# Open Data Platform in Korea

**Global Science experimental
Data hub Center**

**September 6, 2018
Seo-Young Noh**

www.kisti.re.kr

# Contents

1. **Open Science:** *Connected Science*

2. **Data-driven R&D Era**

3. **Data Infrastructure in Korea**

4. **Linking Data Repositories:** *Practical Implementation*

5. **Summary:** *Responses to Requests*

# Open Science:
## *Connected Science*

- **OECD Principles and Guidelines for Access to Research Data from Public Funding (2006-07)**

- **Initial discussion of Open Science at CSTP in 2011**

- **Many Open Science related activities on-going (PSI, open gov data, open educational resources, MOOCS…)**

- **OECD produced the first Open Science report, mainly focusing on Open Access, Open Collaboration and Open Data (2015)**

- **Several expert groups in GSF have been formed to build advisory policy for Open Science: Research Infrastructure, Data Infrastructure for Open Science**

OECD Principles and Guidelines for Access to Research Data from Public Funding

**Acknowledgment on the importance of open**

OECDpublishing

Please cite this paper as:
OECD (2015), "Making Open Science a Reality", *OECD Science, Technology and Industry Policy Papers*, No. 25, OECD Publishing, Paris.
http://dx.doi.org/10.1787/5jrs2f963zs1-en

OECD Science, Technology and Industry Policy Papers No. 25

**Making Open Science a Reality**

OECD

OECD

Open science is *more than open access to publications or data*; it includes many aspects and stages of research processes.  [...]

Open science is *a broader concept* that includes
- the interoperability of scientific infrastructure
- open and shared research methodologies (such as open applications and informatics code)
- and machine-friendly tools allowing, for example, text and data mining.

Source: POLICIES TO PROMOTE OPEN SCIENCE: EVIDENCE FROM OECD COUNTRIES, Giulia Ajmone Marsan

# Key features of Open Science

■ **Main goal of Open Science:**
- ➲ **Provides cost-effective access to digital research data from public funding**
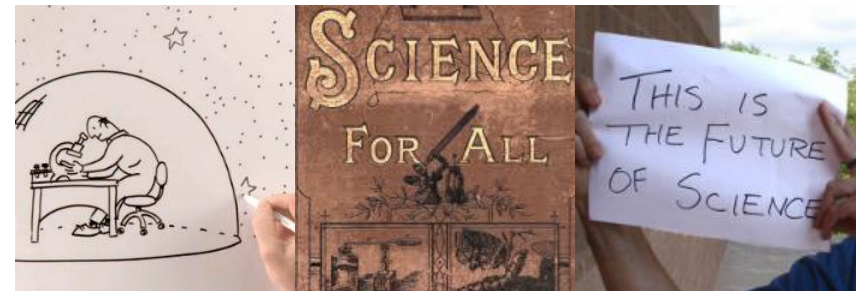- ➲ **Enhances utilizations of research data to scientific communities as well as societies including corporate sectors**

**Science cannot exist in a bubble…**



**should work with public communities…**
**The Future of Science**

■ **Benefits**
- ➲ **Easy Research** → Efficiency, Removing Redundancy → **Solving Contemporary Problems**
- ➲ **Tackling Big Problems** → Enabling Big Science → **Solving Problems of Humankind**
- ➲ **New Value Creation** → Enabling Convergence → **Solving Unknown Problems**
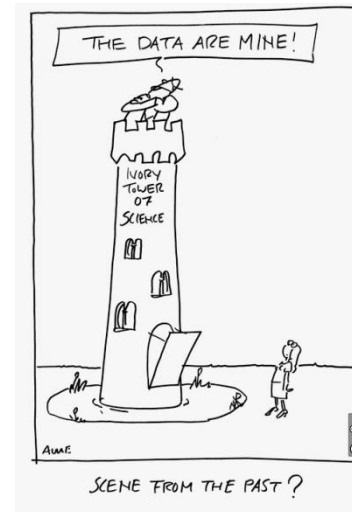
■ **Ways and Means**
- ➲ **Open Data, Open Access and Open Collaboration** through Information and Communication Technology
- ➲ "Open Access" and "Open Collaboration" look straightforward, but *"Open Data" is not so simple, requiring deep understanding the features of data*
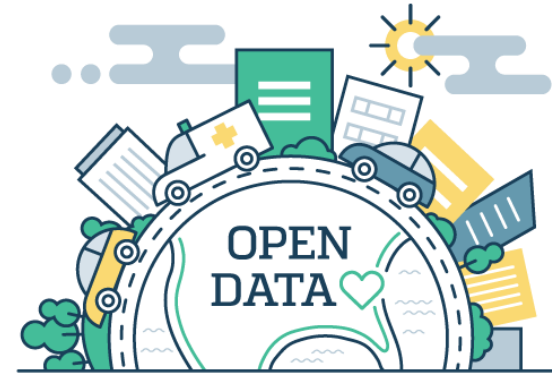
## Open data is the key part of Open Science:

