Digital RI for Open Science: The National Practice in CSTCloud

Jianhui LI (lijh@cnic.cn)

Computer Network Information Center, CAS

CBAS deputy Director General

CODATA VP

2023.6.29



CNIC Introduction

 Computer Network Information Center (CNIC) is a special research institute of Chinese Academy of Sciences (CAS) to power scientific research and innovation, advance scientific research management, and promote public science outreach by advanced technologies. CNIC is the cradle of China's digital transformation and open science with data, network, computing and public outreach initiatives.





Hello Science | Hello World



Data Center Located In Huairou District Beijing, China

Main Campus Located In Haidian District Beijing, China

Requirements for E-infrastructure in CAS

Data-intensive & big data driven Science



• Decades investment in cyber-infrastructure development and e-science in CAS, totally over 2 billion RMB

E-Infrastructures in CAS



CSTNet has 12 Branches across China, and provides services to approx. 1MM. end users from 300 institutions. Her domestic bandwidth is 131G and international bandwidth is 113G.Below shows the CSTNet Global Interconnection.



COMPUTING

Chinese National Grid

CNGrid has connected 19 supercomputing centers, the capability of clustering computing has reached 200PF, supported for thousands of national science and technology program and key engineering. 40+ independent high performance computing software.

CAS Supercomputing Grid environment creates a three-layer architecture environment and has integrated more than 15PF computing resources from 35 CAS institutions



ADS transmutation system



titanium alloy structure



f Simulation of the Atlantic ocean circulation

GRID

20 National Scientific Data Centers are launched by the Ministry of Science and Technology,P.R.C as the pilot to improve research data sharing following the national rules, "Measures for Managing Scientific Data". Among all of which, 11 are from the Chinese Academy of Sciences.

No.	Nati	onal Sci	entific Data Center		
1	National high en	ergy ph	ysics science data center		
2	National genom	e scienc	e data center	•	
3	National microbial science data center				
4	National space s	cience o	lata center		
5	National astrono	omical se	ciences data center		
6	National earth (No.	National Scientific	Data Center	
7	National polar s	11	National glacial frozen deserts scie	ence data center	
8	National Qingha	12	National metrological science data	a center	
9	National ecolog	13	National earth system science data	a center	
10	National materi	14	National population health science	e data center	
	data center	15	National basic sciences public science data center		
		16	National agricultural science data	center	
		17	National forestry and grassland sci	ience data center	
		18	National meteorological science data center		
		19	National earthquake science data center		
		20	National Marine science data center		

OPEN RESEARCH DATA

NETWORK

China National Policy

- Scientific Data Management Regulations issued by the State Council in March 17, 2018
 - intended to clarify the responsibilities of officials and scientists who regulate and use the information.
 - improve the management, security, accuracy and openness of scientific data,
- Scientists have to submit data to related National Scientific Data Center for archiving and open services
- Scientific data sets will also better identify their origins and researchers, allowing clearer citations and stronger protection of intellectual property

Established 20 National Scientific Data Centers

MOST launched 20 National Scientific Data Centers in 2019



Responsibilities of National Scientific Data Centers

- -Integration of scientific data
- -Classification, processing and mining of scientific data
- Open and share scientific data
- -Strengthen exchanges and cooperation in scientific data nationally and internationally

Integrate distributed e-infrastructure for Open Science

 Open science infrastructure features: Federated, Accessible, Internationally Interconnected, Interoperable (UNESCO, 2021).



