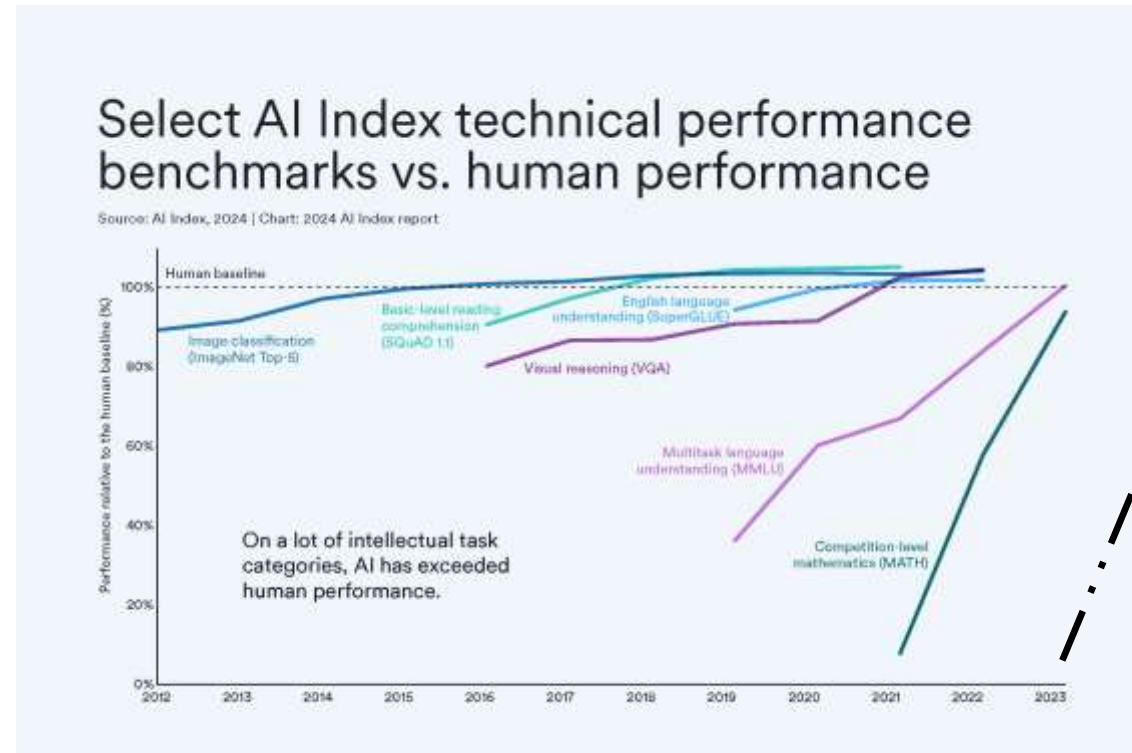# LLM-based physics analysis agent (Dr. Sai) at BESIII and exploration of future AI virtual scientist

Ke Li (like@ihep.ac.cn)

on behalf of Dr. Sai working group

# Outline

- Motivation
  - what is LLM
  - why we need LLM
- Dr. Sai project
  - architecture
  - training
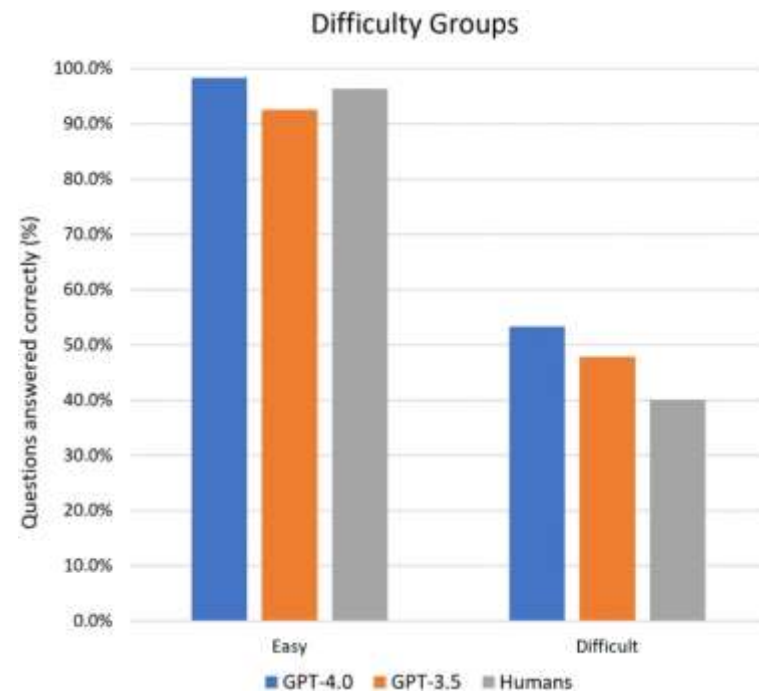- Status and Prospects



2024 AI index report

Analysis at High energy physics ?