- ⮞ **Transparency** in experimental methodology, observation, and collection of data

- ⮞ **Public availability and reusability** of scientific data

- ⮞ **Public accessibility** and transparency of scientific communication

## Data should be valuable

- ⮞ **Reusable data** (ex: observation data, knowledge database)

- ⮞ **Data requiring long time** for data accumulation (ex: pathology tracking data, climate change tracking data)

- ⮞ **Data requiring big budget** for data acquisition (ex: large equipment-based experiments)

- ⮞ **Data requiring huge computing power** for data generation (ex: simulation data)

**Why open data is matter?  ... many terminologies, properties**

*Transparency, Public Availability, Public Access, Reusable, Redistribution, Universal Participation*

# Interoperability

**Interoperability makes it possible for diverse research groups inter-operate scientific data, intermixing different datasets,** *leading to open a new way for unveiled values*

**Connected Scientific Data**

**Connected Science**

**Value Creation**

**Social Responsibility**

**Connected Data**

**Connected Devices**

**Value Creation**

~~**Social Responsibility**~~ **Business**

# Open S/W – Good Example

**Linus Tovalds opened free operating system (Open S/W),**

**later combined with GNU project, influencing great impact**

From: torvalds@klaava.Helsinki.FI (Linus Benedict Torvalds)
Newsgroups: comp.os.minix
Subject: What would you like to see most in minix?
Summary: small poll for my new operating system
Message-ID: <1991Aug25.205708.9541@klaava.Helsinki.FI>
Date: 25 Aug 91 20:57:08 GMT
Organization: University of Helsinki

Hello everybody out there using minix –     **Hobby work**

I'm doing a (free) operating system (just a hobby, won't be big and
professional like gnu) for 386(486) AT clones. This has been brewing
since april, and is starting to get ready. I'd like any feedback on
things people like/dislike in minix, as my OS resembles it somewhat
(same physical layout of the file-system (due to practical reasons)
among other things).

I've currently ported bash(1.08) and gcc(1.
This implies that I'll get something practic
I'd like to know what features most people
are welcome, but I won't promise I'll impl

Linus (torvalds@kruuna.helsinki.fi)

PS. Yes – it's free of any minix code, and i
It is NOT protable (uses 386 task switchin
will support anything other than AT-hard

Sharing is good, and with digital technology, sharing is easy

– Richard Stallman

©The Kolaveridi

**Richard Stallman**

*Software* **to be distributed in a manner such that its users receive the** *freedoms to use, study, distribute and modify* **that software**