CSTCloud: A national research e-infrastructure

- CSTCloud supports multidisciplinary open scientific researches with integrated cloud services for the discovery, usage and delivery of S&T resources.
- Continually funded by CAS Informatization Program(2017-2020,2020-2025)



CSTCloud: Targeted users and tailored services

CST

Cloud



- Service discovery
- One-stop access
- Open and Use

Service Provider

- Service Registration
- Demonstration & Promotion
- Management & Transaction

For the **100,000** scientific researchers in CAS, providing a cloud service environment that supports innovation in multidisciplinary fields, supporting the scientific discovery of big data and big computing

Resource Owner

- Service monitor
- Quality statistics
- **Comprehensive Evaluation**

CSTCloud Service Catalogue and System



CSTCloud: Online Services



Storage Resources: 150PB

Software Resources: 1000+

Network Resources: 100G



Service Catalogue



CSTCloud Service Application and Analysis



- CSTNet
 - About 300 institutes, 1 million users
 - 100G domestic and 100G international interconnection
 - The best quality in China Network Carriers
- HPC
 - 200PFlops Storage 100GB
 - >600 users
 - >1.7 million jobs

- Cloud Computing and storage
 - >1000 cloud hosts
 - 100PB Storage
 - 12 data centers federated
- CAS Passport/Check-in
 - >1200 applications and 1.32 million registered users
- eduroam
 - 400 institutes and universities

CSTCloud: Open Data services



CAS Data Cloud

CAS Data Cloud integrates resources and services from the perspective of three major types of services: infrastructure, data resources, and application platform. Up to now, the cloud platform has integrated nearly **960.14 TB** of shareable scientific data, with **132,500** users, **178,863,200** cumulative online visits, and **22,867,200** cumulative data downloads.



Data set: 18369 Online Resources: 3.11 PB Data Center: 35 Total number of visits: 116 million



RESOURCE CATEGORY:

CSTCloud: Research Data repository



http://www.scidb.cn/en

Science Data Bank is a public and non-profit generalist data repository which is developed and maintained by the CNIC, implementing the FAIR principles of data sharing. The repository accepts submissions from the entire scientific community across disciplinary boundaries and file formats, and publishes all kinds of research data with no charge of fee.



CSTCloud for big science research





📐 the way to new energy

iter





CSTCloud for SDGs

- The SDG Big Data Platform provides rich data resources and cloud services for SDGs Decision Making and science discovery, that
 - With **10PB data** covering multiple subjects, such as geography, remote sensing, ground monitoring, and social statistics;
 - Characterize and Profile Scientific Workflows featuring lifecycle Big Earth Data management, including massive data storage, curation, computation, analysis and visualization;
 - Integrate over one hundred algorithms and tools for advanced data analysis and management;
 - Provide one-stop cloud services through the bilingual portals to support the UN Sustainable Development Goals.



Data Analysis As Service: EarthDataMiner

- Online interactive data analysis environment
 - Integrated with DataBank: upload models, select data, and process data products through instruction operations
 - Algorithm & Model library: more than 150 algorithm developed and provide cloud service: FAAS(Function As A Service)
 - Web IDE: supporting users to write data analysis code (Python) online



SDGs indicator on-demand computing

1. Upload data to the EarthDataMiner system



3. Search radar data from Databank develop by CASEarth



2. Write python code to preprocess data



4. Compute the SDGs Indicator and Visualizaiotn



SDGs indicators on demand analysis Tools



SDG13.2.2 Annual average CO2 concentration evaluation



SDG15.1.1 Forest coverage detection



SDG6.6.1 Surface water change over time indicator



SDG13.1.1 Natural disaster impact evaluation



SDG11.3.1 Global urbanization index monitoring



Converged Virtual Research Environment : SDG Workbench

- Data be accessed by
 - applications/machine Transparently
- SDGs Tools
 - ✓ Integrated Tools
 - ✓ Data Analysis Tools
 - ✓ Data Products Tools

Open-source Tools

- ✓ Developing
- ✓ Machine Learning
- ✓ Data Visualization
- Creating / Using / Releasing Spark
 Cluster On-demand
- Cloud-Native DevOps CI/CD
- Virtual Collaboration
 - Setup a virtual team or virtual organization



Example: GGW-BDF & SDG15.3.1 LDN



 provide high spatial resolution datasets on land cover, land productivity and soil carbon change to better track the progress of LDN.



CSTNet Global interconnection



Connecting CSTCloud with the World

• A cross-continental federated e-infrastructure and virtual research environment for global cooperation and open science.