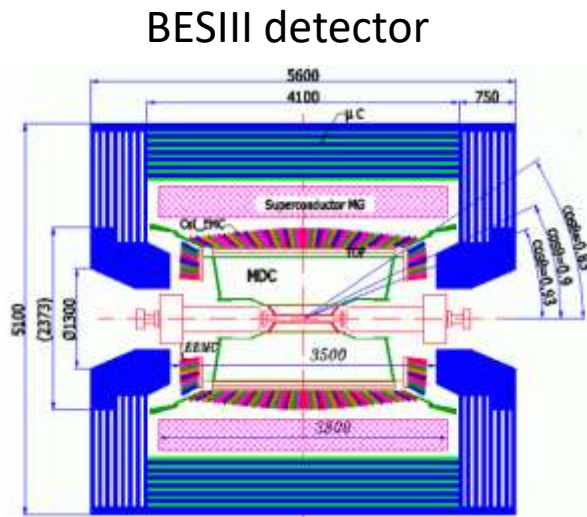# What is Large Language Model (LLM)

- Large language models (LLMs), normally build on transformer architecture, has demonstrated impressive performance in **text/code generation**

  - GPT4o, Gemini, LLaMa3...

  - Could be used for HEP studies

  - Game changer

- **For us, open-source foundation model + higher level model for HEP + fine tuning for BESIII**
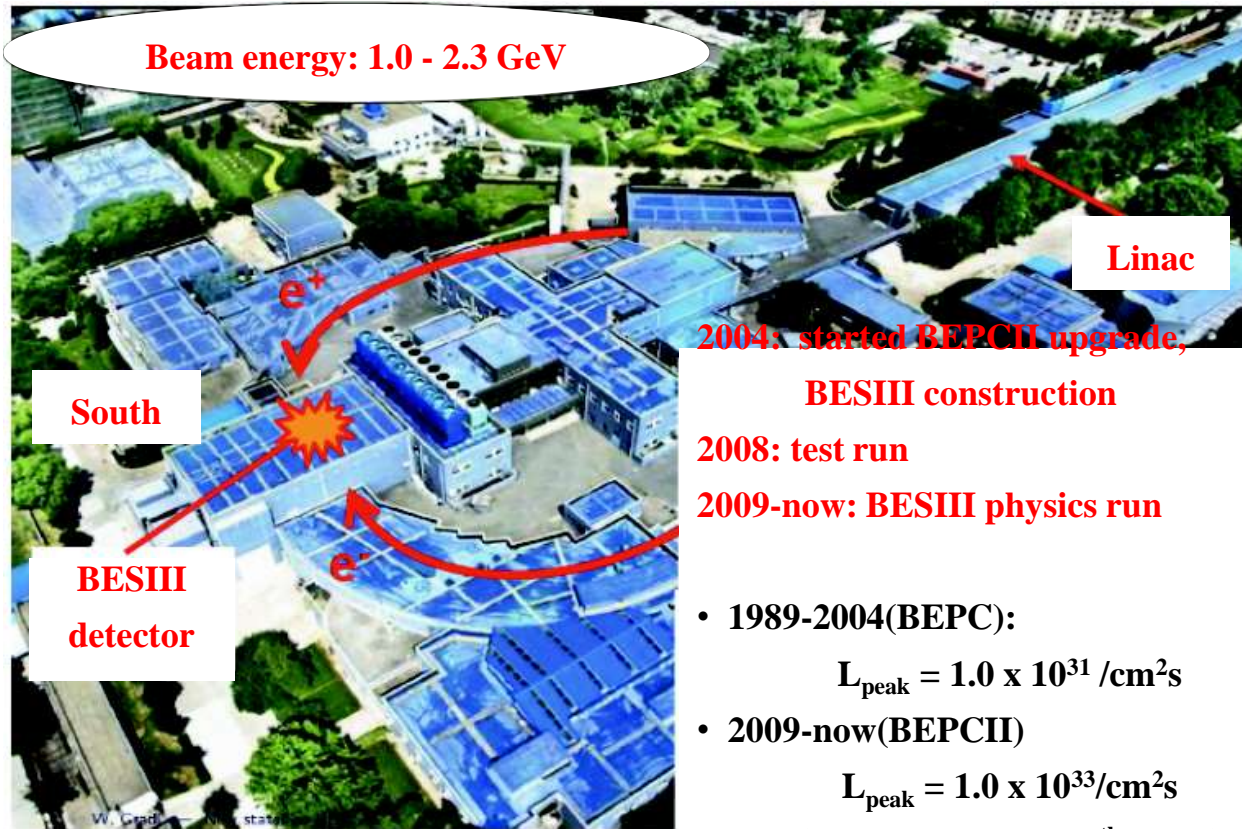


