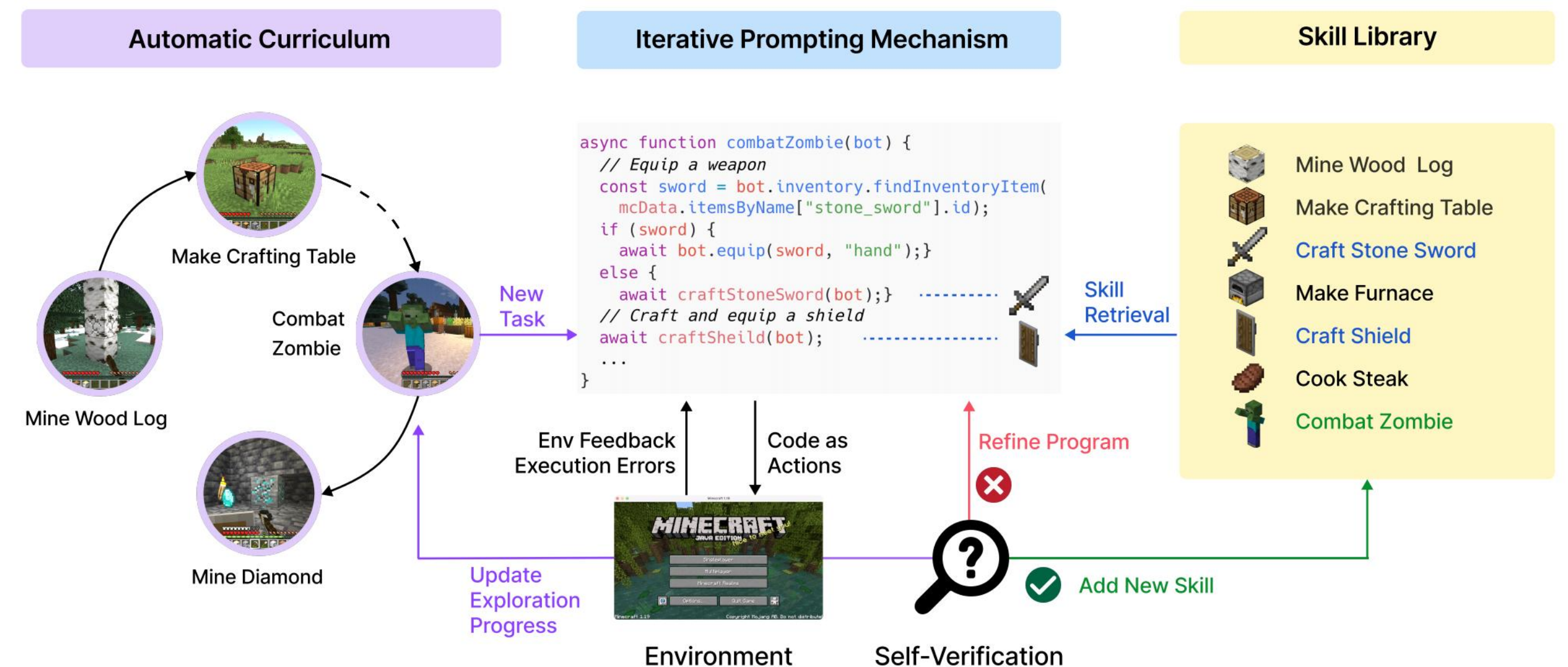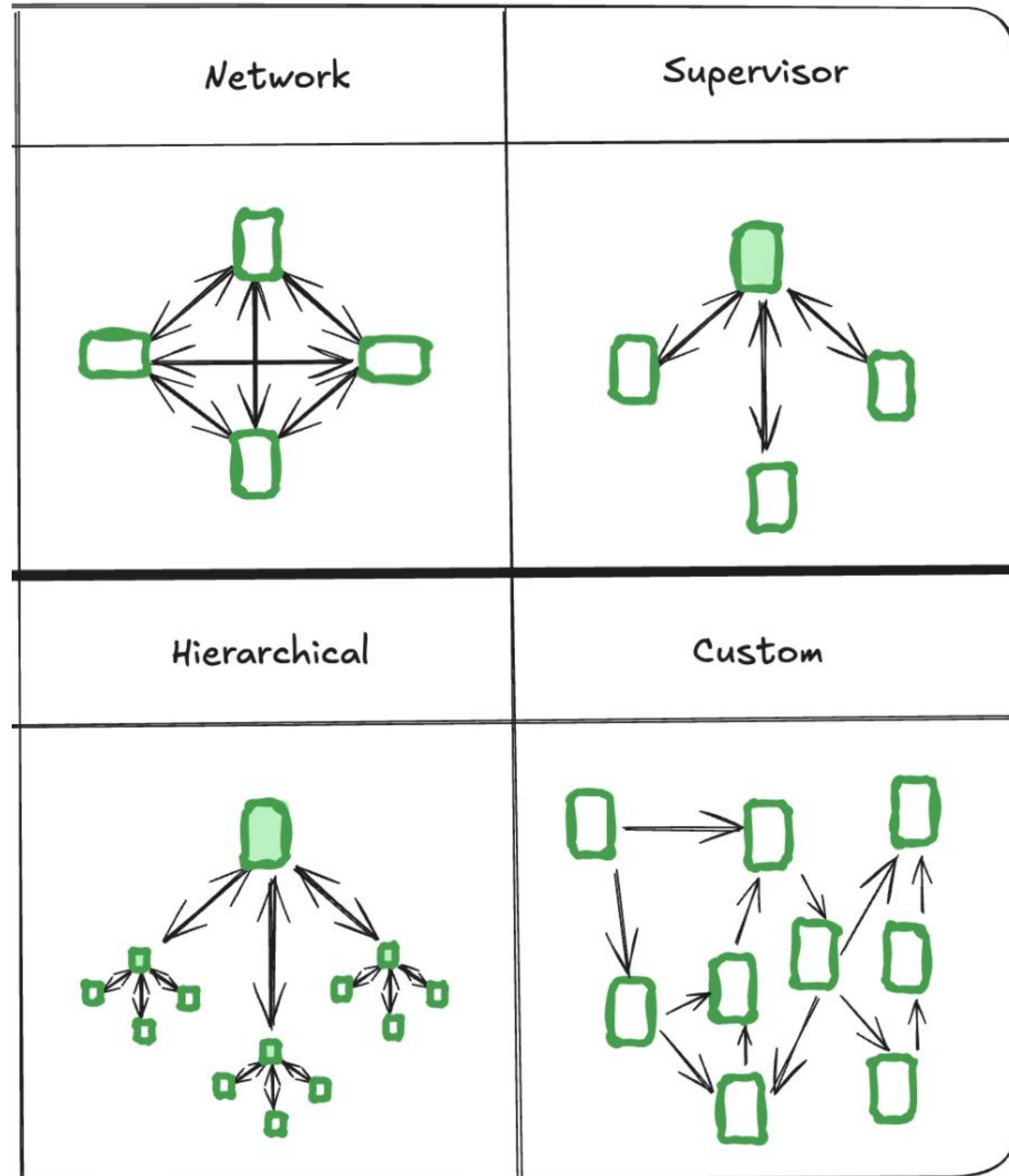- Memory: manage the old behavior and reply history



\* <u>Figures cited</u>

Figure 2: VOYAGER consists of three key components: an automatic curriculum for open-ended exploration, a skill library for increasingly complex behaviors, and an iterative prompting mechanism that uses code as action space.