**Young Linus Tovalds**

**GNU Project since 1984**

GPL v3
Free Software
**Free as in Freedom**

**GNU Public License**

# Data-driven R&D Era
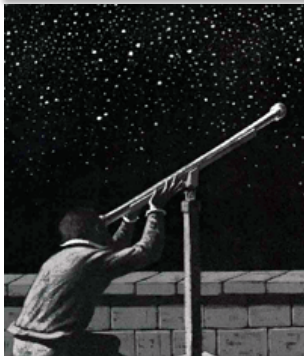
# Data & Infrastructure are Key in Scientific Discovery

**Describing natural phenomena based on Observation**

**Modeling and Theory**

**Computing Simulation**

**Data Analysis of tremendous data produced from large experimental facilities**

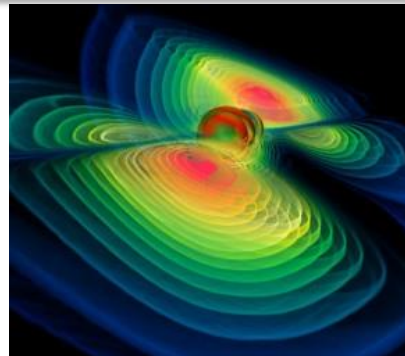**Research Paradigm Shift to Data Intensive Scientific Discovery**

| | | | |
|---|---|---|---|
| **1st Generation: Observation** | **2nd Generation: Theory** | **3rd Generation: Simulation** | **4th Generation: Data** |
| Galileo's telescope | Higgs Theory | Black Hole Simulation | **CERN's CMS and ATLAS experiments → Higgs discovery** |

## More chance to do research with advanced equipment, higher chance to get Nobel prize

87% of Nobel prizes have been given to researchers who produced outstanding scientific discoveries using advanced experimental equipment since 1914.

**Source: The Fourth Paradigm**

- **CERN** <u>**noticed a signal like a new particle in CMS & ATLAS experiment**</u> in December 2015.

- The **750 GeV diphoton excess** in particle physics was an anomaly in data collected at the Large Hadron Collider(LHC) in 2015, <u>**which could have been an indication of a new particle.**</u>

- However, <u>**the anomaly was absent in data collected in 2016**</u>, suggesting that the diphoton excess was <u>**a statistical fluctuation.**</u>

- In the interval <u>**between the December 2015 and August 2016 results**</u>, the anomaly generated considerable interest in the scientific community, including about **500 theoretical studies**.

**We are in data-driven science era!!!**

<u>**Our trust is in data**</u>

The 750 GeV 'thing'

- The plan is to have th
  mode' while in paralle
- Special attention will

1. $S \to Z$
2. $S \to W$
3. $S \to hh$
4. $S \to t\bar{t}$

**CMS Spokesperson
April 2016 CERN Resource Review Board**

**"We are about to open a new page in HEP history"**

- The hints of a possible new state of matter have excited the HEP community: in the experiments our trust is in data ... and we will get a lot more data in 2016 which should tell us if we are about to open a new page in HEP history

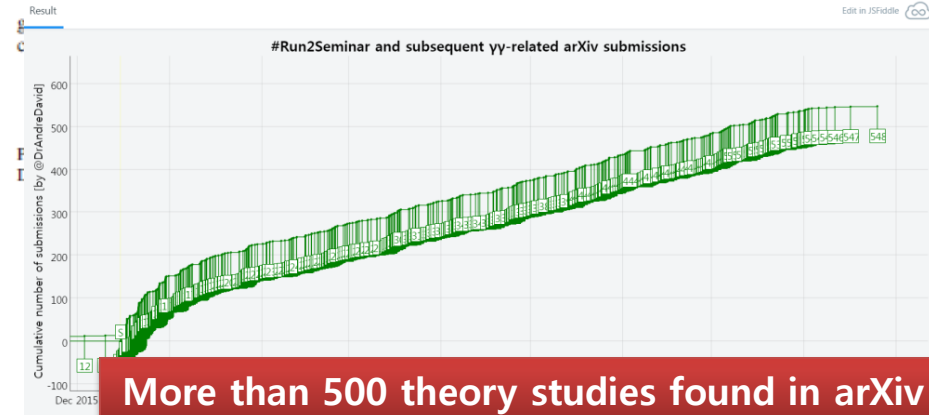PRL **116**, 150001 (2016)     PHYSICAL REVIEW LETTERS     week ending 15 APRIL 2016

**Editorial: Theorists React to the CERN 750 GeV Diphoton Data**

Last December, the ATLAS and CMS Collaborations at the Large Hadron Collider reported
p... [1] ...
have recently reanalyzed their data [3,4], and the signal has become slightly stronger. Though the results are extremely intriguing, more data are required to establish if the excess is real, or a statistical fluctuation.

