Computing Center, IHEP, CAS
National HEP Data Center

# AI Progress in IHEP

Zhengde Zhang, on behalf of IHEP-CC
Beijing
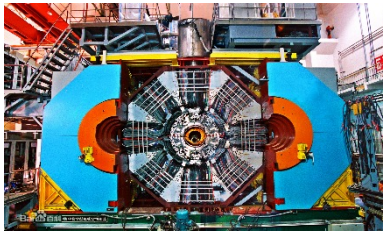Dec, 4, 2024

# Outline

# Why AI for HEP?

## From Large Scientific Facilities to Big Data, then AI becomes a necessity.
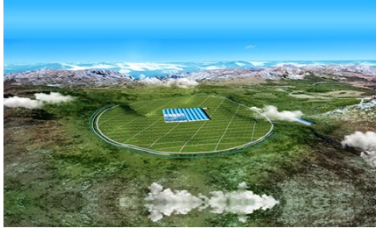
### LFS of IHEP:

- Beijing Electron Positron Collider (BEPC)
- Beijing Spectrometer (BESIII) Experiment
- Jiangmen Underground Neutrino Observatory (JUNO)
- China Spallation Neutron Source (CSNS)
- Large High Altitude Air Shower Observatory (LHAASO)
- Beijing Synchrotron Radiation Facility (BSRF)
- High Energy Photon Source (HEPS)

### HEP Big Data:

- The National HEP Data Center was established in 2019.
- Currently, over **40PB** of scientific data has been stored. (1PB=1024TB)
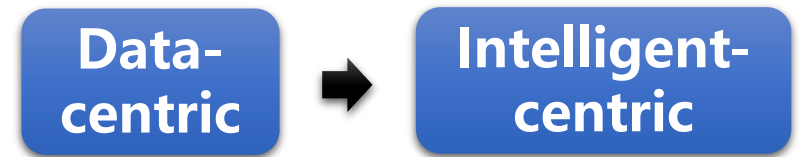- HEPS is expected to generate 800TB/day of raw data after completion in 2025



**BESIII**　　**LHAASO**　　**CSNS**　　**HEPS**

**Scientific Research Paradigm:**

Data-centric → Intelligent-centric

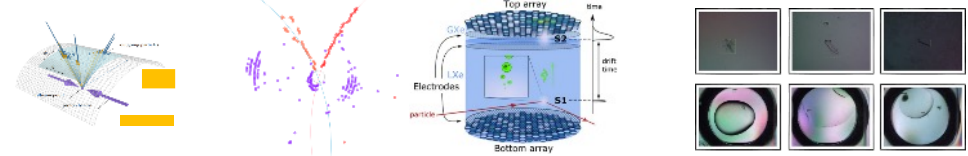$$\Omega(n^3) \rightarrow \Omega(n)$$

**AI offers powerful solutions for unlocking physical potential in massive datasets.**
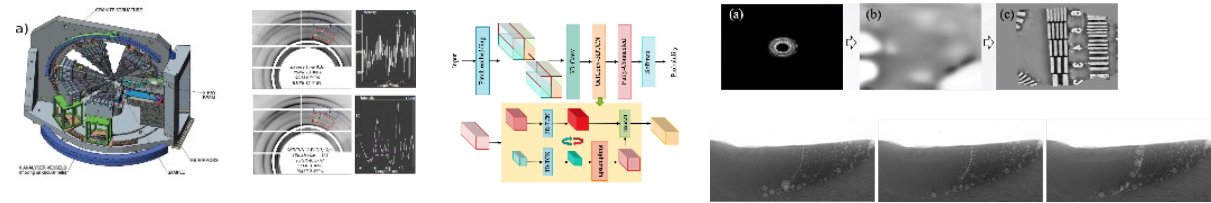
# Overall Strategy for High Energy Physics AI

## The 14th Five-Year Plan of the IHEP:

- Advanced computing technologies including AI, QC make more powerful modeling and simulation possible, which are indispensable means for major breakthroughs in HEP.
- It is necessary to strengthen the application of ML in **real-time processing**, **simulation**, **reconstruction, analysis**, etc., to enhance the ability of high energy physics exploration and new discoveries.
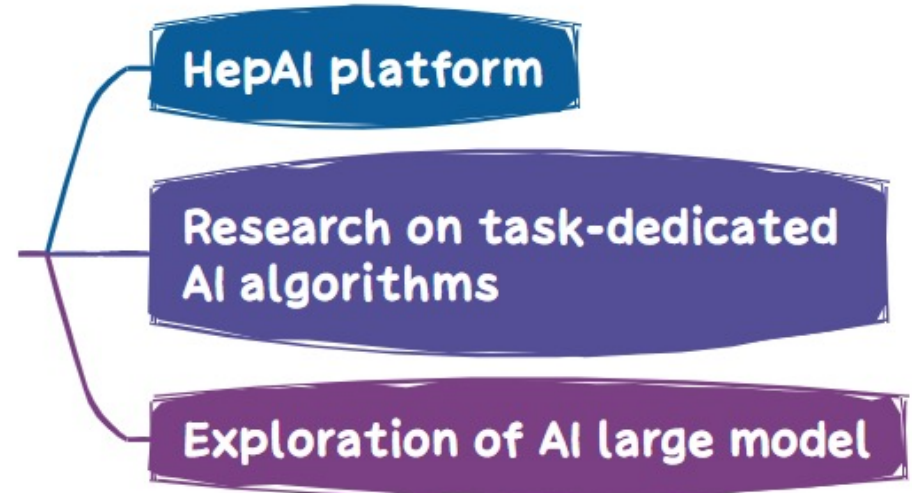
## Our Exploration experiences:

- Prioritizing the development of **AI-Ready scientific datasets** is of utmost importance.
- Insist on the concurrent development of **software frameworks** and **specific AI applications**.
- Investigate the **large models** to propel new discoveries of HEP.

AI4HEP

HepAI platform

Research on task-dedicated AI algorithms

Exploration of AI large model

# Part. 01

HepAI platform

# The HepAI Platform

# HepAI Platform

- Accelerate scientific research in multi-disciplinary scenarios.
- Simplify model iteration and flow.
- Serve as a **common basic infrastructure** for the development of AI algorithms and applications.
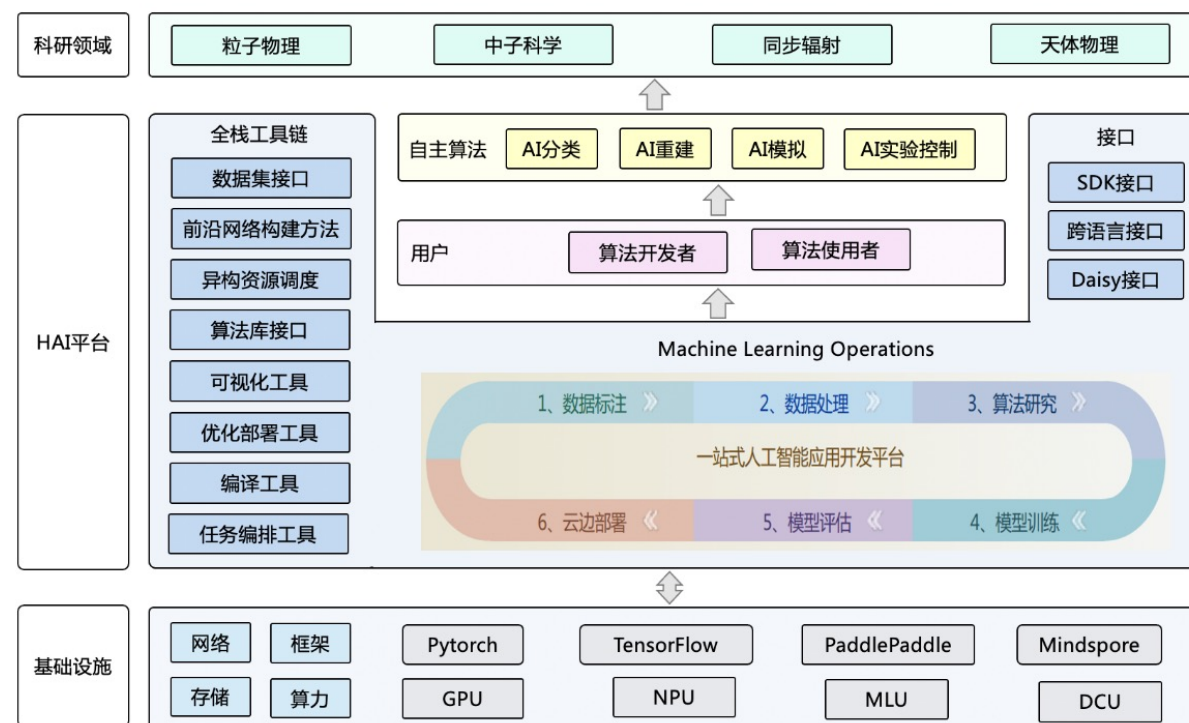
Three Key Elements of AI:

| Data | Algorithms | Comp. Power |
|------|-----------|-------------|

- HepAI Core Codes and Framework (80%)
- 10 AI algorithms.
- 4 AI datasets.
- Heterogeneous computing resources including GPU, NPU, and DCU.

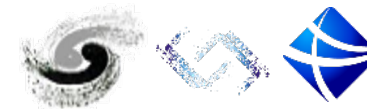Goal: Make AI4HEP **simpler and more advanced!**

Portal site: https://ai.ihep.ac.cn
Open source: https://code.ihep.ac.cn/zdzhang/hai



The architecture of HepAI platform
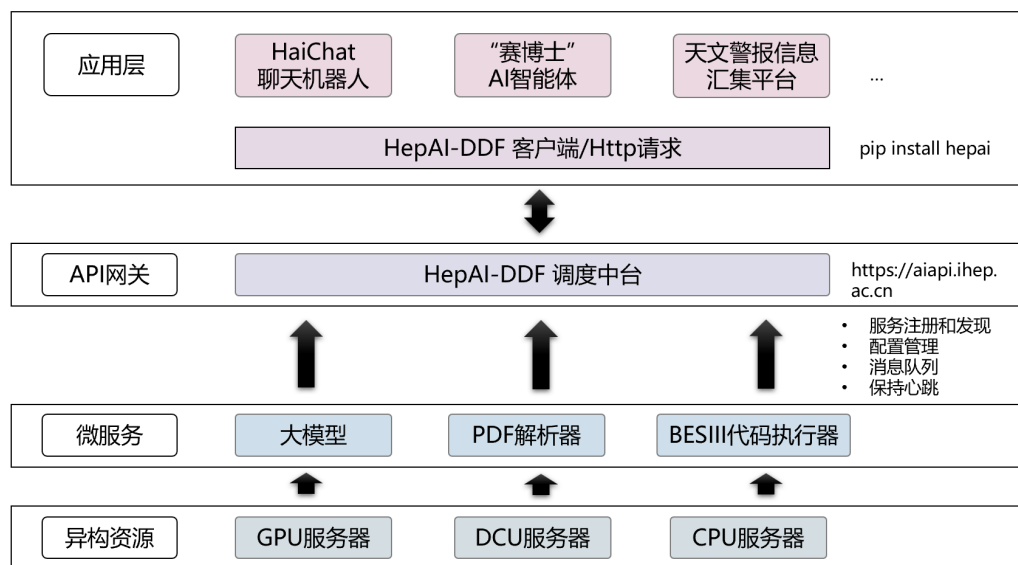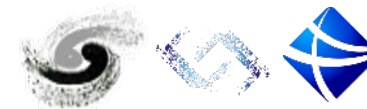
# Key Core Tech. HaiDDF Upgrade to v2

## HepAI Distributed Deployment Framework

HaiDDF2, short for High-energy AI Distributed Deployment Framework, is a key core technology of HepAI platform. It provides crucial technology for unifying access to services such as AI models, scientific tools, vector databases, and knowledge graphs deployed on heterogeneous resources.