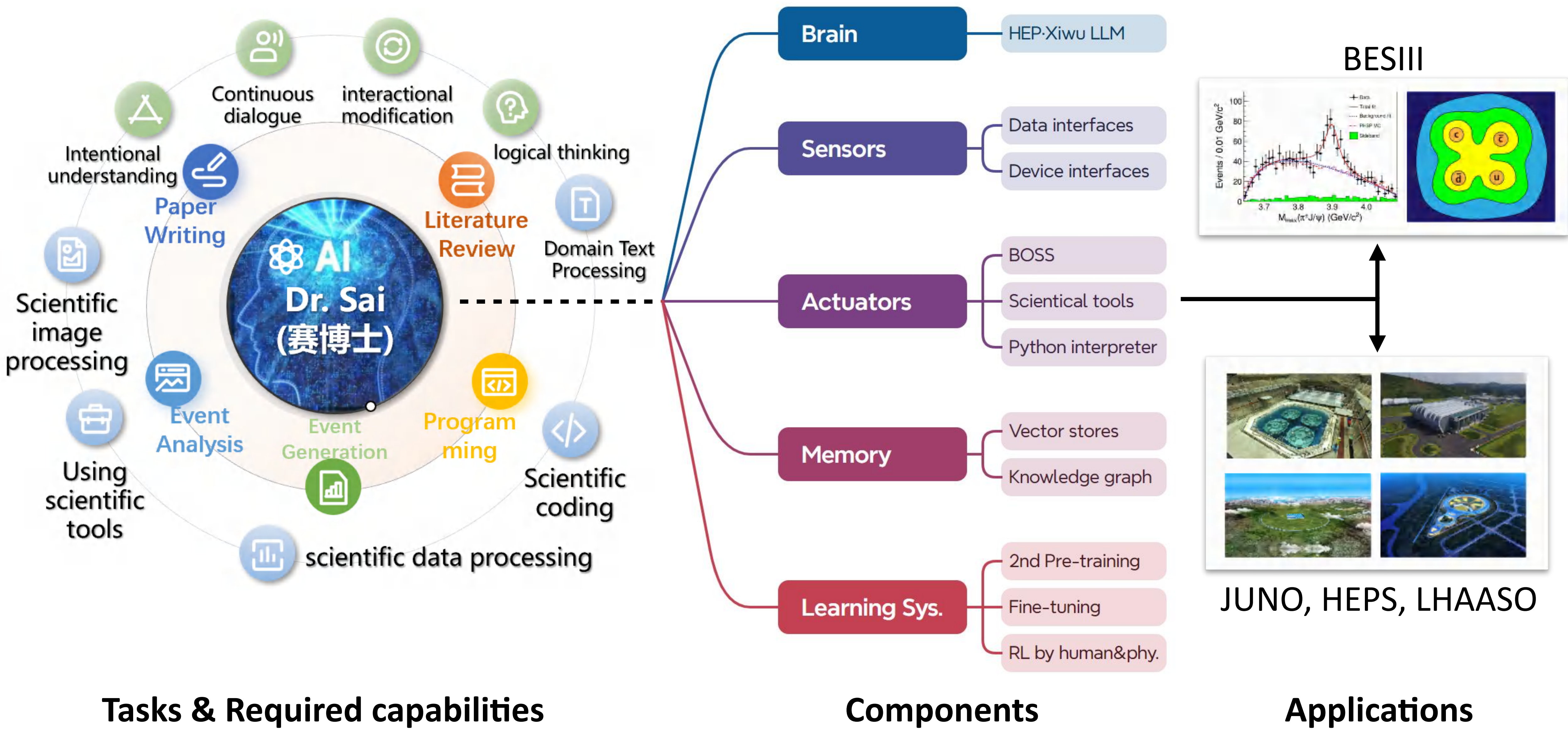\* <u>VOYAGER</u>: an agent of Minecraft

# Multi-Agent System (MAS)



For different types of tasks, agents will have different prompts and tools, so a multi-agent system (MAS) is needed.
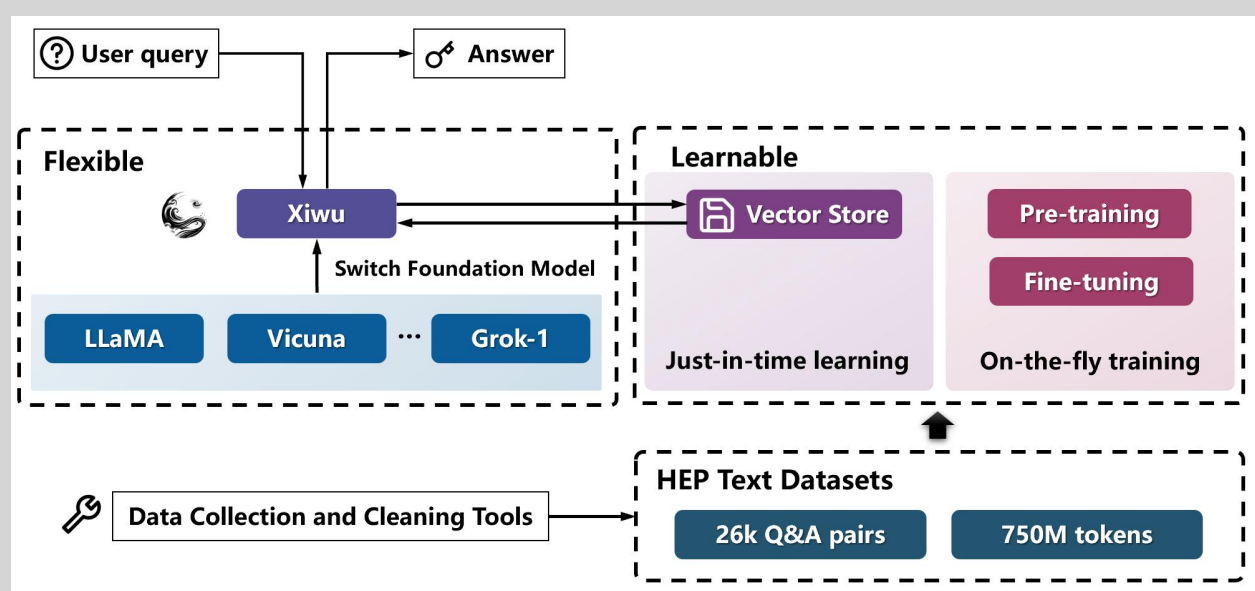
Some architectures of MAS:

- **Round-table**: fixed-order speaking
- **Network**: free connection among agents
- **Supervisor**: each agent is connected with a supervising agent, who decides to call others
- **Hierarchical**: Supervising agent as intermediate layer
- **Custom**: each agent talks with designated agents
- **Talk-reasoner**: Supervising agent interact with human

# Overview of Dr.Sai project



**Tasks & Required capabilities**

**Components**

**Applications**
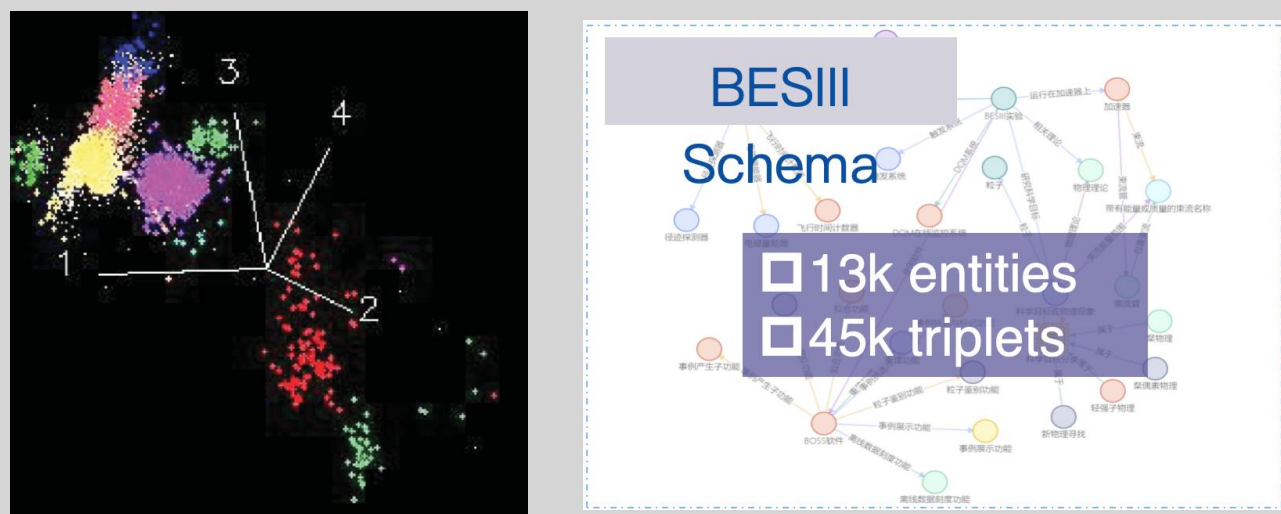
# Dr.Sai single agent

## Brain



HEP Xiwu LLM

- Based on Llama3-8B, 70B, and DeepSeek
- Trained on A800 and DCU from Sugon
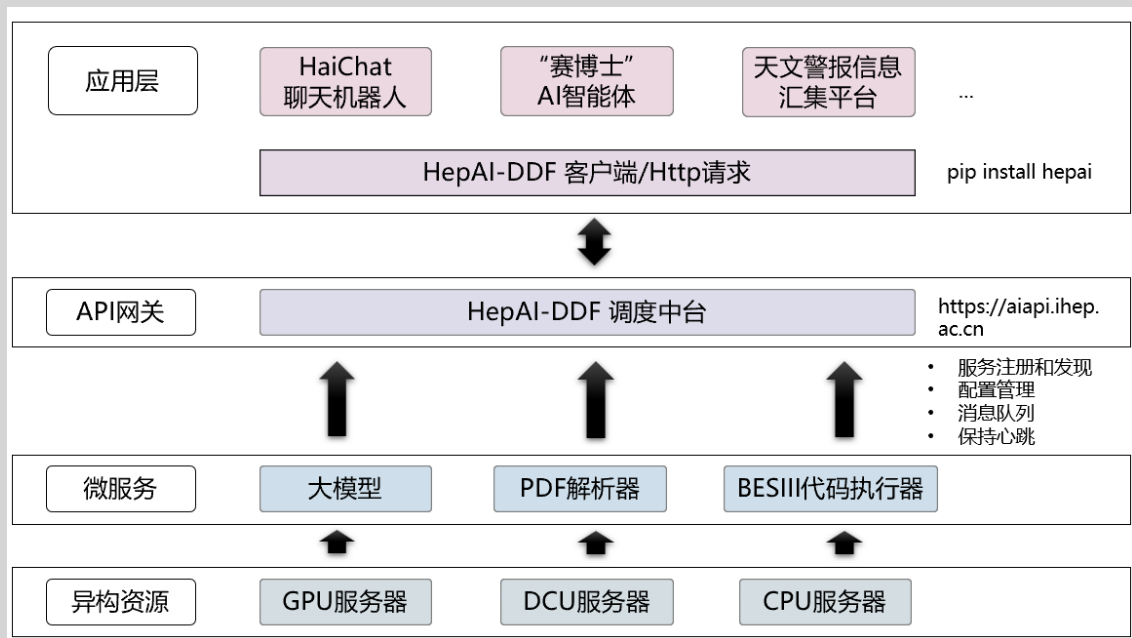- Enhanced domain capabilities

