# Data interoperability in China: practice and challenges

Lianglin Hu (HULL@cnic.cn)

CNIC.CAS/NBSDC.CN/GOSC DataIO WG

2022.10.10

# Abbreviated terms
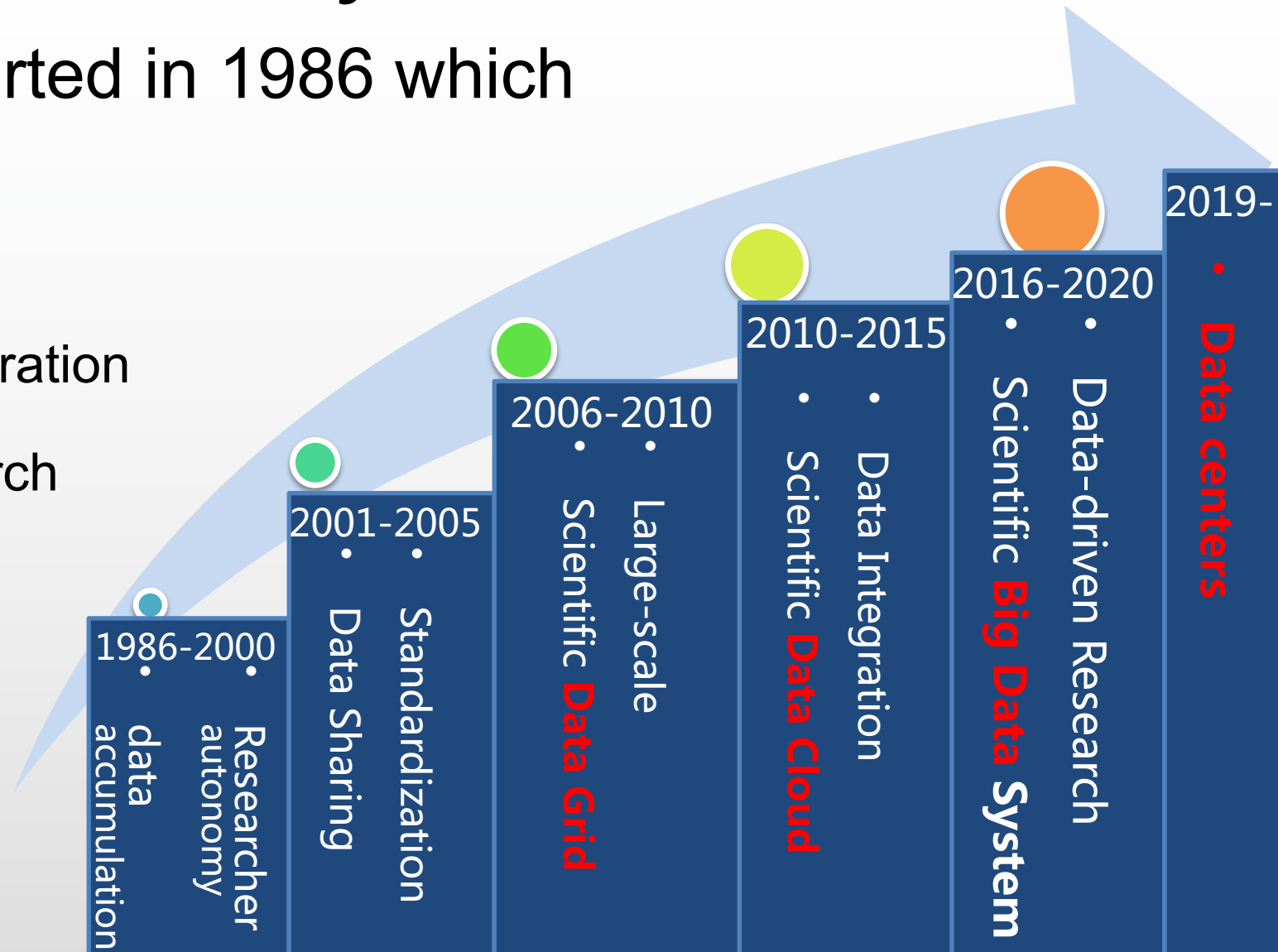
| Abbreviated term | All words term |
|---|---|
| CAAS | Chinese Academy of Agricultural Sciences |
| CAS | Chinese Academy of Sciences |
| CAS Project | Scientific Data Project in CAS, started in 1986 |
| CNIC | Computer Network Information Center, CAS |
| CSTCloud | China Science and Technology Cloud, building and running by CNIC |
| DataIO | Data interoperability |
| EIF | European Interoperability Framework |
| GOSC | Global Open Science Cloud |
| MOST | Ministry of Science and Technology, China |
| MOST Project | National Scientific Data Sharing Project, started in 2001 by MOST |
| National Measure | Measures for Scientific Data Management, China |
| NBSDC | National Basic Science Data Center, China |