Features :
- Unified API Gateway with API-KEY authentication
- Support for unified access to models and algorithms deployed on heterogeneous resources like CPUs, GPUs, and DCUs
- Support for services including AI models, scientific tools, vector databases, and knowledge graphs
- Automatic heartbeat monitoring and load balancing scheduling
- Supports both streaming and non-streaming outputs
- **Remote Models**! Deploy models on a remote server and use the get_remote_model method to obtain remote model instances, allowing you to call them just like local models.

| 应用层 | HaiChat 聊天机器人 | "赛博士" AI智能体 | 天文警报信息 汇集平台 | ... |

HepAI-DDF 客户端/Http请求    pip install hepai

| API网关 | HepAI-DDF 调度中台 | https://aiapi.ihep. ac.cn |

- 服务注册和发现
- 配置管理
- 消息队列
- 保持心跳

| 微服务 | 大模型 | PDF解析器 | BESIII代码执行器 |

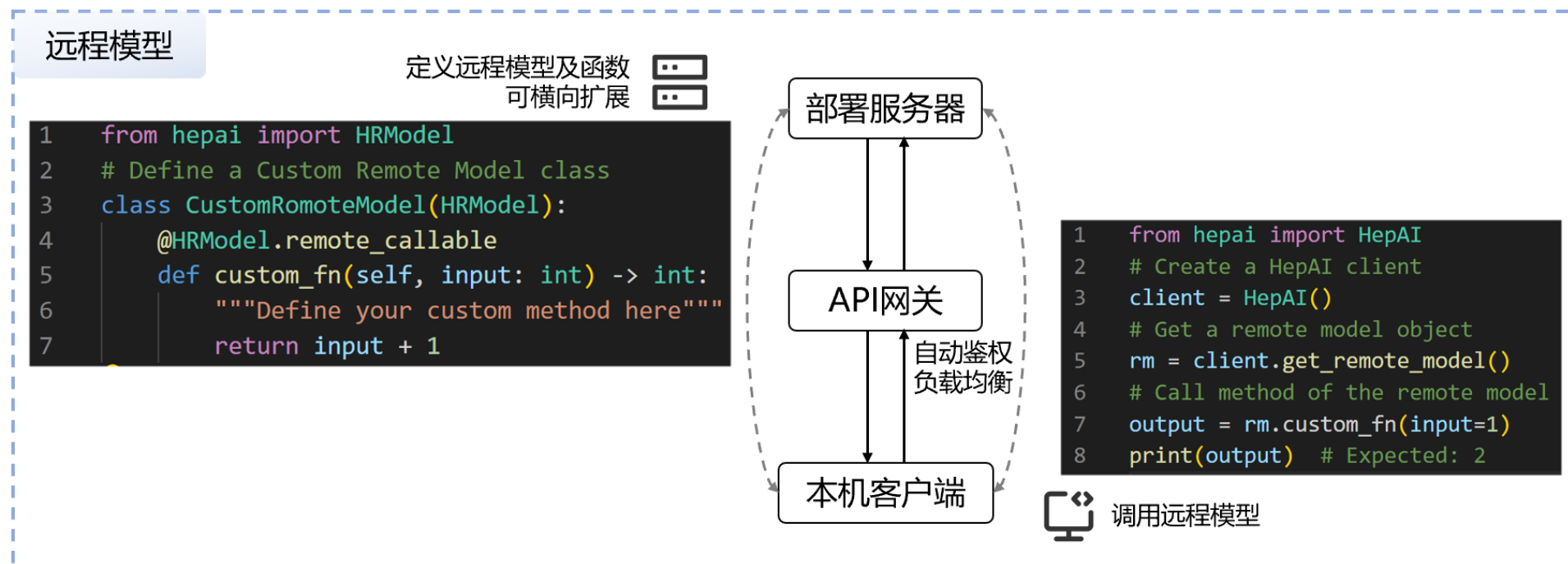| 异构资源 | GPU服务器 | DCU服务器 | CPU服务器 |

# HaiDDF2

**Admin UI**



- Admin
- App Admin
- Database
- Authentication
- Remote Models

# HaiDDF2

## Python Remote Model !

By deploying models and any other software programs onto cloud servers via DDF Workers, and pairing them with the HepAI client, it enables low-latency, distributed invocation of remote models and more.

pip install hepai -U

# HaiChatv2 Upgrade to v3

● Generative AI system improve the daily work efficiency.

Old mode: When encountering a problem, researchers search for information or ask a mentor.
New mode: Researchers consult the LLM first.

- Translation, proofreading, and providing solutions, etc.
- Guiding algorithm and software development

Available LLMs:
- HEP·Xiwu
- Qwen
- Baichuan
- LLaMA
- ChatGPT
- GPT-4

HaiChatv2

HaiChatv3

**New Features:**
- System Prompt Configuration
- Text-to-Speech Functionality
- Like Button Feature
- Multi-Model Output Comparison
- Message Editing
- Continue Generation
- Automatic Generation of Titles with Emoji
- Chat History Download
- Chat Title Search

☐ Total users: 3022 → 4850
☐ Daily Active Users: 800 (de-duplicated)

# Preliminary Dr.Sai v1.0 was online

Dr.Sai v1.0.0 was available on July 1, 2024    https://drsai.ihep.ac.cn

(Based on Chainlit)

- Four parts：
- BESIII AI Assistant
  - Equipped with BESIII internal knowledge
  - Automatic physical analysis by task decomposition, code/text generation, calling BOSS to execute code
- Personal Assistant
  - Allow individuals to upload knowledge
- Image Generation
  - Drawing based on AI
- Chatbot
  - Pure LLM: Xiwu, LLaMA3, GPT etc
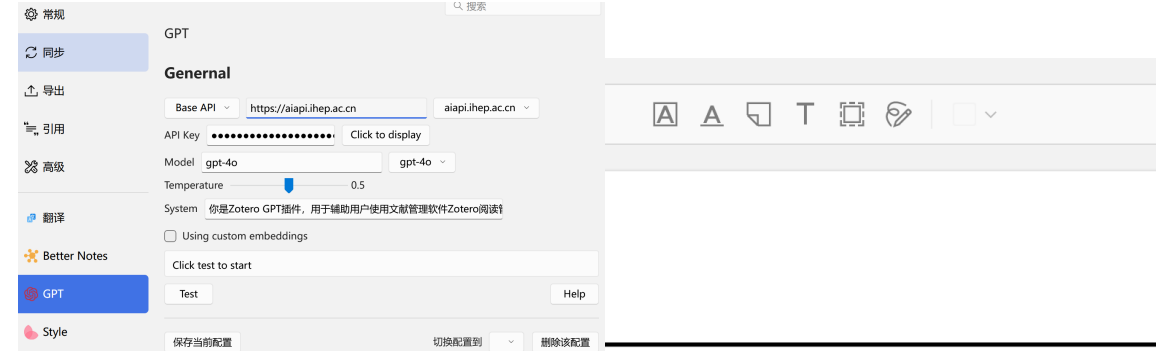  - Supports PDF and image input



WebUI

# Intergrated Applications

## HaiAcademic



## Support Zotero GPT Plugin



## 60+Deep Learning Paper Codes, Implementations



XIWU: A BASIS FLEXIBLE AND LEARNABLE LLM F
ENERGY PHYSICS

Zhengde Zhang[1], Yiyu Zhang[1], Haodong Yao[1], Jianwen Luo[2], Rui Zhao[1], Bo Huang[1], Jia
Yipu Liao[1], Ke Li[1], Lina Zhao[1], Jun Cao[1], Fazhi Qi[1,*], and Changzheng Yuan

这篇论文讲了什么

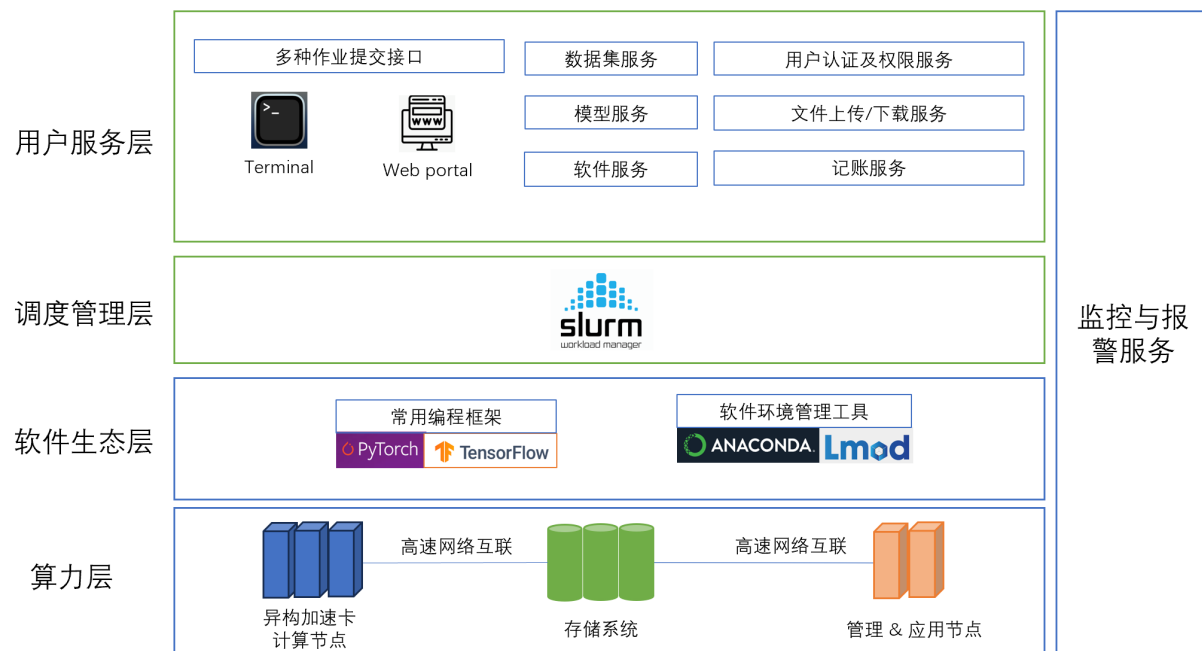🪐 AskPDF    🌟 Translate    ✨ Improve writing    Clipboard    Annotations

Large Language Models (LLMs) are undergoing a period of rapid updates and changes, v
of-the-art (SOTA) model frequently being replaced. When applying LLMs to a specific

# AI Computing Power

New AI Computing Power Infrastucture

☐ Initial investment of 9 million RMB (2024.03)