Scientific Reports volume13, Article number: 18562 (2023)

3

# BESIII at Beijing Electron-Positron Collider II

**A double-ring collider with high luminosity**

BESIII detector



**Beam energy: 1.0 - 2.3 GeV**

**Linac**

**South**

$e^+$

$e^-$

**BESIII detector**

**2004: started BEPCII upgrade, BESIII construction**

**2008: test run**

**2009-now: BESIII physics run**

- **1989-2004(BEPC):**

    $L_{peak}$ = 1.0 x $10^{31}$ /cm²s

- **2009-now(BEPCII)**

    $L_{peak}$ = 1.0 x $10^{33}$/cm²s

    **(Achieved on Apr. 5th, 2016)**

Beam energy: 1.0 -2.3GeV
Luminosity: $1 \times 10^{33}$ cm⁻²s⁻¹
Optimum energy: 1.89 GeV
Energy spread: $5.16 \times 10^{-4}$
No. of bunches: 93

**Zoom in IP**

IP

- Preparing for upgrade (2024-2025)
- More data will be collected

# Why we need LLM



- More data will be collected after BEPCII-upgrade
- \>500 physics results from ~500 people in the past 14 years
  - One result normally took **~3 years**
- We need a more efficient workflow in order to achieve the goals in BESIII white paper !
  - E.g. Test Lattice-QCD,  light/charm hadron spectroscopy, charm quark weak decay

# Why we need LLM

- Major effort in BESIII analysis is spent in writting/testing/updating code/text
  - **LLM is good at code/text generation !**
- Key problems for LLM at HEP
  - how to make sure the outputs are reliable?
  - how to avoid hallucinations ?
  - Current solutions:
    - **More accurate and good quality data for training**
    - **More tests and validations**
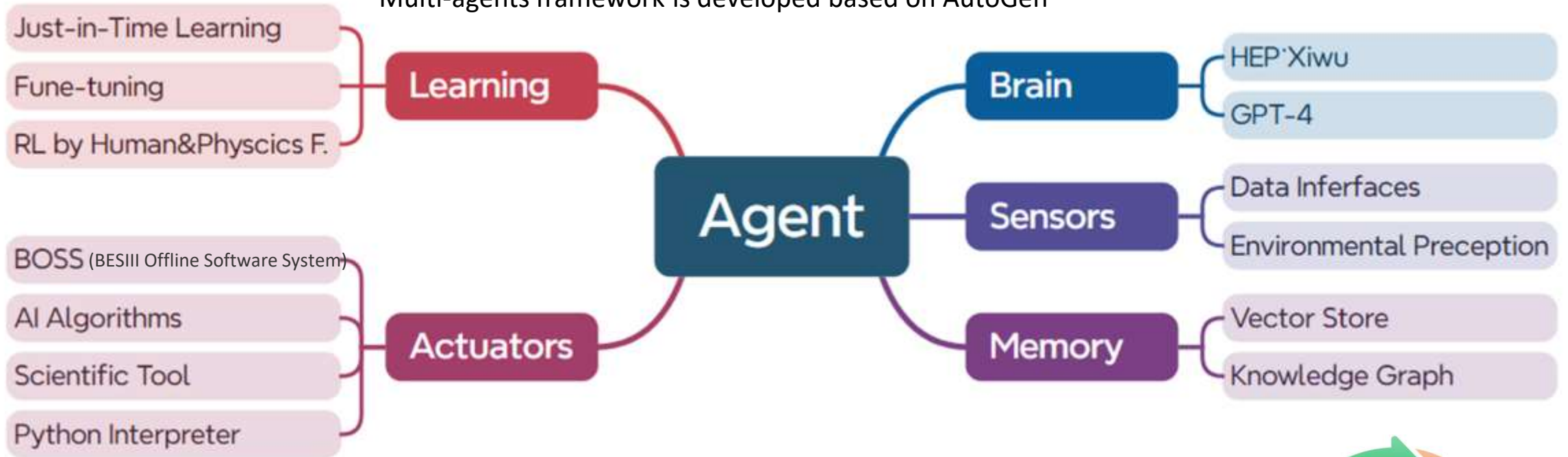    - **More proper architecture**

# Dr. Sai project for BESIII/HEP

- AI Agent: AI tools capable of autonomously performing complex tasks
    - LLM = brain  -> AI agent = human
- AI agent based on **Xiwu** model (LLM for HEP)
    - based on Llama 2/3, will train with BESIII internal data, e.g. memo, source code, Q-As from internal review
- One milestone: **AI assistant**, It can help scientist on data analysis, e.g. MC generation,  signal extraction, and a navigator inside BESIII
    - target at **June 2024** !
- Goal: **AI scientist**, it can analyze the data automatically like an expert