# Contents

- ➢ Two projects

- ➢ DataIO practice

- ➢ Challenges

# Scientific Data Project in CAS

- A Long-term mission started in 1986 which funded by CAS

  - 60+ institutes involved

  - long-term, large-scale collaboration

  - data from research, for research

**1986-2000**
- data accumulation
- Researcher autonomy

**2001-2005**
- Data Sharing
- Standardization

**2006-2010**
- Scientific **Data Grid**
- Large-scale

**2010-2015**
- Scientific **Data Cloud**
- Data Integration

**2016-2020**
- Scientific **Big Data System**
- Data-driven Research

**2019-**
- **Data centers**

# National Scientific Data Sharing Project

- **China has long attached importance to the management and sharing of scientific data**

  - In 2001, "implement the scientific data sharing project and strengthen the national science and technology innovation capability" was put forward. In December, the Ministry of Science and Technology launched the first scientific data sharing pilot for meteorological science data sharing.

  - In December 2002, the Ministry of Science and Technology launched 8 other shared pilots: agriculture, forestry, hydrology and water resources, earthquakes, surveying and mapping, earth system science, sustainable development, and rural modern science and technology information.

  - In 2003, with the support of the Ministry of Finance, the Ministry of Science and Technology set up the National Science and Technology Infrastructure (NSTI). The scientific data sharing project was included as an important component in the construction of the NSTI.

  - In 2009, the first batch of scientific data sharing projects were accepted and transferred to the operational service phase. So far, eight scientific data sharing platforms have been identified by the NSTI.

  - In 2017, the "13th Five-Year National Science and Technology Innovation Base and Special Conditions for Capacity Building" integrated the Science Data Centers into the National Science and Technology Innovation Base of Basic and Conditional Support.

# Summary

1. **Both are long-run project**

   - CAS project - 36 years
   - MOST project - 21 years

2. **Both are national level indeed**

   - CAS project was approved by the State Planning Commission in 1986

3. **Both support the best scientific data working teams in China**

   - Some institutions of CAS project are also participating body of MOST project
   - MOST project includes a number of institutions out of CAS
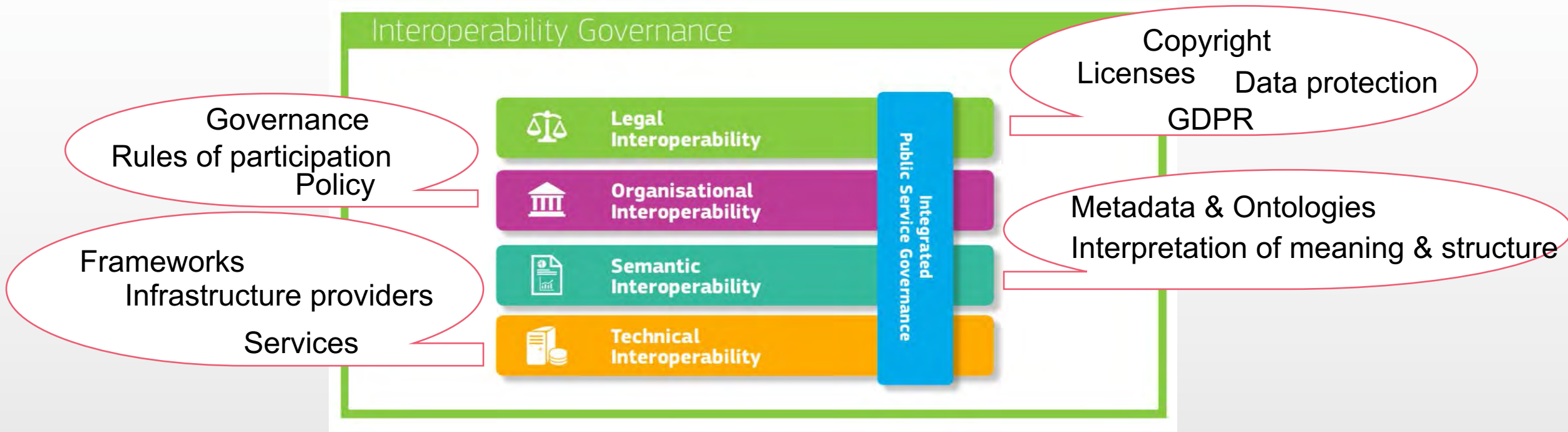
4. **Both lay a good foundation**