☐ Raise funds from other channels



**Accelerator Cards**
- 2 nodes of NVIDIA GPU
  - 8 * A800 (80GB)
  - 8 * L40 (48GB)
- 4 nodes of Sugon DCU (Domestic)
  - 32 * K100AI (64GB)
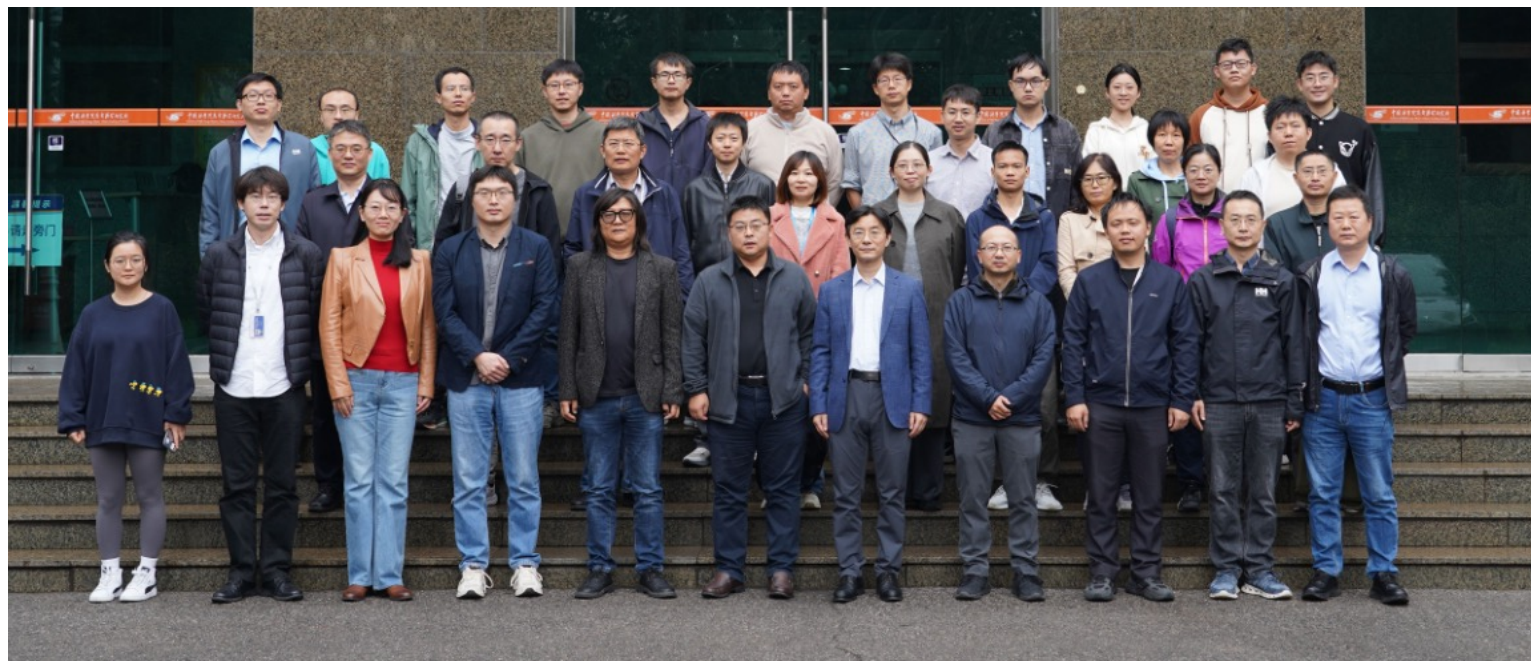
**Storage**
- 200TB flash memory

# ML Collaboration Group

□ Establishment of the Machine Learning Collaboration Group of IHEP

□ From October 16 to 18, 2024, the IHEP, CAS, held the Machine Learning Symposium and the inauguration ceremony of the Machine Learning Collaboration Group. Over 100 researchers from various departments within the institute participated in the meeting either online or in person. They engaged in in-depth discussions and exchanges on research related to machine learning and quantum computing from diverse disciplinary perspectives. During the conference, the first charter of the Collaboration Group was also discussed and approved.



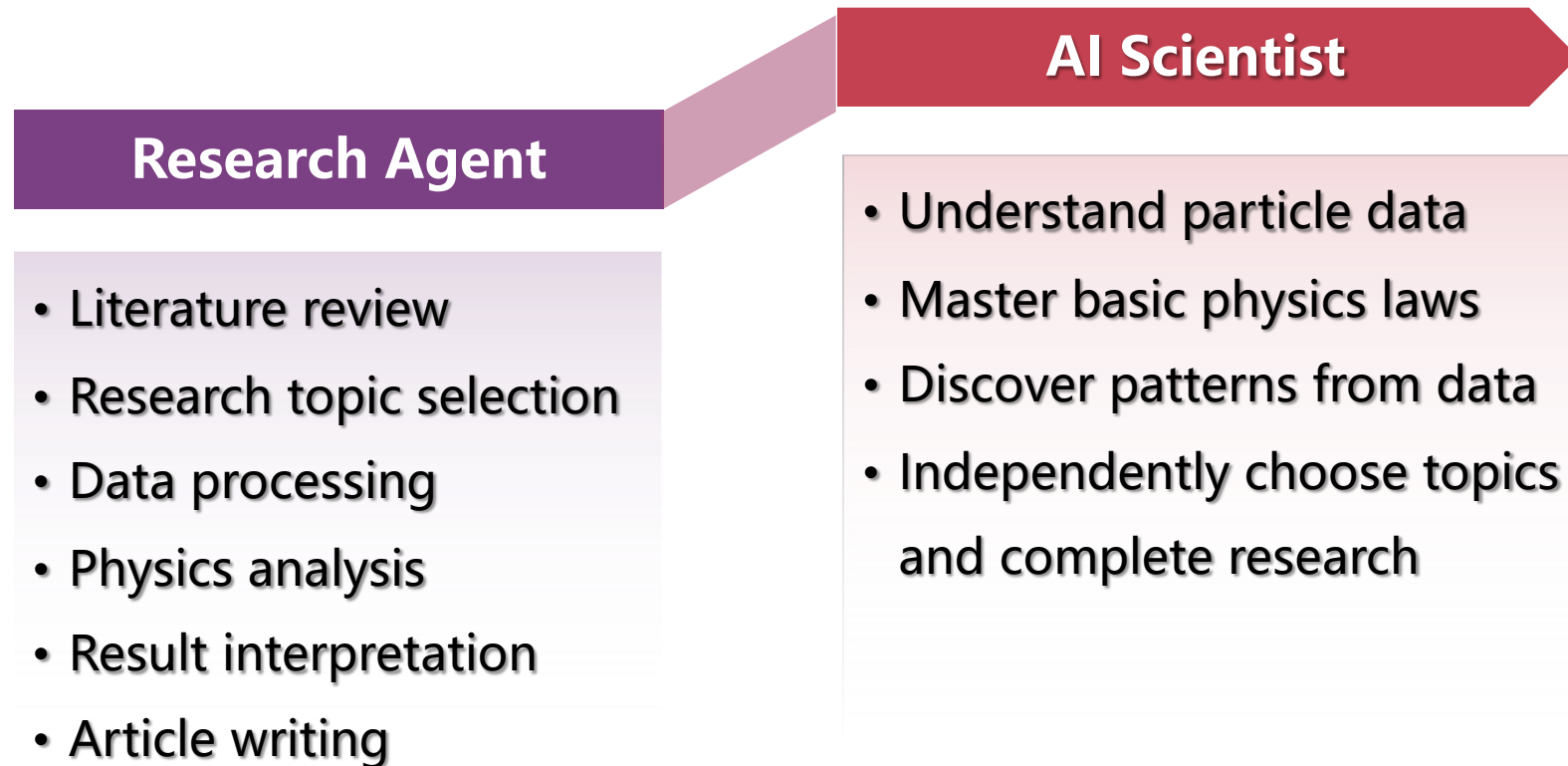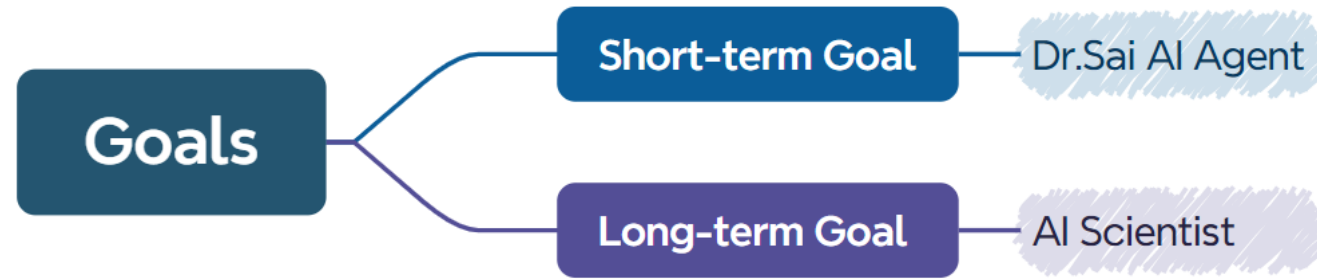高能所机器学习合作组成立

文章来源：计算中心 2024-10-23                    【大】【中】【小】

　　10月16至18日，中国科学院高能物理研究所机器学习研讨会（以下简称研讨会）暨高能所机器学习合作组（以下简称合作组）启动仪式召开，来自所内各个研究部门的共一百余位研究人员通过线上或线下的方式参加了会议，从不同学科的角度对机器学习和量子计算相关的研究展开了深入探讨和交流。会议期间讨论并通过了首个合作组章程。

# Part. 02

**Exploration of AI large model**

# Dr.Sai: An AI agents system for BESIII experiment

**Goals**
- Short-term Goal — Dr.Sai AI Agent
- Long-term Goal — AI Scientist

**Research Agent**

- Literature review
- Research topic selection
- Data processing
- Physics analysis
- Result interpretation
- Article writing

**AI Scientist**

- Understand particle data
- Master basic physics laws
- Discover patterns from data
- Independently choose topics and complete research

# What is AI Agent?

An AI agent refers to a **system** or **software** that can make autonomous decisions or perform actions on behalf of its users based on its knowledge, programming, environment, and inputs.
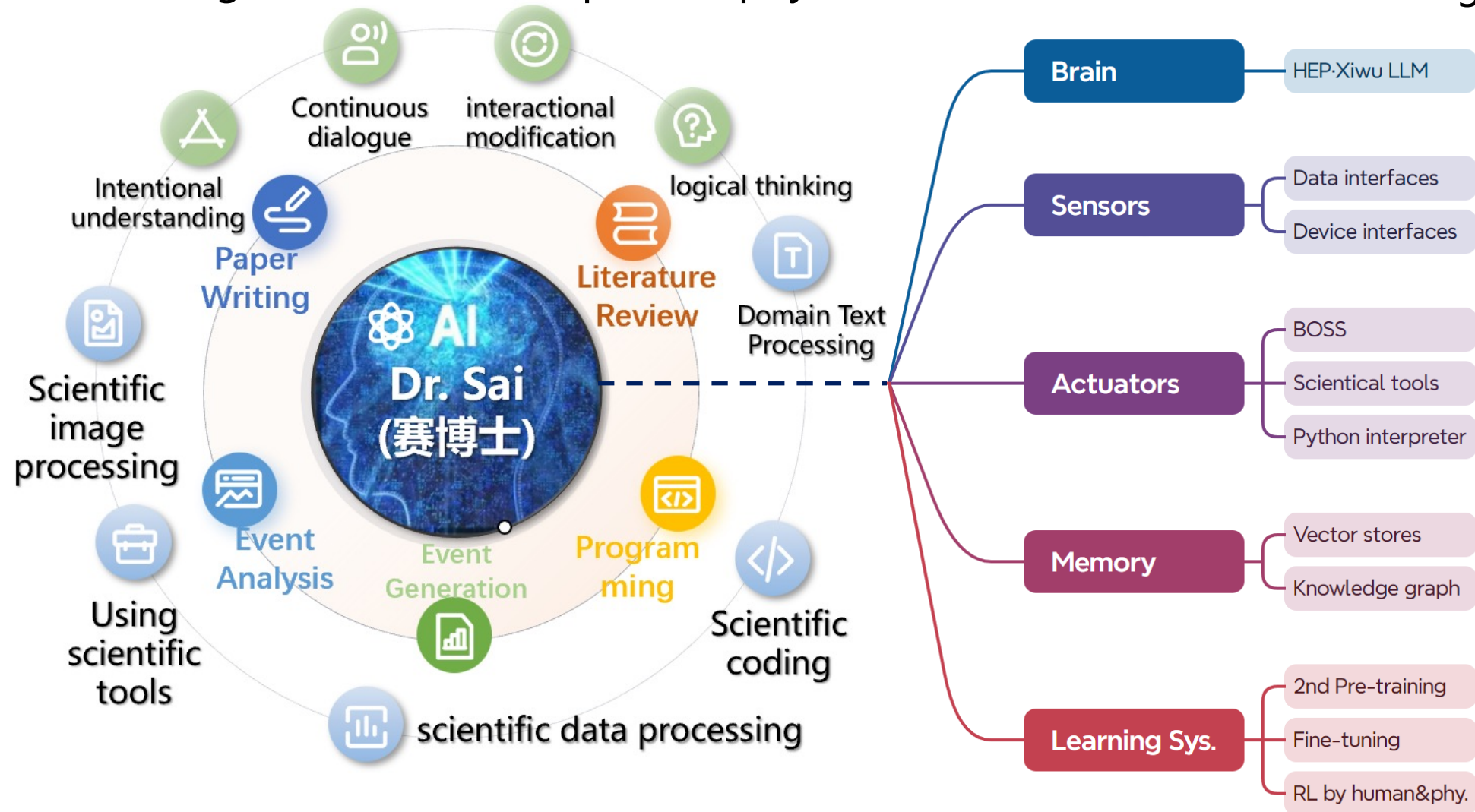