~20 people from IHEP, UCAS, LZU and JLU,  lots of fun stuffs, **welcome to join !**

# Dr. Sai

Multi-agents framework is developed based on AutoGen
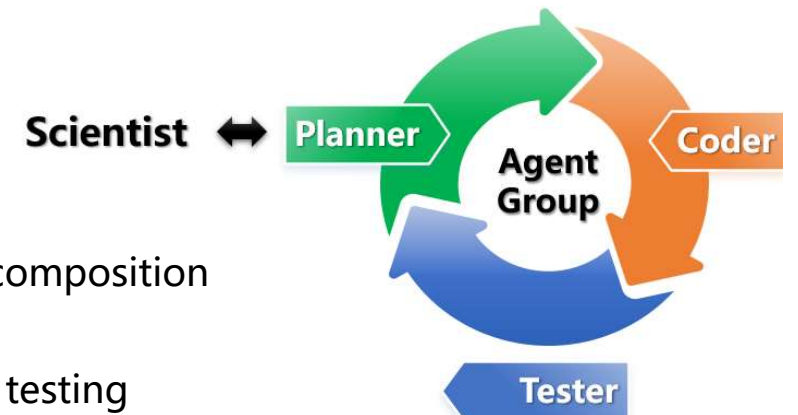


Key of this project:

**make the results from AI more reliable**

- New architecture
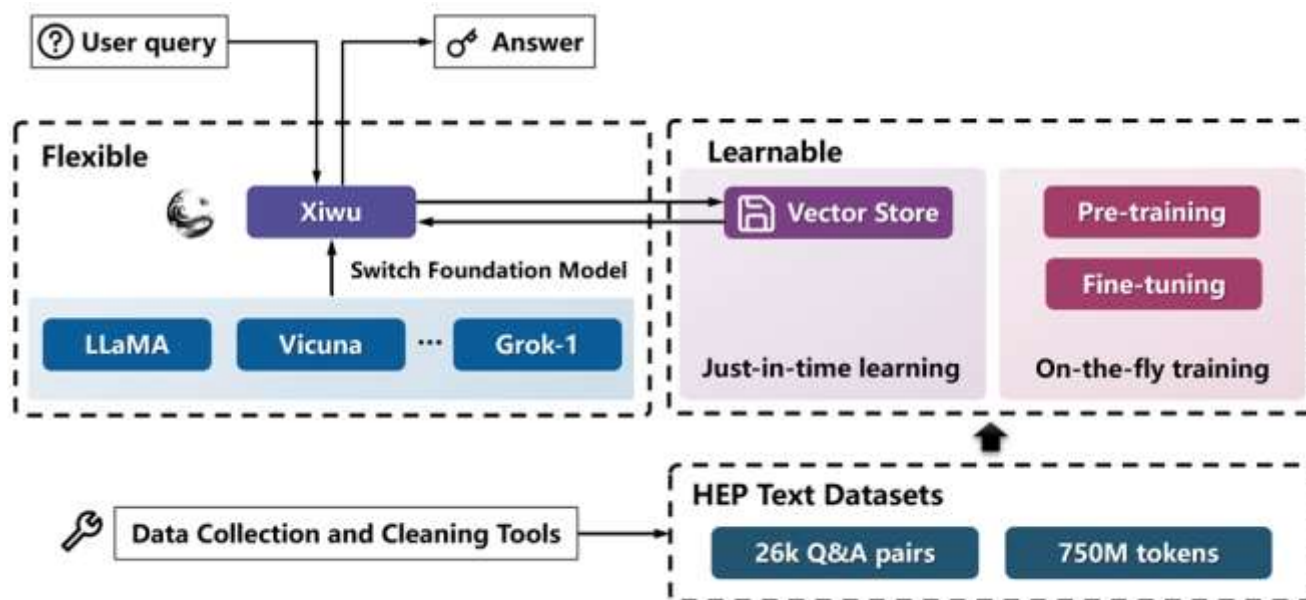- Good quality data
- In-the-fly validation and test

Agents:
- Planner: Planning and tasks decomposition
- Coder: Write code/text
- Tester: Using scientific tools for testing

Human can interact via HumanProxy

# The brain of Dr. Sai – Xiwu model

- Xiwu: a basis flexible and learnable LLM for HEP

- First version release at April (refer to arXiv:2404.08001 for more details)