## Memory



Vector store
Knowledge graph

- Vector store
  - BAAI BGE-M3
  - Llama-Index
- Knowledge graph
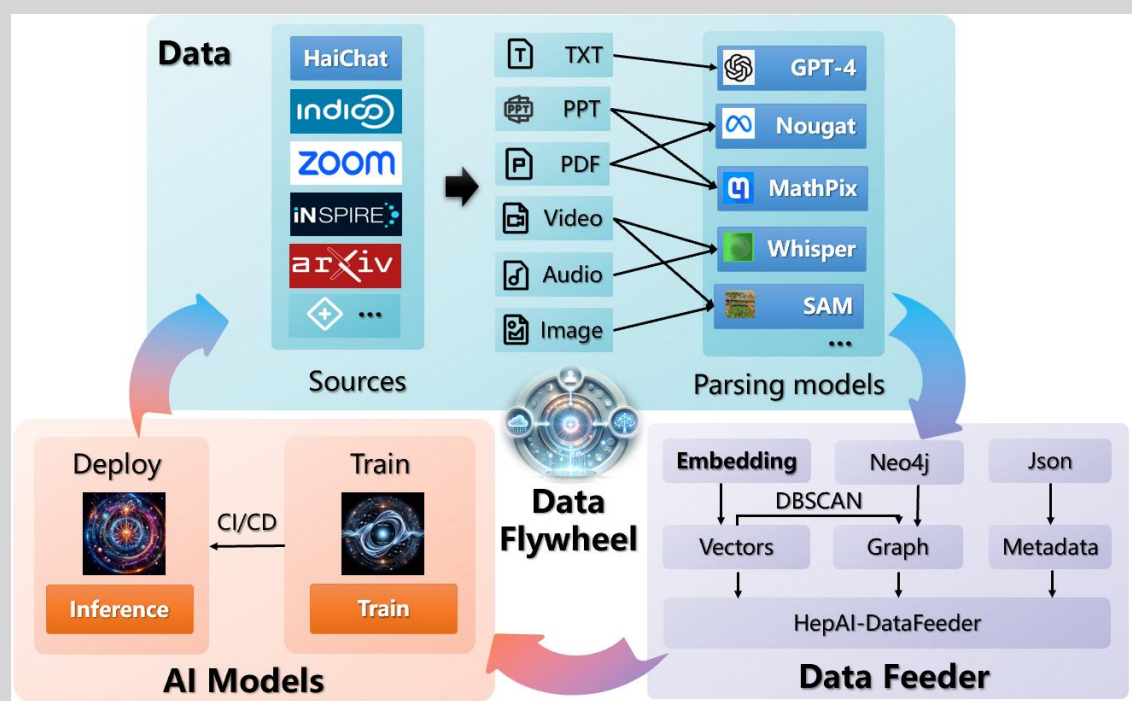  - LightRAG

## Actuators



HepAI DDF v2.0

- BESIII actuator
  - BOSS worker
- Daisy framework
  - HEPSCT, ptycho
- Unified access of models, tools
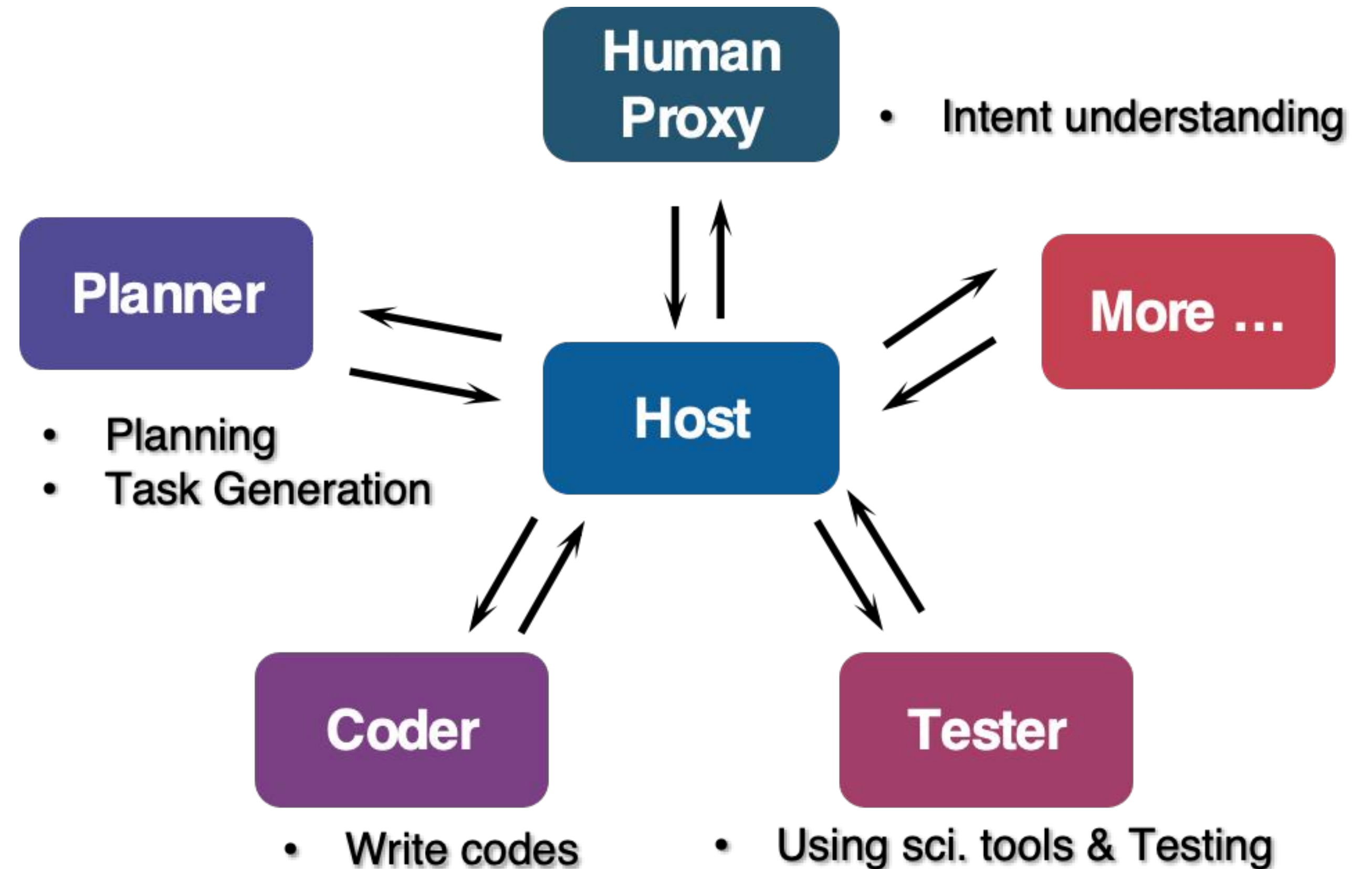
## Sensors



Data Flywheel System

- Data flywheel
  - HaiNougat
  - HaiChat
  - Whisper, SAM
  - arXiv, indico
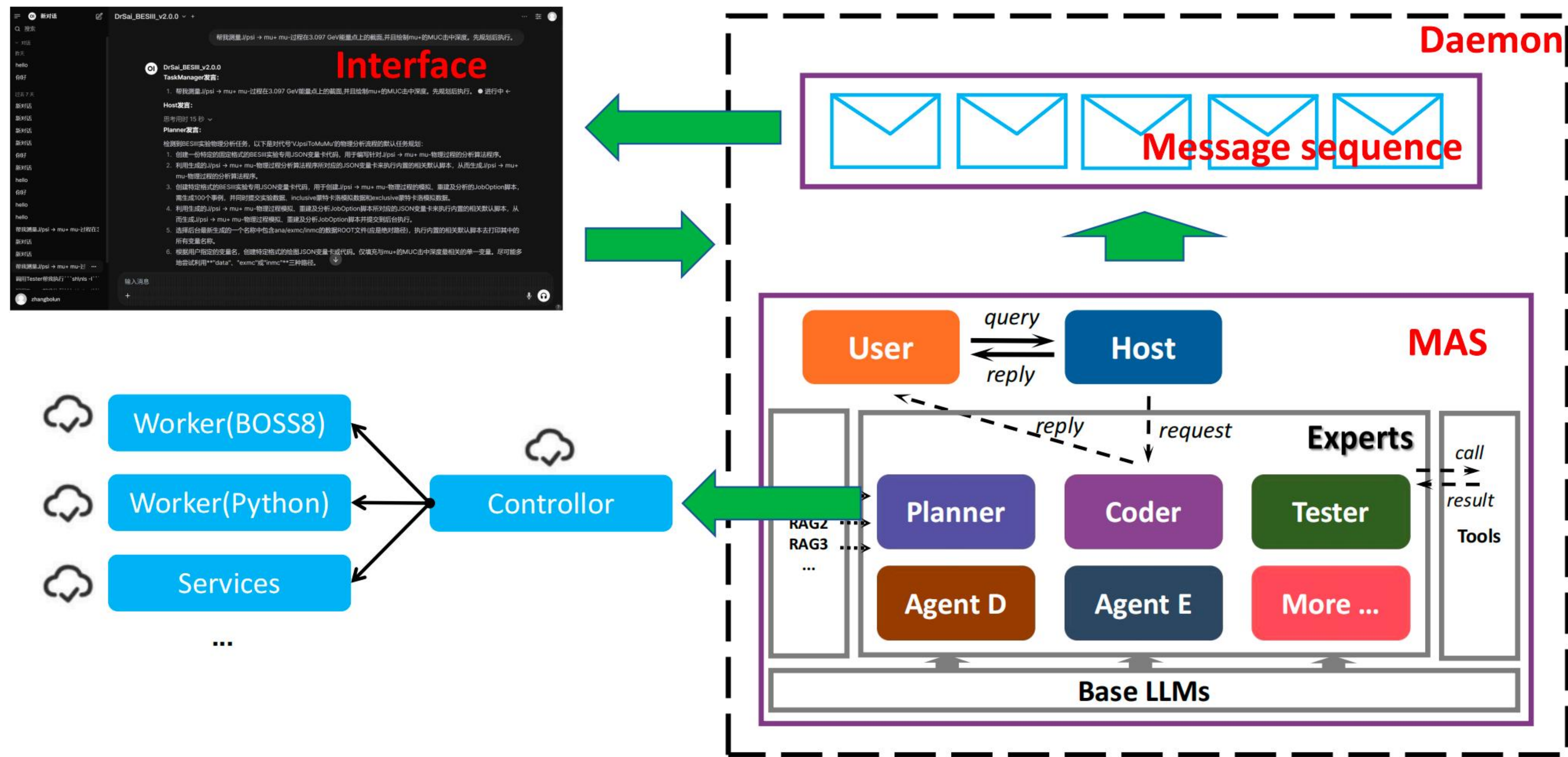- DOMAS framework
  - Scientific metadata