A LLM ≠ A person

**A agent ≈ A person**
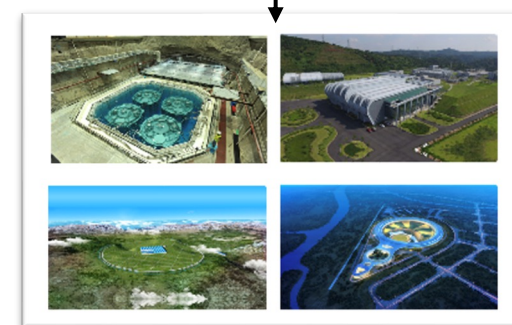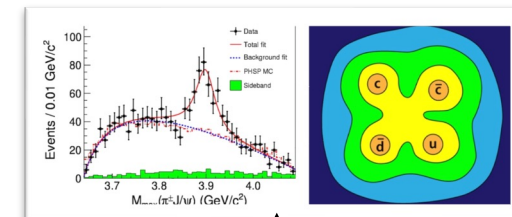
The LLM is the **processing core** (the Brain) of agent.

# Overview of Dr.Sai Agent

Let the large model conduct particle physics research

Essential: Modeling the research process.
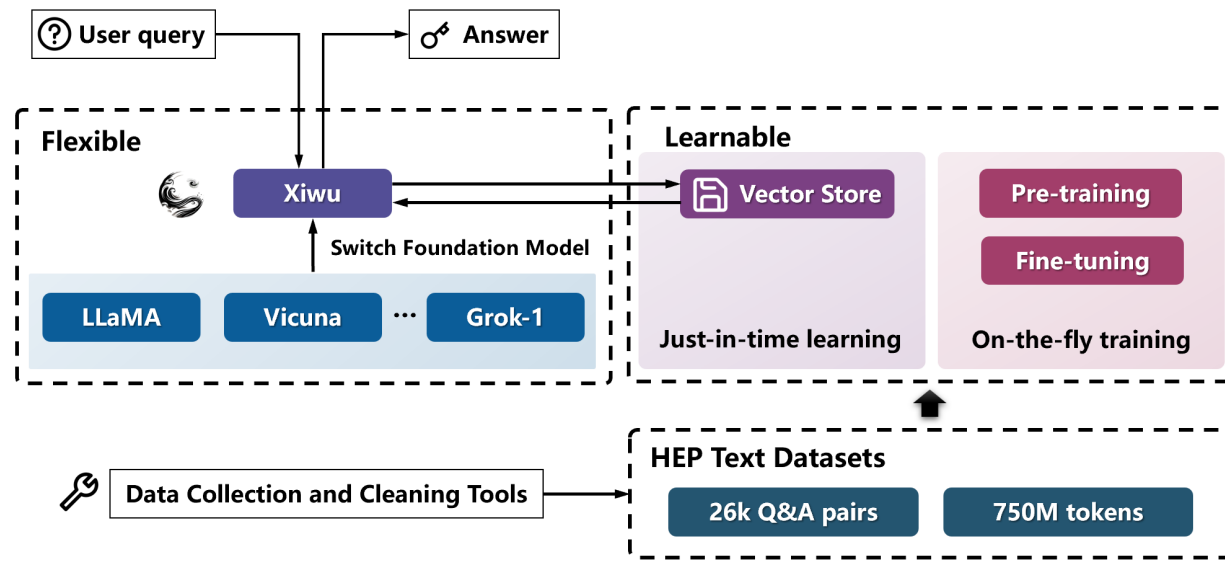


**Tasks & Required capabilities**
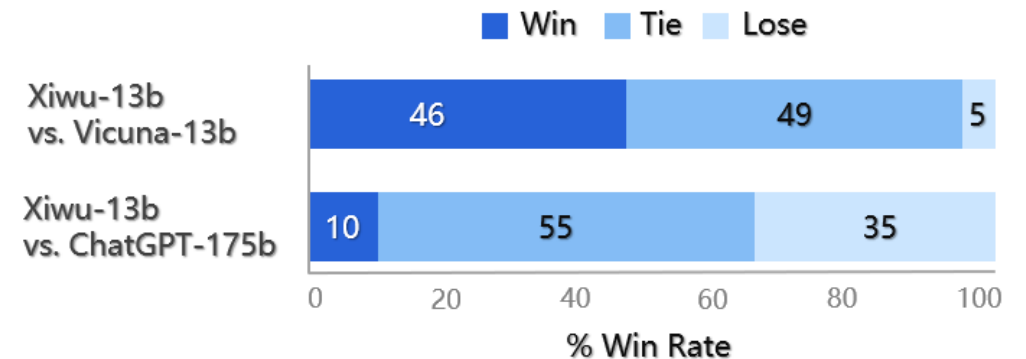
**Components**

**Applications**

## Xiwu(溪悟): A Basis Flexible and Learnable LLM for High Energy Physics [1]

- Currently, Xiwu is based on LLaMA3-8B. Historically, LLaMA, LLaMA2 (7B, 13B) etc.
- Secondary pre-training and fine-tuning.
- Significantly better than the base model in HEP Q&A and internal code generation.
- New version based on LLaMA3-70B and Qwen2-72B is in training.
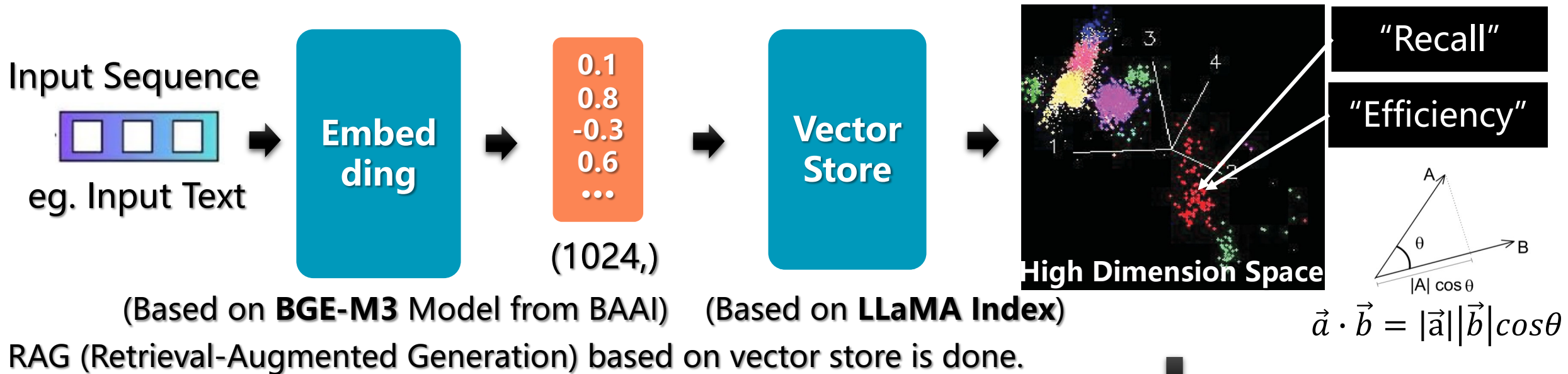


Test Results

[1] arXiv:2404.08001

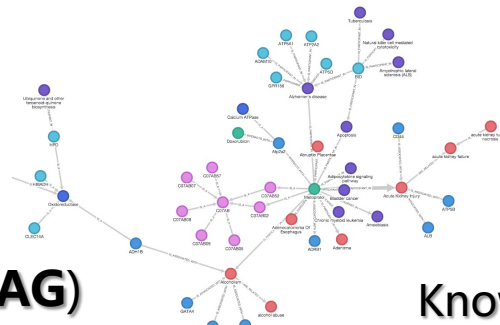https://github.com/zhangzhengde0225/Xiwu

- Save HEP knowledge into **Vector Store** and **Knowledge Graph**.
- Retrieve them to enhance the model's generation accuracy and reduce hallucination.
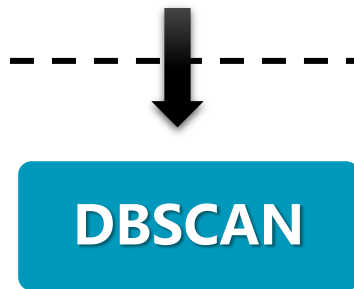
Input Sequence

eg. Input Text

**Embed ding**

0.1
0.8
-0.3
0.6
...

(1024,)

**Vector Store**

"Recall"

"Efficiency"

**High Dimension Space**

$$\vec{a} \cdot \vec{b} = |\vec{a}||\vec{b}|cos\theta$$

(Based on **BGE-M3** Model from BAAI)     (Based on **LLaMA Index**)

RAG (Retrieval-Augmented Generation) based on vector store is done.

- self-reflection could be realized.

**DBSCAN**

- Unsupervised clustering

Knowledge Graph is studing (Based **GraphRAG**)
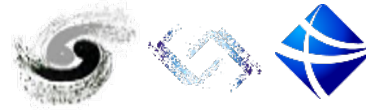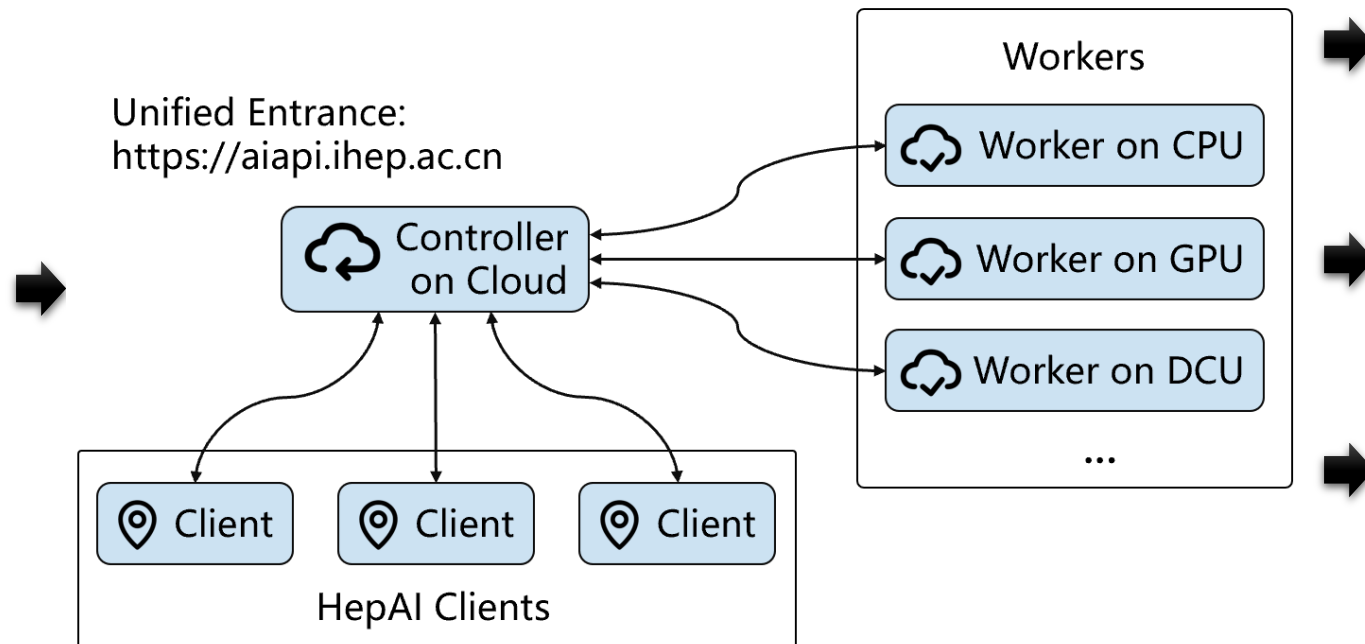
Knowledge Graph

# The **Actuators** of Dr.Sai – HepAI DDF

- The Distributed Deployment Framework (HepAI-DDF) is developed.
  - Featured with flexibility, cross-language, cross-platform, hererogeneous.
  - Allows Dr. Sai to easily scale its actuator components.
- BOSS (BESIII Offline Software System) actuator enables BESIII code execution and result retrieval.
- No longer a Q&A assistant; It can perform operations.