  - high level model based on open-source foundational LLM, e.g. LLaMa
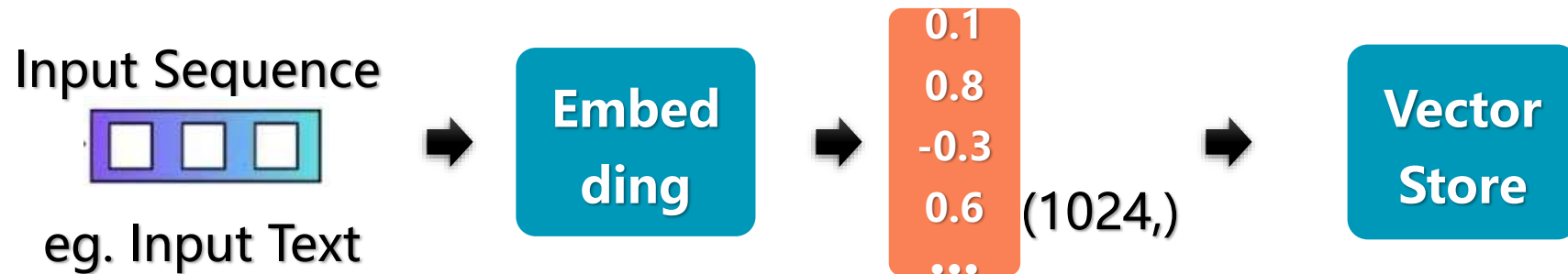
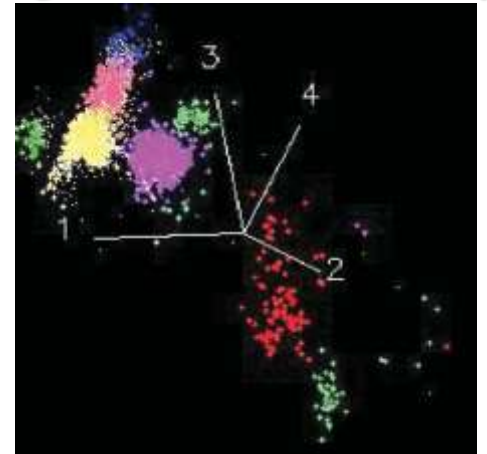  - **First LLM for HEP, version 2 is on-going**

# The Memory of Dr. Sai - RAG

- Retrieval-Augmented Generation (RAG)

  - Most promising solution to avoid hallucinations

  - Goal: store private data so no need for retraining

  - Current approach: vector store

    - Embeddings (BGE-M3 model):

      - Convert input data into vectors of a multi-dimensional space

  - Usage: store BESIII internal data

    - user send BESIII related questions
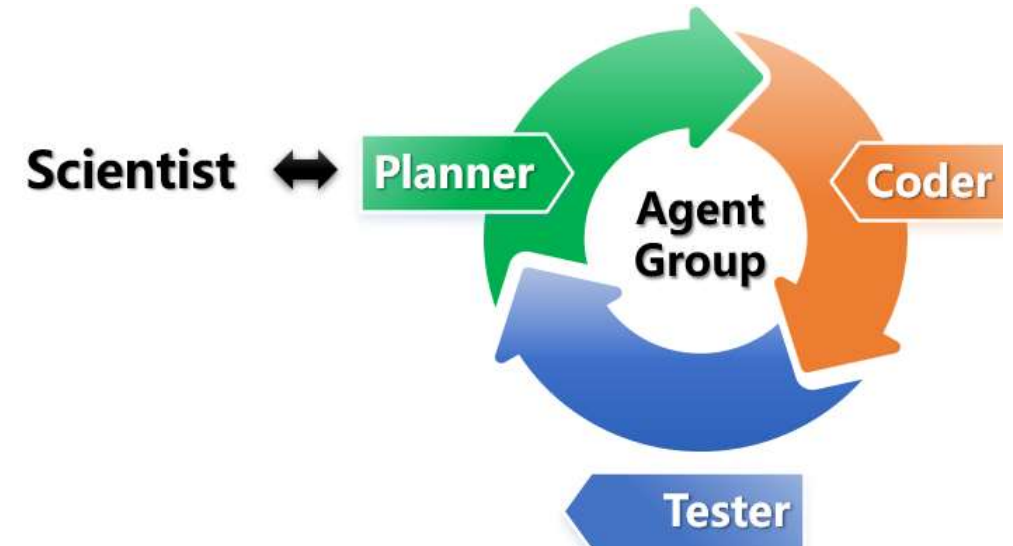
    - RAG return  question + BESIII internal data to LLM

**High Dimensional Space**

**Input Sequence**

eg. Input Text

**Embed ding**

0.1
0.8
-0.3
0.6
... (1024,)

**Vector Store**

# Multi-agents managment system

- Developed based on AutoGen framework

- Normally one agent is dedicated for one task, HEP data processing is very complicated, impossible to build one common agent

- Multi-Agents:
    - **Planner**: Planning and tasks decomposition
    - **Coder**: Write code/text
    - **Tester**: Using scientific tools for testing
    - Common tools: plotting et.al.
    - ...