Over 250 theory papers have appeared following the December announcement, and a number of them were submitted to us. We found it appropriate to publish a small sample of them. To maximize the coherence and fairness of our choices, we obtained informal advice from several experts.

**More than 250 theory papers submitted to PRL**



#Run2Seminar and subsequent γγ-related arXiv submissions

**More than 500 theory studies found in arXiv**

**Science relies on data, requiring infrastructure for data.**

**Data is getting more important and growing fast.**



**Data Infrastructure is the one of key factors for successful science and tackling big problems of humankind.**
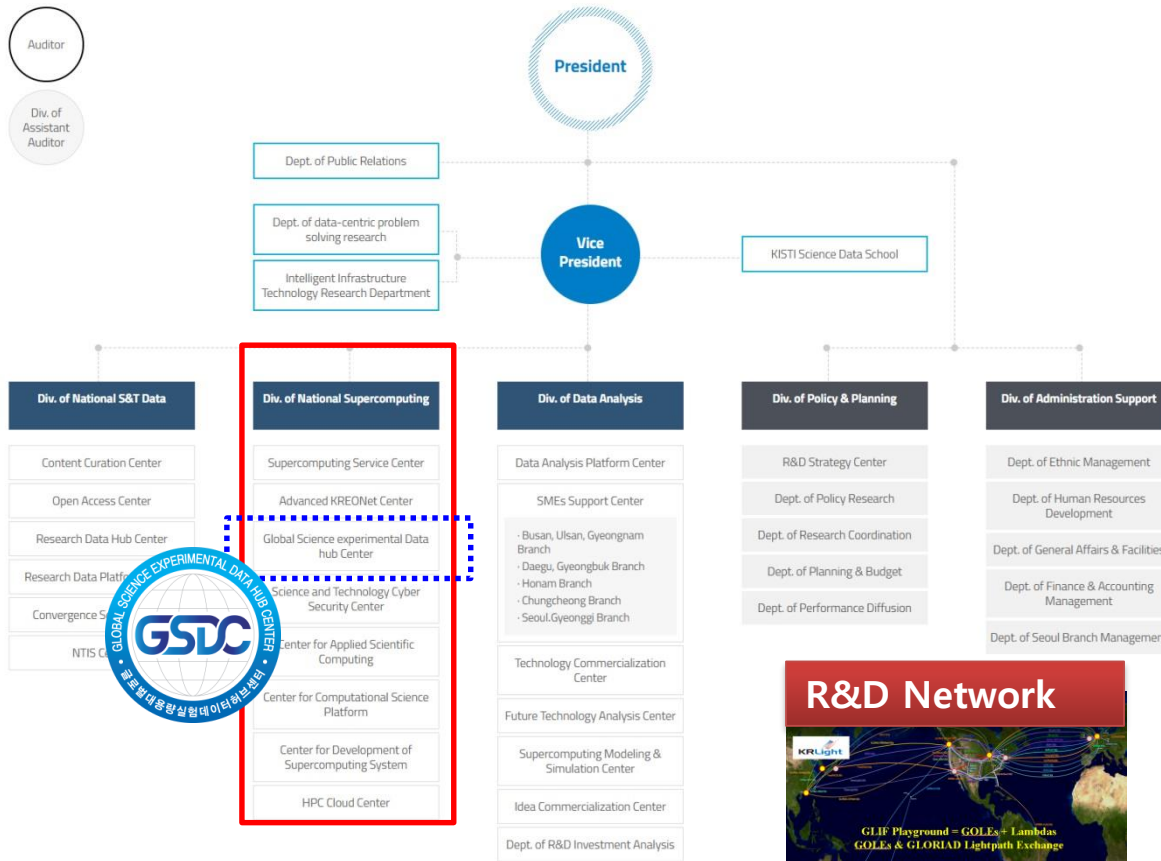