AI
**Dr. Sai**
**(赛博士)**

Unified Entrance:
https://aiapi.ihep.ac.cn

Controller on Cloud

Client Client Client
HepAI Clients

Workers

Worker on CPU
Worker on GPU
Worker on DCU
...

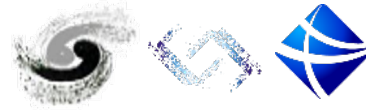**BOSS 7.1.0**
For executing physical analysis code.

**Daisy**
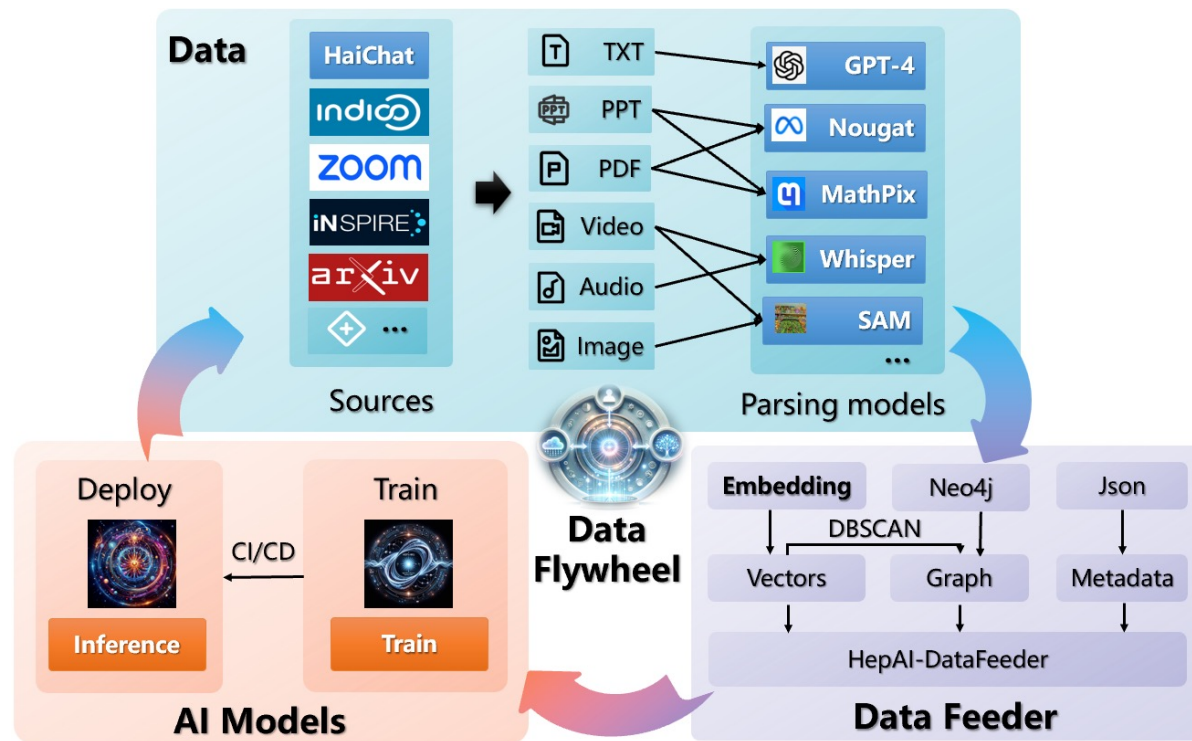For invoking scientific tools

**Python Interp.**
For interpreting general Python code.

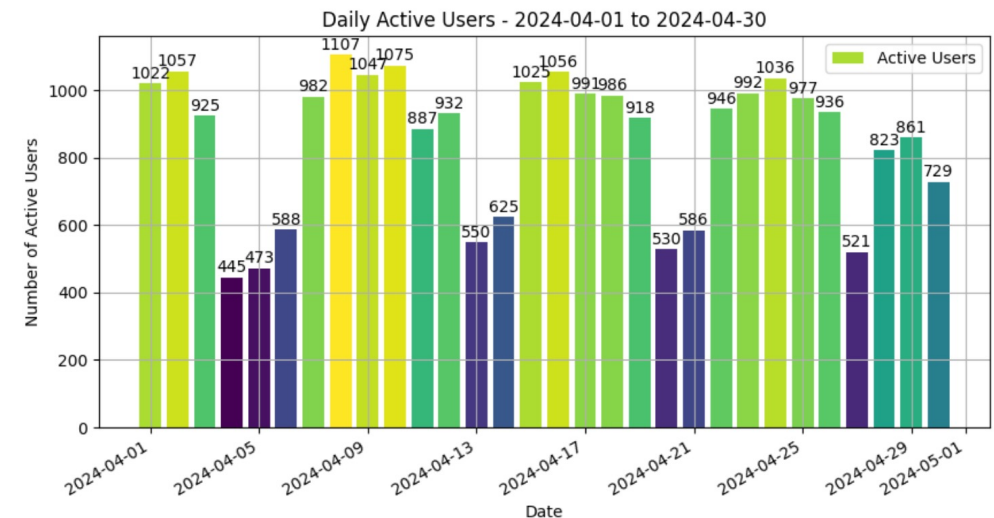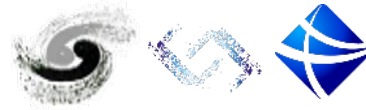# The Sensors of Dr.Sai – Data Flywheel

- The "data flywheel" enables continuous model iteration and evolution.
  - Data flywheel, i.e. **data-driven flywheel effect**
  - Improves models by constructing circular data pathways
  - Attracts more users, generating more data, further enhance the models.

- **HaiChat** service
  - AIGC service based on LLM
  - 4000+ users
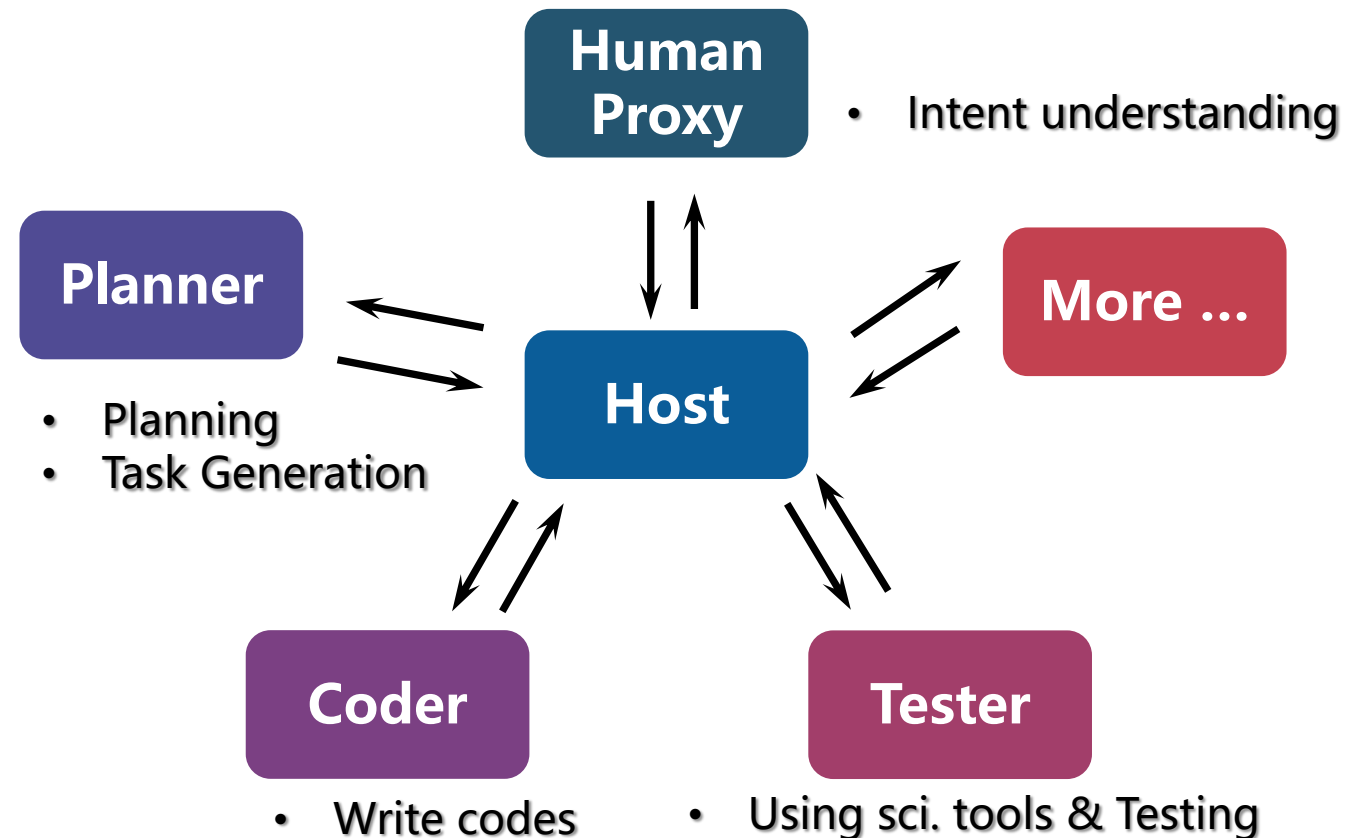  - Daily active users exceed 900.
  - Real needs from HEP users

# Multi-Agent Collaborative System

Dr. Sai's multi-agent collaborative system - handling complex tasks
- Based on AutoGen framework.                                    Autogen: arXiv.2308.08155)
- 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 intrdroduced to allow humans intervene at any time.

