

# Recent developments in the data analysis integrated software system of HEPS

Yu Hu

Institute of High Energy Physics, CAS

ISGC 2024/03/26



# Outline

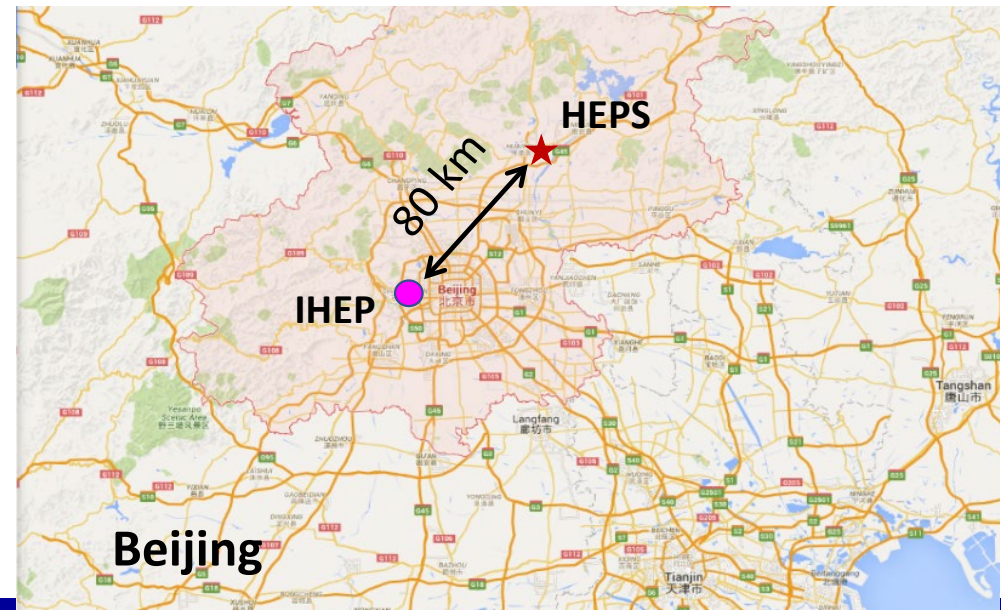
---

- 1. Introduction**
- 2. Demand and Challenges of scientific data and software system**
- 3. The architecture and design of the framework**
- 4. The recent Progress of the system**
- 5. Summary**

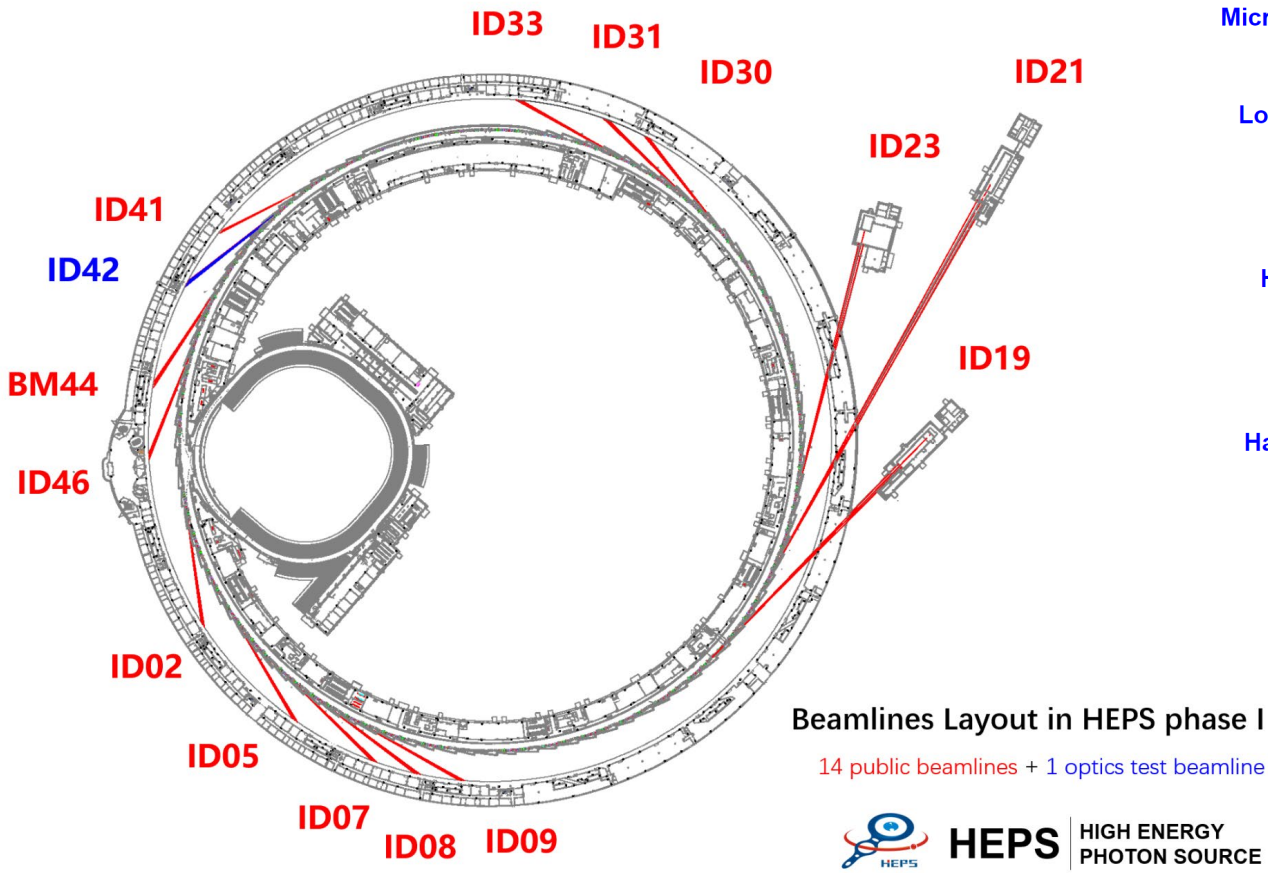
# High Energy Photon Source (HEPS)

- New light source in China — High energy, high brightness
- Located in Beijing - about 80KM from IHEP
- Officially approved in Dec. 2017
- The construction was started in the middle of 2019
- The whole project will be finished in mid-2025

Main parameters	Unit	Value
Beam energy	GeV	6
Circumference	m	1360.4
Emittance	pm·rad	< 60
Brightness	phs/s/mm <sup>2</sup> /mrad <sup>2</sup> /0.1%BW	>1x10 <sup>22</sup>
Beam current	mA	200
Injection		Top-up



# Beamlines in HEPS phase I



- Microfocusing X-Ray Protein Crystallography-ID02 Beamline
- Low-Dimensional Structure Probe Beamline-ID05
- Engineering Materials Beamline-ID07
- Hard X-Ray Coherent Scattering Beamline-ID09
- Pink Beam SAXS Beamline-ID08
- Hard X-Ray Nanoprobe Multimodal Imaging-ID19 Beamline
- Hard X-Ray Imaging Beamline-ID21
- Structural Dynamics Beamline-ID23
- ID30-Transmission X-Ray Microscopic Beamline
- ID31-High Pressure Beamline
- ID33-Hard X-Ray High Resolution Spectroscopy Beamline
- BM44-Tender X-Ray Beamline
- ID41-High Resolution Nanoscale Electronic Structure Spectroscopy Beamline
- ID42-Optics Test Beamline
- ID46-X-Ray Absorption Spectroscopy Beamline

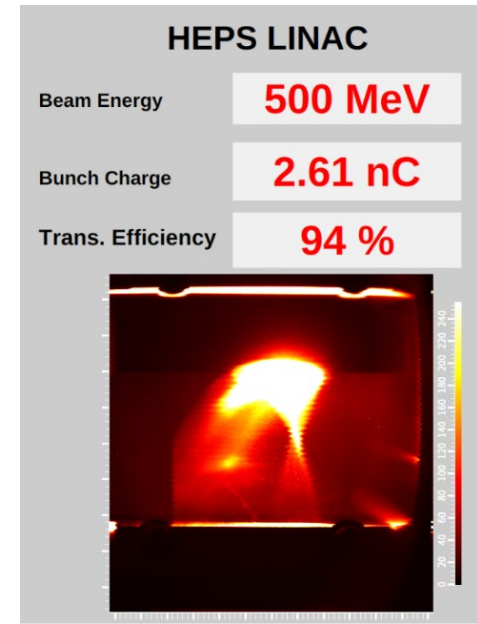
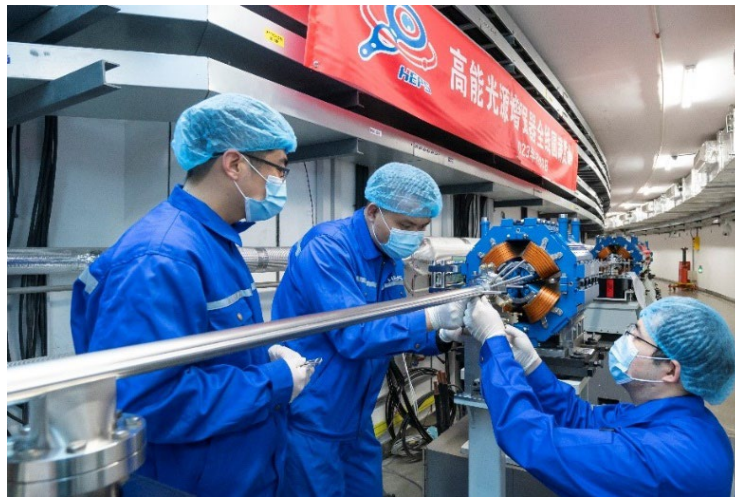
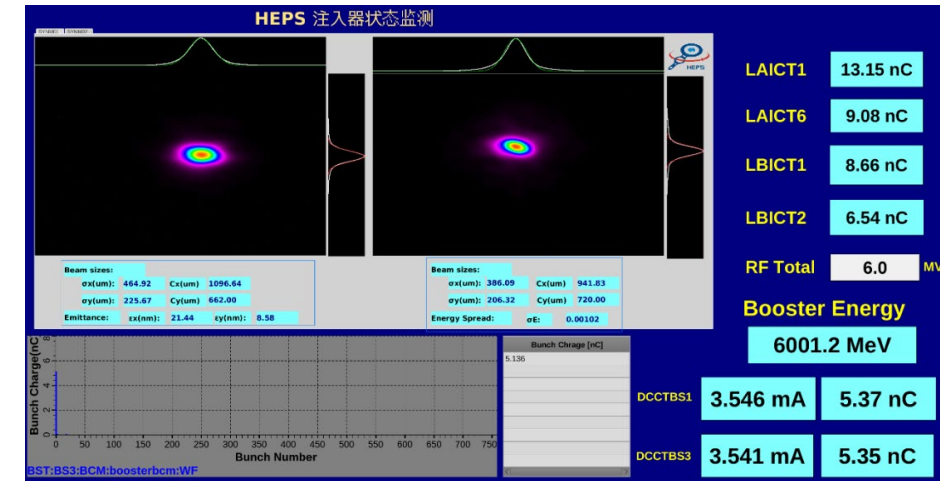
14 public beamlines + 1 optics test beamline in Phase I  
 Can accommodate over 90 beamlines in total