# Dr.Sai-BESIII

**Dr.Sai-BESIII** is an application of Dr.Sai framework, focus on the automated workflow of the **BESIII** experiment
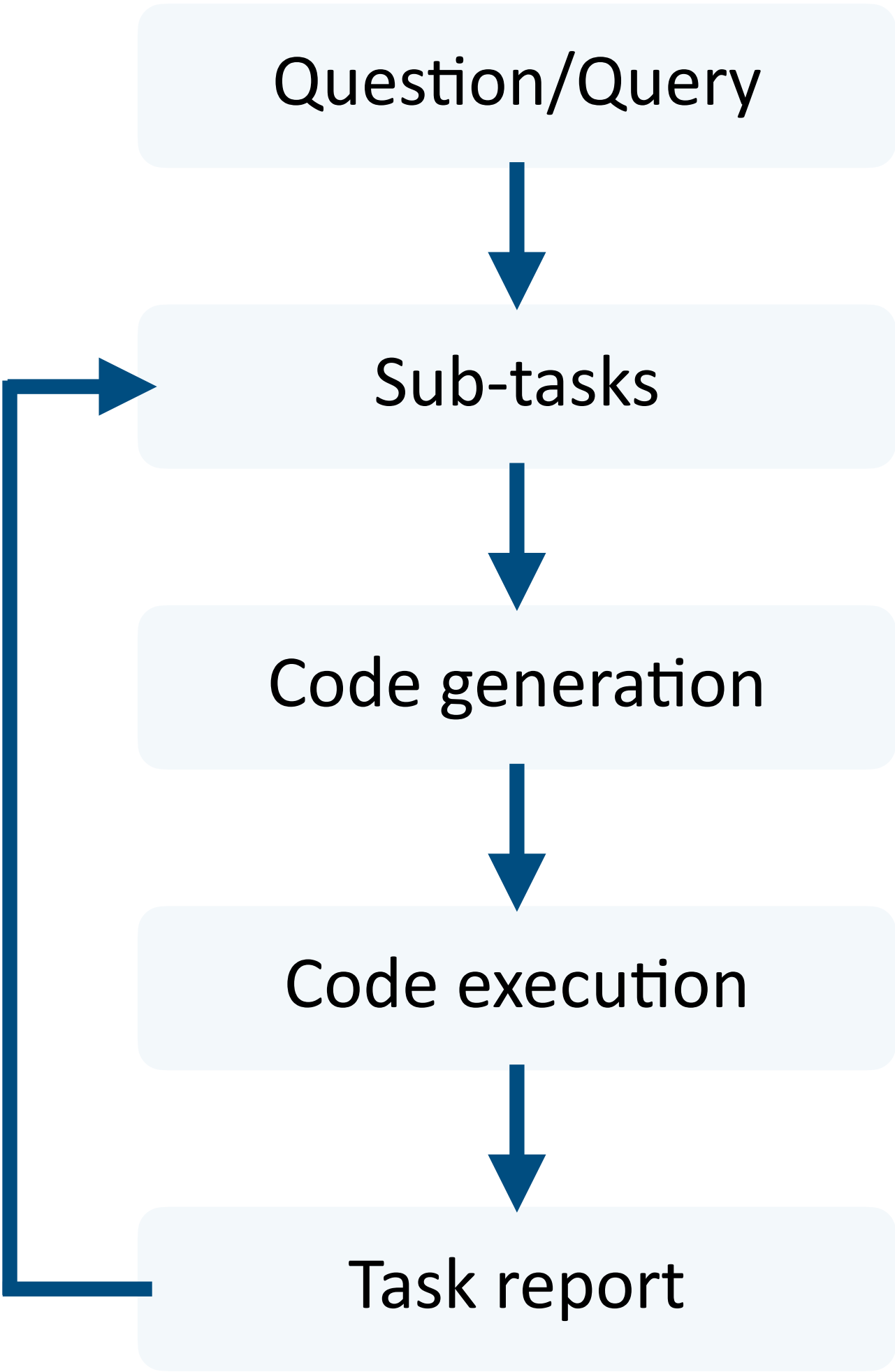
- Talker-Reasoner architecture based on AutoGen framework
- Each agent is equipped with specific knowledge, tools, and LLM
- A Host agent is introduced to manage group chats, making it easy to expand with more agents.
- A human proxy agent is introduced to allow humans intervene at any time.