**Human Proxy**
- Intent understanding

**Planner**
- Planning
- Task Generation

**Host**

**More ...**

**Coder**
- Write codes

**Tester**
- Using sci. tools & Testing

Workerflow:
Msg → host → **LLM**
Planner → **RAG** → **LLM** → host
Host → **LLM**

# User Interface of Dr.Sai

Dr.Sai v1.0.0 was available on July 1, 2024     https://drsai.ihep.ac.cn

(Based on Chainlit)

- Four parts：
- BESIII AI Assistant
  - Equipped with BESIII internal knowledge
  - Automatic physical analysis by task decomposition, code/text generation, calling BOSS to execute code
- Personal Assistant
  - Allow individuals to upload knowledge
- Image Generation
  - Drawing based on AI
- Chatbot
  - Pure LLM: Xiwu, LLaMA3, GPT etc
  - Supports PDF and image input



WebUI

# Example of Single Agent

**Planner**

- Decompose complex task into multiple sub-tasks



Example

## Coder

- Generate internal BESIII codes

Example



generate a ROOT code to draw a Gaussian function

(auto select "Coder")

To generate a ROOT code to draw a Gaussian function, you can use the following script. This script will create a canvas, define a Gaussian function, and then draw it.

```cpp
ROOT

// filename: draw_gaussian.C

{
    // Create a canvas
    TCanvas *c = new TCanvas("c", "Gaussian Function", 800, 600);

    // Define a Gaussian function
    TF1 *gaus = new TF1("gaus", "gaus", -10, 10);
    gaus->SetParameters(1, 0, 1); // Set the parameters: amplitude, mean, and sigma

    // Draw the Gaussian function
    gaus->Draw();

    // Save the canvas as a PNG file
    c->SaveAs("gaussian_function.png");
}
```

copy the code block

Note: Demonstrated the ability to generate ROOT code, using BOSS code internally.

This script will create a Gaussian function with an amplitude of 1, a mean of 0, and a standard deviation (sigma) of 1. The function will be drawn on a canvas and saved as a PNG file named "gaussian_function.png".

# Example of Single Agent

## Tester

- Excute analysis code in BESIII environment
- Run on remote HPC server
- Retrieve the results

Example

code type

backticks

file name

command

Execute this code:
```sh
# filename: test.sh

du -sh ~
```

(auto select "Tester")

There are 1 code blocks in the message.
The output are:
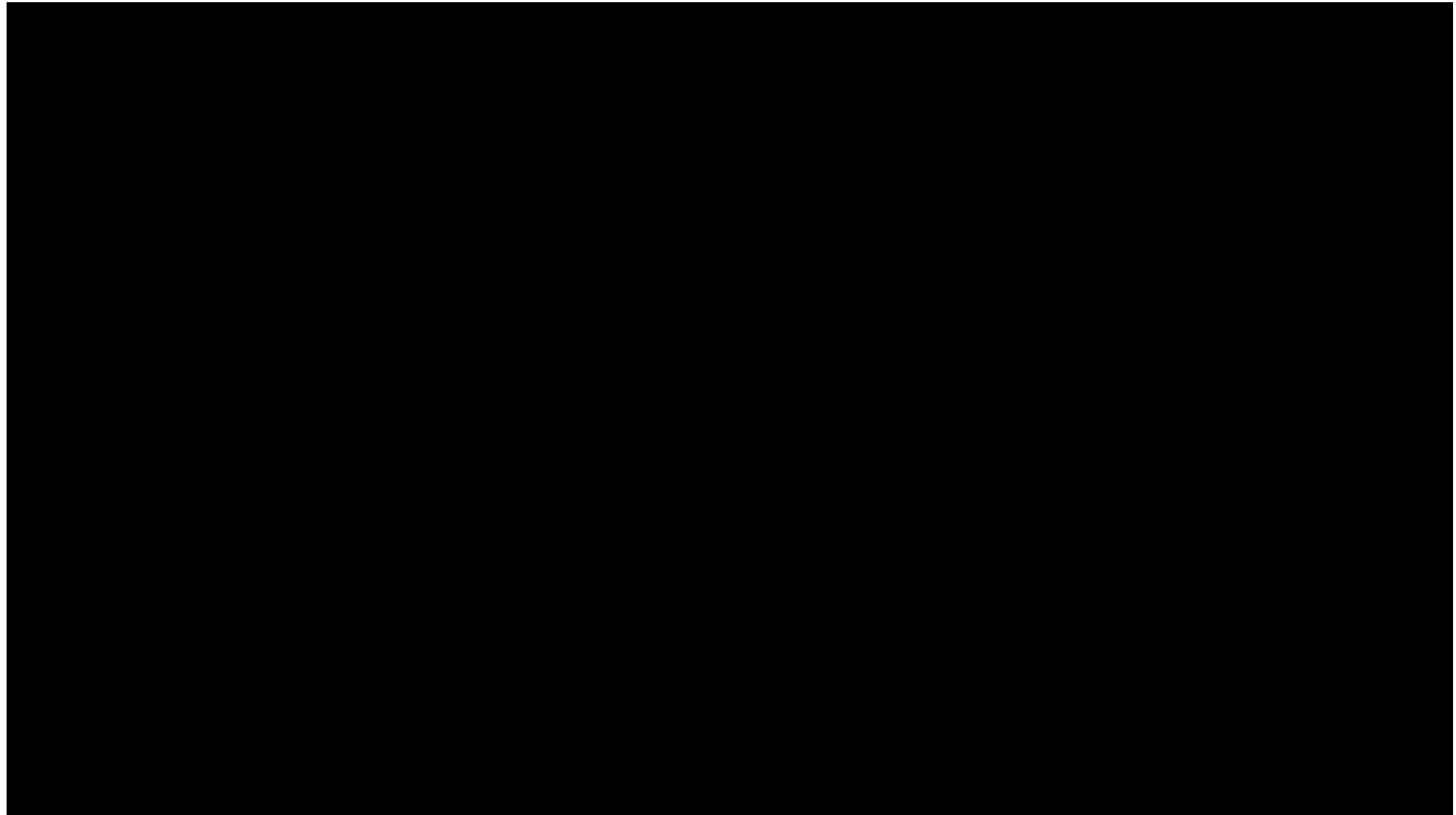
exitcode: 1
output: 140M  /afs/ihep.ac.cn/users/z/zhangbolun
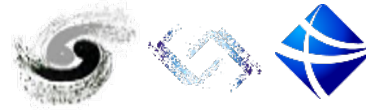error:

The generated files are:
<pic: None >
<pdf: None >

# Example of Multi-Agent System

- **Host** agent select suitable speaker
- **Coder** generate domain code
- **Tester** Call BOSS (BESIII Offiline Software System) to excute
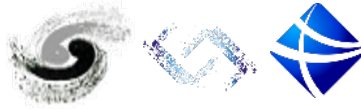- Draw a signal histogram

Example

# Summary

- **HepAI Platform**

  - The core component **HaiDDF**v1 → v2. General Remote Model!

  - Generative AI system **HaiChat**v2 → v3.

  - New AI Apps. (**HaiAcademic**, **Zotero Plugin, DL Codes tutorials**) have been intergrated

  - New AI Computing Power Infrastructure (**48 GPU&DCU cards**)

  - IHEP **ML Collaboration Group** has been officially established.

- **AI Agents based on LLMs**

  - The **HEP•Xiwu** LLM has been trained to enhance its domain-specific capabilities.

  - Initial version of  **Dr.Sai Agents** for physics analysis based on LLM was launched.

  - The feasibility study has been successfully concluded, and we have now entered the formal research phase.

# Contact

- **We are seeking more collaborations.**
- **Recruiting international students, postdoctoral fellows, staff, and talents through recruitment programs.**

- **Welcome to discuss any interests related to AI for High Energy Physics!**

- Fazhi Qi (齐法制)
- qfz@ihep.ac.cn
- Computing Center, Institute of High Energy Physics, CAS, Beijing
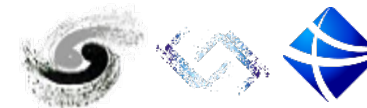
- Zhengde Zhang (张正德)
- zdzhang@ihep.ac.cn

# Backup

# The limitations of general LLMs

| Question1: What is HEPS? | |
|---|---|
| **[GPT-4 answers incorrectly]:** HEPS stands for Harmonised European Prevalence of Infection Study. It is ... | **[Xiwu answers correctly]:** HEPS stands for High Energy Photon Source. It is ... |

| Question2: Please provide an example of simulate particle decay on BOSS. | |
|---|---|
| **[GPT-4 answers nonsensically]:** | **[Xiwu answers sensibly]:** |

**[GPT-4 answers nonsensically]:**

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 answers sensibly]:**

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:

```
1   Decay  psi(2S)
2   1.0    J/psi   pi+ pi-   JPIPI;
3   Enddecay
4
5   Decay  J/psi
6   0.5    e+     e-          PHOTOS VLL;
7   0.5    mu+    mu-         PHOTOS VLL;
8   Enddecay
9
10  End
```