**Experiment Hall**

# Progress of the HEPS project

- ❑ The construction of the civil structure completed. Now at the equipment installation stage
- ❑ 2023.01, HEPS booster installation completed
- ❑ 2023.02, Start installation of storage ring
- ❑ 2023.03, HEPS achieved the first electron beam accelerated to 500 MeV
- ❑ 2023.11, HEPS electron beam ramped up to 6 GeV
- ❑ 1st SR X-ray to be emitted in 2024



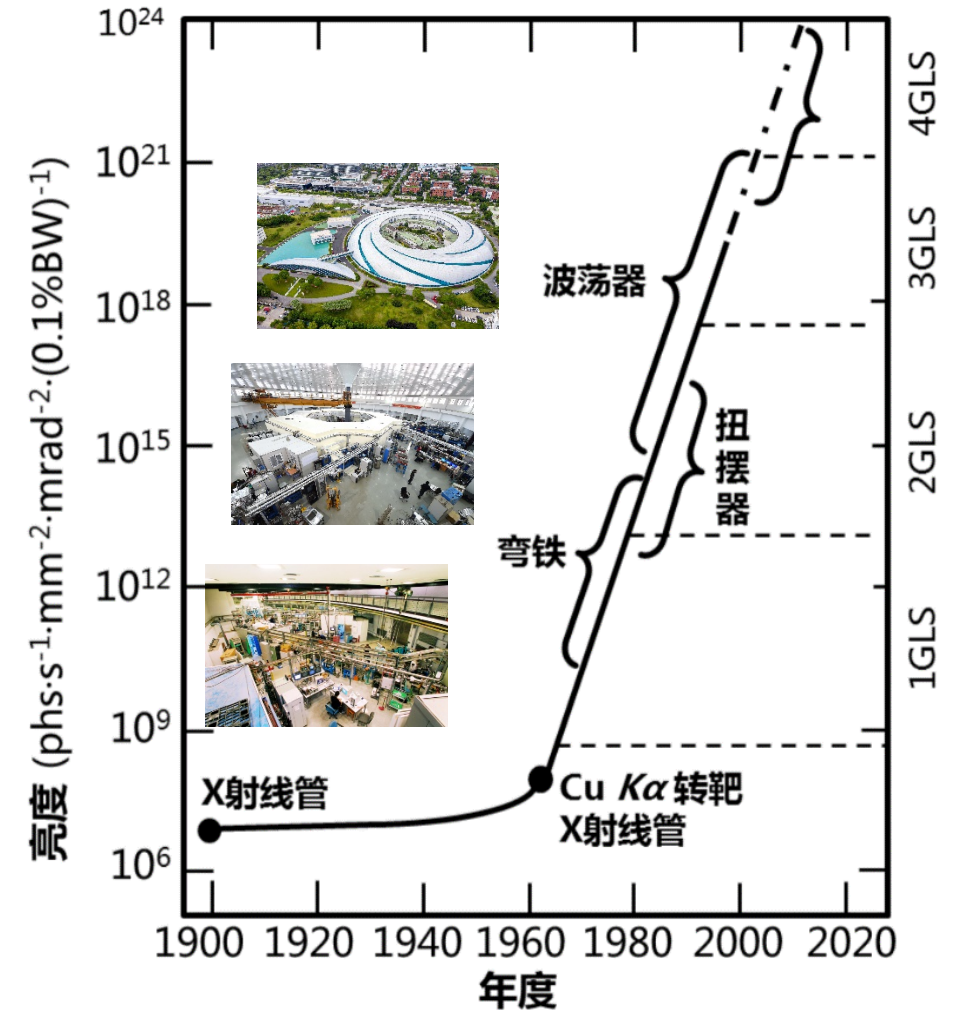
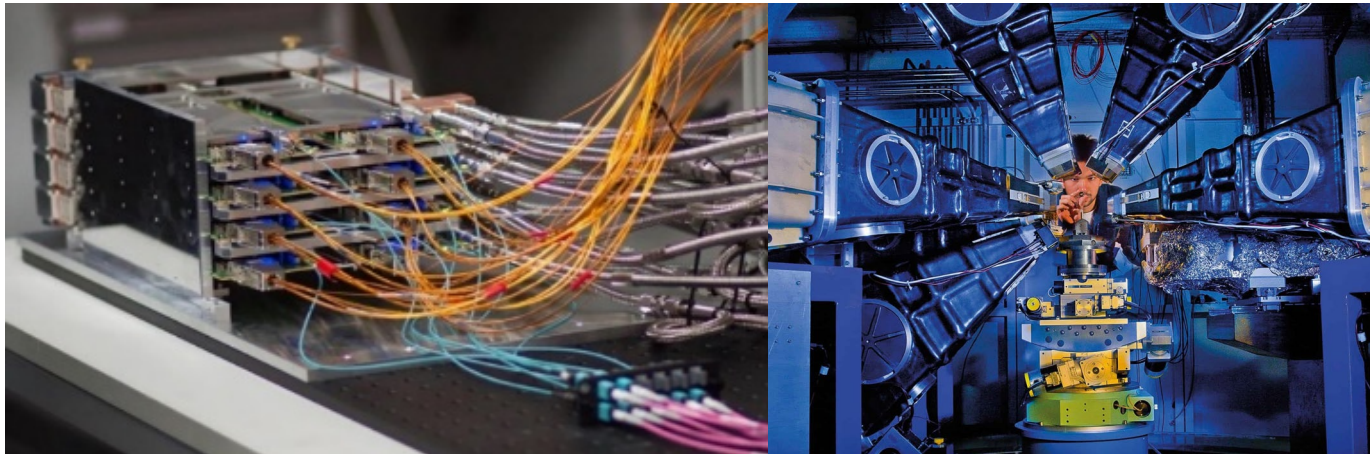
# Outline

---

1. Introduction
- 2. Demand and Challenges of scientific data and software system**
3. The architecture and design of the framework
4. The recent Progress of the system
5. Summary

# Data Challenges

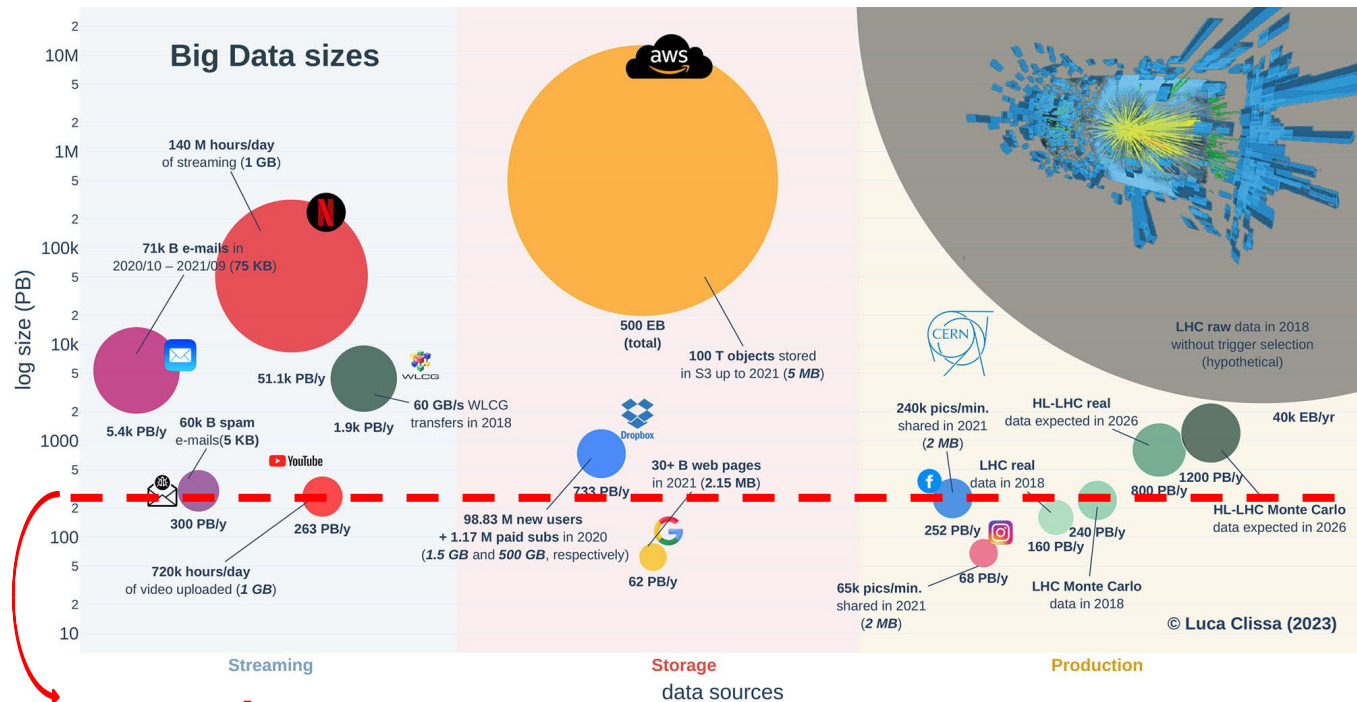
- Increased source brightness/ luminosity
  - More raw data in greater detail and less time
- Detector capabilities constantly improving:
  - Increased dynamic range, faster readout rates, larger pixel arrays
  - Bigger frames, higher frame rates
  - => more raw data