# Workflow of Dr.Sai-BESIII
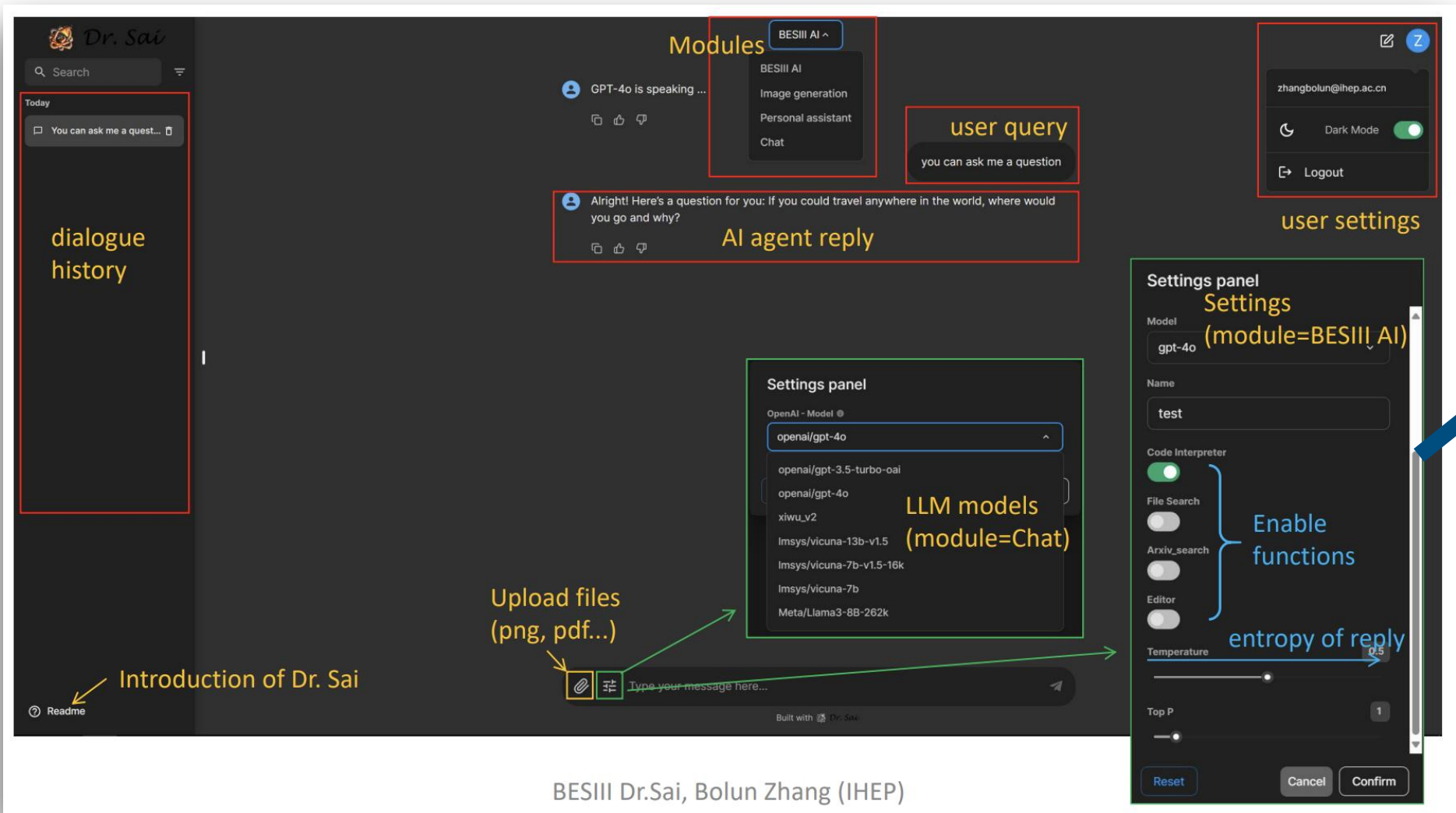
# One example of Dr.Sai-BESIII



- Multi agent collaboration
- Professional domain abilities
  - BESIII Knowledge
  - Physics analysis planning
  - BOSS coding
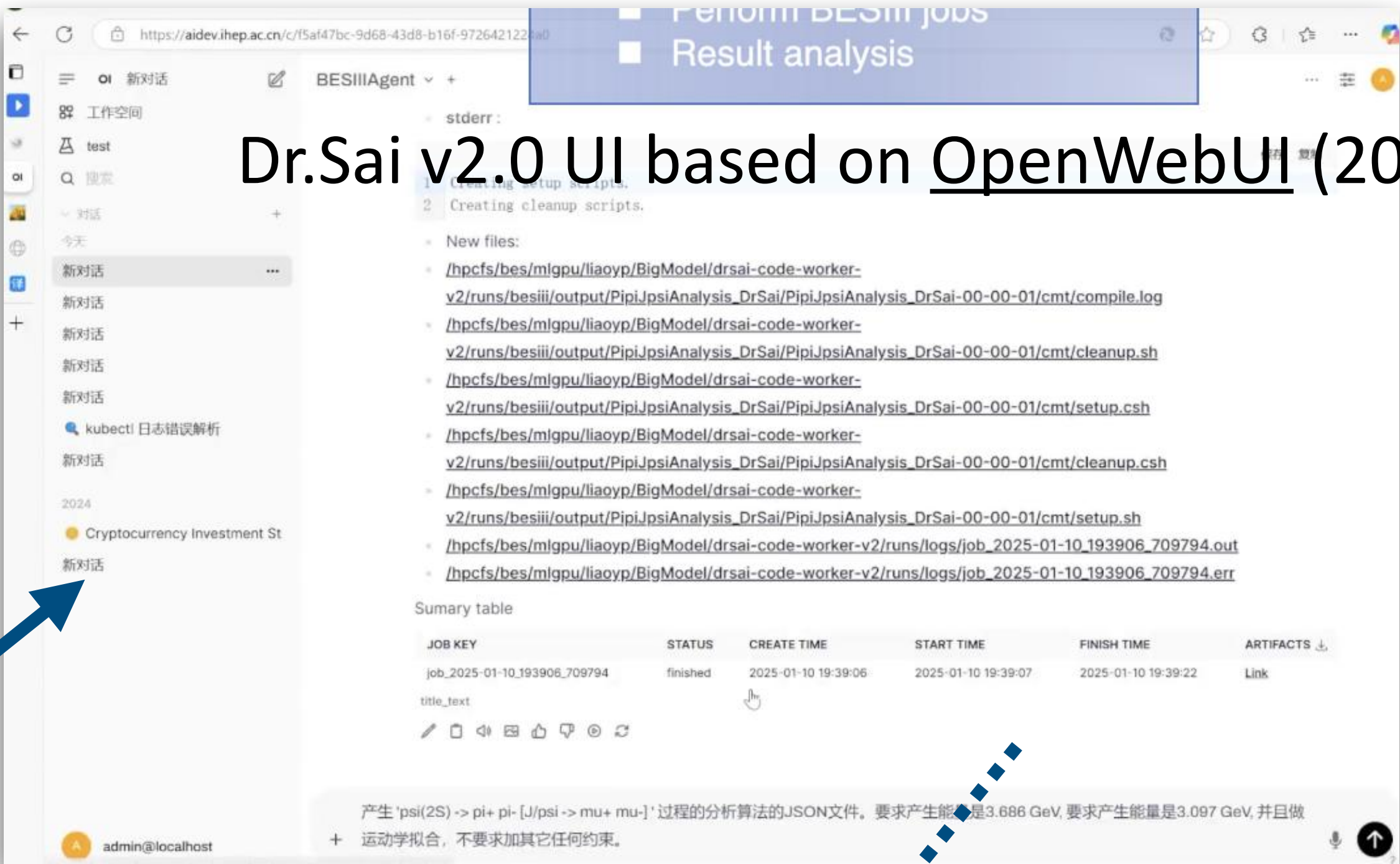  - Perform BESIII jobs
  - Result analysis

Question/Query → Sub-tasks → Code generation → Code execution → Task report

# Interface of Dr.Sai-BESIII

Dr.Sai v2.0 UI based on OpenWebUI (2025.07)



Dr.Sai v1.0 UI based on Chainlit (2024.07)



Dr.Sai v2.x (Planned) UI based on Magentic-UI (2025.xx)

# Prospect of Dr.Sai-BESIII



**Dr.Sai v1**

Test feasibility

**Now**

**Dr.Sai v3**

Add the selection optimization and simple signal extraction

**Dr.Sai v5**

Automate the whole workflow

Start at 2023.06

2024.07

2025.07

2026.12

2027.06

2029.06

**Dr.Sai v2**

Automated workflow up to the preliminary event selection

**Dr.Sai v4**

Preliminary results from the first application at BESIII

# In the end

⭐ LLM could be very helpful for HEP, and also for other domains!

   ⭐ Not just generate draft code/text, but also can be used to **automate the analysis**

⭐ **Dr.Sai** - a multi-agent collaborative system has been developed, and one of the application **Dr.Sai-BESIII** preliminarily demonstrates its ability to automate physics analysis processes
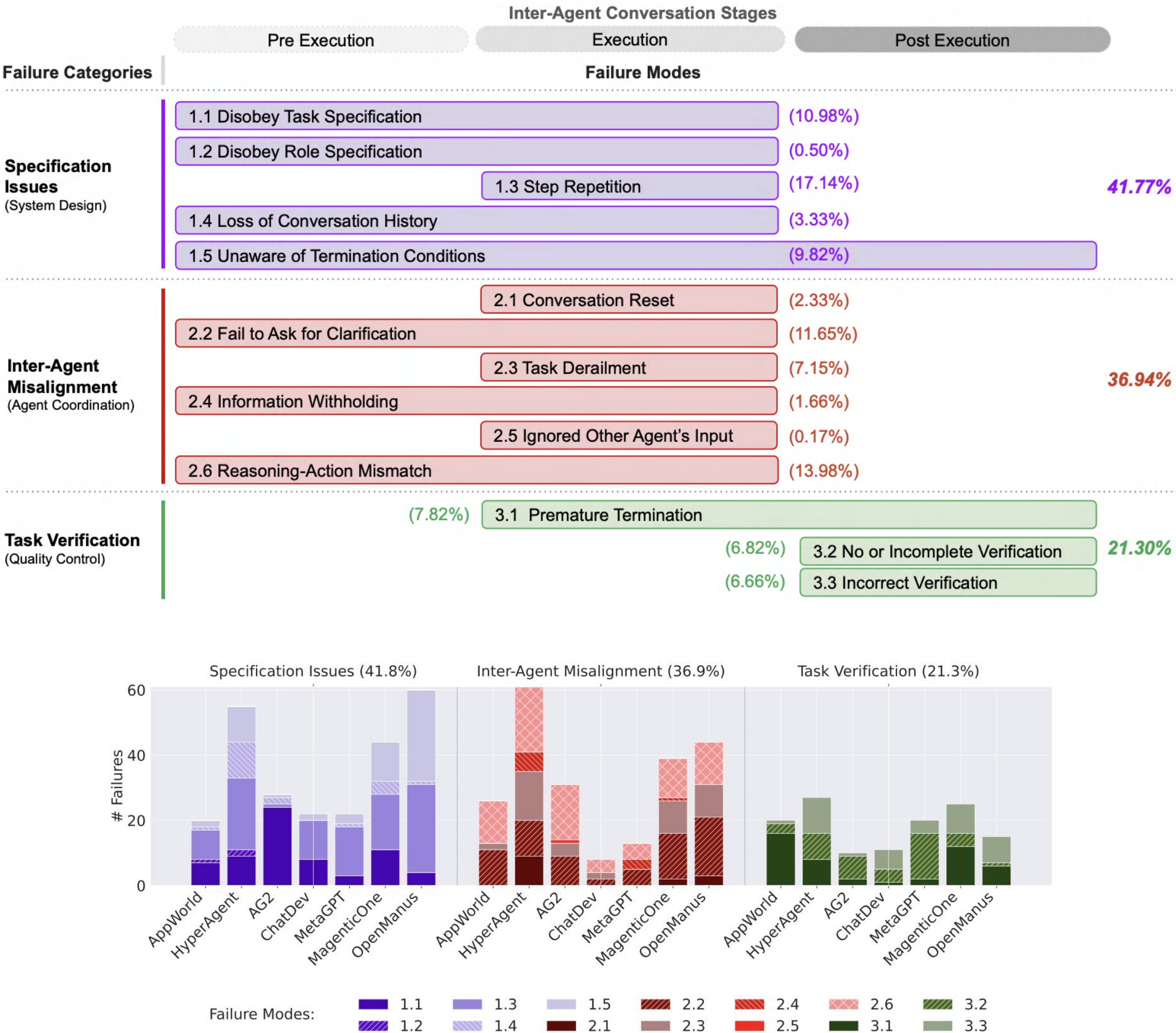

⭐ Dr.Sai is a new system, which still needs further improvement

   ⭐ DSL development; Benchmark design; Better HEP-LLM …

   ⭐ In the future, use Dr.Sai to **quickly explore hundreds of physical processes**, which is expected to be used to search for new hadron states or new decays