   - data, standards, softwares, platforms, services and so on
   - policy basis for National Measures

5. **Both focus on data centers now**

# Contents

➢Two projects

➢DataIO practice

➢Challenges

# European Interoperability Framework



Governance
Rules of participation
Policy

Frameworks
Infrastructure providers
Services

Copyright
Licenses    Data protection
GDPR

Metadata & Ontologies
Interpretation of meaning & structure

*Source: The European Interoperability Framework four levels of interoperability*

# L1 : National policy - Measures for Scientific Data Management

- drafted by the Ministry of Science and Technology
- issued by the General Office of the State Council on March 17, 2018
- establish a guideline for China's scientific data management and sharing

**1** Deeply follow the development trend of scientific data in the era of big data. Fully draw on advanced experience and mature practices at home and abroad to strengthen the life cycle management of scientific data.

**3** Vigorously promote the open sharing of scientific data resources. Adhere to the principle of "openness being the norm and non-openness being the exception", especially the sharing and opening of the scientific data supported by the state fund.

**2** Prioritize ensuring data security. Scientific data sharing should be based on security and controllability, and strengthen the data resources protection capacity.

**4** Focus on the weak links in the work of scientific data in China. Propose measures from the main responsibility, intellectual property rights, and the aggregation mechanism.

# Policy interoperability

Academy Implementation Rules

Provincial Implementation Rules

Ministerial Implementation Rules

Measures for Scientific Data Management

# Policy interoperability

- **Provincial Implementation Rules (14+1)**
  - Shaanxi province
  - Heilongjiang province
  - Gansu province
  - Yunnan province
    - kunming city
  - Hubei province
  - Anhui province
  - Inner Mongolia Autonomous Region

- Jilin province
- Guangxi province
- Chongqing city
- Jiangsu province
- Hainan province
- Shandong province
- Sichuan province

# Policy interoperability

- **Academy Implementation Rules**
  - Chinese Academy of Sciences
  - Chinese Academy of Agricultural Sciences

- **Ministerial Implementation Rules**
  - Ministry of Transport

# L2 :
# Organisational interoperability

# L2 : Organisational interoperability

## ●National science data center

| No. | National Scientific Data Center | No. | National Scientific Data Center |
|---|---|---|---|
| 1 | National High Energy Physics Science Data Center | 11 | National Cryosphere Desert Data Center |
| 2 | National Genomics Data Center | 12 | National Metrology Data Center |
| 3 | National Microbiology Data Center | 13 | National Earth System Science Data Center |
| 4 | National Space Science Data Center | 14 | National Population Health Data Center |
| 5 | National Astronomical Data Center | 15 | National Basic Science Data Center |
| 6 | National Earth Observation Data Center | 16 | National Agricultural Science Data Center |
| 7 | National Polar Science Data Center | 17 | National Forestry and Grassland Data Center |
| 8 | National Qinghai Tibetan Plateau Data Center | 18 | National Meteorological Data Service Center |
| 9 | National Ecosystem Science Data Center | 19 | National Earthquake Science Data Center |
| 10 | National Corrosion & Protection Data Center | 20 | National Marine Data Center |

