# HEPS Data Service System

Fazhi QI

on behalf of  HEPS CC group

2019.05.30

# Overview

- 战略研讨会：<span style="color:red">希望HEPS做到科学实验过程和数据处理的自动化….</span>

- 暂时纯粹从<span style="color:red">需求</span>出发的设计，侧重于功能性设计

- 希望做到"大"IT的框架设计

  - <span style="color:red">用户视角：</span>用户实验全过程的信息化、自动化、便利化

  - <span style="color:blue">设施运行视角</span>：设施运行的数据化/数字化、尽量的自动化，高效、可靠、协同

  - IT任务视角：软硬件架构和功能模块化，功能丰富、接口及数据标准化；数据格式标准化；综合服务服务化；

# HEPS

- High Energy Phonon Source (HEPS)
- An major Infrastructure project of 13th Five-Year Plan
- The most brilliant source worldwide
- High energy, low emittance
- Investment of 5 billion RMB
- Construction started from Jul. 2019, operation expected in 2025
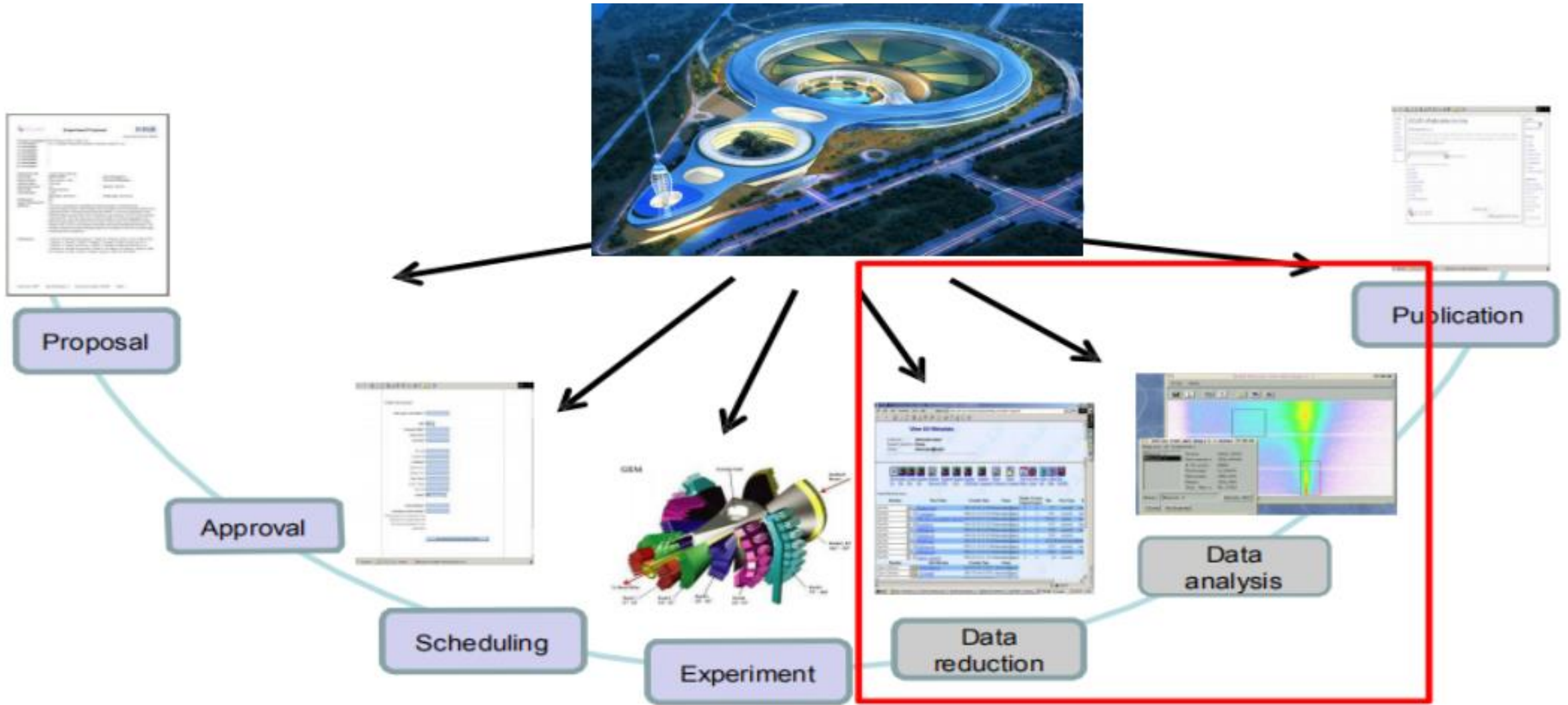



BSRF


NSRL


SSRF


TPS

**Table:** Estimated data rates of HEPS beamlines.

| Beamline | Average output (Bype/day) | MAX burst (Byte/day) |
|---|---|---|
| B1 | 820G | 4T |
| B2 | 14T | 20T |
| B3 | 112.5G | 1.4T |
| B4 | 2T | 2T |
| B5 | 6.8G | 10G |
| B6 | 20G | 50G |
| B7 | 95T | 520T |
| B8 | 10.32G | 30G |
| B9 | 5G | 10G |
| BA | 35T | 35T |
| BB | 91.08G | 165.6G |
| BC | 1T | 1T |
| BD | 275M | 500M |
| BE | 11.2T | 25T |
| TOTAL | 159.27T | 608.67T |

| Beamline | Bandwidth demanded |
|---|---|
| B1 | 40Gbps |
| B2 | 100Gbps |
| B3 | 1Gbps |
| B4 | 40Gbps |
| B5 | 1Gbps |
| B6 | 10Gbps |
| B7 | 100Gbps |
| B8 | 100Mbps |
| B9 | 100Gbps |
| BA | 10Gbps |
| BB | 100Mbps |
| BC | 10Gbps |
| BD | 1Gbps |
| BE | 40Gbps |
| TOTAL | 553.1Gbps |

| Beamline | Data format |
|---|---|
| B1 | tiff |
| B2 | hdf5 |
| B3 | tiff、SIF、dat |
| B4 | hdf5 |
| B5 | hdf5、spec、mca、tiff |
| B6 | tiff、hdf5 |
| B7 | |
| B8 | txt、mca |
| B9 | txt、tiff |
| BA | hdf5 |
| BB | tiff |
| BC | pxt |
| BD | hdf5、mca、txt |
| BE | binary和tiff |
| TOTAL | Tiff、SIF、dat、mca、pxt、hdf5、spec、binary |

Proposal

Publication

Approval

Scheduling

Experiment

Data reduction

Data analysis

# Research (user view)

- **Sample prepare**
- **Safety sheets/tests**

- **User Register**
- **Proposal**
- **Approval**
- **Scheduling**

**Synthesis**

**Prepare**

**publication**

**Experiment**

- **Scientific Database**
- **Article/Certification**

- **Data Processing**
- **Data Analysis**
- **Metadata Cat**
- **Data Storage/Mgr**

# Data lifecycle



Proposal

Experiment

Sample

Raw Data Metadata

Data Analysis

Publication

**User**

Pre-Experiments | Experiments Duration | Post-Experiments

## Facility Management System

- ✓ Facility Database
- ✓ Facility Operation & Monitoring System

## User Service System

- ✓ User Database
- ✓ User Safety Training System
- ✓ Proposal System
- ✓ Sample Management System
- ✓ Single-Sign-On
- ✓ Import proposal details from user portal: Experiments, samples details, etc.
- ✓ Create users, groups and setup permissions
- ✓ Create folders structure for storing scientific data

## Data acquisition layer

- ✓ Set of data aggregator running on cluster of high-performance computers
- ✓ Collect, monitor and store experiment data
- ✓ Disseminate data to online cache and analysis pipeline

## Run management service

- ✓ Coordinates PC layer activities related to data

## Metadata catalogue (MDC)

- ✓ Organize and keep track of experiment data in a homogeneous and consistent way

## Control Layer

- ✓ ......