**KISTI has been in preparation for big data research era. Our mission is gradually expanding to national role for data intensive research.**

# Data Infrastructure:
## *KISTI, Korea*

# KISTI...providing powerful ICT infra. service



**New Supercomputing Building**

**25.7PFlops and 20PB ranked 11th in Top 500**

**Supercomputer**

**R&D Network**

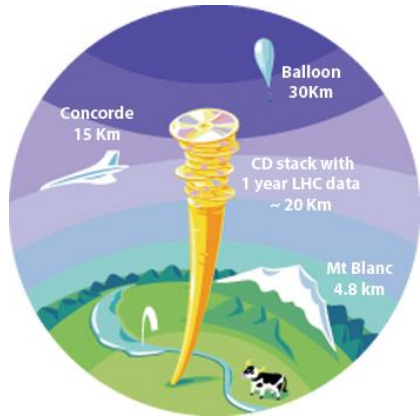**Unique organization in Korea, our mission is to support various R&D activities through ICT infra. service**

**Fast and secure network**, providing domestic researchers with a **constraint free collaborative research environment** through KREONET (locally) and GLORIAD (globally)

# Global Science experimental Data hub Center

**Collaboration with global laboratories**

**Large-scale Scientific Data:**
20Km CD stack with data produced per year in CERN

**Data from large and high-valued research equipment**

⮌ **(Global)**
**Asia representative Data Hub**

⮌ **(Domestic)**
**Scientific data management and analysis platform service**

# Growing every year by ~1,200 cores and ~2PB

**Best equipment procured every year**

**Centralized data repository model**, but **distributed model** in near future

**Interlocking with existing systems done by 100% KISTI experts**

25 Storage Racks with 6 Different Models

Storage Area Network 5.6 PB

Tape Storage 3 PB

60 Switches with 10 Different Models

Backbone Switch

FDF

Network Attached Storage 4.5 PB

ALICE

LIGO

550 Server with 14 Different Models

CMS

BELLE  RENO  Admin  TEM  Genome  NEW  Testbed

# Data Repository Services (1)

**High Energy Physics**

1. **ALICE(A Large Ion Collider Experiment)**
   - To generate similar conditions that have existed a fraction of the second after the Big Bang
   - 1,655 scientists from 159 institutes, 41 countries

2. **CMS(Compact Muon Solenoid)**
   - To investigate a wide range of physics, including the search for the Higgs boson, new physics
   - 3,000 scientists from 172 institutes, 40 countries

3. **Belle/BelleII(Japan KEK)**
   - To investigate CP-violation effects and new physics
   - 428 scientists from 67 institutes, 20 countries

**Astro Physics**

4. **LIGO(Laser Interferometer Gravitational Wave Observatory)**
   - To detect cosmic gravitational waves and to develop gravitational-wave observations
   - 1,000 scientists from 60 institutes, 16 countries

**Particle Physics**

5. **RENO(Reactor Experiment for Neutrino Oscillation)**
   - To measure a limit on the neutrino mixing matrix parameter(Yeonggwang Nuclear Power Plant)

**Medical Science**

6. **Genome Research**
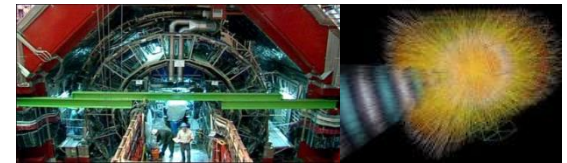   - Genome data analysis for the next generation of personalized treatment: 50 Korean researchers