# Organisational interoperability

●**CAS scientific data center**

# Organisational interoperability

- **Gansu provincial science data center**

```
                                    ┌─────────────────────────────┐
                                    │  Ecological and Environmental│
                                    │     Science Data Center      │
                                    └─────────────────────────────┘

┌──────────────────────┐           ┌─────────────────────────────┐
│ General Scientific   │           │      Climate Change          │
│    Data Center       │           │    Science Data Center       │
└──────────────────────┘           └─────────────────────────────┘
National Cryosphere Desert Data Center
                                    ┌─────────────────────────────┐
                                    │     Natural Disaster         │
                                    │    Science Data Center       │
                                    └─────────────────────────────┘
```

# Organisational interoperability

**National science data center 20**

- Data centers National High Energy Physics Science Data Center
- National Genomics Data Center
- National Microbiology Data Center
- National Space Science Data Center
- National Astronomical Data Center
- National Earth Observation Data Center
- National Qinghai Tibetan Plateau Data Center
- National Ecosystem Science Data Center
- National Earth System Science Data Center
- National Basic Science Data Center
- National Cryosphere Desert Data Center

**CAS science data center 1+18+X**

**Gansu provincial science data center**

**1+3**

# L3 : Semantic interoperability

- Metadata
  - Data about data



dataset  metadata

dataset  includes data and its'  metadata

# Semantic interoperability -- dataset metadata

- **National Standards for dataset**

  - GB/T 30523-2014,Science and technology infrastructure

    —**Resource core metadata**

  - GB/T 31073-2014,Science and tech

    —**Core metadata of service**

  - GB/T 39913-2021,Science and tech

    —**User metadata**

- Science and technology infrastructure
  - scientific instruments and equipments
  - biological germplasm resource
  - experimental material resource
  - standard substance resource
  - human genetic resource
  - specimen resource
  - scientific data
  - S&T achievements
  - S&T literature

●**Resource core metadata elements**

- ● ResourceIdentifier
- ● ResourceTitle
- ● ResourceCategory
- ● Keyword
- ● Description
- ● AccessConstraints
- ● DateOfUpdate
- ● PointOfContact
- ● OnlineAddress

☐ scientific instruments and equipments
☐ biological germplasm resource
☐ experimental material resource
☐ standard substance resource
☐ human genetic resource
☐ specimen resource
☐ **scientific data**
☐ S&T achievements
☐ S&T literature

**All 20 national science data centers publish datasets based on this metadata standard.**

● **CAS Scientific Data Core metadata @ CAS project**



DataSet:
- DatasetDescriptionInfo — 16 elements
- DQInfo — 2 elements
- DistributionInfo — 7 elements
- MetadataReferenceInfo — 4 elements
- ServiceReferenceInfo — 2 elements
- StructureDescriptionInfo — 3 elements
- Coverage — 3 elements
- Contact — 4 elements

# Semantic interoperability
## -- NBSDC Data



- **Data(set) resources**

  - Physics

  - Chemical

  - Material

  - Plant

  - Zoology

  - Information Science

  - ......