## Accelerator & Beamline Optimization

- ✓ Running & Operation Log Collection and Analysis
- ✓ Feedback to Facility

## Data management online -> offline

### During measurement (run)

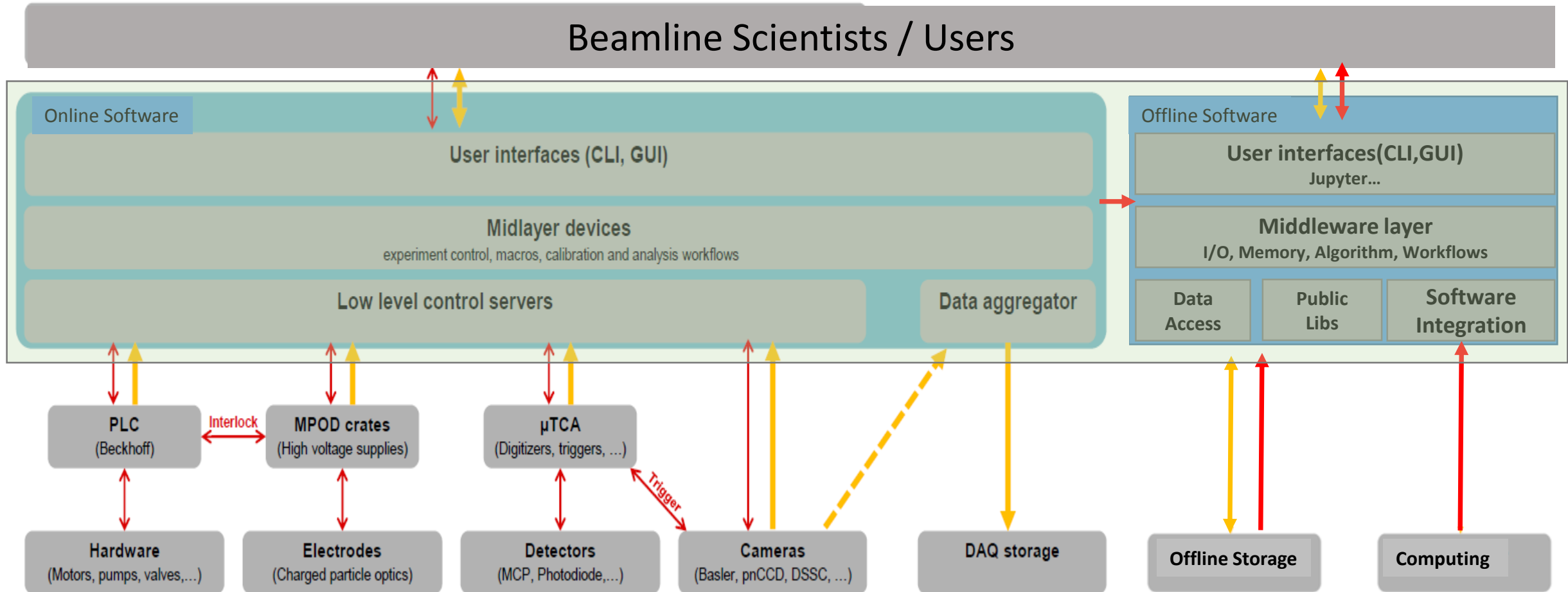- ✓ Calibrated and raw data available in hutch (GUI, online)

### Data migration after each run

- ✓ Data manager decides on quality of the data: "good", "unclear", "not interesting"  "good" and "unclear" data transferred to "Offline cluster"
- ✓ Migration triggers computation of calibrated data at online cluster
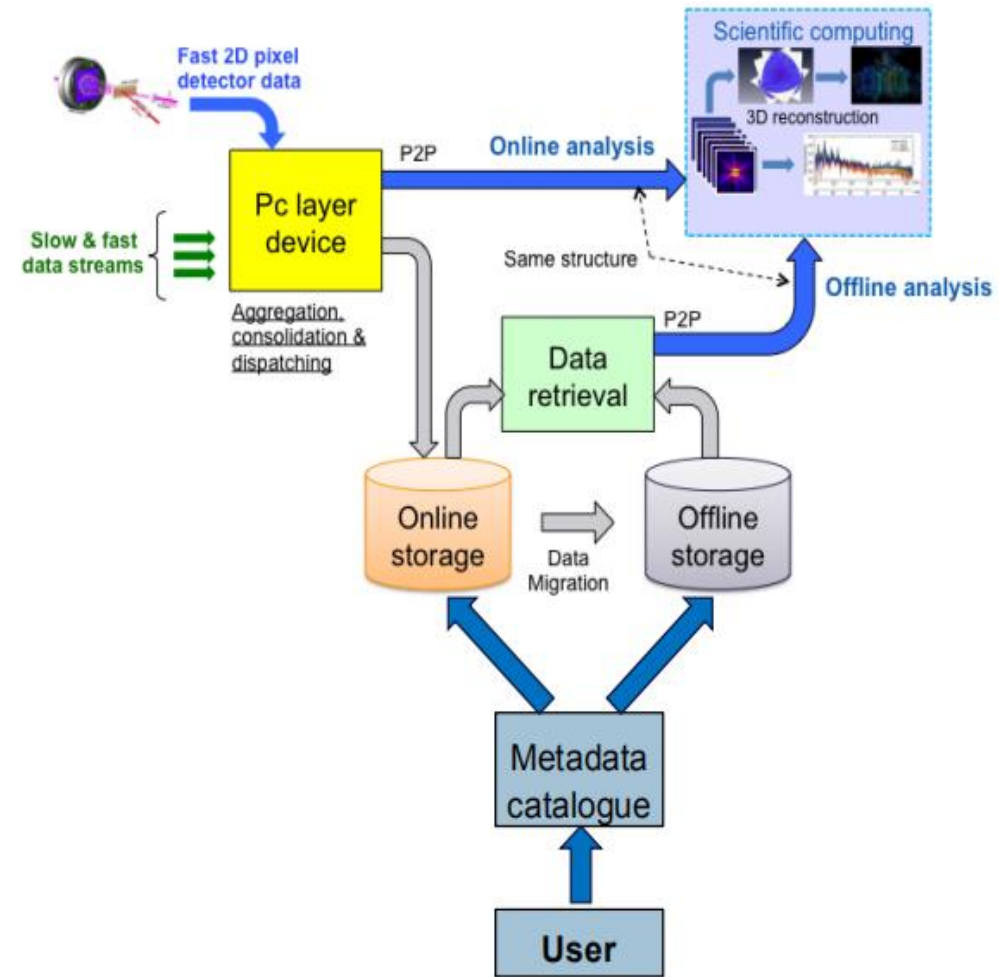
### Post experiment

- ✓ Raw and calibrated data is available
- ✓ Data catalog is available
- ✓ Data service is available for search , access , download , interface to remote analysis
- ✓ Analysis resources / "Offline cluster"

Distributed Infrastructure and Services: Storage , Computing , Database , Middleware , Software / Cyber and Data Security services
Policy and Rules

# Core Software

# HEPS Data Policy (Refer to Eur.-XFEL)

- Raw data and metadata is stored, curated, and arch[ived] by HEPS
- Access to data through searchable metadata catalo[gue]
- <span style="color:red">Raw and metadata will be open access after embar[go] period</span>
- Electronic logbook will be provided for documentat[ion]
- Access to metadata catalogue and computing infrastructure through HEPS user account
- Online analysis and offline analysis should be provid[ed]
- Related system
  - Control / DAQ / Offline / … /User Office

# Control and Data taking

**New setup/experiment**

**New run**

- ## Setup
  - Retrieve experiment and sample metadata for the given proposal
  - Receive and process run configuration
  - Dispatch this metadata to Aggregators

- ## Data taking
  - Instruct start and stop of runs

- ## Publish data
  - Build run summary and register data files

# Scientific data collection/DAQ



Configuration input to MDC

**Proposal (beam time)**

| Experiment 1 (sample 1) | Experiment 2 (sample 2) | Experiment 1 (sample 3) |
|---|---|---|

User submit/edit proposal

store
monitor
ignore

Enter data taking mode
Enter monitor state
Start run
Stop run
Start run
Stop run
Leave monitor state
Enter monitor state
Start run
Stop run
Start run
Stop run
Leave monitor state
Leave data taking mode

User performing experiment(timeline)

Conf input

| Experiment 1 (sample 1) | Experiment 2 (sample 2) | Experiment 1 (sample 3) |
|---|---|---|

Runtime input to MDC

| Run 1 | Run 2 | Run 3 | Run 4 | Run 5 | Run 6 | Run 7 |
|---|---|---|---|---|---|---|
| N data files | N data files | N data files | N data files | N data files | N data files | N data files |

Logical structure of experiment and runs

# Data aggregators

- Functions
  - Correctly receiving data from fast and slow sources
  - Consolidating data streams from all data sources
  - Monitoring, fast-feedback analysis, rejection and reduction
  - Formatting and recording to HDF5 data files
  - Data formats for FPGA binary, Control and Instrument data

- Inputs and outputs
  - DAQ/ Control/ User system
  - Rawdata / Metadata catalogue