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

```
1   //*************job options for generator
        (KKMC)**************
2   #include "$KKMCROOT/share/jobOptions_KKMC
        .txt"
3   KKMC.CMSEnergy = 3.686;
4   KKMC.BeamEnergySpread=0.0009;
5   KKMC.NumberOfEventPrinted=1;
6   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!

- **Hallucination：**
Produce content that is nonsensical or untruthful in relation to certain sources[3, 4], which is particularly detrimental in scientific fields
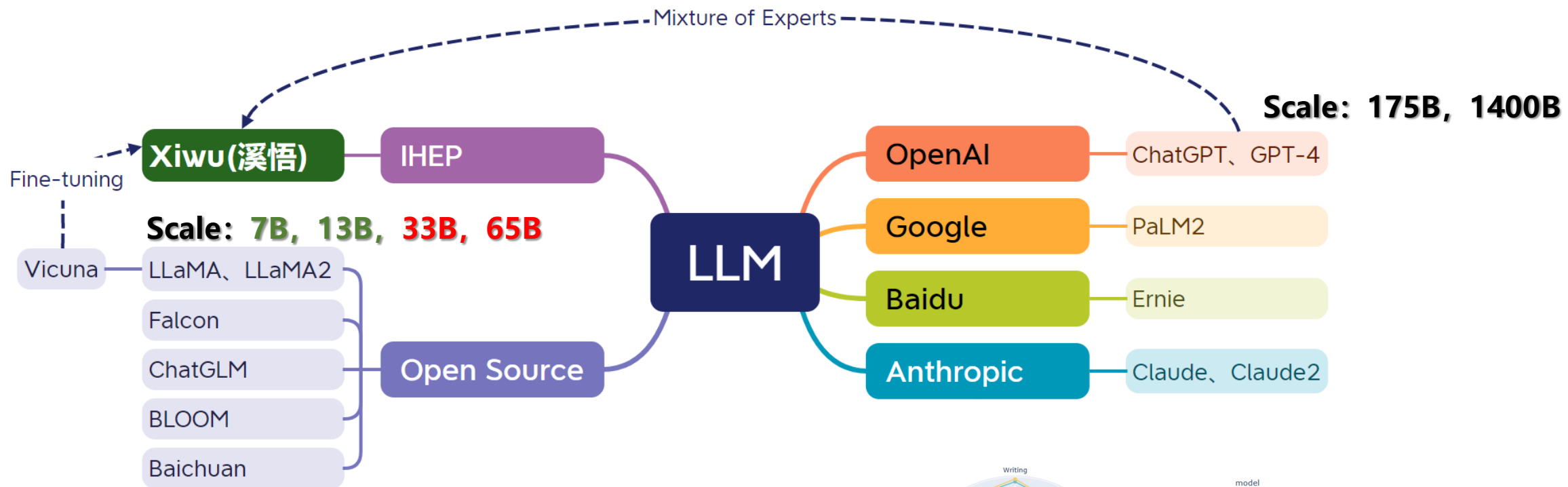
- **Lack of rapid learning ability**

- **Solutions**
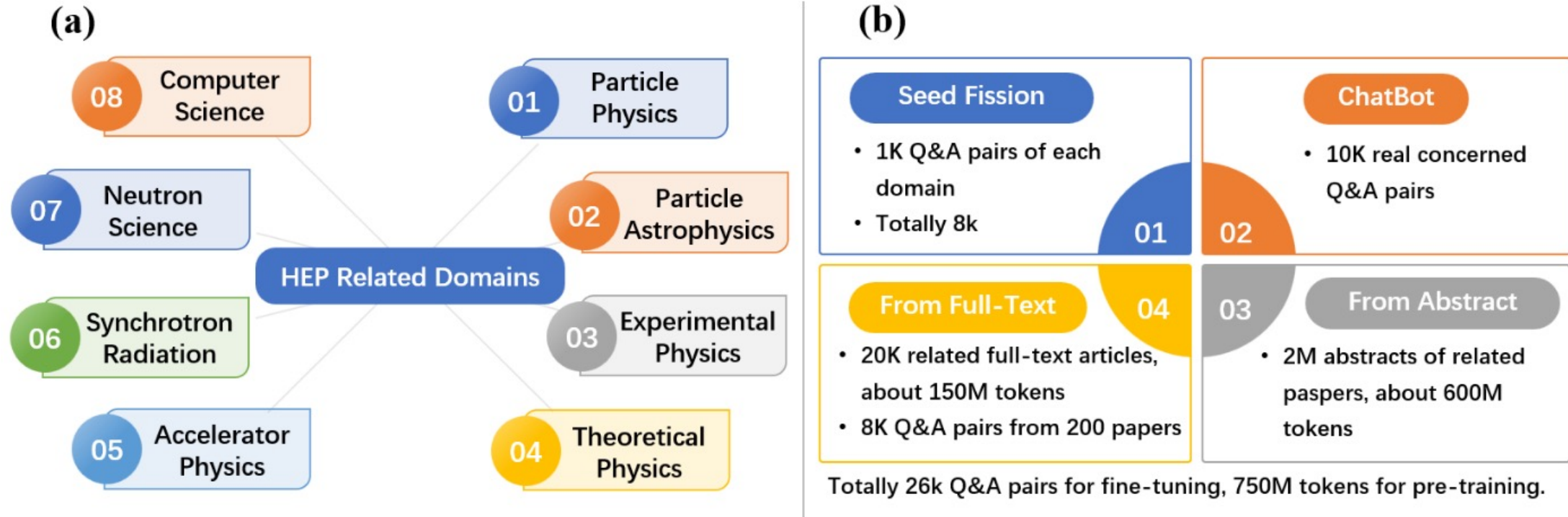- No.1 **Fine-tuning LLM**
- No.2 **External Knowledge**

## Xiwu: a customized LLM for High Energy Physics



**Scale：** 175B，1400B

**Scale：** **7B，13B，33B，65B**

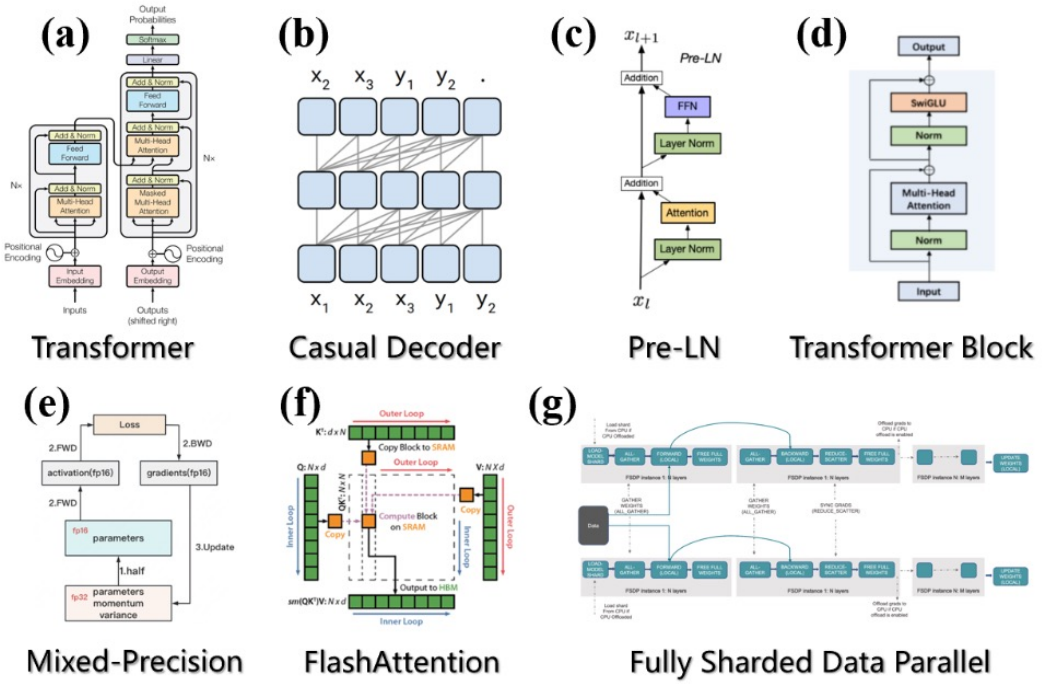Evaluation Matrix： Writing, Coding, Extraction, Math, Reasoning, Humanities and STEM.

(a) HEP Related Domains; (b) Four methods employed to collect and clean up data

Collected, cleaned and trained dataset on July 2023: **5K Q&A pairs**
Collected, cleaned and **NOT** trained dataset on Sept. 2023: **21K Q&A pairs + 750M tokens**

New source: open Indico-Conferences-2168, pdf, ppt & mp4, about 100GB on disk.

# HEP·Xiwu: Training techniques and Results



**The illustration algorithm components and training technologies**

- 8*A100(40G) GPU，Xiwu-7B & 13B √
- 40*V100(32G) GPU，Xiwu-33B，Failed

**Advanced training technologies for LLM：**

- Transformers with Self-Attension
- Casual Decoder
- Pre-LN
- Transformer Blocks with multi-head attesion
- Mixed-Precision for fast training
- FlashAttention improves the utilization rate of the GPU's FLOPs.
- Low-Rank Adaptation (LORA) significantly reduces the need for fine-tuning.

**High Energy Physics 100 Question Test Results：**



**The Xiwu significantly outperforms the baseline model!**

# Just-In-Time Learning System

Pure LLM cannot handle dynamically changing information and are difficult to train on unprecedented data. ➡ Vector store as a memory module

# Xiwu vs GPT-4

- Q: **How to use BOSS to process BES data?**

  - Ans. from Xiwu :

    Correct ➡

- Ans. from GPT-4 :

  ↓ Nonsensical



- Xiwu excels in HEP-specific question answering and generating scientific codes.
- Xiwu-MoE combines domain knowledge with general capabilities for a more comprehensive approach.
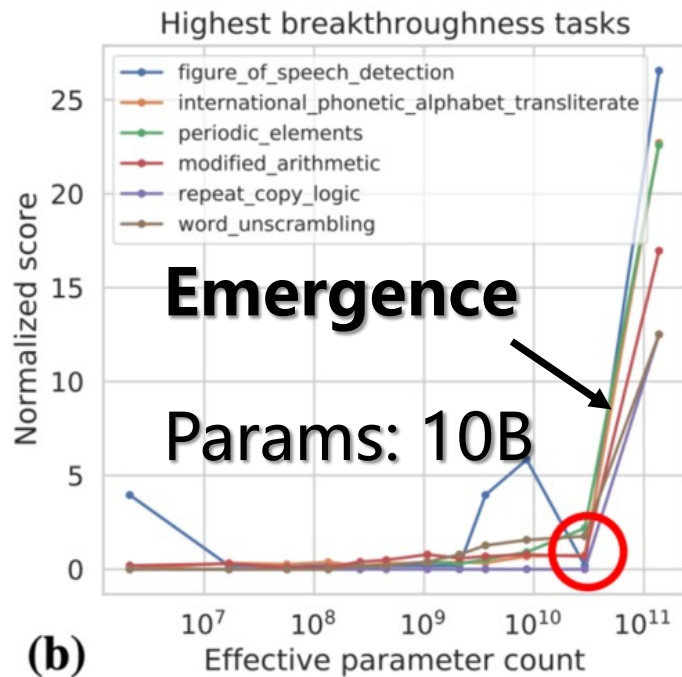
**Part. 03**

Exploration of AI large model

# The first pathway:
# Research Agent based on Xiwu

# Particle Physics Large Model

GPT-4 has surpassed dedicated AI models in all downstream tasks of natural language processing (NLP), indeed showing sparks of artificial general intelligence (AGI).



**Emergence**

**Params: 10B**

Why and how LLM achieves such remarkable intelligence?

- The large amount of data forces neural networks to learn generic and useful "**neural circuits**", while the large size of models provide enough **redundancy** and **diversity** for the neural circuits to specialize and fine-tune to specific tasks.
- The huge size of the model could have several other benefits, such as making gradient descent more effective by connecting different minima or by simply enabling smooth fitting of high-dimensional data.

**Idea:** **40PB HEP data** ➡ Leverage large-scale unsupervised pre-training methodologies to enable AI to grasp the **global rules** from all data, and employ Reinforcement Learning by Physical Feedback (**RLPF**) to facilitate its emergence.
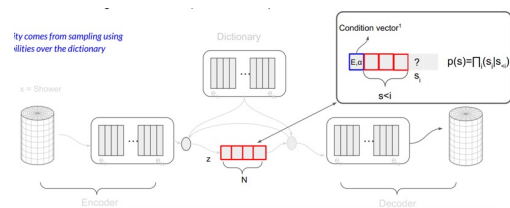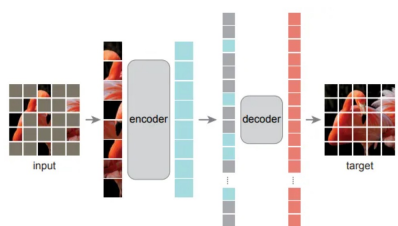