  •Human can interact via HumanProxy



- If test failed, the feedback will be used to improve the prompts at next iteration.
- For each task, we have multi-unittests

# Training data

- Recent papers on arXiv

  - PDF files parser: HaiNougat, advanced iteration of the Nougat model

- Good quality chat history from IHEP-AI platform

  - The data is cleaned by human or AI (GPT4)

  - 180k Question-Answer pairs in 3 months

- **BESIII internal data**

  - internal memo, parsered by HaiNougat (>5000 PDF files)

  - Q-A pairs from internal paper review ( >50k Q-A pairs)

  - BESIII Offline Software System (BOSS) source code

  - BESIII public webpages and internal webpages

  - The data (agenda, slides) on indico will be used later

- All the BESIII internal data sets are stored in RAG or used in training and fine-tuning

# Simple test: internal navigator at BESIII

- Same with the chATLAS project at ATLAS
  - Navigator and assistant to replace the simple 'search'
  - BESIII internal data at websites (bes3.ihep.ac.cn)
    - Not public yet
  - In general, better performance than I expected
    - E.g. Question 'where is the XXXX MC sample',
    - Answer 'The path of the sample is in XXXXXXXX'

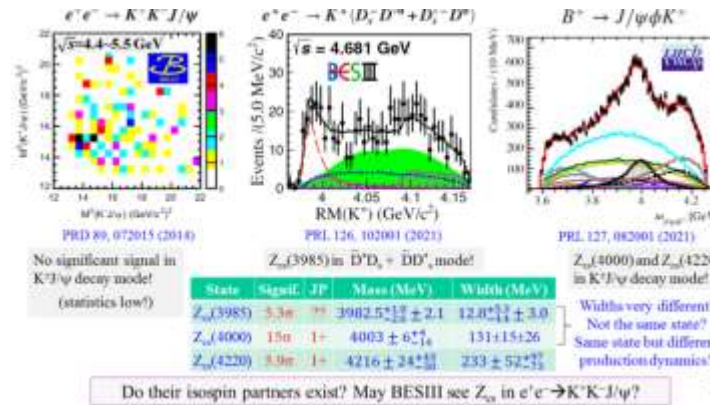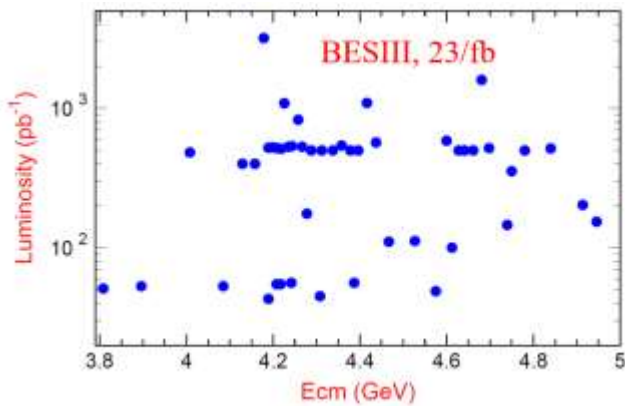question 1 = "What's the link to BESIII offline software group's main page?"

https://docbes3.ihep.ac.cn/~offlinesoftware/index.php/Main_Page

question 2 = "In the KKMC generator, up to which order are the electroweak corrections calculated?"
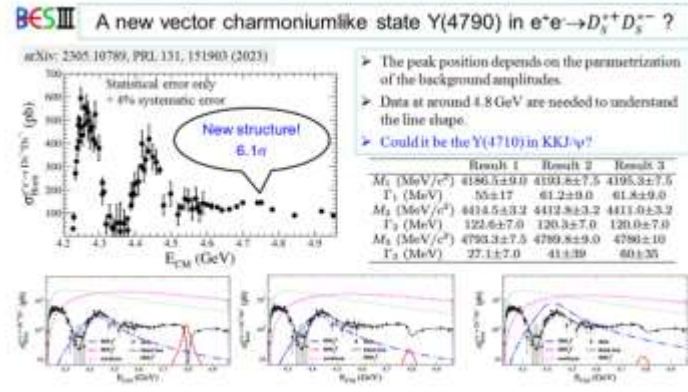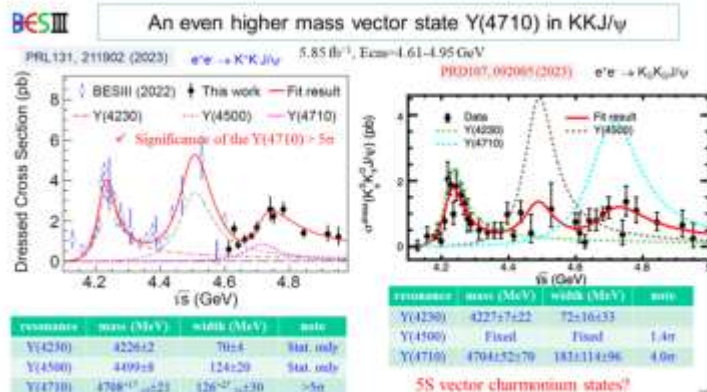
In KKMC, the most important features are the ISR-FSR interference, the second-order subleading corrections, and the exact matrix element for the two hard photons. Effects due to photon emission from the initial beams and the outgoing fermions are calculated in QED up to the second order, including all interference effects, within the Coherent Exclusive Exponentiation (CEEX), which is based on Yennie-Frautschi-Suura exponentiation. Electroweak corrections are included in the first order, with higher-order extensions, using the DIZET 6.21 library. Final-state quarks hadronize according to the parton shower model using PYTHIA. Decays of the [math]\displaystyle{ \tau }[/math] lepton are simulated using the TAUOLA library, taking into account the spin polarization effects as well. The code and more information on the program are available at the KKMC web page.