Data sources

Data aggregators

Metadata catalogue

Online storage

# DAQ & Online Analysis

- FPGA / CPU / GPU based
- Distributed Resources
  - Storage and computing
  - Locates in beamline
  - Central management for infrastructure and public software
  - User defined application software
- Central Resources
  - Locates in Computer-Center
  - Virtualization for beamlines
  - Central management for infrastructure and public software
  - User defined application software
- Interfaces
  - Common data interface with analysis algorithms
- Rapid feedback through GUI
  - Scientific data
  - Control data
  - Facility operation logs
  - ……

# Data Services
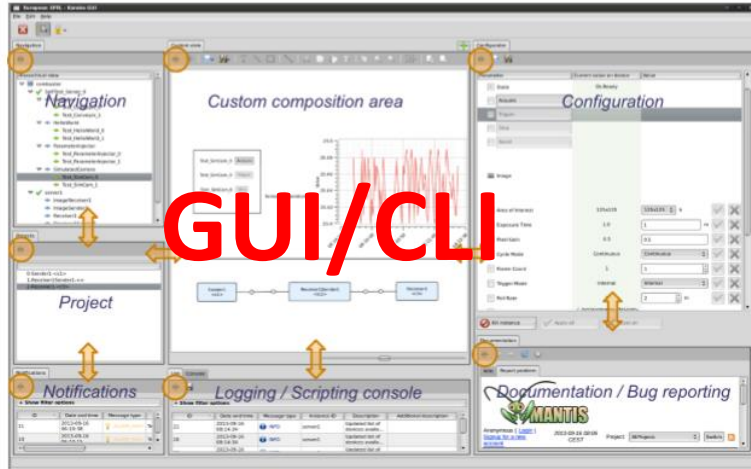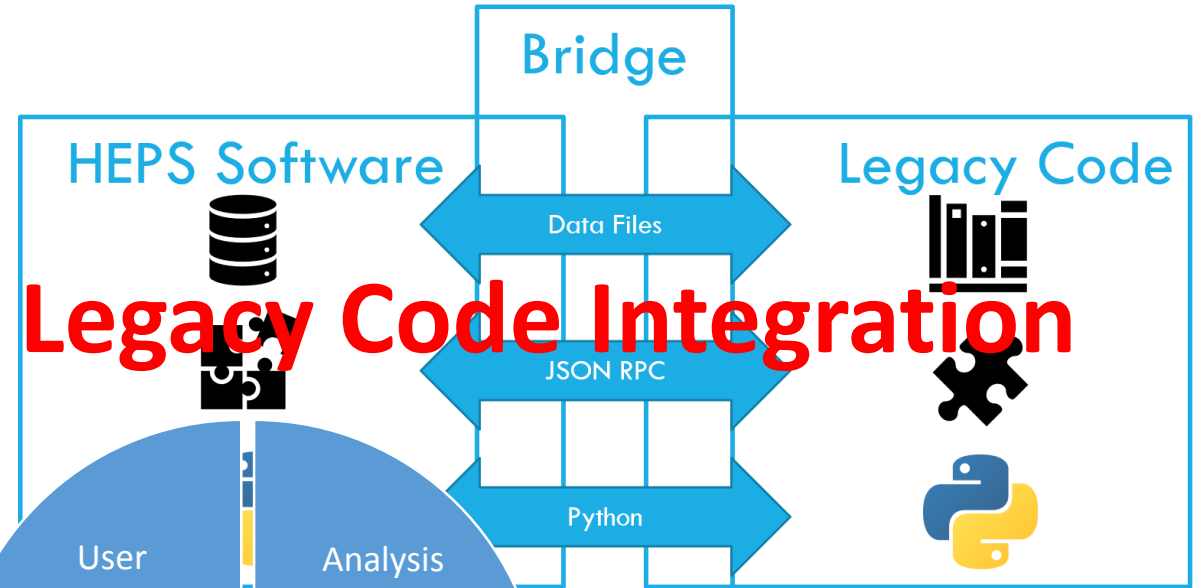
# Key Technologies and Services

# Software Framework



GUI/CLI

Bridge

HEPS Software

Legacy Code

Data Files

Legacy Code Integration

JSON RPC

Python

User Interface

Analysis Tools

SNiPER For HEPS

Distribution System

| Data Analysis Module/Tool | Data Analysis Module/Tool | ... | Data Analysis Module/Tool | Data Analysis Module/Tool |

Modularized Extensible Framework    Kernel    Modularized Extensible Framework

Near Real-time Feedback    DAQ Data Stream    Data Files    User Interface

**Online**    **Offline**

MQ & RPC

GUI/CLI Master

Web App

Event loop

Event loop

Event loop

DAQ

Work flow

# Data management

# Raw Data Management / Raw Data Format

- Raw data: data files collected after each run

- Data format
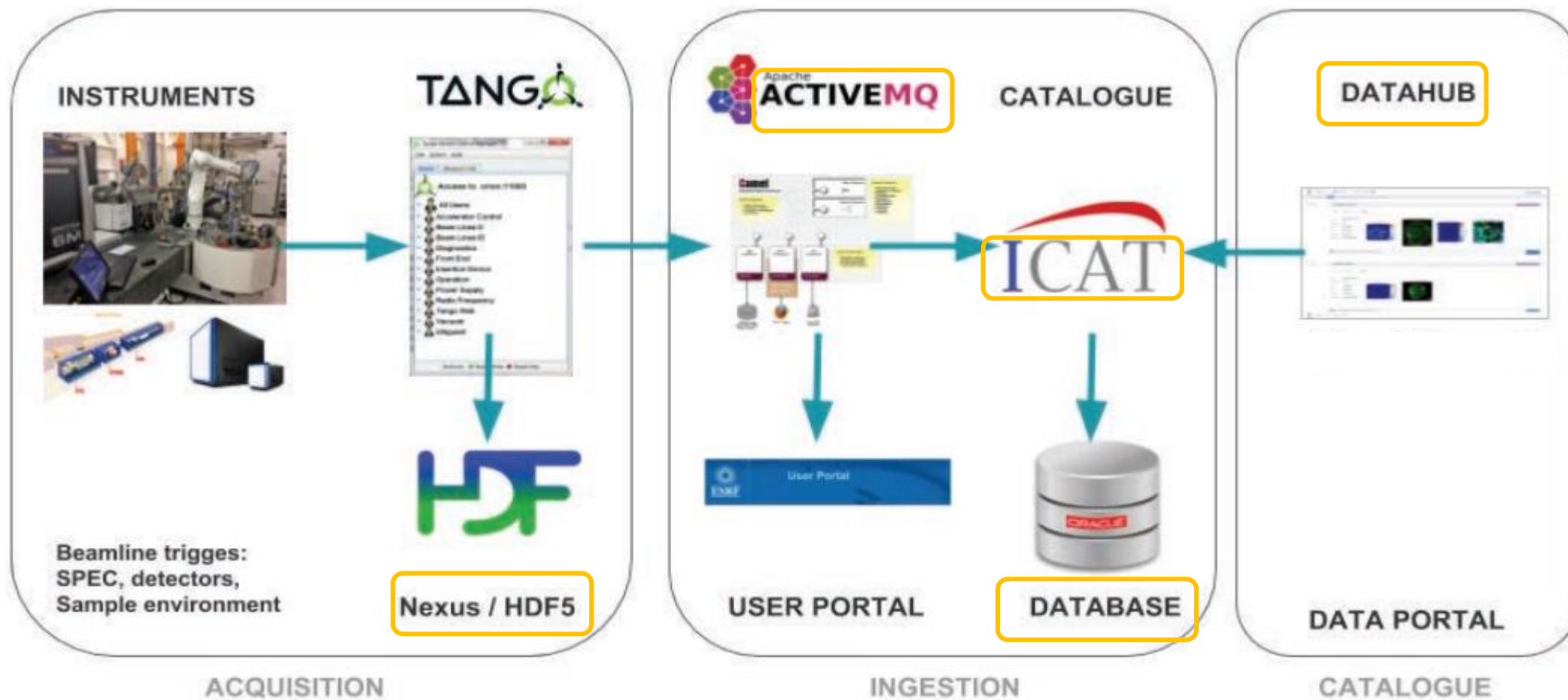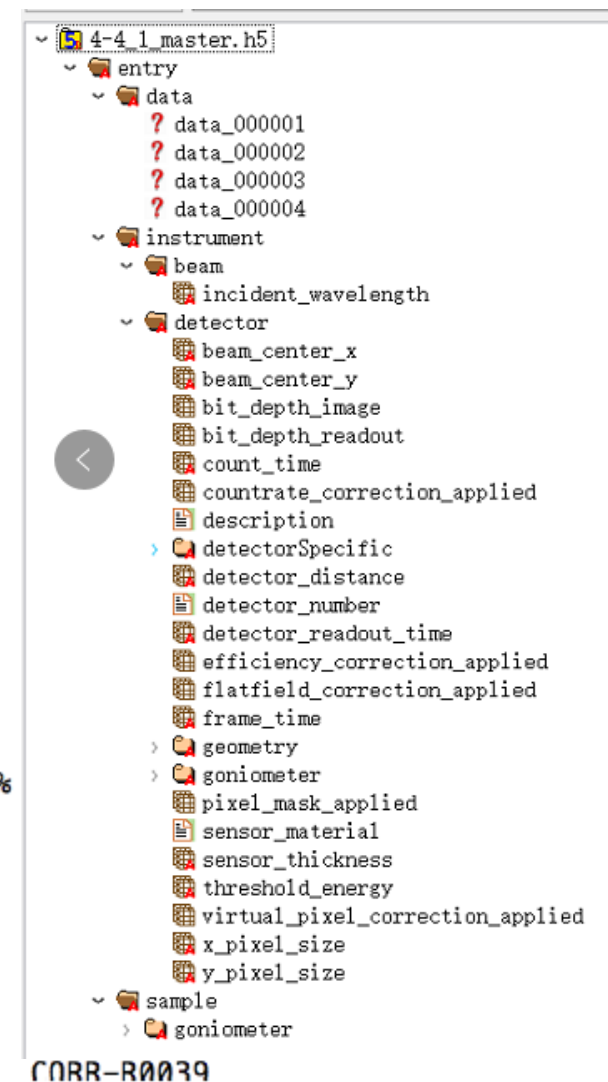
  - HDF5: a file format for managing any kind of data