**Biology**

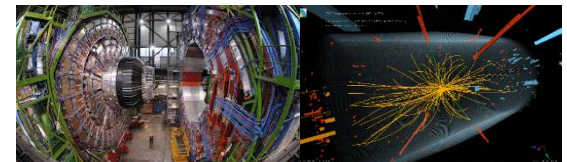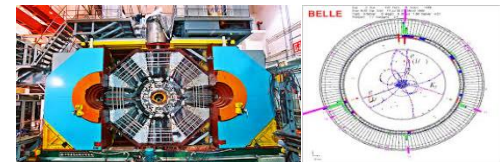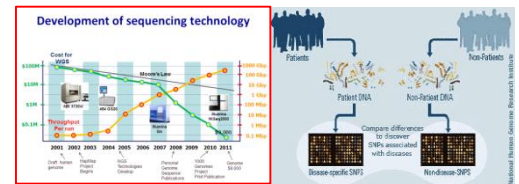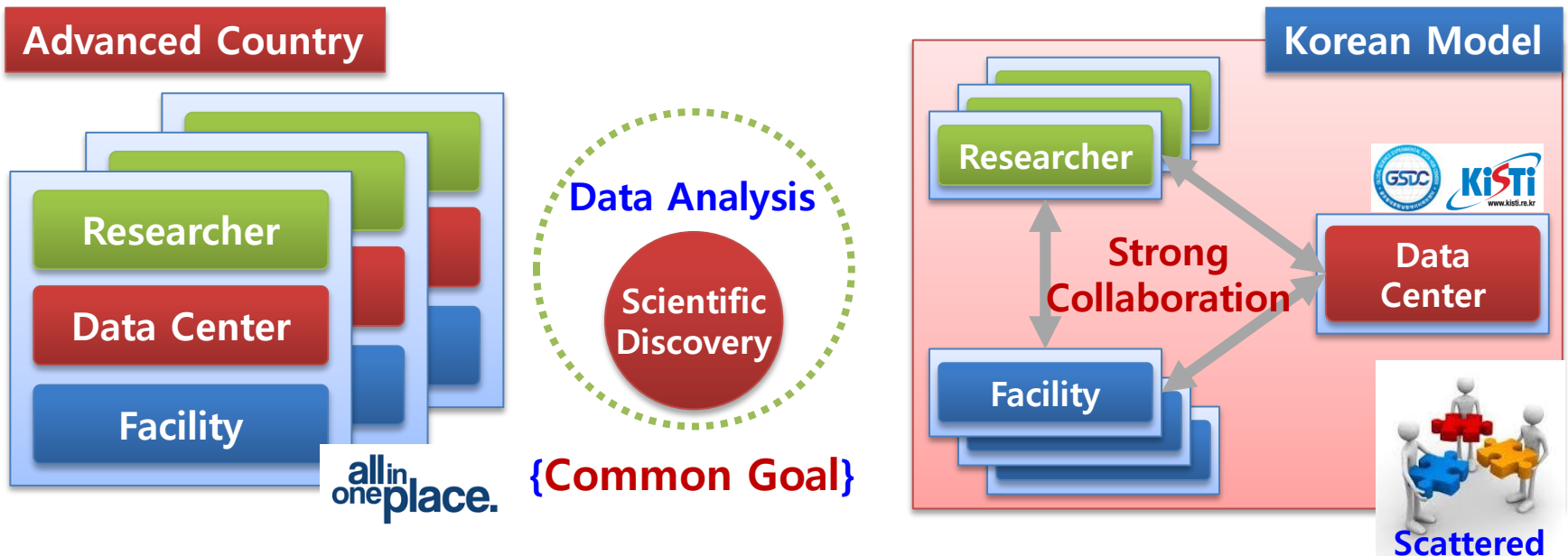■ **Structural Biology**
- ➲ Data management and analysis service for Transmission Electron Microscope-based research

**Safety & Disaster**

■ **Volcanic Disaster Prevention**
- ➲ Data repository service for volcano simulation
- ➲ Helping to build a policy for volcanic disaster response scenarios

**High Energy Physics**

■ **CMS Tier-2 Service**
- ➲ Unified WLCG service in Korea
- ➲ A new Tier-2 center targeting 2018 service

**Astro Physics**

■ **KAGRA(Kamioka Gravitational Wave Detector, Japan, 2017)**
- ➲ A project of the gravitational wave study group
- ➲ Data transmission, repository and data analysis environment in conjunction with LIGO data analysis environment