# Status and prospects for Dr. Sai

- XYZ studies at BESIII:
  - Many (intermediate) resonants, many decay channels -> interferences
  - Difficult to retrieve useful information correctly



**How to understand them? what is the physics behind them ?**

# Status and prospects for Dr. Sai

- Exotic hadron (XYZ) studies at BESIII:
  - Many (intermediate) resonants, many decay channels -> interferences
  - Difficult to retrieve useful information correctly
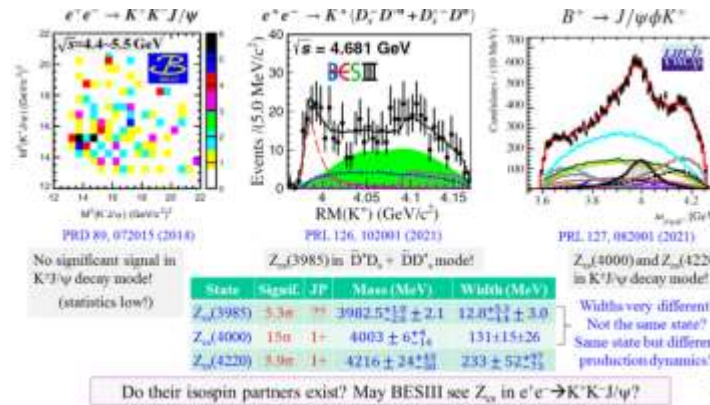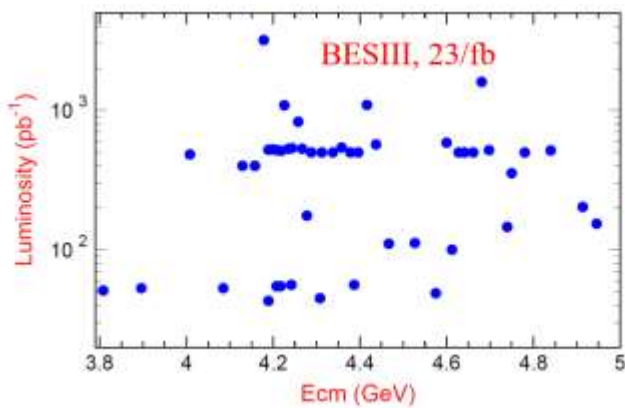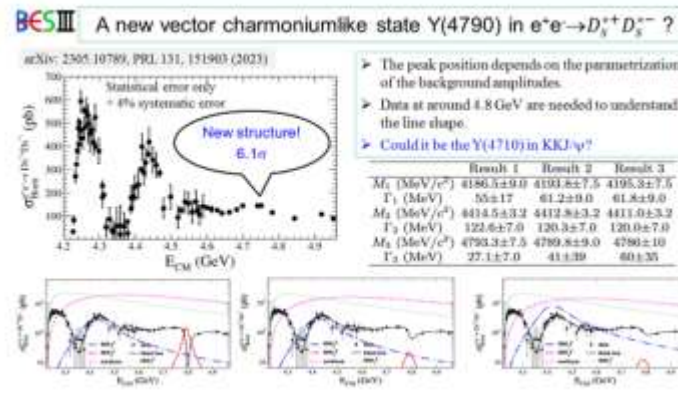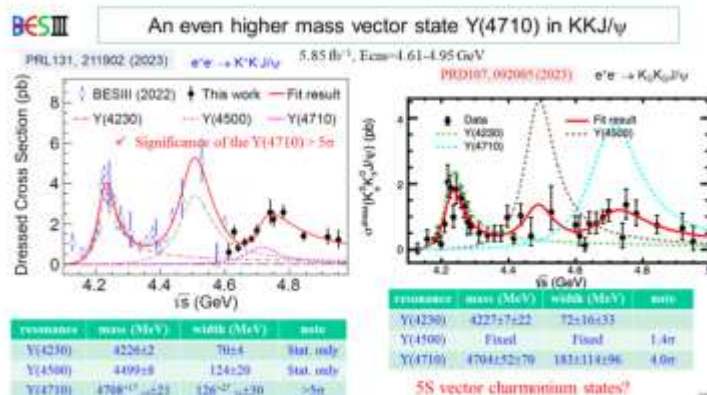