- Data format standardization

  - convert other data format (TIFF、SIF、MCA、TXT…) to HDF5

- Raw data repository:

  - File system storage

  - TIFF、SIF、MCA、TXT

```
[haufs@max-exfl014]/gpfs/exfel/exp/SPB/201701/p002038/proc%
CORR-R0039-AGIPD00-S00000.h5    CORR-R0039-AGIPD05-S00002.h5
CORR-R0039-AGIPD00-S00001.h5    CORR-R0039-AGIPD05-S00003.h5
CORR-R0039-AGIPD00-S00002.h5    CORR-R0039-AGIPD06-S00000.h5
CORR-R0039-AGIPD00-S00003.h5    CORR-R0039-AGIPD06-S00001.h5
CORR-R0039-AGIPD01-S00000.h5    CORR-R0039-AGIPD06-S00002.h5
CORR-R0039-AGIPD01-S00001.h5    CORR-R0039-AGIPD06-S00003.h5
CORR-R0039-AGIPD01-S00002.h5    CORR-R0039-AGIPD07-S00000.h5
CORR-R0039-AGIPD01-S00003.h5    CORR-R0039-AGIPD07-S00001.h5
```

```
4-4_1_master.h5
  entry
    data
      data_000001
      data_000002
      data_000003
      data_000004
    instrument
      beam
        incident_wavelength
      detector
        beam_center_x
        beam_center_y
        bit_depth_image
        bit_depth_readout
        count_time
        countrate_correction_applied
        description
        detectorSpecific
        detector_distance
        detector_number
        detector_readout_time
        efficiency_correction_applied
        flatfield_correction_applied
        frame_time
        geometry
        goniometer
        pixel_mask_applied
        sensor_material
        sensor_thickness
        threshold_energy
        virtual_pixel_correction_applied
        x_pixel_size
        y_pixel_size
    sample
      goniometer
CORR-R0039
```

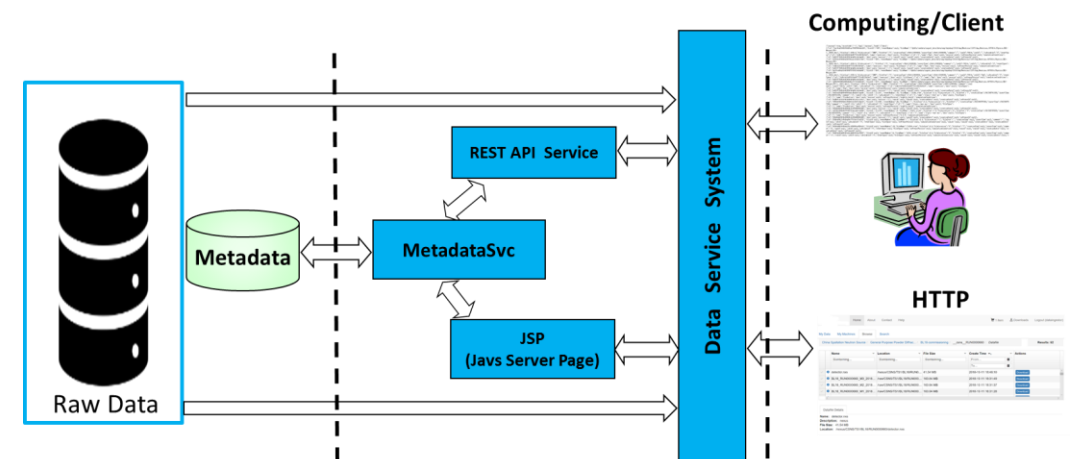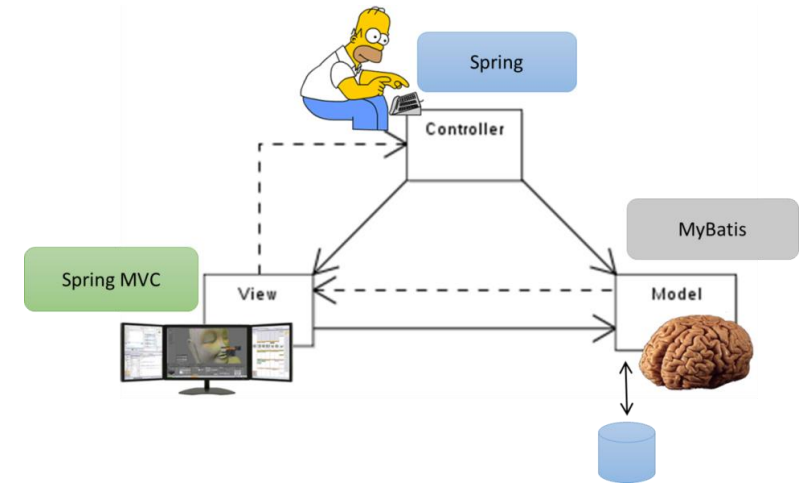# Scientific Aggregators and Converter Module

- Aggregation
  - raw data file + scientific metadata
  - Scientific metadata: sample, beamline and experiment parameters
- Conversion
  - to HDF5