The Global Open Science Cloud (GOSC) Initiative

- The Global Open Science Cloud is a robust network to bridge trusted research einfrastructures among which different research resources are connected and all stakeholders are linked for innovative science discovery in the dynamically evolving global open science environment.
- Important early project as part of ISC CODATA Decadal Program with seed funding from CAS.



GOSC: Working Groups

Strategy, governance and sustainability



Qunli Noorsaadah Han Abdul Rahman



Working Group (WG) to understand the expectations, interest areas, resour

Your answers will be kept confidential within the WG and may be anonyr purpose of summary analyses

While the questions are not required fields, we encourage members to rea





Technical Infrastructure



Jianhui Mark Dietrich

Li



- Continued work in Sub Working Groups to develop technical reference frameworks for key functions
 - AAI (Authentication, Authorisation and Identification)
 - Data Access/Storage Authorisation
 - Service and Data Discovery 0
 - Remote Compute Execution/Authorisation

Thu 31 March

- Interactions with sister WGs to identify common or overlapping challenges
- Developing an Overall Technical Framework for GOSC to guide creation of testbeds





Happy

Sithole







Natasha Simons



Case Study template

Case Study:	ICSM ANZ Metadata Working Group (MDWG) https://www.icsm.gov.au/what-we-do/metadata-workin g-group
Participant name, affiliation	Irina Bastrakova (ICSM / ANZ), Melanie Barlow (ARDC), Lesley Wyborn (ANU), Rowan Brownlee (ARDC)
Participant email address	Irina.Bastrakova@ga.gov.au Melanie.Barlow@ardc.edu.au Lesley.Wyborn@ardc.edu.au Rowan.Brownlee@ardc.edu.au
Date	Wednesday 4 May 2022

Policy and legal



Lili Sarah Zhang Jones

ODATA Research progress - overview

Strand 1: A light-weight policy review - RoP

- EOSC Rules of Participation
- Open Storage Network Acceptable Use Policy
- Australian Research Council Data Management Policy

- Legal & Economic laver
- · Community Layer (social/cultural/ethical)

Strand 3: OSCs metrics

Data Characteristics Affecting Levels of Openness

Strand 4: Worldwide OSCs stories





- Community outreach Monthly Member Meeting Bio-weekly subgroup meeting
- · Conferences and events

Strand 2: Implementation guidelines for OSCs alignment

- · Resources Layer (information resources/digital objects)



MOSAIC example (Hans Pfeiffenberger)

OSN example (Melisa Cragin, TBD)

Pascal Heus





Ojsteršek

GOSC: Case Studies

Incoherent scatter radar data



Ingemar Häggström



Biodiversity and Ecology



Joe Mille



Zhishu Xiao

Diffraction Data



John Helliwell



SDG-13 climate change and natural disasters Sensitive data in population health



Gensuo Jia

Monthip Sriratana



Bapon Fakhruddin



Lei

Liu

Jildau Bouwman

Lauren Maxwell

Francisca Oladipo

Incoherent Scatter Radar Data Fusion and Computation

GOSC

This CS supports the international collaboration of the radar community, particularly with a focus on **data and technical interoperability**. The management of large-scale radar data provides an excellent scenario to validate the technical maturity of the GOSC testbed.



Integrating Machine-derived Camera-trap Data into the Global Shared Science Networks

Through a collaborative platform for global camera-trap data sharing and analysis service, this CS will contribute to global biodiversity research, especially focusing on **distributed big data management, intelligent analysis, and cloud computing** for high-quality integration and optimization for camera-trap data management.







This CS would lead to **single, definitive, protein models** derived from their raw diffraction data sets, which is important to the crystallographic community and the broader research community for drug discovery, especially for the pandemic crises that COVID-19 poses for society.



Sensitive Data Federation Analysis Model in Population Health

This CS seeks to demonstrate better ways of sensitive data sharing. Commonly agreed FAIR implementation profiles will be created based on the **FAIR data points** established in various 'GOSC' regions. The feasibility of distributed analytics over datasets held in various regions will then be explored and demonstrated. **Community-accepted standards** will be used throughout to facilitate the implementation of this CS.