# Semantic interoperability -- Data

● **Example： Chemical Database @NBSDC**



化学主题数据库
Chemical Database

首页　搜索》　服务》　数据资源　技术支持》　帮助》

基本性质
热化学数据
相变数据
估算性质
安全与处置数据
毒性数据
质谱数据
药物和天然产物
精细化工产品
危险品
化学试剂

基本性质

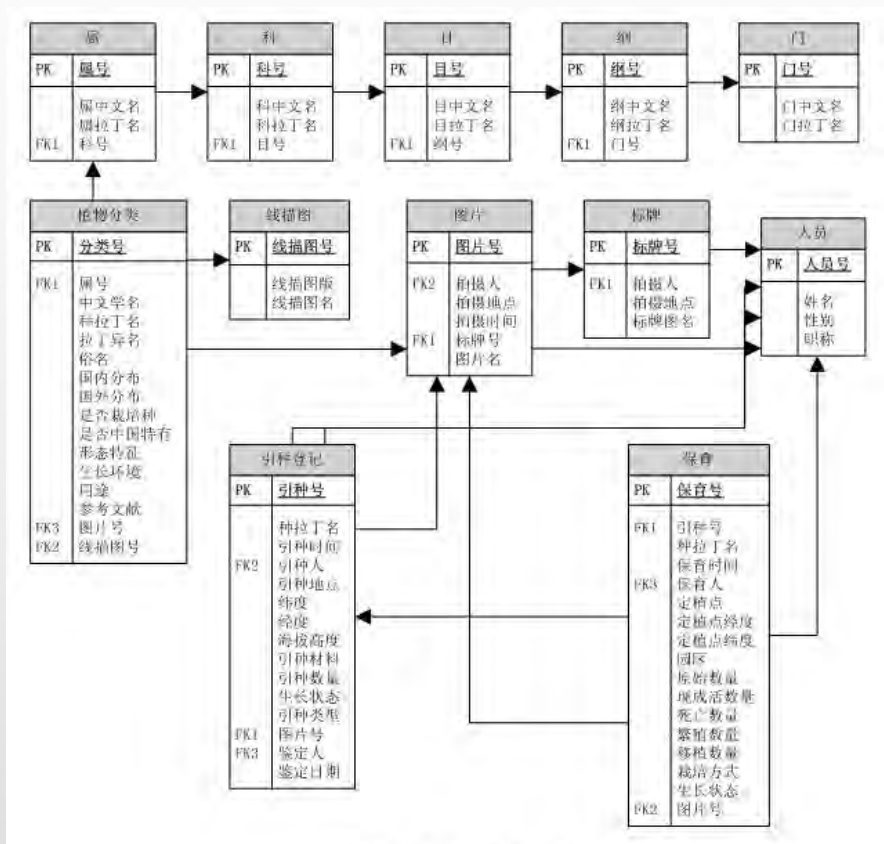化合物名称： (1R,2R,3R,4R
化学式： C6H6Cl6
分子量： 290.82984
CAS号： 319-84-6

● **Identification Mapping table**

| IDtype | ... | DataType | |
|--------|-----|----------|---|
| CAS RN | ... | varchar(10) | |
| InChIKey | ... | varchar(27) | |
| SRN | ... | Long(10) | |
| id | ... | varchar(50) | |