- Control +  DAQ + Offline

# Metadata Management/ Data catalogue

- Meta data
  - administrative : data management lifecycle, ownership, filecatalog
  - scientific: describing the sample, beamline and experiment parameters relevant for the users data analysis
- Data model:
  - define common generic (fixed) meta data
  - allow completely flexible and rich structured (beamline, instrument specific) scientific meta data
- Storage
  - MongoDB backend
  - NoSQL -no fixed structure and common data format
  - Map/Reduce queries and the option to use JavaScript as the query language
  - Powerful indexing and support for file storage
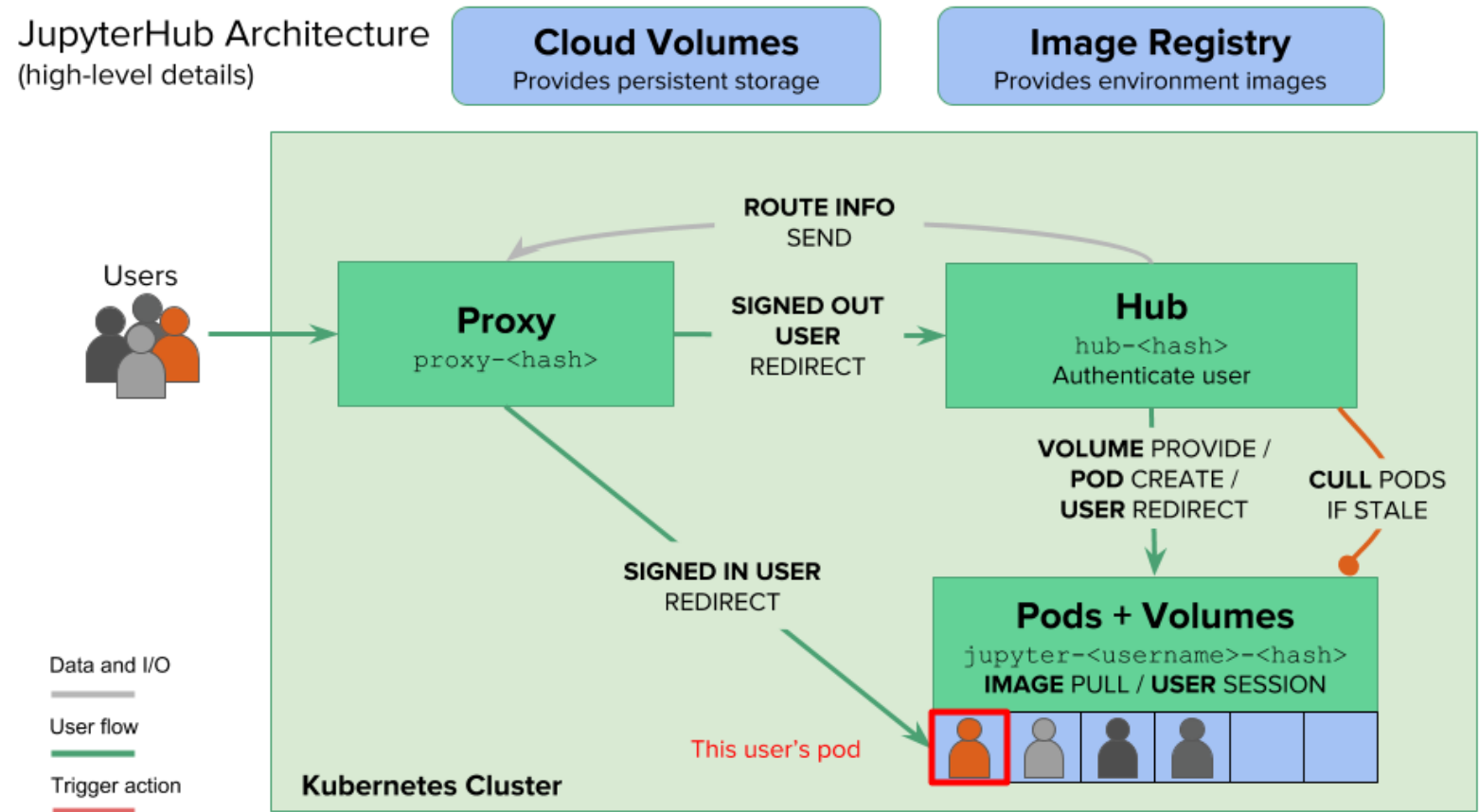  - Fault tolerant and drivers for most languages

# Data Service System

- Web based frontend
  - Find data based on the meta data
  - Download data based on  data permissions
  - Data can be linked to proposals and samples
  - Data can be linked to publications
  - Helps keeping track of data provenance
- API Services
  - Computing: Local & Cloud
  - Client

# Software –tools platform

- Gitlab/GitHub
- Matlab
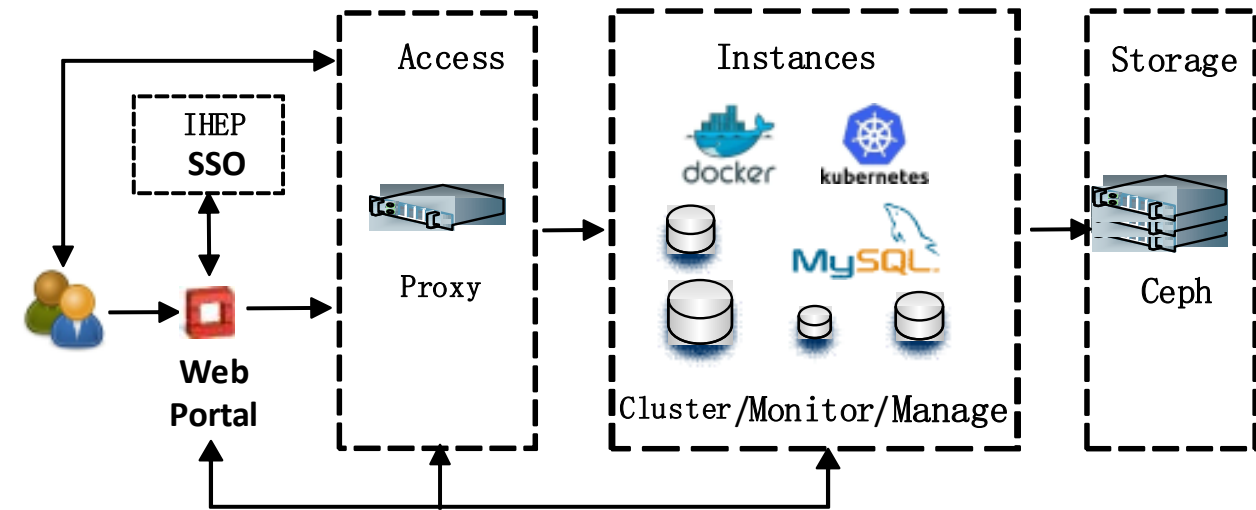- Ansys
- JuputerHub
- ……

# Database Services System

- Services for Database
  - Public databases
  - Facility Database
  - Operation Database
  - Scientific Database

- High performance and stability

- RD and NoSQL

No access to underlying hardware
No DBA support or application support

# Facility Database

- 独立的固定资产管理
  - 采购、报销、入库、位置、状态
  - 属性管理
- 资产关系管理（Relationship Management)
  - 父-子关系
  - 兄-弟关系
  - 树状结构
  - 举例：一个设备由多个部件构成，资产管理数据库能够反映出其部件组成关系及其关联的财、人关系