From Raw Biodiversity Data to Operational Indicators through the Essential Biodiversity Variables (EBV-OSC)

The objective of this CS is to operationalize EBV indicators by targeting the highest levels of FAIRness (Findable, Accessible, Interoperability, Reusable) for both **data and source code implementation**, so that data and tools can be widely shared and reused.



SDG-13 Climate Change and Natural Disasters

CASEarth for Sustainable Development Goals (CASEarth4SDGs) is a platform system of data sharing and online computing for monitoring, measuring, and evaluating SDG indicators. Supported by this system, this CS mainly focuses on climate change and natural disasters in the SDG-13 field, addressing technical, semantic, and policy interoperability to support decision-making.

GOSC: International Programme Office (IPO)

- Expanding GOSC cooperation network and maintaining membership
- Facilitating conferences, training workshops and other events to increase social visibility
- Coordinating community outreach, day-to-day operation and future development activities of GOSC

GOSC IPO CNIC Team





Lili ZHANG CNIC

CNIC



Simon Hodson CODATA



CODATA

GOSC IPO CODATA Team



institutional, national, and regional Open Science clouds and platforms to create a global virtual environment for globalized research and innovation. So far, four thematic Working Groups and five exemplar scientific Case Studies

First IPO in CODATA's History



International Office located at CNIC



IPO Launch Event (10.12)

CSTCloud Federated with EGI

- Technique solution has been verified, ready to provide service
- Shared resource is limited, need more resource...



A Typical Use Case – 3D Radar Data Federation analysis



• EISCAT Data Access Portal (DIRAC)

- User access via Check-In/Perun
- Access token is passed to DIRAC File Catalogue (by DIRAC Client)
- DIRAC File Catalogue return token with user information
- (DIRAC Client) access DIRAC Storage Element with token
- DIRAC Storage element enables search of EISCAT data (in the EISCAT storage)
- DIRAC Server submits Jobs to Cloud
 - TUBITAK (10core, 30TB) + CSTCloud (30Core, 100GB)
 - Token access to the Cloud is ongoing
- EISCAT Compute Center
 - Jupyter notebook
 - User access via Perun
 - Matlab enabled
 - Not yet access to DIRAC Storage Element cannot perform search, and have lower performance
 - EISCAT Storage: EISCAT data
- EISCAT Notebook Container
 - Enable EISCAT user access
 - Notebook+DIRAC image
 - dirac.egi.ei cvmfs repository is configured alongside eiscat cvmfs repository (eiscat.egi.eu)
 - DIRAC client is available via cvmfs
 - pass user access token to the DIRAC File Catalogue to get user information
 - Access the DIRAC Storage Element to search and retrieve data from the EISCAT storage
 - Will include Matlab
 - Token to use EISCAT license?
 - Will include EISCAT Tools

Analysis interoperability

- Analysis (GUISDAP)
- Internally radar independent
- EISCAT wrapper
 - Single beam setup
- SYISR wrapper
 - Multi beam setup
 - Test also for EISCAT_3D

From the Use case group Report on the 2AHM meeting by Ingemar



results



SYISR a

Based on Jupyter notebooks



- 3. Jupyter Notebook landing page
- •Jupyter Notebook Desktop workspace with GUISDAP application analysis running

EISCAT User Manual PITHIA-NRF CROLHL project Anders Tjulin Carl-Fredrik Enell Ingemar Häggström Mária Miháliková April 2022 User Manual

EISCAT Scientific Association

Notebook is in production: https://notebooks.egi.eu/ https://jupyter.eiscat.se/

https://radardata.fed.cstcloud.cn/

From the Use case group Report on the 2AHM meeting by Ingemar

Call for International Partnerships

- Open mind, open framework, open collaboration for open data & open science.
- Coming together is a beginning. Keeping together is progress. Working together is success. Henry Ford (1863 1947, American industrialist)



Thank you!