Analytical Chemistry

Environmental chemistry

Applied chemistry

Identification of compounds

Physical chemistry

Organic chemistry

Pharmaceutical chemistry

● Data from 3 institutes

● Collaborate services on the same page

热化学数据
相变数据
估算性质
安全与处置数据
毒性数据
质谱数据
药物和天然产物
精细化工产品
危险品
化学试剂

●**Example： Plant Collection DB @NBSDC**
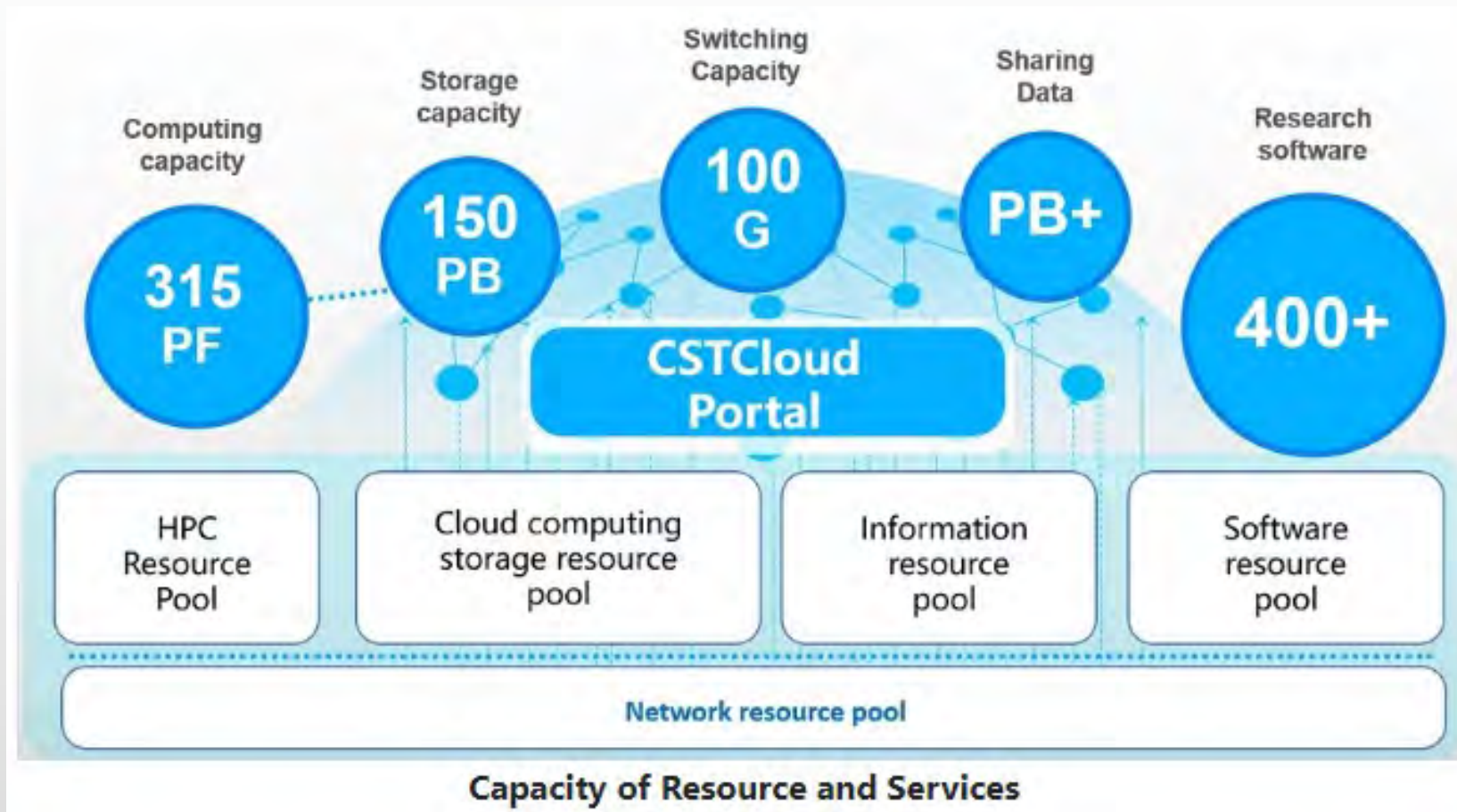


common data model

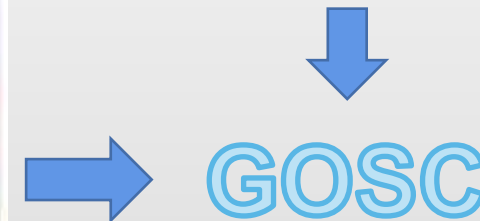●**Scientific Data Softwares @ CAS project**

# Technical interoperability -- cyber infrastructure

● **China Science and Technology Cloud（CSTCloud）**



- Network resource service
- Cloud storage service
- Object storage service
- Cloud backup service
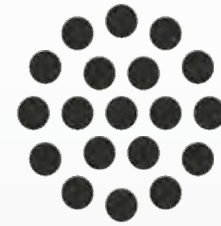- Computing service
- ......

# Contents

➢Two projects

➢DataIO practice

➢Challenges

# Challenges

- **The ecology of open and shared scientific data should be improved urgently**
  - The number of open and shared scientific data is still insufficient compared with the scale of data resources generated in China.
  - The availability and usability of scientific data that have been opened and shared are not high, the quality is uneven, and the overall level still needs to be improved.
  - The ecological environment for open and shared scientific data needs further improvement, including further policies and regulations, and exploration of innovative forms of open and shared data.

# Challenges

- **The comprehensive utilization of scientific data is still far from enough**
  - The management mode of scientific data center is difficult to meet the high efficiency requirements of multiple applications.
  - The cross-disciplinary application is still in the demonstration stage.
- **The scientific data centers are lacking International influence**
  - Although many data centers have been formed, but their overall strength is not prominent, especially lacking international influence.
  - Data exchange and interoperation with international data centers are also severely missing