**General Purpose Accelerator**

■ **PAL(Pohang Accelerator Laboratory, Pilot Project, 2018)**
- ➲ Data repository and analysis service for 4th generation accelerator
- ➲ **Trial service for data convergence produced from heterogeneous large equipment for structural biology (Government Supported Funding)**

**Medical Science**

■ **Brain Research (Under discussion)**
- ➲ Platform service for nationwide management, repository, transmission, and analysis of brain data
- ➲ Extending knowledge of HEP data analysis platform service to brain data

**High Energy Physics**

■ **RISP(Rare Isotope Science Project, RAON, under discussion)**
- ➲ Exploration of the origins of chemical elements, structural study of new isotopes and applied medical research.
- ➲ Discussion with Government for data center role of RAON

# Korean Model for Data Intensive Research

**Focusing on a centralized model for data repository at the beginning**

**by fully utilizing ICT specialized institute like KISTI/GSDC**



**Advanced Country**

Researcher

Data Center

Facility

all in one place.

**Data Analysis**

Scientific Discovery

**{Common Goal}**

**Korean Model**

Researcher

**Strong Collaboration**

Data Center

GSDC  KiSTi  www.kisti.re.kr

Facility

**Scattered**

- ➲ **Large scale research group**
- ➲ **Large scale research facility**
- ➲ **Dedicated data center**

**VS.**

- ➲ **Small scale research group**
- ➲ **Small or medium scale facility**
- ➲ **Not easy to have a dedicated data center (in size and experts)**

# Unified Data Analysis Environment (Centralized Model)

Equipment A

| Equipment E | Equipment F | Equipment Z |
|---|---|---|

| Group Z | Group E | Group A |
|---|---|---|

**Fast Experimental Data Transmission**

**User Service[Analysis] Data Flow**

Science Gateway

**Science Neutral Unified Analysis Environment**

**Logically Unified Single Resource**

**Analysis S/W**

/swA
/swB
/swC

/swZ

**Fast Data Access**

Data A    B    C    Data Z

**Science Data DMZ**

Group C

Group A    Group B    Group Z

**Automatic S/W Connection On-Demand**

## 【Advantages】

1. **Pluggable Science** → **Supports in unified way** for various groups and equipment

2. **Data Infra. Sharing** → **Reuse and full utilization** of infra. **saving tax-payer's money**

3. **Simple R&D Process** → **Fast results** from data acquisition to data analysis

# Linking Data Repositories:
## *Practical Implementation*

# R&D Open Data Task Force: Open Data Platform

**A Task Force** has been setup in Ministry of Science and ICT, focusing on ...

1. **Fostering** data-driven R&D communities

2. **National level** regulation for data management

3. **Scientific program development** for open data



**Public Hearing - Dec. 18, 2017**

**Main activities ...**

1. **Analyzing best practices** including European Open Science Cloud in EU, Big Data Hub Program in U.S

2. **Defining priority and categorization** of data intensive research fields and **new program development**

3. **Developing R&D Open Data Platform,** helping seamless data sharing, data accessing, data analysis, data linking across disciplines

## Accelerating R&D productivities
## through ICT-based R&D e-Transformation

- Solving unknown problems

**【Expect ation】**

**Social Benefits**

⟷ **Improving R&D Productivity**

⟷ **New Value Creation**

⟷ **R&D Convergence**

⟷ **R&D Reliability**

- Public participation

- Easy R&D
- Fast R&D

- Connected R&D

- Reproducibility
- Error probing

**【Driving Wheels】**

**AI· Big Data· ICT**

**(Technology)**

**R&D Data Open· Sharing**

**Culture: Data Open & Sharing**

**(Policy)**

- R&D e-Transformation
- ICT utilization

- Collaboration
- Overcoming barriers

**【Base】**

**S/W Platform**

**Infrastructure**

**Fostering Community**

**Data Management**

**Systemization**

- Unrestricted data access

- Data acquisition
- Digitalization

- Main players

- Data tracking

- Regulations
- Programs

# Stepwise Implementation: from part to all



➲ **Covering entire data management, linking data repositories, providing easy and seamless access across various R&D**