**How to understand them? what is the physics behind them ?**

**As a experimental people, I don't know.**

**But the cross section measurements of ALL channels should be one necessary condition.**
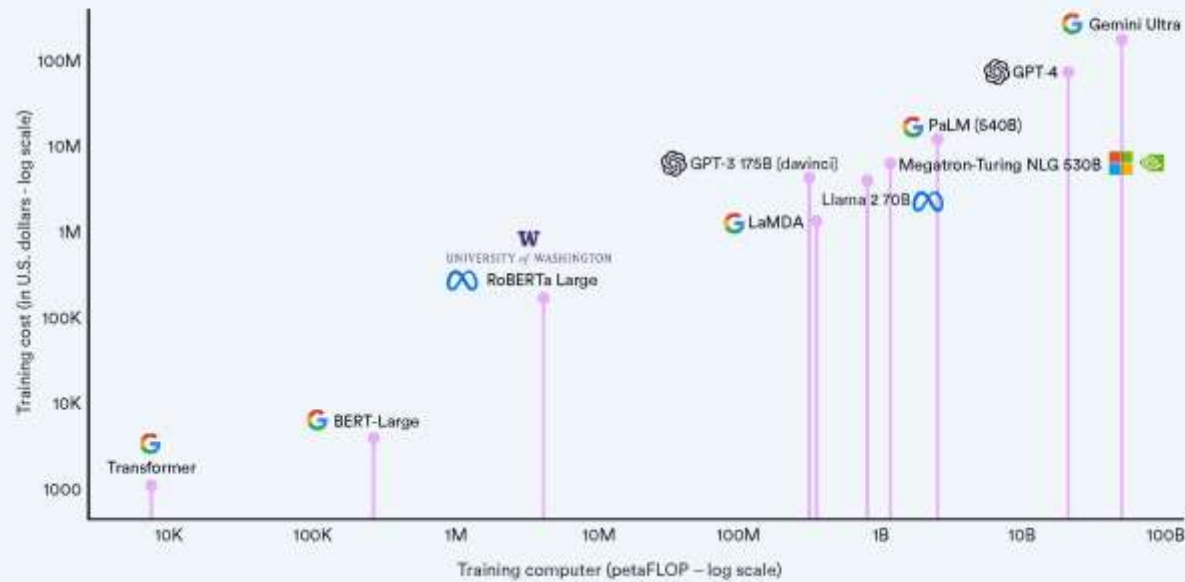
# Summary

- AI agents for HEP - **Dr. Sai** is build for BESIII experiment
  - Assist people on multiple tasks, e.g. data analysis
  - Internal testing and debugging
  - Beta version at June 2024, stable version at the end of 2024
  - Personpowers are more than welcome
- Next:
  - Xiwu2: train with more data, multi-modal, e.g. slides on indico, experts' chat history at IHEP AI platform https://ai.ihep.ac.cn/
  - AI virtual scientist: automate full chain of physics analysis

# back-up

Estimated training cost and compute of select AI models

Source: Epoch, 2023 | Chart: 2024 AI Index report

This is a C++ code for a class called `Gam4pikp` which is used to analyze data from the BaBar experiment. The class contains several methods for filtering and sorting data, as well as outputting results.

This is a C++ code for a data analysis algorithm called Gam4pikp. The algorithm is designed to analyze data from high-energy particle collisions and identify specific patterns of particles. The code appears to be a part of a larger program that is used to analyze data from the Large Hadron Collider (LHC) at CERN.

This is the implementation of a ROOT-based algorithm called `Gam4pikp` which is used to analyze particle physics data. The algorithm is designed to identify and reconstruct events containing four-pion ($4\pi$) and four-kaon ($4K$) resonances, and to perform various calculations and data analysis tasks.

This is a C++ code for a particle physics analysis tool, specifically a program that analyzes data from the Belle II experiment at the SuperKEKB collider. The code is designed to identify and reconstruct particles produced in high-energy collisions, and to study the properties of these particles.

# Status and prospects for Dr. Sai

- **Under construction and testing, plan to release the first version (two AI agents) at June 2024**
  - one dedicated for BESIII and another for public, **stay tune**
- One application: software and training
  - BOSS (C++ code) upgrade
    - step 1: simple improvements using new C++ features, e.g. array to vector
    - step 2: re-structure the code for each file individually
    - step 3: AI-assisted update on algorithms
  - Outreach and training:
    - Train junior graduated students to understand BESIII and data analysis better