# Facility Lifecycle Management

- 采购—报销—入库—上线—运行（巡检）—维修（升级）--下线—退库—报废

统计报表

灵活定制可视化图形报表，一键生成统计分析；

资产使用状态报表、资产分类报表、资产走势图；

多维度生成资产明细分析、汇总报表，支持按条件导出；

全生命周期管理

跟踪资产启用直至资产退出，对整个资产全生命周期进行全程跟踪管理；

为管理者提供详实的资产数据汇总，可视化图形报表，管理变得更轻松；

入库
领用
退库
调拨
借用归还
维修
处置
盘点
全生命周期管理

资产盘点

盘点平台多样化，支持微信公众号、安卓平台、IOS平台以及WINDOWS工业平台采集设备和手机端APP；

可通过盘点计划、盘点批次、个人核查等多方式建立盘点任务；

盘点数据自动提交，快速自动生成盘点报表，实现多人同时参与盘点，高效提升工作效率；

采集设备　　安卓平台　　IOS平台

日常管理

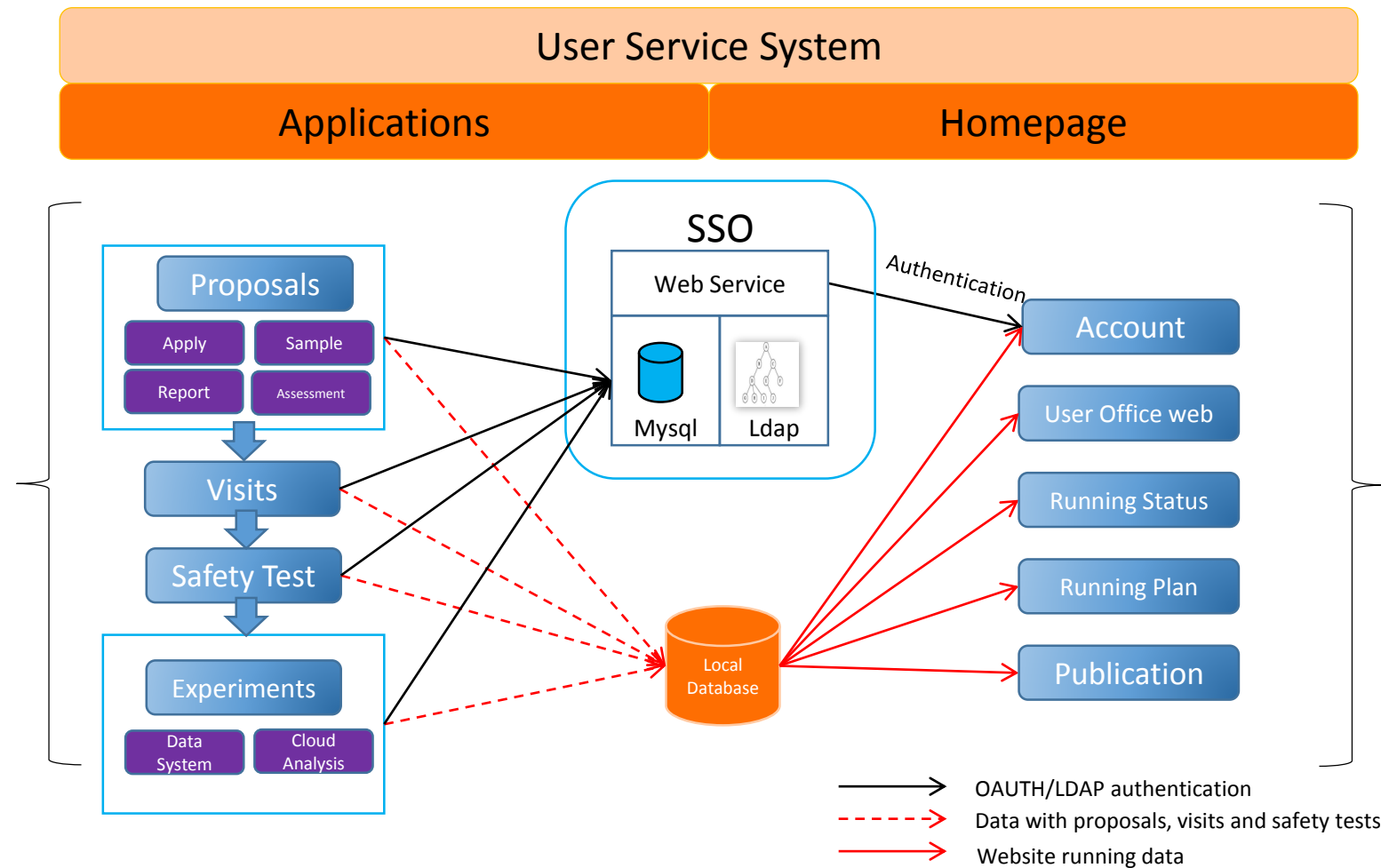资产入库、盘点、领用、变更、维修、调拨、报废清理等全生命周期管理；

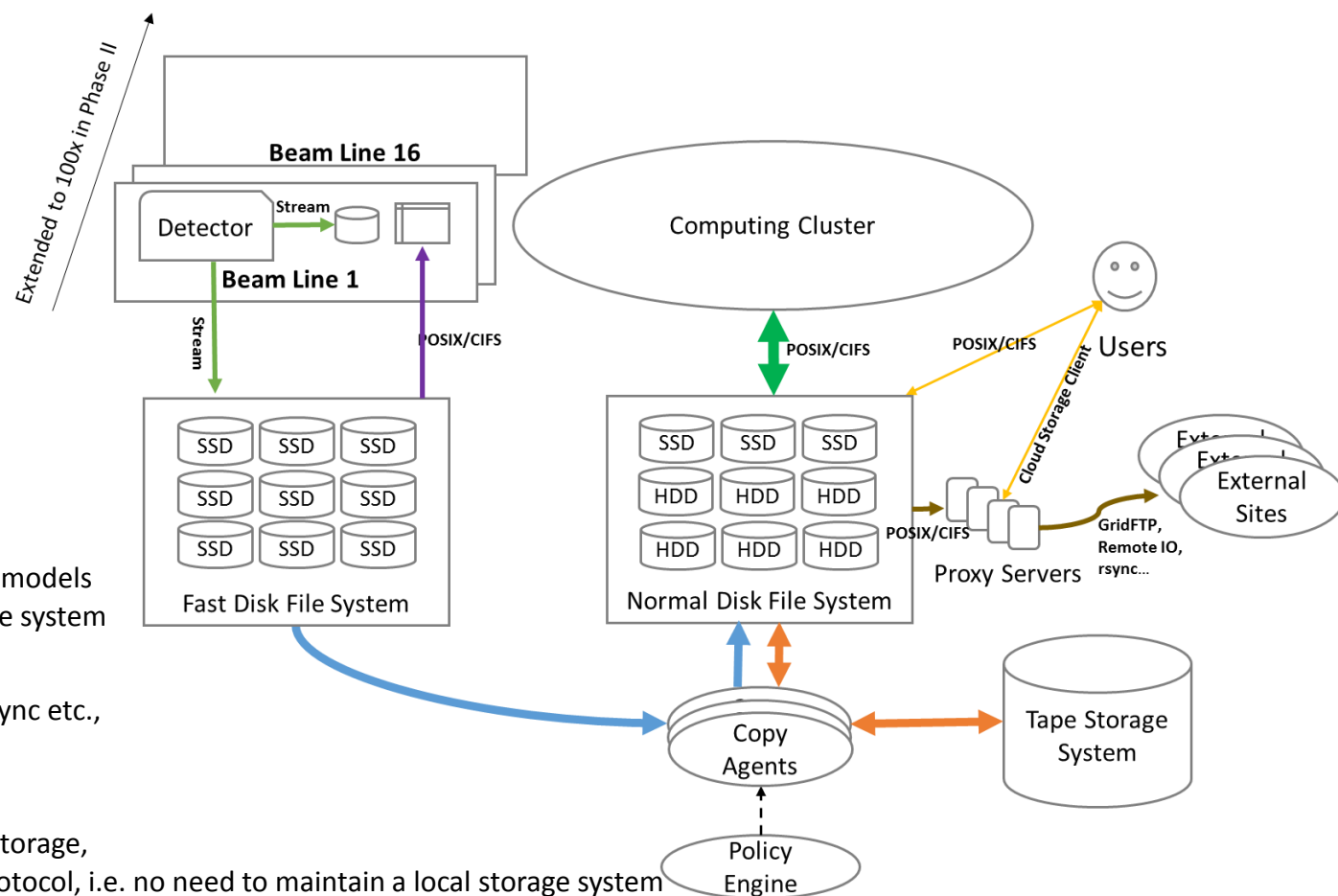告别传统的纸记手抄和Excel模式，全面进入无纸化资产信息管理方式；

# User & Proposal Service System

- **One-stop Service System**
  - Online experiment solution for scientists and students
  - One IHEP SSO Account for all services

- **System Function**
  - Proposals Application
  - Experimental risk assessment
  - Expert review
  - Safety Tests
  - Visits Arrangements
  - Running status statistics

# Distributed / Central Storage Resources

- Fast Disk File System (managed by Lustre file system)
  - for fast on-line analysis,10GB/s write, 20GB/s read
  - Sequential read and write ,1 day data life,160TB Capacity
- Normal Disk Storage (managed by Lustre file system)
  - for batch jobs, interactive analysis, backups and data transfers
  - Mixed I/O patterns:  sequential read and write, random read
  - 10GB/s write, 50 GB/s read, 6 month data life, 30 PB  Capacity
- Tape Storage (managed by CERN CASTOR )
  - for data preservations and backups
  - Sequential Read and Write
  - 5 year  data life, 300 PB Capacity
- Policy engine (developed by IHEP-CC )
  - make decision of data copies by  static rules or prediction of AI models
  - scheduling and monitoring data copy  requests between storage system
- Data transferring System (SPADE)
  - Transfers data to remote site before data analysis by gridFTP, rsync etc.,
  -  multi-stream, presumable, schedulable transfers
- Data Federation (developed by IHEP-CC )
  - Users at remote sites can see an identical namespace of HEPs storage,
  - Remote sites can access data at HEPS storage by remote I/O protocol, i.e. no need to maintain a local storage system
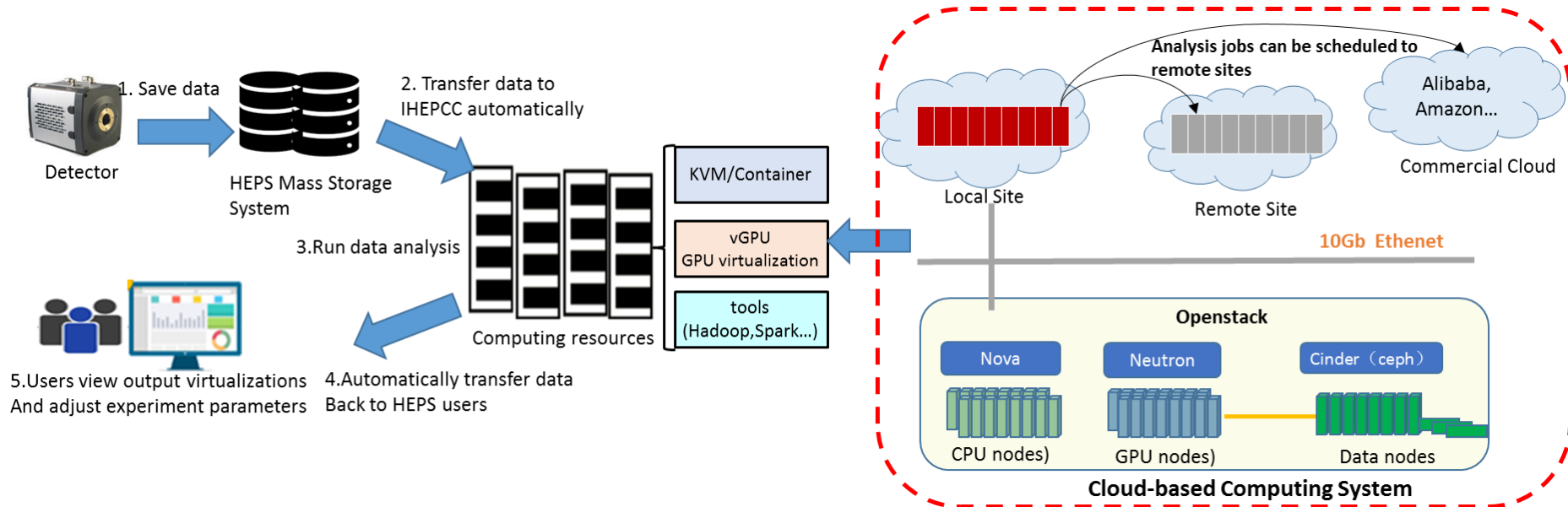
# Distributed / Central Computing Resources