Development of synchrotron radiation light source

# Data Challenges

- ❑ >200 petabytes of raw data per year for Phase I of HEPS (15 beamlines)
- ❑ More than 90 beamlines volume in total
- ❑ Data volumes to reach exabytes in the near future



## HEPS Phase I

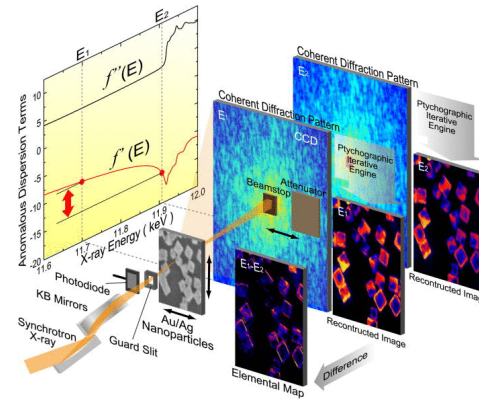
## Data volume of HEPS Phase I Beamlines:

Beamlines	Burst output (TB/day)	Average output (TB/day)
Engineering Materials	600.00	200.00
Hard X-ray Multi-analytical Nanoprobe	500.00	200.00
Structural Dynamics	8.00	3.00
Hard X-ray Coherent Scattering	10.00	3.00
Hard X-ray High Energy Resolution Spec.	10.00	1.00
High Pressure	2.00	1.00
Hard X-Ray Imaging	1000.00	250.00
X-ray Absorption Spectroscopy	80.00	10.00
Low-Dimension Structure Probe	20.00	5.00
Biological Macromolecule Microfocus	35.00	10.00
pink SAXS	400.00	50.00
High Res. Nanoscale Elec. Struc. Spec.	1.00	0.20
Tender X-ray beamline	10.00	1.00
Transmission X-ray Microscope	25.00	11.20
Test beamline	1000.00	60.00
<b>Total average:</b>		<b>805</b>

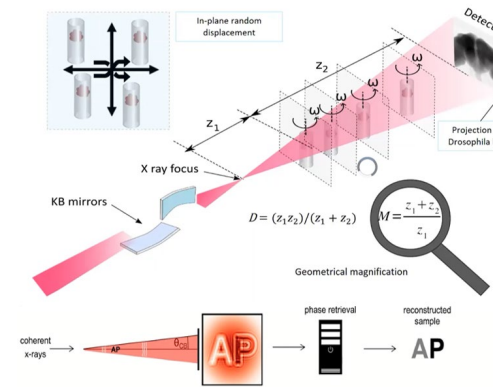


# Data Challenges

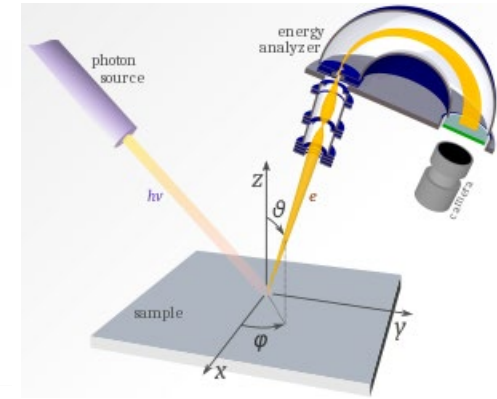
- ❑ New and more complex experiments
- ❑ Multi-modal experiments require combining data from multiple samples, techniques, and facilities
- ❑ In situ and in operando experiments require real-time feedback and autonomous control
- ❑ Data throughput and volume vary greatly with experiments and scientific goals
- ❑ New users from a wide variety of backgrounds and domains



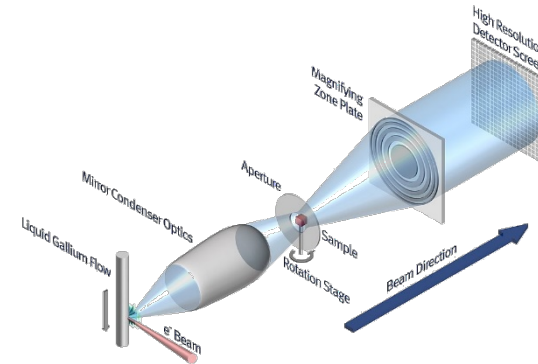
Fluorescence mapping



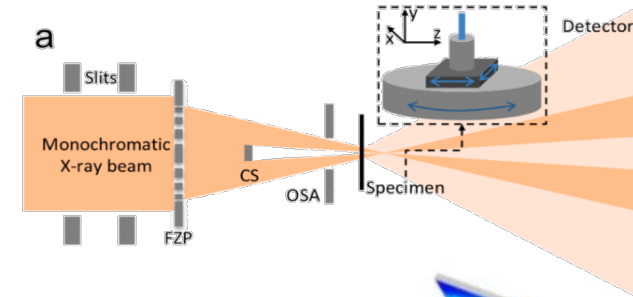
Nanoholotomography



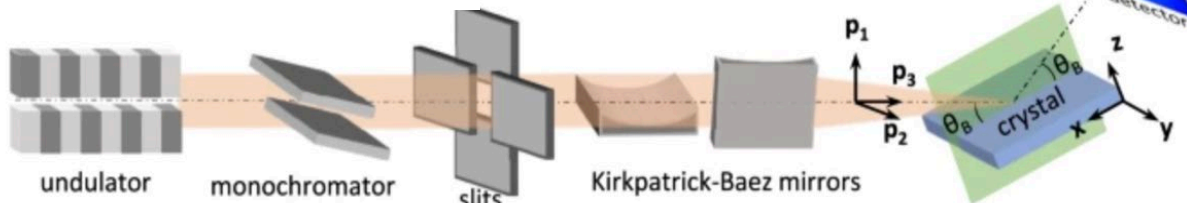
ARPES



Nano CT



Ptychography CT



Bragg ptychography

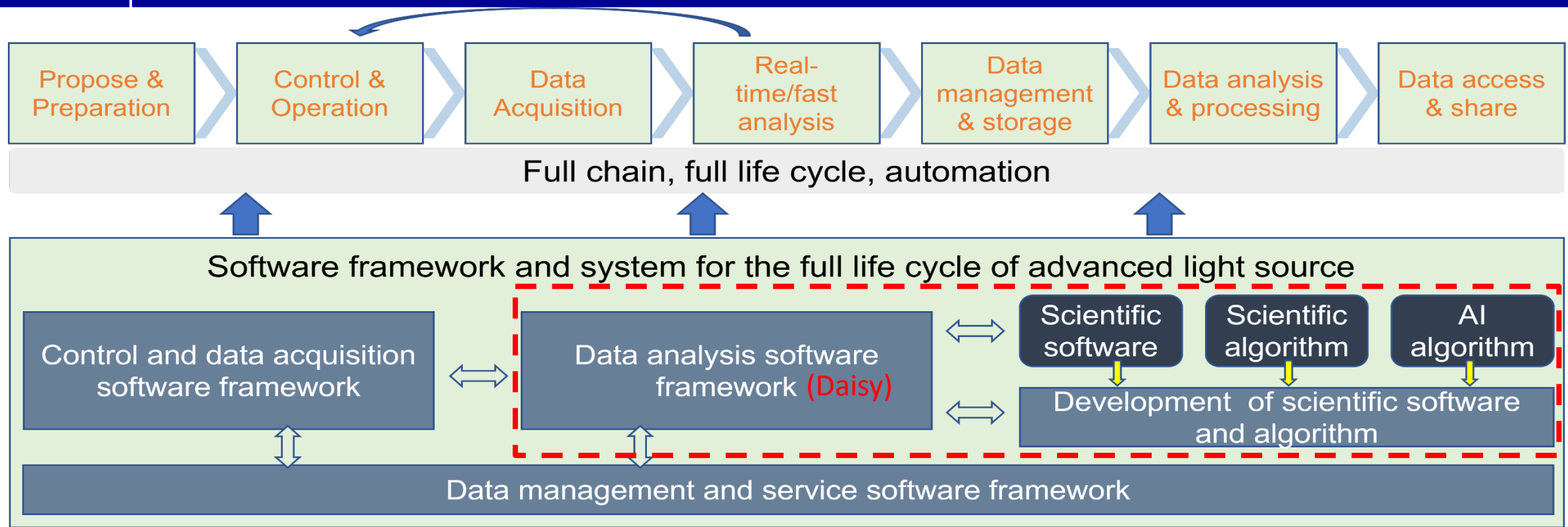
# Data Challenges

- Analysis and management of large datasets at advanced scientific facilities is becoming progressively more challenging
- Development and integration of advanced analysis and management tools is needed
  - Provide storage, organization and management of massive scientific data
  - During the experiment, provide real-time analysis and fast feedback to guide the experiment steering and optimize the data acquisition
  - After the experiment, process the massive offline data, accelerate the scientific discovery
  - Provide the scalable distributed heterogeneous computing power, meet the diverse computing requirements of different scientific goals
  - Make things easy and faster for users