Yipu Liao (廖一朴)
IHEP, CAS, Beijing
liaoyp@ihep.ac.cn

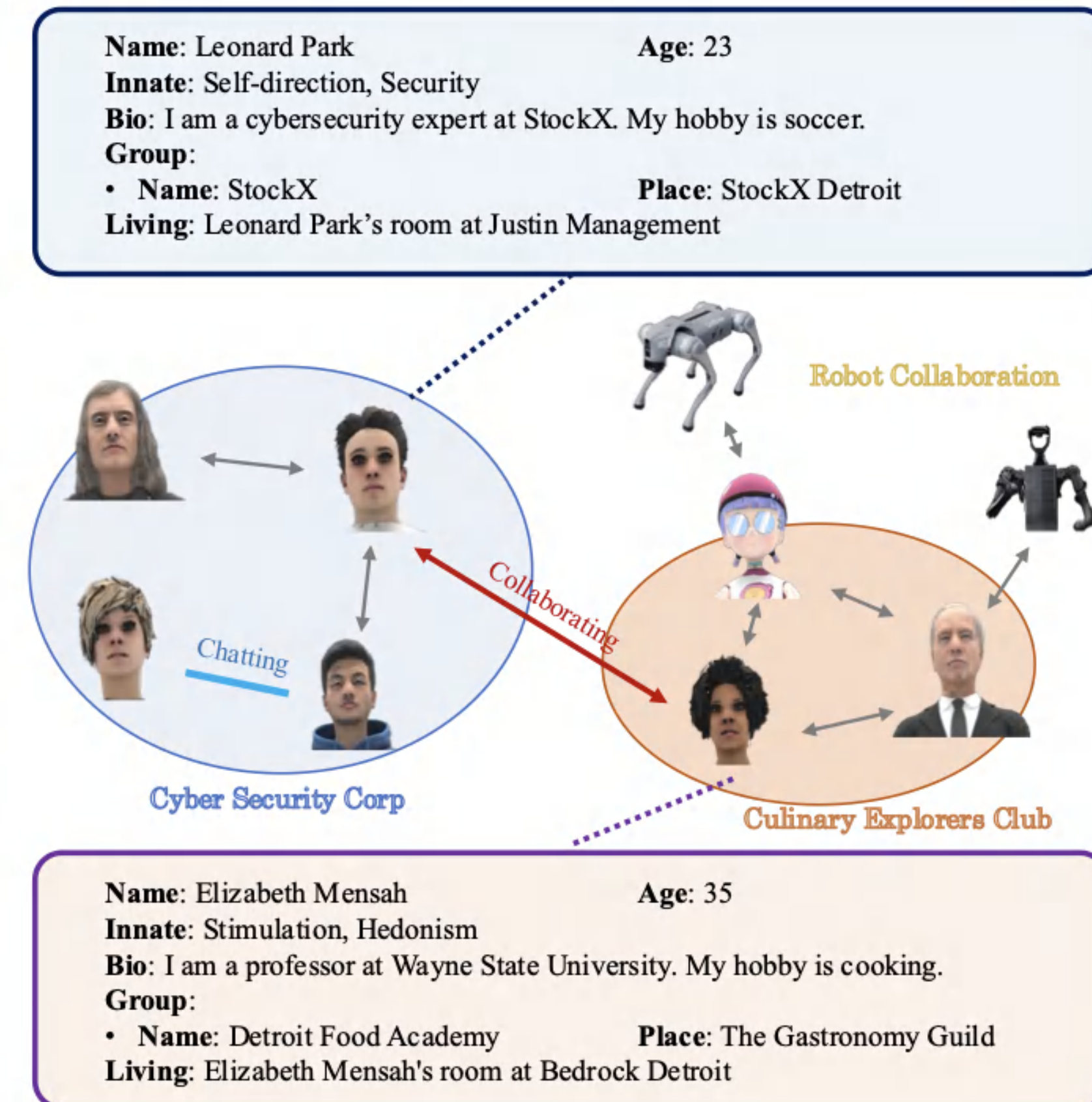# Thank you for listening!

# LLM nowadays: Multi-Agent System

However, a lot of questions appear in the MAS (Still developing):

1. **Specification issues** (System design)

   Wrong system design decisions, and poor or ambiguous prompt specifications.

2. **Inter-agent misalignment** (Agent coordination)

   Incorrect or redundant reasoning of agents

3. **Task verification** (Quality control)

   Difficult to determine whether the task result is correct

\* <u>Why Do Multi-Agent LLM Systems Fail?</u>

21

# Some applications based on MAS

- Real geospatial data and physics engine
- Agents with different identities and Robots to form virtual communities with World Model (3D LLM)
- Focus on community interaction and task execution (such as campaigning and voting)

# Some applications based on MAS



- A Multi-Agent AI (USTC) + automated Lab for Chemistry
- Chemical focus
- Reader + Planner + Robot Operator + Executor + DSL

# Xiwu (溪悟): a HEP-specific fine-tuned LLM

- The **brain** of Dr.Sai to enhance the domain capabilities
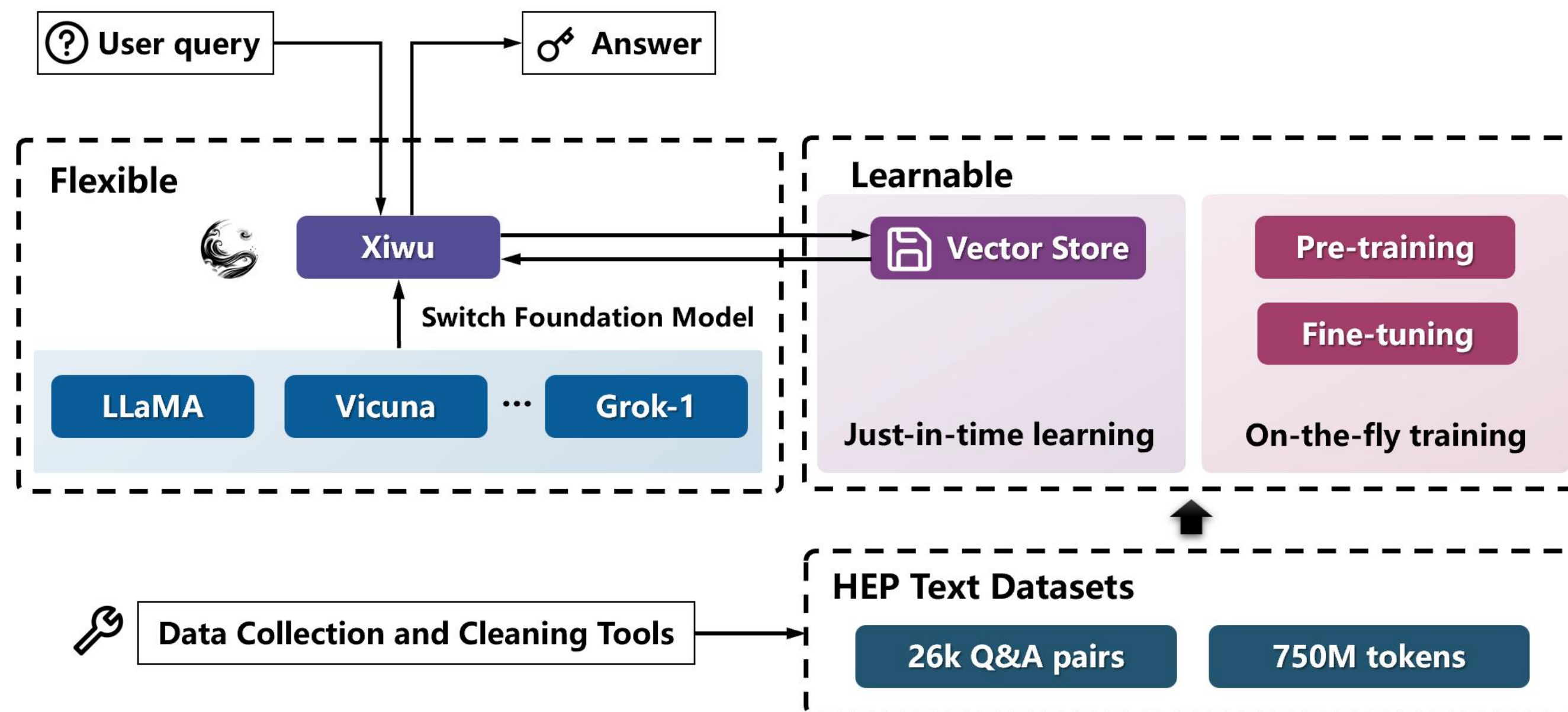- Currently, Xiwu is based on LLaMA3, will be upgrade to **DeepSeek**, Historically, LLaMA, LLaMA2
- Secondary pre-training and fine-tuning techniques
- Trained on GPU-A100 previously, now trained on DCU-K100 (GPU from Sugon)
- Significantly better than the base model in **HEP Q&A** and **internal code generation** (at that time)



Test results

* arXiv:2404.08001        * *Preliminary trial work*        24        https://github.com/zhangzhengde0225/Xiwu

# Xiwu (溪悟): a HEP-specific fine-tuned LLM

**?** Please provide an example of simulate particle decay on BOSS    *BOSS: BESIII Offline Software System*

**GPT-4:**

BOSS (BesIII Object-oriented Simulation Software) is a software package used for event simulation at the BESIII experiment. Here is a simplified example of how a particle decay simulation might be conducted using BOSS:

**Initialization**: Before the simulation can begin, the software must be initialized. This includes setting the parameters for the experiment, such as the type of particle to be simulated, the energy levels, etc.

**Generation of Initial Particles**: The software generates a set of initial particles based on the parameters set in the initialization step. These particles are then "fired" into the simulated detector.

**Tracking**: Each particle's path through the detector is tracked and recorded. This includes any interactions the particle has with the detector material, such as scattering or decay.

**Decay Simulation**: If the particle is unstable (i.e., it can decay into other particles), the software simulates this decay process. This involves randomly determining when the particle will decay (based on its known lifetime) and what particles it will decay into (based on known decay modes).

**Data Recording**: All data from the simulation, including the initial and final states of each particle, are recorded for later analysis.

**Analysis**: After the simulation is complete, the recorded data can be analyzed to learn more about the particle's properties and behavior. This might include things like determining the particle's lifetime, measuring its decay modes, etc.