- Cloud-based computing system
  - Deploy the various computing platform in cloud like batch farm, HPC farm and also develop with a rich collection of tools to work for Daas
  - Integrate distributed computing resources as well as commercial clouds to expand computing scale based Openstack+HTCondor
  - Using GPU virtualization technology to provide the HPC service in cloud
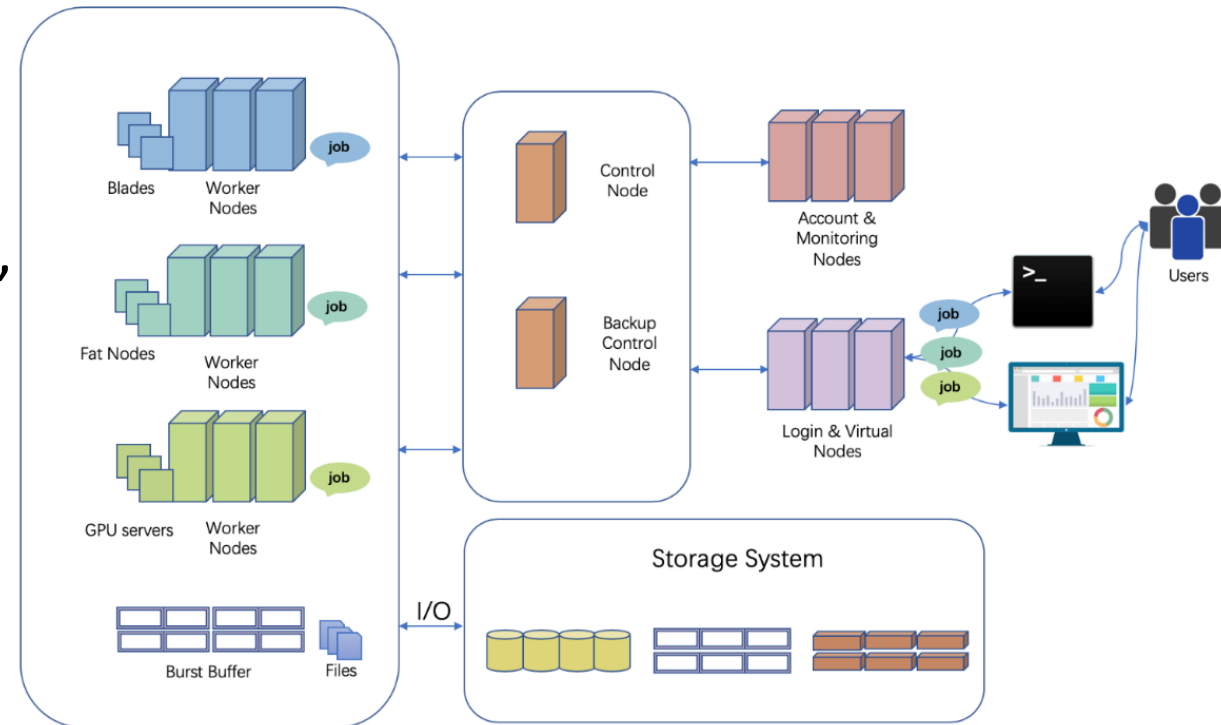  - Kubernetes+Docker to provide virtualized resources for WEB-based analysis
- Requirements and Design
  - Online data analysis(Streaming data services):developed a rich collection of tools with the Hadoop and YARN including Spark for real-time analyzing streaming data after data acquisition
  - Offline Data analytics: as implemented with the HEPS software to provide Data analysis as a service for users

# High Performance Computing

- Flexible and elastic scale expansion design (designed by IHEPCC)

- Convenient and unified user job interface (developed by IHEPCC)

- Heterogeneous resources in cluster including CPU nodes, GPU nodes, fat nodes, with rich scheduling policies (Managed by Slurm)

- Various Operating systems and softwares provided to the job based container technology and transparent to user (developed by IHEPCC)

- Burst storage integrated with cluster providing high IO throughput during the job life time (Integrated by IHEPCC)
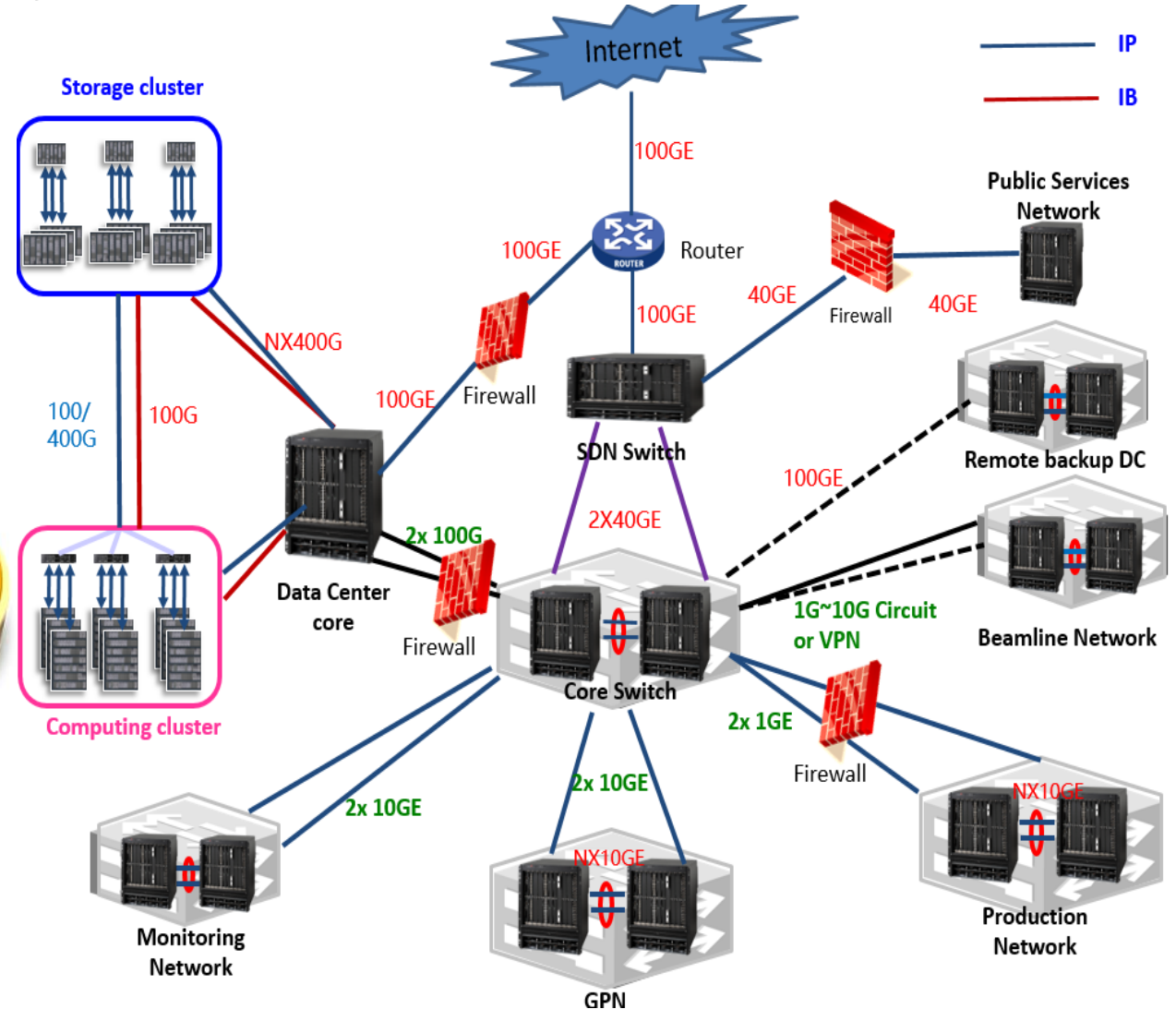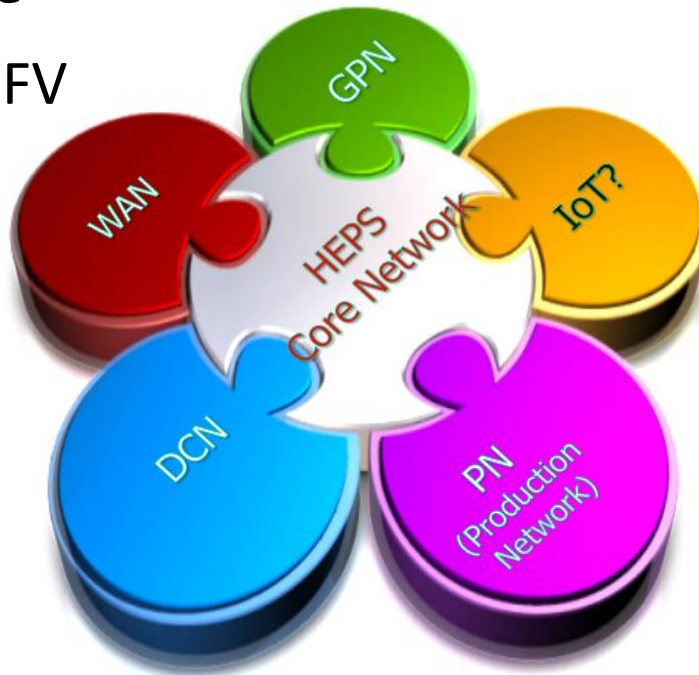
# ML Platform