# Outline

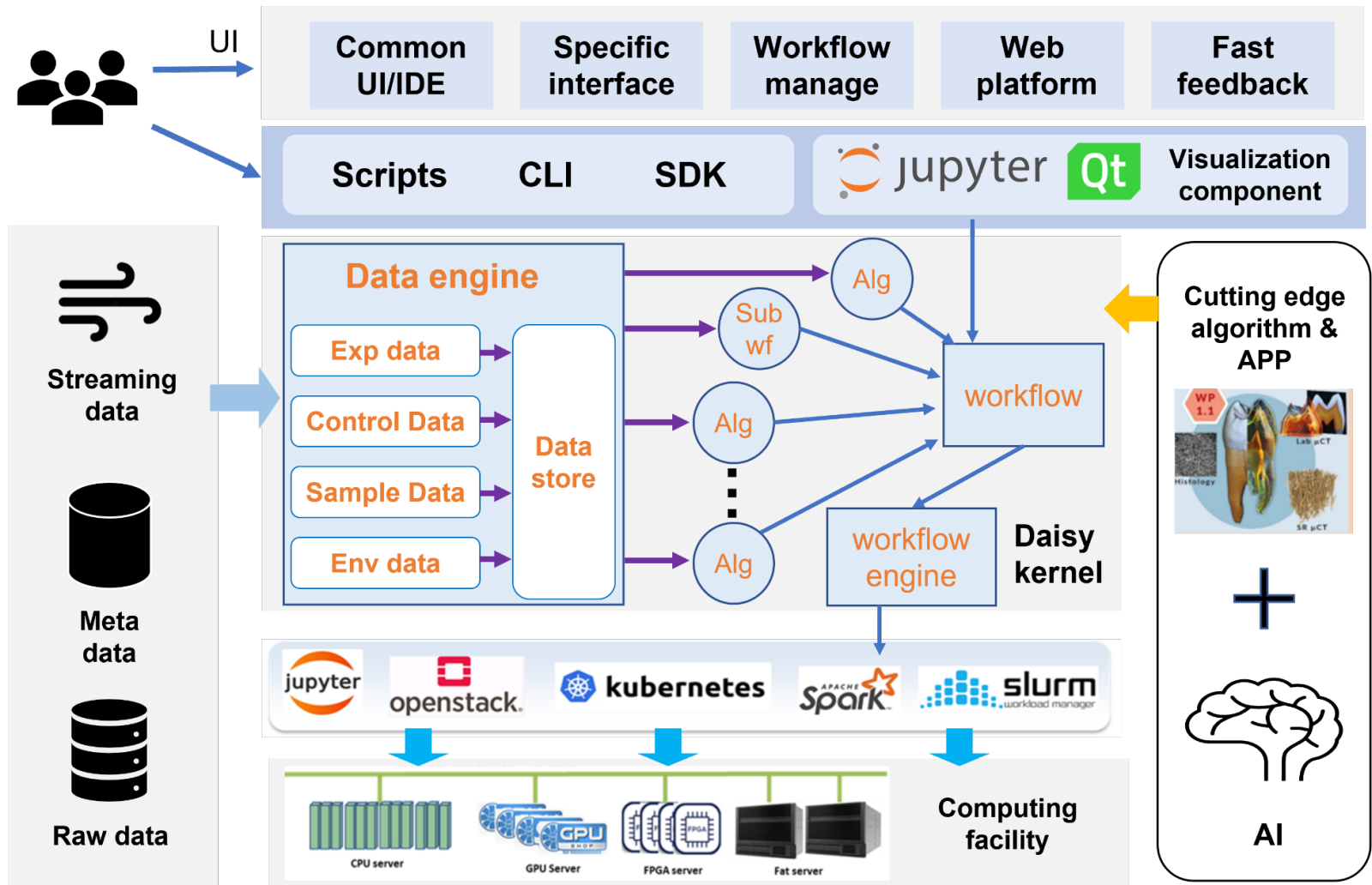
1. Introduction
2. Demand and Challenges of scientific data and software system
- 3. The architecture and design of the framework**
4. The recent Progress of the system
5. Summary

# Full data lifecycle software system



- ❑ Software framework and system for the full data life cycle of advanced light source
- ❑ Implement the tracking and management of scientific data throughout the entire process
- ❑ Support the development of new advanced data analysis methods, as well as the integration of existing software into the framework
- ❑ Designed for light source at first, but also suitable for other facilities

# Data analysis software framework—Daisy



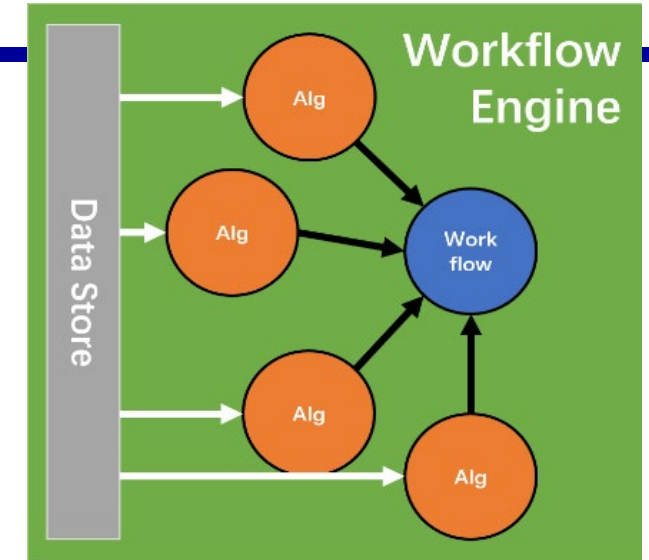
- Kernel of the framework
- Derivative technology modules to meet the data processing requirements of advanced scientific facilities
  - Data object management module for high-throughput data I/O, multimodal data exchange, and multi-source data access.
  - Scalable cluster computing power support for data processing with different scales, different throughputs, and low latency
  - Interface and developing environment for scientific software integration and development
- Domain specific App and flexible general workflow management system based on the framework

# Kernel of the Daisy framework

Extract domain models independent from technology, and establish relationships between models to form a domain architecture

## Four core modules are provided:

- **Algorithm:** The smallest unit in framework, defining the domain model, basic data processing module, support integration of third-party libraries.
- **Workflow:** Defines the domain architecture, execute processing tasks by calling a series of algorithms, supporting nesting.
- **Workflow Engine:** Manages the runtime environment and the distribution of the algorithm modules. Uncouple the process task from the computing environment.
- **Datastore:** Manages the creation and transmission of data objects between algorithms.



### Algorithms

- Input Data Processing
- Output Data Defined

### Workflow Engine

- Handle Data Store
- Running Time Management

### Business Domain

- Algorithms
- Workflow

### Running Time

- Workflow Engine
- Data Store

### Workflow

- A sequence of Algorithm
- Workflow is also an algorithm

### Data Store

- Data Object Management

# Outline

---

1. Introduction
2. Demand and Challenges of scientific data and software system
3. The architecture and design of the framework
- 4. The recent Progress of the system**
5. Summary

# Daisy graphical user interface

Daisy Workbench

File View Interfaces Help

Workspaces

Load Delete Clear

Sort Save

name

tooth

Interfaces

Help

Integration

MaxMin

PyFAI calib

XRF Batch Fitting

Spectra Matching

tooth

	1	2
2	[27008.75 ...	[27098.75 ...
3	[27051.75 ...	[26986.25 ...
4	[27192.75 ...	[27107.5 ...
5	[27208. ...	[27020.25 ...
6	[27181.75 ...	[26995. ...
7	[27190. ...	[27033.5 ...
8	[26869.75 ...	[27169.25 ...
9	[27142.25 ...	[26977.75 ...
10	[27407.75 ...	[27282.5 ...

Algorithms

Excute AlgMatrixTransp...

Daisy

- LoadHDF5
- LoadTIFs
- SaveHDF5
- SaveH5VDS
- AlgMatrixTransp...