Please note that this is a simplified description. Actual simulations would be much more complex and require a deep understanding of particle physics.

**Xiwu:**

In BOSS (BESIII Offline Software System), a common particle simulation is:

$$e^+e^- \to \psi(2S) \to J/\psi\pi^+\pi^-$$

The decay card for this example is as follows:

```
Decay psi(2S)
1.0  J/psi  pi+ pi-  JPIPI;
Enddecay

Decay J/psi
0.5  e+  e-     PHOTOS VLL;
0.5  mu+ mu-    PHOTOS VLL;
Enddecay

End
```

In addition, it is necessary to modify the control file of the KKMC generator, such as:

```
*************job options for generator (KKMC)*************
#include "$KKMCROOT/share/jobOptions_KKMC.txt"

KKMC.CMSEnergy = 3.686;
KKMC.BeamEnergySpread=0.0009;
KKMC.NumberOfEventPrinted=1;
KKMC.GeneratePsiPrime=true;
```

In this example, the entire simulation process is conducted at a collision energy of 3.686 GeV and uses the KKMC generator.

This simulated particle decay process is commonly used in high-energy physics experiments to study the properties and behavior of particles. I hope this example is helpful to you!

**Comment:** GPT-4 generated an incorrect concept of BOSS, offering a solution that seems correct but is useless, while Xiwu correctly explained the concept of BOSS, providing the correct decay card and example code for generator.
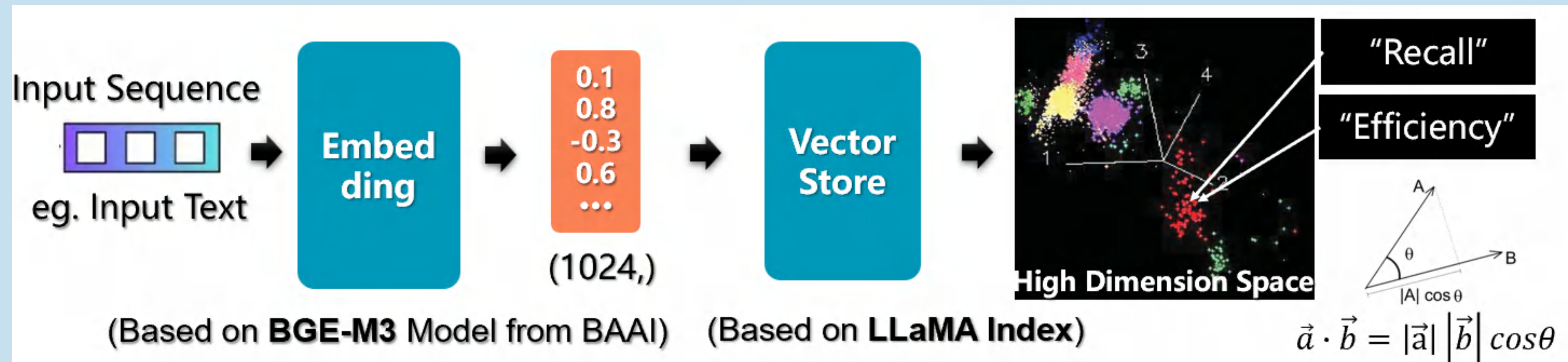
# Current training data for Xiwu (溪悟)

- Cleaned and evaluated papers arXiv
  - PDF files parser: *HaiNougat*, advanced iteration of the Nougat model
- Good quality chat history from HepAI platform
  - The data is cleaned by human or AI (GPT4)
  - 180k Question-Answer pairs in 3 months
- **BESIII internal data**
  - Internal memo/notes, parsed by HaiNougat
  - Question-Answer pairs from internal paper review
  - BESIII Offline Software System (BOSS) source code
  - BESIII public webpages and internal webpages

https://github.com/ai4hep/hai-nougat

# Memory component of Dr.Sai

Memory component is based on **vector store** and **knowledge graph**



- Embeddings (BGE-M3 model from BAAI)
- RAG system (LlamaIndex)
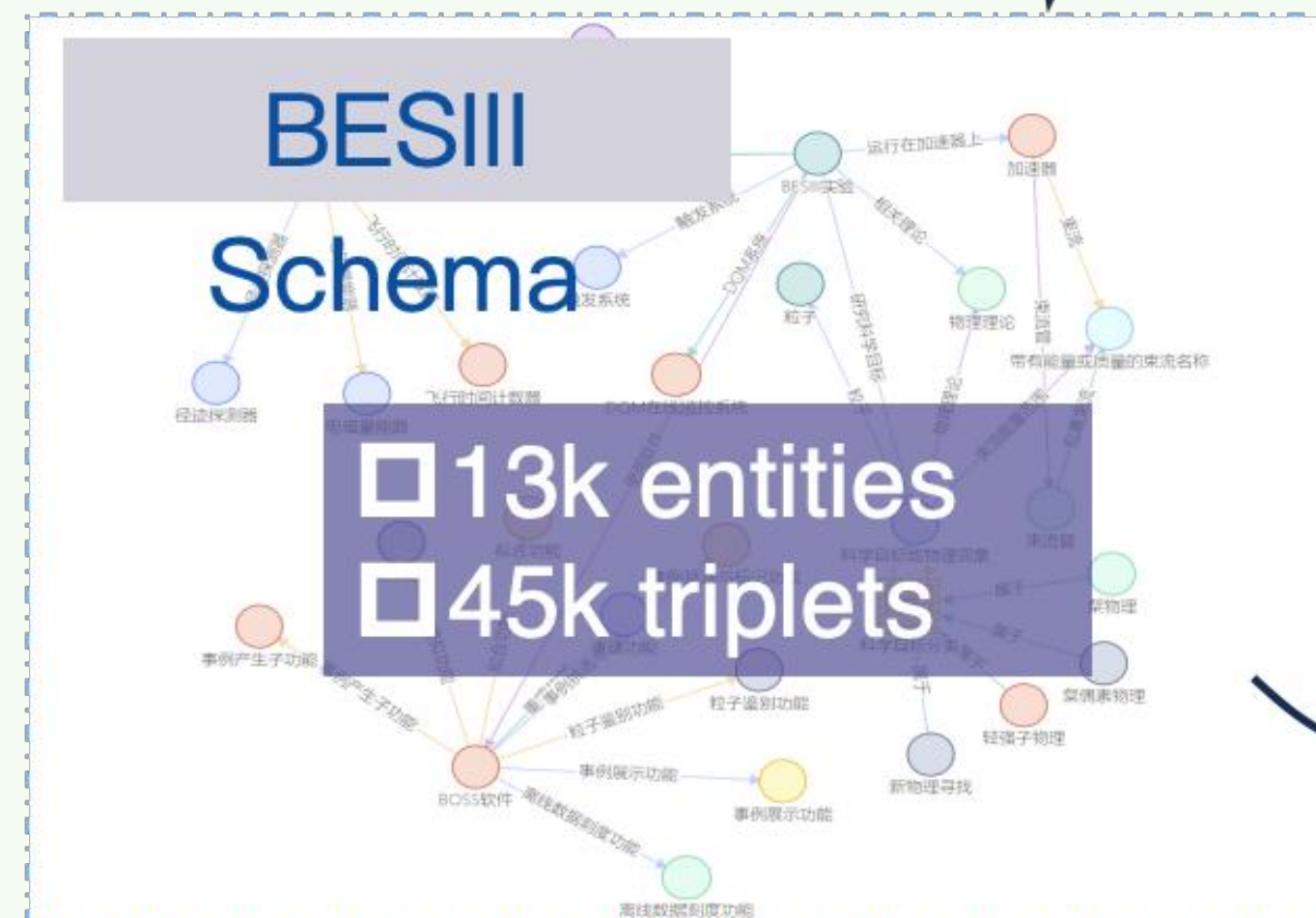- Unified metadata and support of fragment segmentation and full-text search

# Memory component of Dr.Sai

Memory component is based on **vector store** and **knowledge graph**

- Knowledge extraction to build BESIII schema
- Knowledge mapping (Entity alignment, relationship chain)
- Graph storage based on Neo4j
- KG-RAG system (LightRAG)



Tell me the infomation about Y(4360)

[Incorrect Output] Y(4360) is …

Verify

[Verified Output] Y(4360) is …

[Corrected Output]  Y(4360) is …

600+ BESIII papers

BESIII Schema

☐ 13k entities
☐ 45k triplets

*Y(4360): an exotic particle with 4.36 GeV/c2 mass*

28

# From answer to action: HepAI-DDF

## High-performance AI Distributed Deployment Framework

- ☑ **Unified Access**: Seamless integration of AI models, scientific tools, vector databases, and knowledge graphs across heterogeneous resources (CPU/GPU/DCU)

- ☑ **Smart Schedule**: Automated heartbeat monitoring + dynamic load balancing

- ☑ **Unified API Gateway**: Standardized service interfaces with API-KEY authentication for enhanced security

- ☑ **Remote Model Invocation**: Support for remote server model deployment with seamless local calling experience