# Difficulty 1: How is scientific data pre-trained?

## Characteristics of HEP data:
- Floating-point numbers with physical meaning.
- Non-sequential.

## Pre-training Methods:
- Next word prediction ❌
- Mask Auto Encoding (MAE) ➖
- Vector Quantized Variational Auto-Encoder (VQ-VAE) ➖
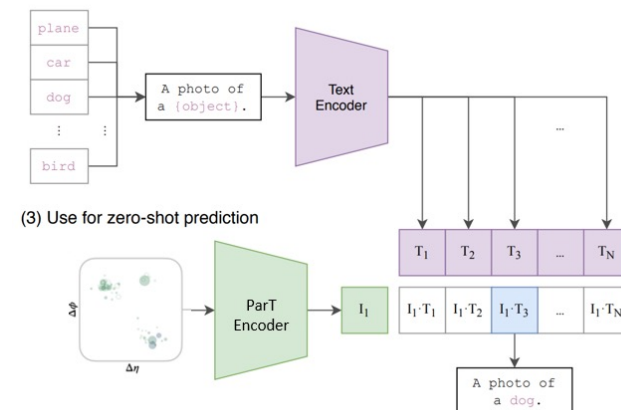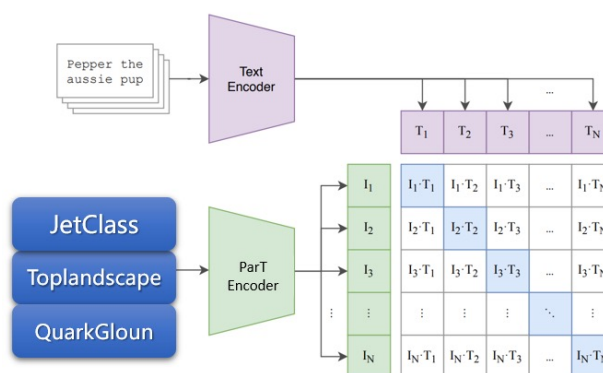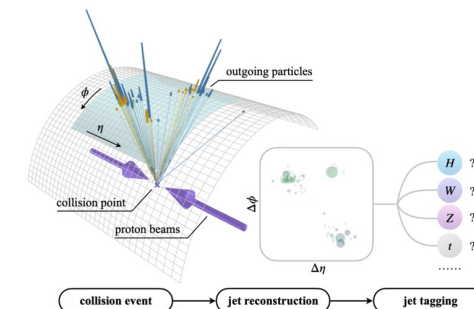- Graph + Transformer ➖
- Cross-modal supervision ➖ May be

**Seeking collaboration**

## Particle Physics Large Model
- Cross-modal self-supervised method i.e. **text as supervisory signal**
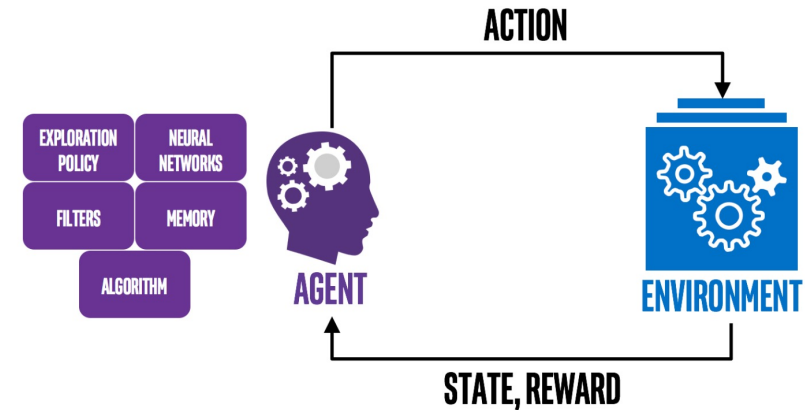- We are trying to develop a model that can handle multiple HEP tasks.

# Difficulty 2: How to align large models with physics?

## Possibility 1: Reinforcement learning with interactive environments.

- A reinforcement learning system that uses the BOSS program as the environment

- Enhance abilities such as scientific code programming in the context of reinforcement learning.



## Possibility 2: Reinforcement learning by Physical Feedback (RLPF)

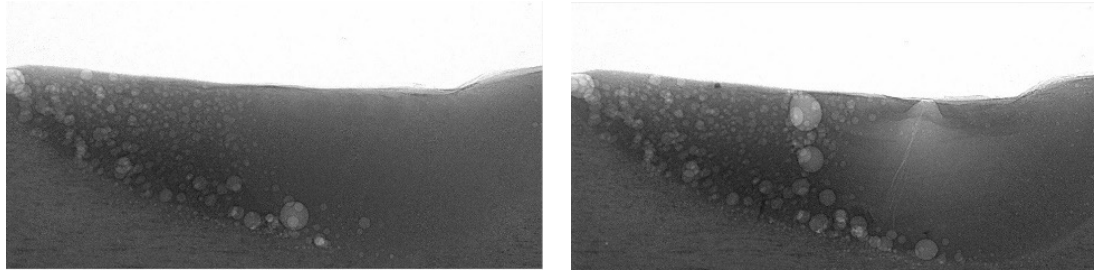- In the evolution from GPT-3 to ChatGPT, Reinforcement Learning from Human Feedback (RLHF) plays a crucial role in alignment. The performance of the 175 billion parameter model without RLHF is even worse than the 7 billion parameter model with RLHF.

- The establishment of **Reinforcement Learning from Physical Feedback (RLPF)** mechanism can achieve alignment between the model and physical principles. (Requires further thought)

# Example 1: Defects real-time feedback

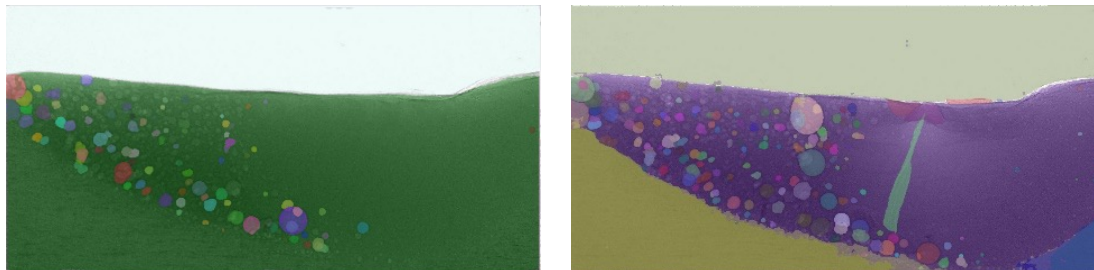## Defects real-time feedback of additive manufacturing

(Dr. B.B.Zhang, Z.D.Zhang)

- Complex
- Low resolution
- dynamic

- Defect Detection and Tracking

- Defects such as bubbles and melt pools directly affect material properties
- Laser intensity and powder feeding rate can control the distribution of bubbles and melt pools

**Old Mode:** The laser intensity and powder delivery rate are determined based on experience.

**New Mode**: Based on AI real-time feedback of defect information, the laser intensity and powder delivery rate are **dynamically controlled**.
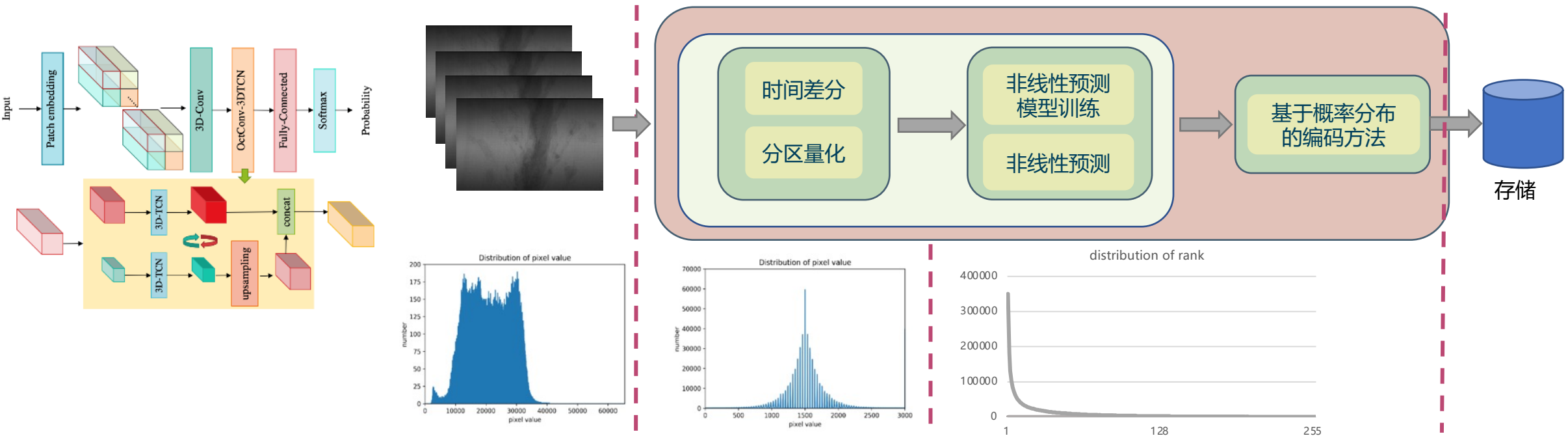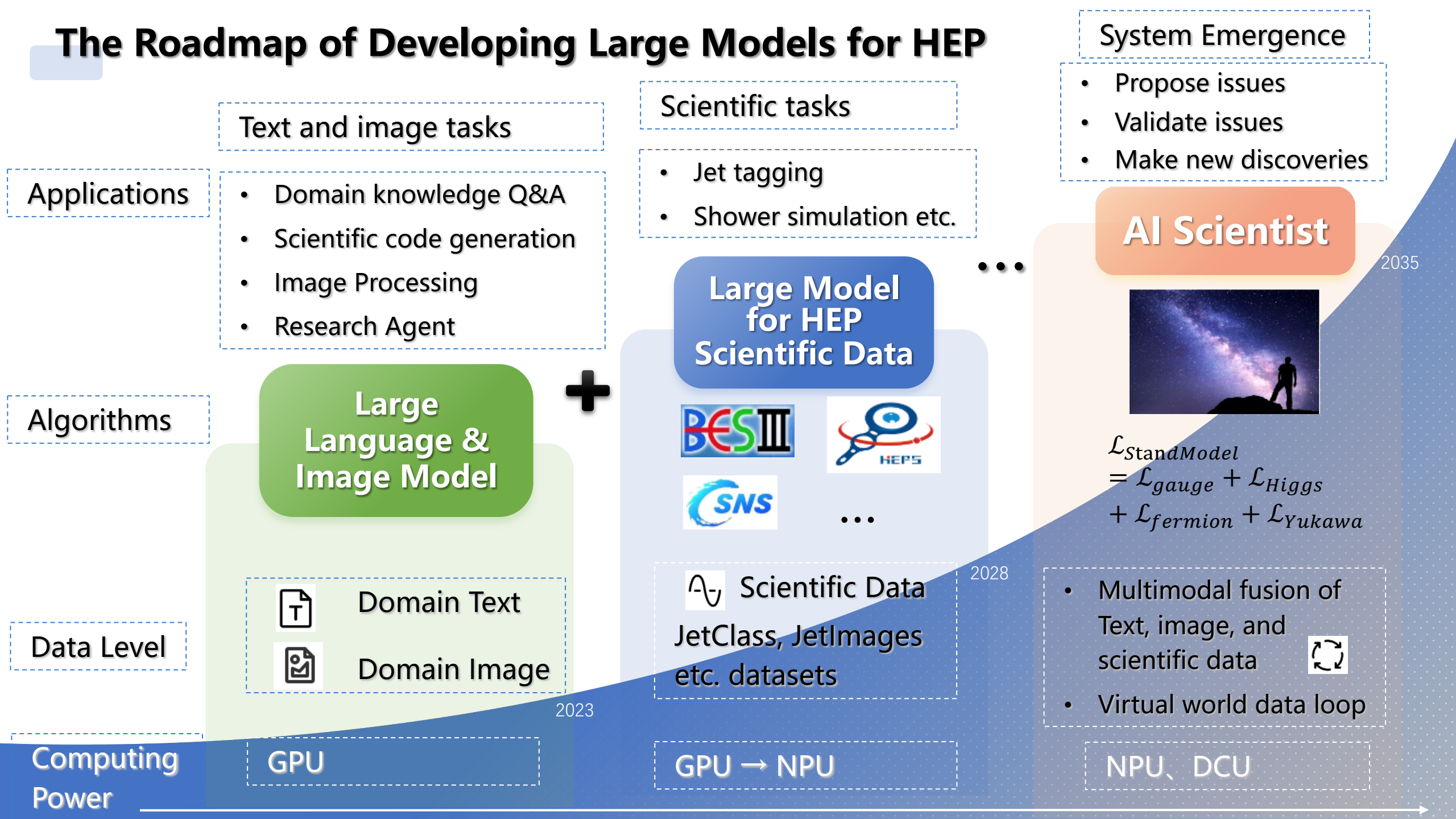
# Example 4: Image data compression

**lossless compression for light source images** *(Dr. S.Y.Fu)*

- Improving compression rate through convolutional neural network prediction probability and Huffman coding

- Saving **58%~66%** storage space

- Equivalent to reducing several hundred PB of data storage per year for HEPS

# The Roadmap of Developing Large Models for HEP

**System Emergence**
- Propose issues
- Validate issues
- Make new discoveries

**Applications**

Text and image tasks
- Domain knowledge Q&A
- Scientific code generation
- Image Processing
- Research Agent

Scientific tasks
- Jet tagging
- Shower simulation etc.

**AI Scientist**

2035

**Large Model for HEP Scientific Data**

**+**

**Algorithms**

**Large Language & Image Model**



$$\mathcal{L}_{StandModel} = \mathcal{L}_{gauge} + \mathcal{L}_{Higgs} + \mathcal{L}_{fermion} + \mathcal{L}_{Yukawa}$$

**Data Level**

Domain Text

Domain Image

Scientific Data
JetClass, JetImages etc. datasets

2028

- Multimodal fusion of Text, image, and scientific data
- Virtual world data loop

2023

**Computing Power**

GPU

GPU → NPU

NPU、DCU