IPython

```
In [4]: load start:tooth
workflow:LoadHDF5.config
INFO: initialized, read
HDF5 File: /root/tooth.h5
workflow:LoadHDF5.execute
INFO: Load data /exchange/
data as tooth from /root/
```

System Memory Usage

1.72/62.76 GB (2%)

Plots

Show Hide Close

Plot Name

Plot

0

100

0 100 200 300 400 500 600

0.011373519897460938s

QCoreApplication::exec: The event loop is already running

## Daisy workbench:

- General-purpose GUI based on PyQt5
- Include data object list, algorithm list, data view/visualization, and IDE for developers
- Interfaces of custom GUIs for a variety of scientific techniques

jupyterhub Home huy Logout

启动已选择的分析环境

## 应用分析环境列表

分析环境1

分析环境2

- CT 3D reconstruction  
CT 3D reconstruction service based on tomopy.
- alphafold-with-40g  
alphafold-with-40g
- cumopy  
cumopy

开发者环境

- 

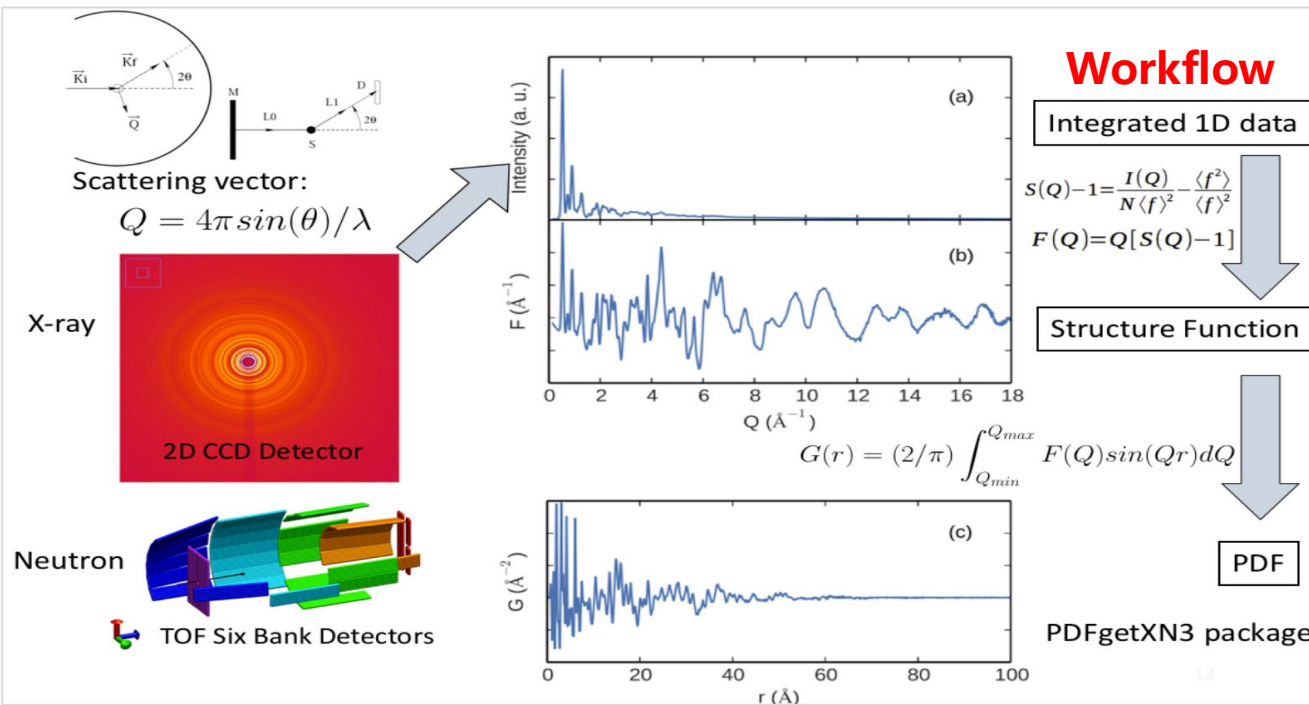
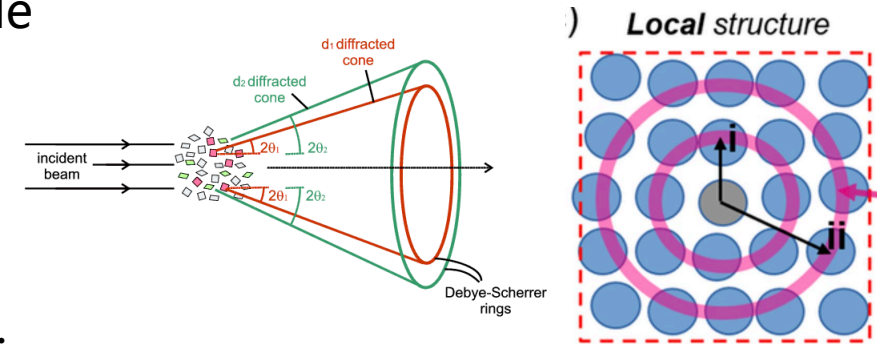
## Web data analysis platform:

- Based on the jupyterlab ecosystem
- Container encapsulates the computing environment
- Scalable computing resource
- Terminal and web scientific APP

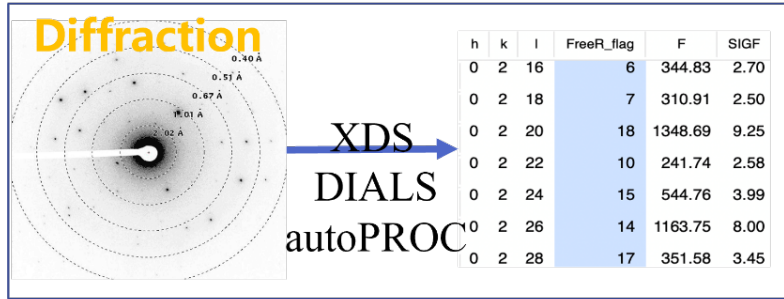


# Application for Pair distribution function(PDF)

- Serve for total scattering experiment, rapid and highly automatable pipeline from raw data to pair distribution functions
- Developed PDFHEPS python package, integrated several X-ray scattering scientific software, such as PyFai, PDFgetX3 and LiquidDiffract
- Web GUI is provided for interactive data processing and visualization



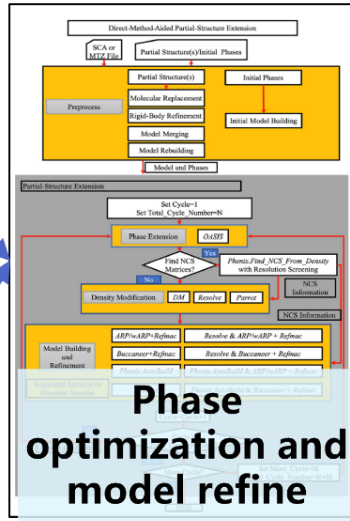
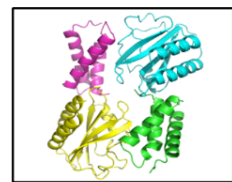
# AI-based application for biological macromolecule



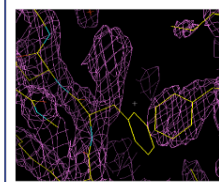
Real-time data processing



Structure prediction based on AlphaFold2



Structure truing based on AI

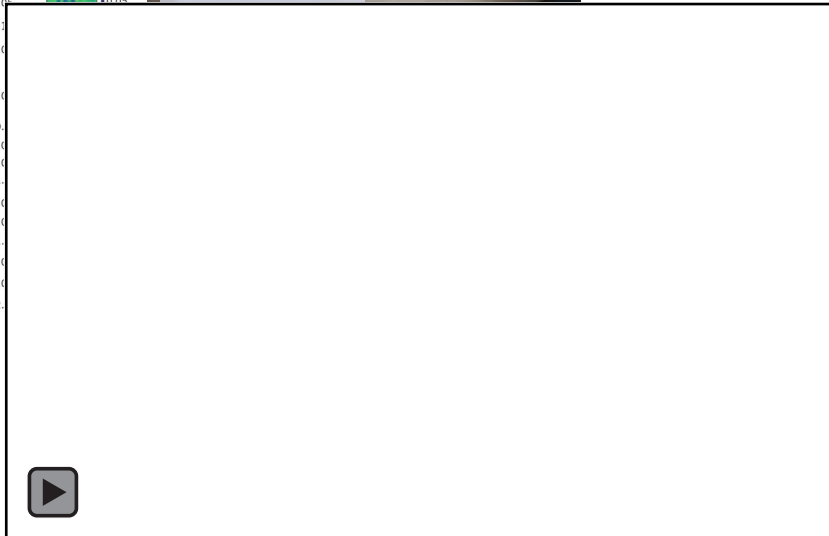
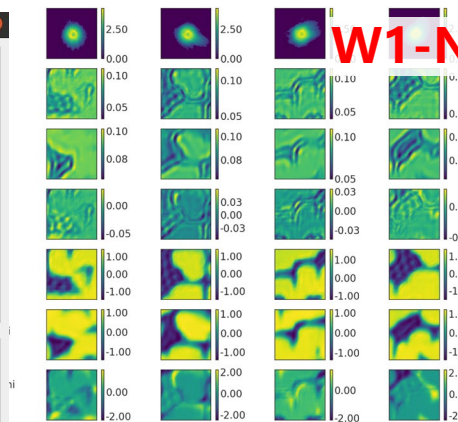
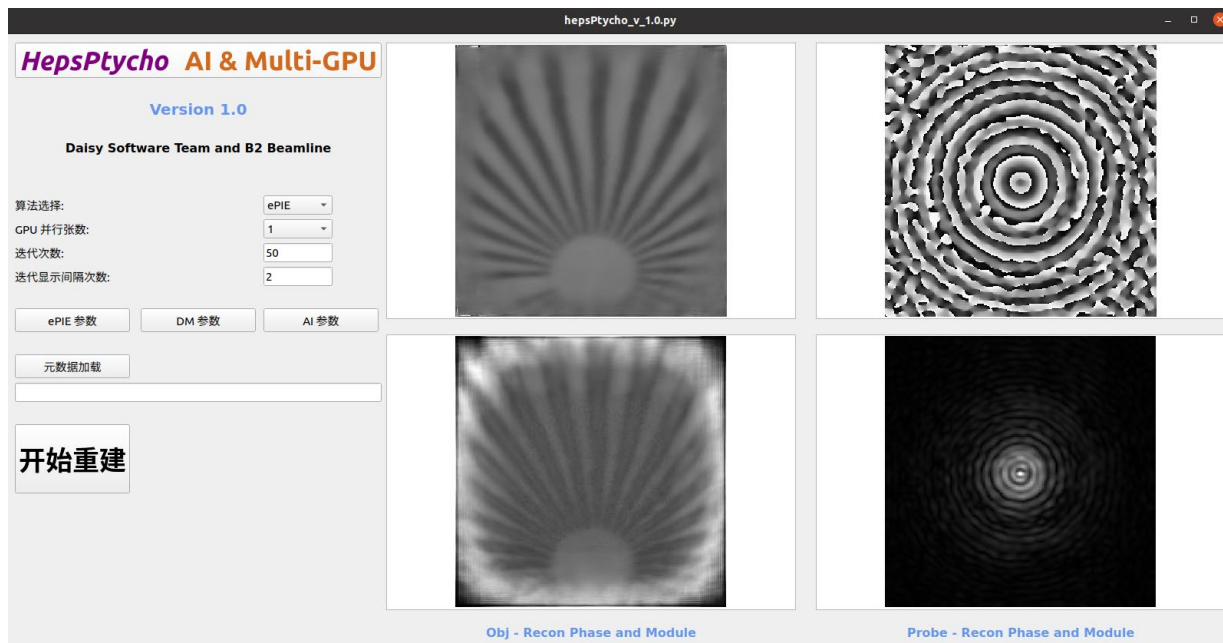