- ☑ **Adapted MCP**: Model Context Protocol, Access to 2600+ MCP Server

*\* Generated by AI*

✦ 10+ models
✦ 2 local tools
✦ 2600+ MCP tools
✦ In one server

# From answer to action: BESIII code worker

**BESIII code worker**　　　　**HaiDDF**　　　　**Dr.Sai**



- I/O Tools (files and data)
- Message interface
- Executor and Job management

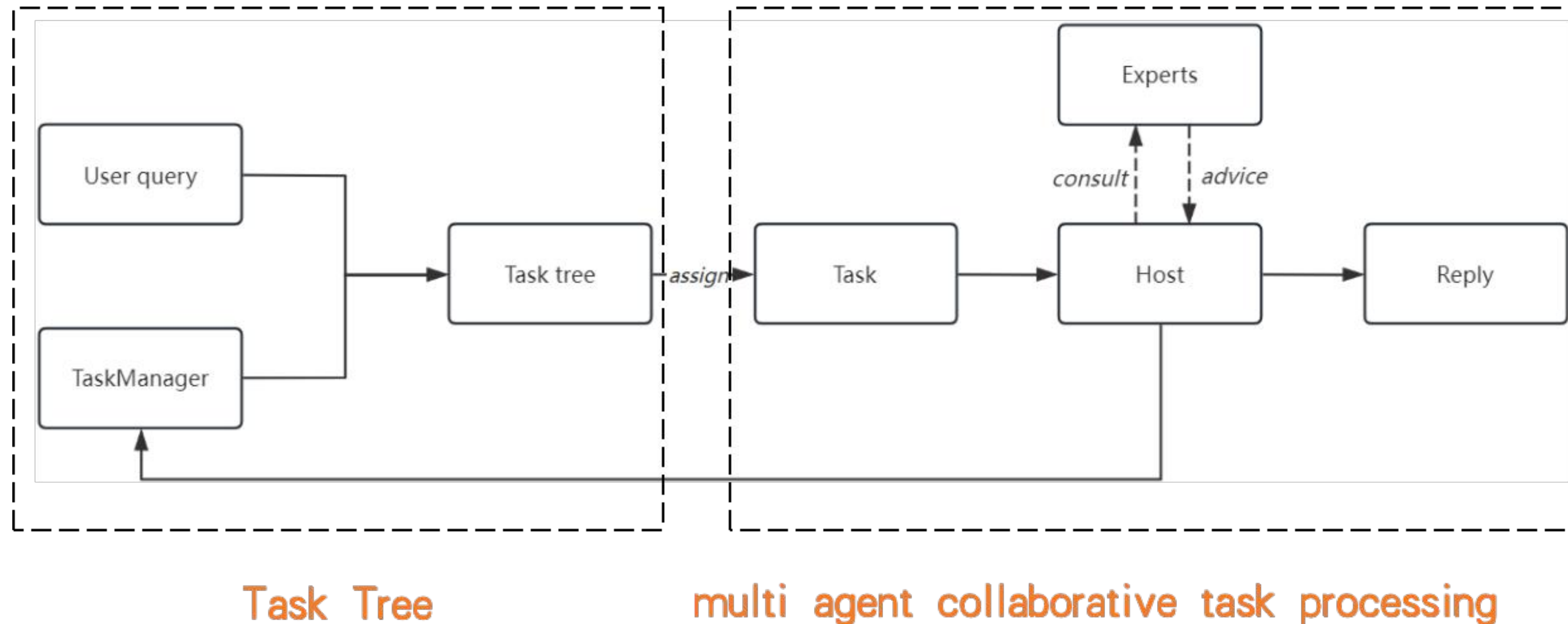From **Q&A chatbot** to **expert assistant**

# Data flywheel system

**The "data flywheel" enables continuous model iteration and evolution**

- Improves models by constructing circular data pathways
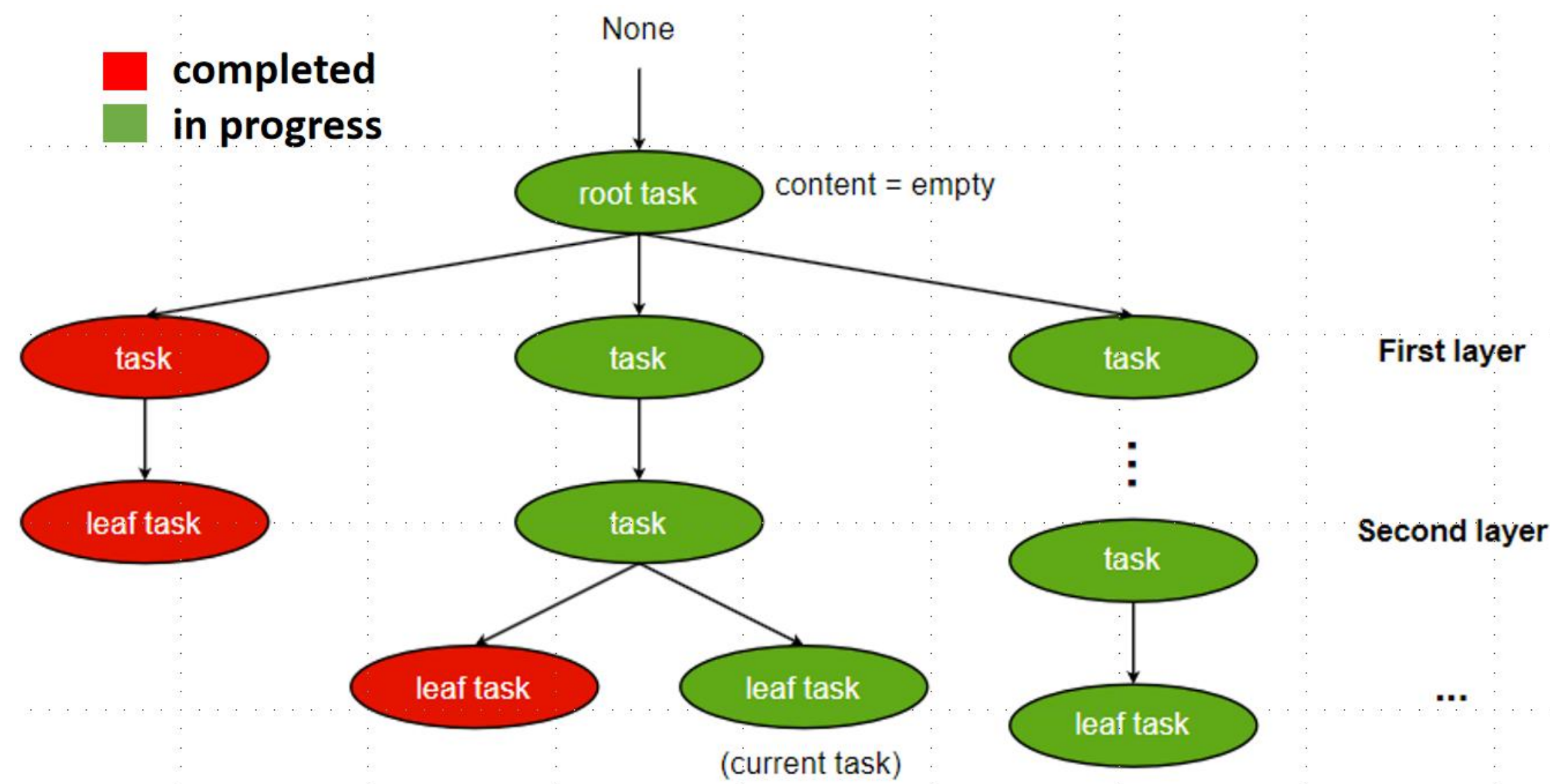- With more collected data for cleaning and labeling, further enhance the models

# MAS optimization: Collaboration mode

- Based **Talker-Reasoner** agent mode
- Multiple rounds of inquiry to achieve task execution results reflection
- Upgrade to **Host** led multi-agent collaboration architecture, with direct response to routine problems and intelligent allocation of complex tasks



Task Tree          multi agent collaborative task processing

# MAS optimization: Task Tree

- Planner divide a user-query (parent task) to sub tasks, and use **Hierarchical** task model
- Adopt the **task tree** to manage complex physics analysis tasks
- Fractal design ensures a closed-loop logic for task backtracking
- Implement task tree addition, deletion, modification, and query



**Task Tree structure**



**Task retrieval logic**

# MAS optimization: Task Tree

• **Also support long-task execution**

*Long-task: tasks with long execution time*



**recover**

根据用户输入进行行动
• "status"：调用Tester查询任务状态
• "continue"：更新任务状态为"completed"，顺序执行旧任务系统中的下一个任务
• other：提示存在未完成长任务，退出groupchat

**pending**

1. 设置task.status="paused"
2. 文本处理：提取jobkey，延展reply内容
3. 更新groupchat
4. 结束groupchat

---

**Expected messages**

Q: <task1>

A: The jobs are submitted, please hold on. Enter key words: ["status", "continue"] to proceed.

>> exit groupchat

Q: status

A: <job status from Tester>

>> exit groupchat

Q: continue

A: <answer to task2> (groupchat continue)

Q: <other query>

A: You have unfinished long tasks!

>> exit groupchat

# MAS optimization: Domain-specific language

- Current LLM do not know the HEP data analysis procedures and do not understand the logics
- We can interpret the analysis to a **Domain-specific language (DSL)**
  - Define each step of analysis in sequence, so LLM can "understand" the procedure
  - 600+ BESIII published paper → transfer to DSL manually
- **DSL is served as a guide to Dr.Sai-BESIII**
  - Stored into RAG; find the DSL for the analysis similar to user's target analysis and take it as reference

# Performance of Dr.Sai-BESIII

## Introduction to AGENTBEENCH

➢ a first benchmark designed to evaluate LLM-as-Agent across a diverse spectrum of different environments.



➢ Decouple framework design:
- Agent controller: 4 agents in Dr.Sai.
- Worker controller: run plenty of workers at the same time.
- Task controller: run different task at the same time.

➢ Supports parallel testing for Multi-Agent and Multi-Task scenarios

➢ Could be used to Dr.Sai evaluation.

# Update: Magentic-UI

Human can access and interrupt at any time based on <u>Magentic-UI</u>