- For Scientific Data Anlaysis
- For Facility Operation Log Analysis
- ……

# Network and Security

- Demand driven network

  - High Performance

  - Intelligent

  - Flexible

  - SDN/NFV

  - IoT

  - 5G

  - ......

# Machine Room

- To be a **GREEN** machine room
- The floor space: $600\,\mathrm{m^2}$ (15 beamlines → 90 beamlines )
    - 100 racks（30 racks for phase I）
    - High-density racks ( Comping resources ) and General racks ( Storage, Network...resources )

- Supporting space: $200\,\mathrm{m^2}$
    - Power
    - Fire Protection
    - ...
- Power Capacity：3300kW

# 计划

- 近期 2019年内
  - 继续深入理解实验及数据处理自动化业务流程: BSRF, SSRF, EXFEL, ESRF, Diamond,PSI ……，HEPS
  - 加入国际类似装置的联盟，发起国内光源类设施相关工作的联盟，开展合作及成果共享
  - 数据管理、数据格式转换、<span style="color:red">软件框架</span>、资源组织方式等设计思路

- 中期 2019-，HEPS光束线建成前
  - 在BSRF上选取若干数线站（如漫散射站、成像站、XAFS、生物大分子站等）对如下工作进行设计、开发、测试
    - 元数据标准制定、抽取、元数据库建设
    - 数据本地存储、数据中心存储、（准）实时数据传输系统搭建
    - 本地数据处理集群建设、云计算处理
    - 软件框架。。。
    - 接口对接、集成等
  - 以上工作如下的调研同步进行（并根据需要进行更新）
    - 线站对数据存储、计算等的要求
    - 线站对软件框架的要求
    - 其它

- 远期
  - 在HEPS线站进行建设、部署，并根据实际情况进行调整
  - 提供实验数据自动化服务

# Reference

- Reference:
  - Eur. XFEL , CSNS , ESRF , ......
- Test-bed
  - BSRF
  - SSRF
  - ...
- Cooperation ?
  - HEPS
  - SHINE
  - SSRF
  - CSNS
  - HSRF
  - ...



Big Data Science at SSRF - A Next Generation Superfacility

Prof. Alessandro Sepe
Head of the Big Data Science Center



**Control, data acquisition, management and analysis**

Thomas M. Baumann
Scientific Instrument SQS

Hans Fangohr
Control & Analysis Software

SQS Early User Workshop
Schenefeld, 12.02.2018

European XFEL



大科学装置数据科学研讨会
2018.10.21 上海科技大学

**Overview of SHINE and Its Data System**

**Ping HUAI （怀平）**
Beamline and Endstations SHINE,
ShanghaiTech University



中国散裂中子源
21 Oct, 2018  Shanghai

Software Architecture of
Data Analysis & Management for CSNS

Junrong ZHANG, Ming TANG, Yakang LI, Fazhi QI
Institute of High Energy Physics

CSNS
CHINESE ACADEMY OF SCIENCES

**Data Acquisition and Management Services at European XFEL**

Djelloul Boukhelef

for ITDM

Karabo 2.0 workshop

DESY, Hamburg 24.01.2017

European XFEL

# 总结及讨论

- HEPS是一个用户类设施...软件和服务层面大量工作需要做，希望能借鉴高能所、高能物理领域的成果和合作模式
- 软件框架
  - 在线分析、离线分析、数据管理？甚至包括控制？
  - 是否有可能统一框架？
- 接口
  - 科学数据接口：仅仅同控制系统接口？同控制系统/DAQ分别有接口？
- 合作
  - 用户类设施 + 高能物理实验？共同性、区别？
  - 人员缺口严重
  - 推荐人才：职工、研究生、博士后

谢谢

# 深化设计报告

- 每一部分可以独立为一个设计报告，能够基于该报告开展合作交流；其它系统、其它设施、协作研发
- 每一部分需要描述出同外部相关系统和部分的接口
- 整合起来之后是一个完整的HEPS IT深化设计方案
- 图文并茂、既有总体设计、也有细节描述，具有可操作性
- 占到用户视角（线站、科学家、用户）和IT提供者视角描述提供的服务
- 尽早发布、寻求合作
  - 评审
  - 得到认可，讨论、合作
  - 合作的原则：共性合作、不共性的共享思路，

# 深化设计报告

- 总体设计报告(Overview)
  - 目标、理念、设计思路、包括的内容
  - 实施计划

- 基础设施
  - 机房
  - 网络：总体设计思路和架构，描述出几张不同的网络之间的关系、关键技术指标、互相之间的接口、各自的功能和承载业务，涉及到的关键技术（甚至包括需要预研验证的）；面向别的系统
  - 存储：总体设计思路和架构，根据业务需求介绍资源组织方式（逻辑集中，物理上可以集中，也可以分离，根据具体使用场景，如线站缓存？数据中心集中存储？长期保存？）和服务场景、性能指标；分别占到用户的角度和IT设施提供者的角度描述资源服务流程；涉及到的关键技术，涉及到的关键技术（甚至包括需要预研验证的）
  - 计算：总体设计思路和架构，根据业务需求介绍资源组织方式（逻辑集中，物理上可以集中，也可以分离，根据具体使用场景）和服务场景、性能指标；分别占到用户的角度和IT设施提供者的角度描述资源服务流程；涉及到的关键技术（甚至包括需要预研验证的）
  - 设施统一安装与管理：总体方案、实现技术、服务场景、关键指标
  - 安全：

- <span style="color:red">科学数据获取：</span>

- 科学数据管理
  - 总体描述科研数据生命周期，每一个过程中数据管理的功能、输入、输出、接口
  - 科学数据服务portal及其后端关联系统及技术
  - 软件实现的思路及每一部分的具体实现方案

- 科学<span style="color:red">数据获取与分析</span>（软件框架及软件集成）
  - 总体描述需求，围绕数据生命周期，数据分析的分类和特点，和数据管理的关系，和其它系统的关系
  - 在线快速分析、离线分析，软件的功能，对计算和存储资源的特性的需求
  - 软件框架与科学软件之间的关系和业务逻辑，不同视角的描述
  - 软件实现的思路及每一部分的实现方案

- 设施运行支撑软件
  - 总体描述需求，围绕实验生命周期，包含的支撑软件
  - 用户管理、权限管理、提案管理、实验过程管理、……
  - 设施资产管理、设施运行数据管理与分析、…
  - 软件实现的思路及每一部分实现方案

# 深化设计报告

- 公共支撑平台
  - 关系型数据库
  - 非关系型数据库
  - 公共软件库管理及服务平台
  - 代码管理平台
  - 智能化运行监控中心（包括运行和安全）
- 预研、测试及合作计划
  - 每一个工作均要涉及
  - 预研目标、指标、计划的线站
  - 合作的单位、要求