Data name	Space group	a	b	c	$\alpha$	$\beta$	$\gamma$	OscWidth	Frames	Resolution	Inner Rmeas	Outer Rmeas	Rmeas	In
8f8u	P 21 21 2	99.91	155.73	58.75	90	90	90	0.1	1800	3.08	0.025	1.129	0.063	-0
puck4-14_1	P 63 2 2	85.35	85.35	106.76	90	90	120	0.5	1440	2.89	0.064	1.295	0.142	0.

- Serve for structure reconstruction of biological macromolecule. Automatic pipeline from diffraction to macromolecule structure
- Web GUI offering real-time data processing status monitoring and result query
- Based on alphafold2, the success rate and accuracy of macromolecular structure reconstruction get improved

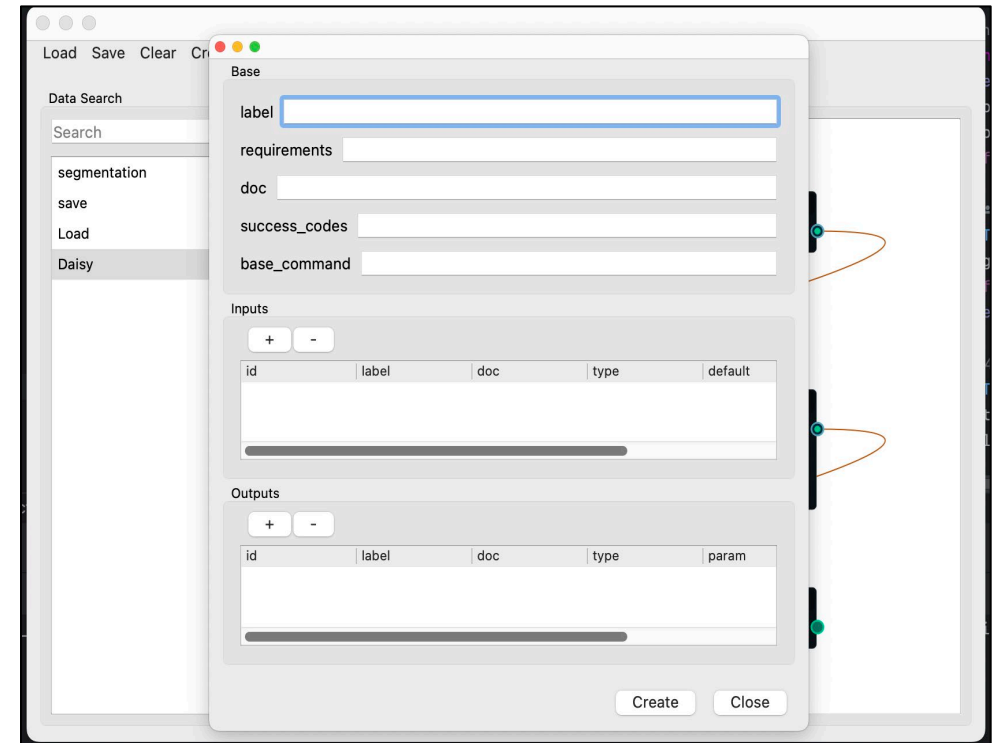
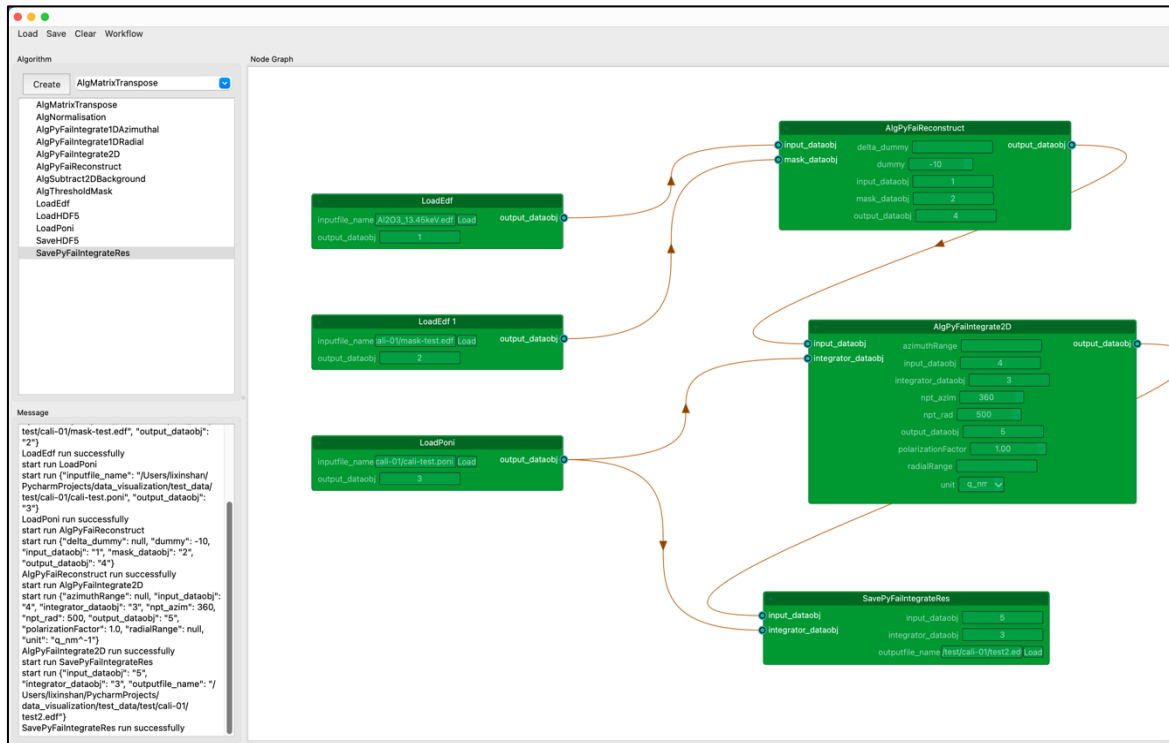
# Application for X-ray ptychography

- For coherent X-ray Imaging. Supports various phase retrieval algorithms such as ePIE and DM
- Supports multi-GPU parallel processing for large-scale data. Migration on Rocm GPU are also in progress
- A fast phase recovery algorithm(W1-Net) based on AI is developed, which is 500 times faster than the traditional method
- W1-Net will be used to optimize the DAQ



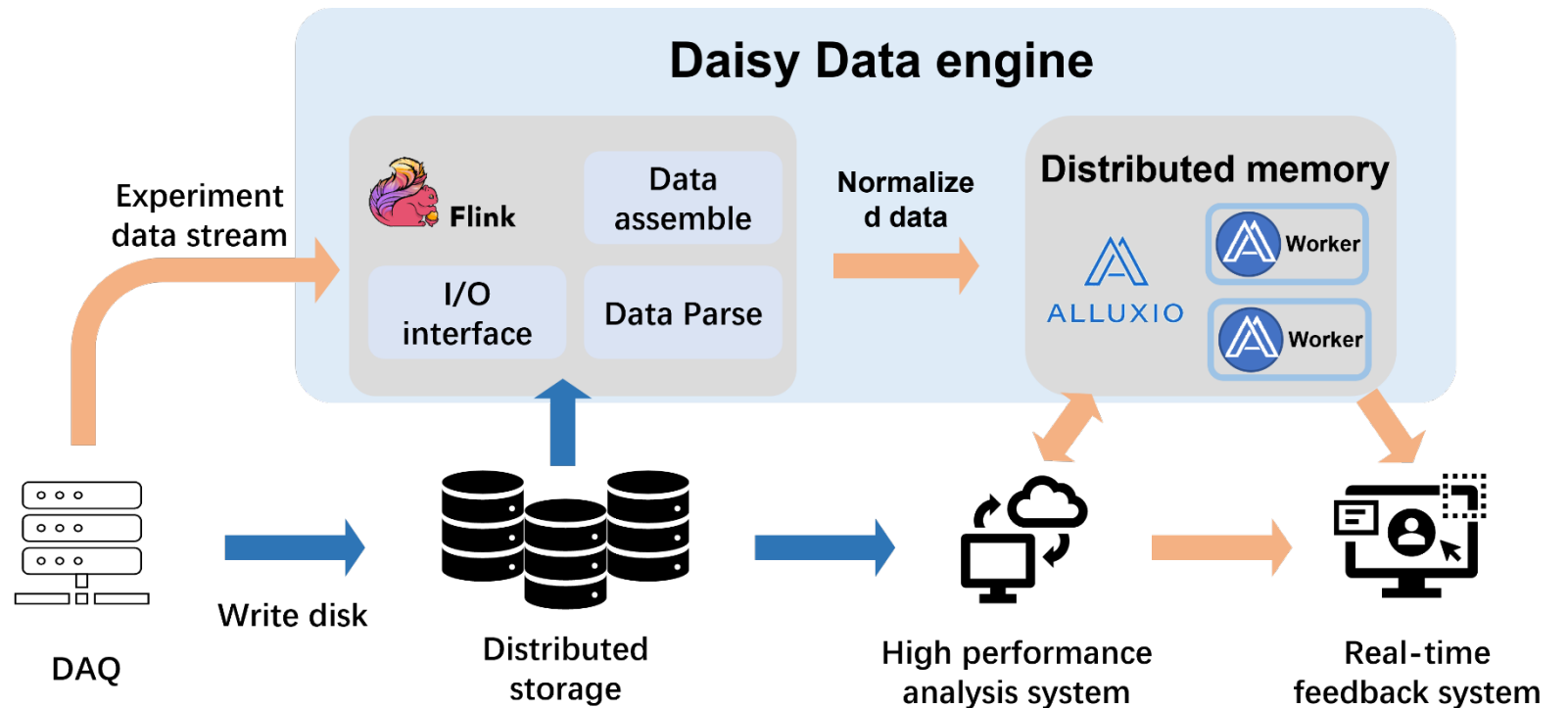
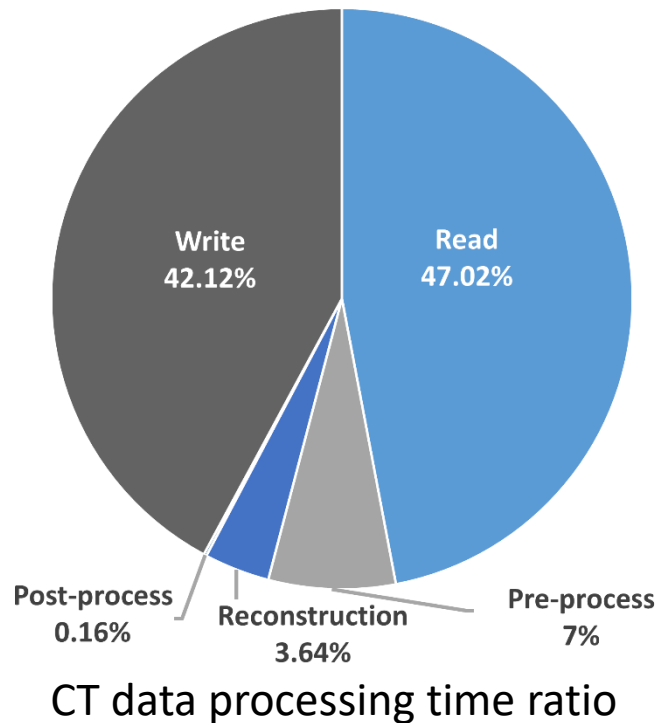
# Workflow management system

- For flexible and general data process task
- GUI support interactive workflow creating, import, export and operation monitoring
- Follow the Common Workflow Language(CWL) standard
- Daisy workflow: automatically parsing algorithms into nodes, execute, monitor, visualization
- Common workflow: create nodes for script commands, can execute on multiple CWL platforms



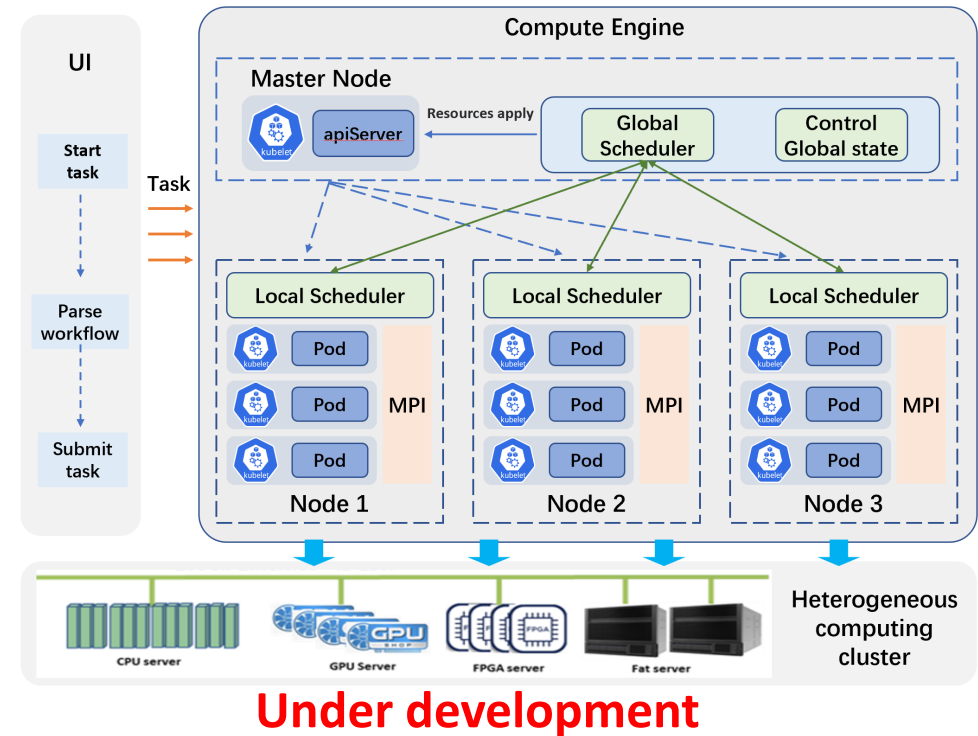
# Data I/O optimization

- I/O is the bottleneck. Employ asynchronous parallel, distributed memory, adaptive storage parameters and compression to optimize the I/O
- For real time, high throughput data process task, the streaming data process method is employed to avoid disk data I/O delay
- Verified in fluorescence and spectral data processing pipeline
- Unified I/O interface to shield the difference of underlying architecture and data structure

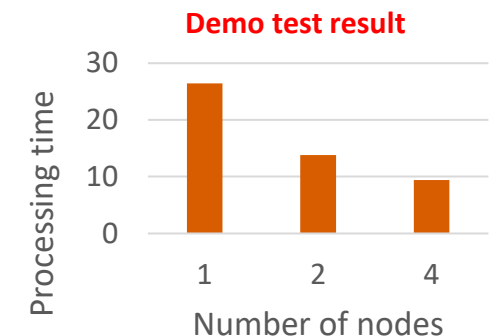


# Distributed data processing support

- A single dataset of HEPS imaging experiment will reach the TB scale
- Scientists expect data processing time at the scale of DAQ time
- A distributed data processing system is designed and developing
- Support heterogeneous distributed computing power
- Provide a unified flexible programming interface API for computing models, to reduce the complexity of parallel programming
- Two layer distributed computing task scheduler to achieve better efficiency



Mode	Detector pixel	Frame rate	Projections number	Data rate	Dataset (TB)	Acquisition time	Daily data (TB/d)	Annual data (PB/y)
Powder CT	6k×6k	19fps@16bit	6k	1.08 GB/s	0.432	6.3 min.	78	9.4
High voxel CT	28kx10k	2pfs@16bit	28k	1.1 GB/s	15.68	240 min.	87	10.4
Fast CT	5k×4k	595fps@8bit	5k	1.7 GB/s	0.1	1 min.	98	5.9



# Support for developers and users

## Version control

- Git for version control, Gitlab hosted project code, connected CI/CD



## Runtime environment

- Container packages the foundational runtime environment
- Harbor manages container images
- CVMFS deploy the pre-compiled software

## Documentation

- User documentation, guide for the developer
- Based on Jupyterbook, Sphinx, readthedocs
- Doxygen generate documentation from source code

DAISY project

Data analysis integrated software system.

Stars 4 | doi 10.1051/epjconf/202125104020

DAISY (Data Analysis Integrated Software System) is a software framework developed using object-oriented technology and programming languages such as C++ and Python. It was originally designed for advanced photon source scientific data processing. During its initial design, it was inspired by some of the world's leading data processing software projects, including DAWN, a data analysis software developed for the Diamond Light Source in the UK; Mantid, a data analysis software framework developed for the ISIS neutron and muon source in the UK; EDNA, an online data processing software framework developed for European Synchrotron Radiation Facility; and Gaudi, a data processing software framework for high energy physics.

The aim of DAISY was to create a versatile and highly extensible basic software architecture. It integrates various methodological algorithms and tools, abstracting away the complexity of the computational architecture and the diversity of computing resources. This framework provides a uniform and simple interface for higher-level application software and users, with additional development of generic components, including desktop tools for data visualisation and analysis, aimed at fostering a rich and thriving software ecosystem.

This documentation is organized into a few major sections.

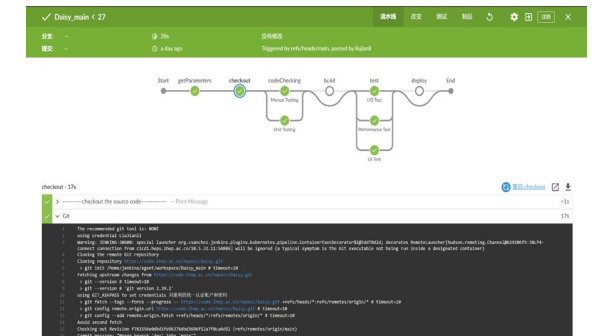
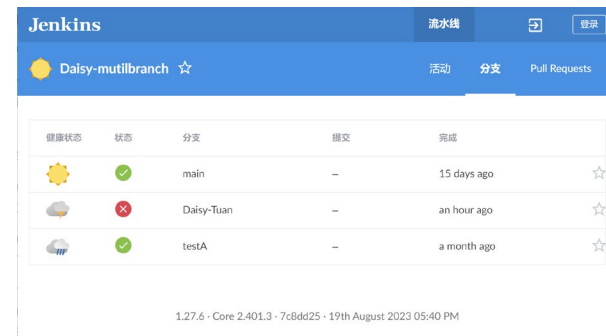
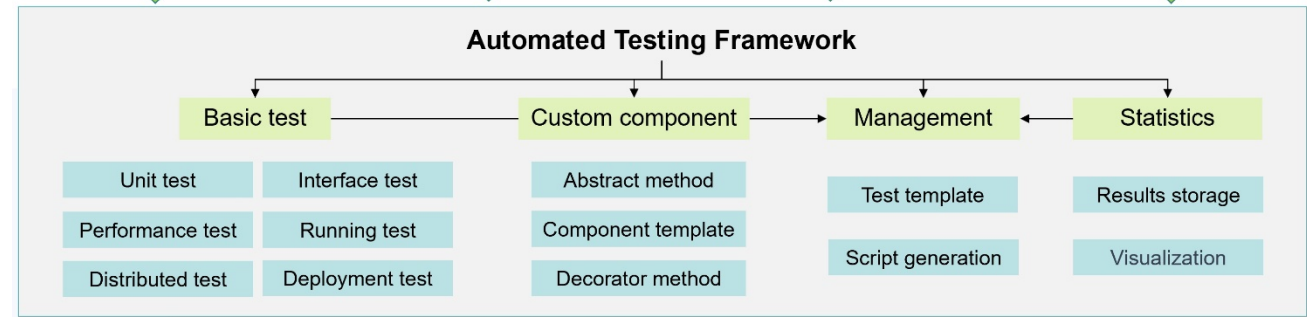
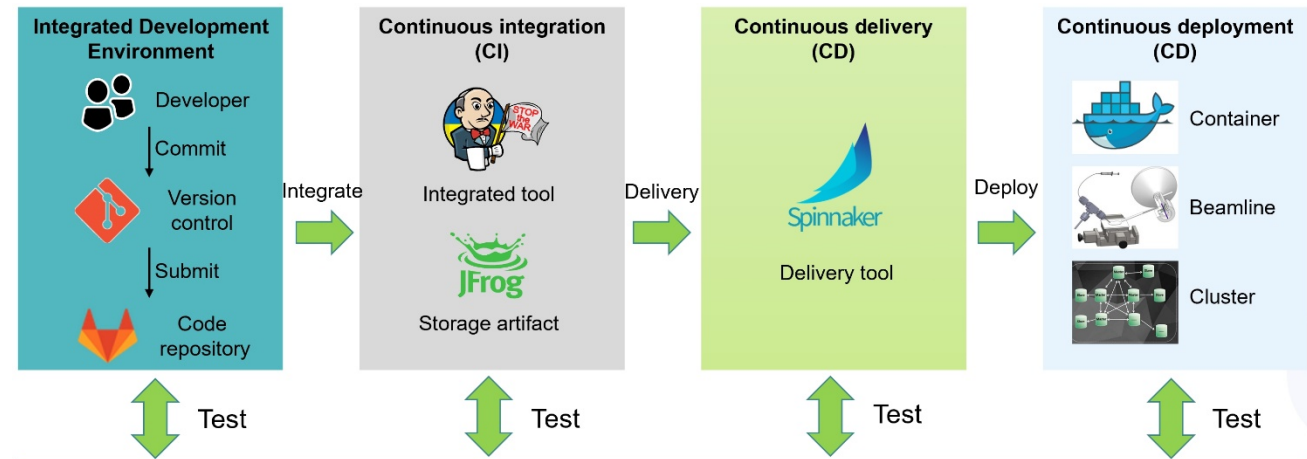
• Overview An overview of Daisy project

<https://daisy.ihep.ac.cn/>

# Support for developers and users

## Continuous integration, delivery, deployment (CI/CD) for software development

- Automated pipeline for software integration, building, test, delivery and deploy
- Continuous monitoring for each stage throughout the lifecycle of software
- Continuous testing to ensure the quality of the codes
- Enables incremental code changes from developers to be delivered quickly and reliably to production
- Based on gitlab, Jenkins, Pytest, PyUnit, and Allure. Some modules are already in production





# Application of Daisy in space astronomy

## Possible application scenarios

- Data processing, analysis, and product generation
- Detector simulation, observation simulation
- Integrate existing software resources to form common software packages

## Web based HXMT data processing platform

- Based on Jupyterlab, Docker, K8s
- Provide data processing environment and services via web browser

## Svom and eXTP data processing software based on Daisy are also in processing

- Integrated the I/O algorithms of fits files
- Some data product generation algorithms have been integrated into the Daisy framework

## Plan to support the new observation plan of HXMT with new algorithm and workflow

1. Search for target data.

(You can search directly for the object name such "Crab", or click on "more search condition" for a more complex search by source position and/or observation time)

Search by object name:  
Object Name: Crab e.g., Cyg\_X-1 Crab (Note use the underline to replace the space in the name. Query more object name, please click: HXMT object name)

... and/or search by date:  
Begin date: 年 / 月 / 日 ... End date: 年 / 月 / 日 ...

... and/or search by coordinates:  
Coordinates: ra: ... dec: ... (degrees, J2000, e.g. 83.633, 22.013)  
Search Radius: 0 (degrees)

... and/or search by proposal:  
Proposal type: ... Proposal number: ... (e.g. P0101299)

... and/or search by timestamp:(Unix timestamp differences between the selected datetime and 2012-01-01T00:00:00, e.g. 179038244, MET, TT time, A Date/Time Conversion Utility)

Start time(MET or MJD): ... Stop time(MET or MJD): ...

Search: Simply search condition Correction clear Duration sort: Default

Select the Observation ID and Exposure ID from the search result:  
Obs. ID: P0101299001 Exp. ID: P010129900101

The selected data files are in the directory: /sdcf/hepdata/AstroHXMT/1/A01/P0101299001/P010129900101-20170827-01-01.E

Please go to the "2.Data processing" tab to process the data!

Search result:  
Search criteria:  
Object name: Crab

index	obsid	tstart	tstop	obsDate	obsEnd	duration	targetid	pi	proposal	target	ra	dec
0	P0101299001	178430732	178517873	2017-08-27T04:05:29	2017-08-28T04:17:50	87141	T023		HXMT...	Crab	83.633	22.0145
1	P0101299002	178797820	179027741	2017-08-31T10:00:17	2017-09-03T01:55:38	230121	T023		HXMT...	Crab	83.633	22.0145
2	P0101299003	179038244	179065022	2017-09-03T04:50:41	2017-09-03T17:50:19	46718	T023		HXMT...	Crab	83.633	22.0145
3	P0101299004	179066854	179144896	2017-09-03T20:46:41	2017-09-04T02:36:03	46364	T023		HXMT...	Crab	83.633	22.0145

<https://sdccompute.ihep.ac.cn/>



# HEPS CC system integration/Test bed/Production

Set up testbed, integrate full data lifecycle software systems to verify the system interfaces, run in the real experimental environment, move to production gradually.

1

**Oct, 2020, BSRF 1W1A**

Simple verification of the data management system

- Network bandwidth is 1Gb/s
- Beamline storage: **2TB** NAS, Dell EMC NX3240, NFS file system
- Central storage: **80TB** disk array, Lustre file system
- Metadata ingest, catalogue, data transfer, data service

2

**July, 2021, BSRF-3W1 test beamline**

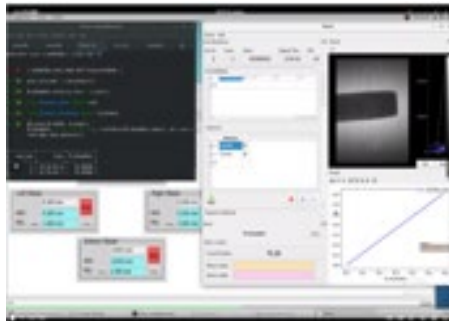
- Network bandwidth updated to 10Gb/s
- Beamline storage & Central storage: **80TB** disk array, Lustre file system
- Integrate **MAMBA, DMS, Daisy, computing system**

3

**July, 2023, BSRF 4W1B/1W1A/4W1A**

**Running in production environment**

- Network bandwidth updated to 25Gb/s
- Beamline storage: Huawei Ocean Store 9950
- Central storage: 80TB disk array, Lustre file system
- **Follow real experiment process, provide Pymca, HEPSCCT to do analyzing**



Data acquisition



Analysis framework Interface



CT reconstruction



Integration test at BSRF

# Outline

---

1. Introduction
2. Demand and Challenges of scientific data and software system
3. The architecture and design of the framework
4. The recent Progress of the system
- 5. Summary**

# Summary

---

- The system design has been finished
- Cooperation with other facilities and community is ongoing
- The basic framework has been stable and tested on the test bed
- Based on the framework, scientific software integration and application development are ongoing
- The development of scientific software ecosystem also needs the support and participation of user community

<https://daisy.ihep.ac.cn/>

# Thank you!





国家高能物理学数据中心

National HEP Science Data Center



高能所计算中心

IHEP Computing Center

**HEPS-CC**