➲ **Enabling Open Science Service by stepwise implementation**



**4 Pilot Projects**
**(launched in 2018)**

1. **Genome Data**
2. **Material Data**
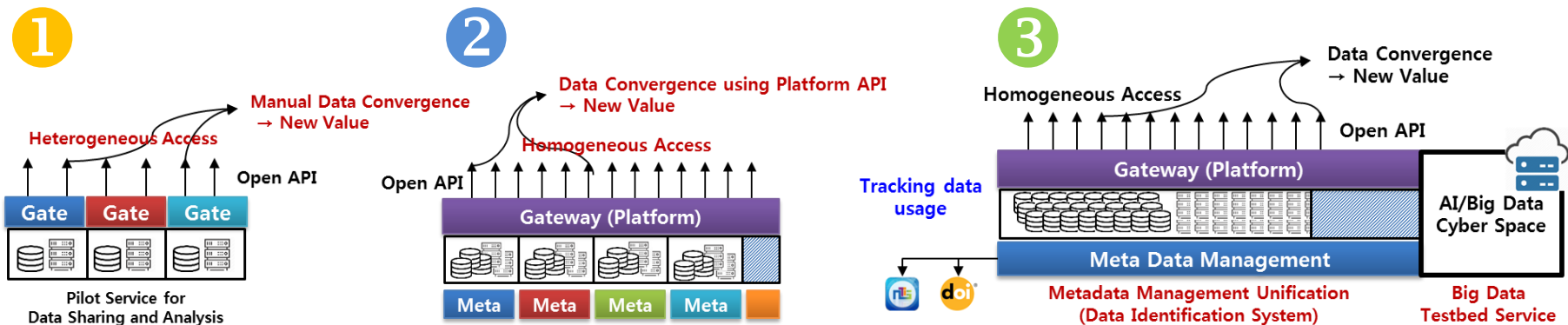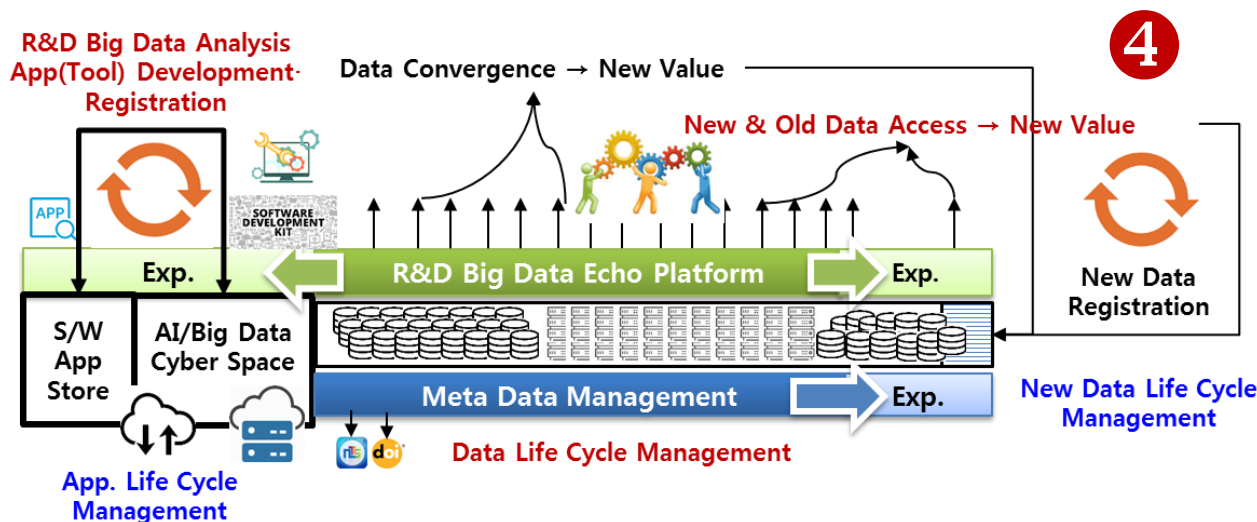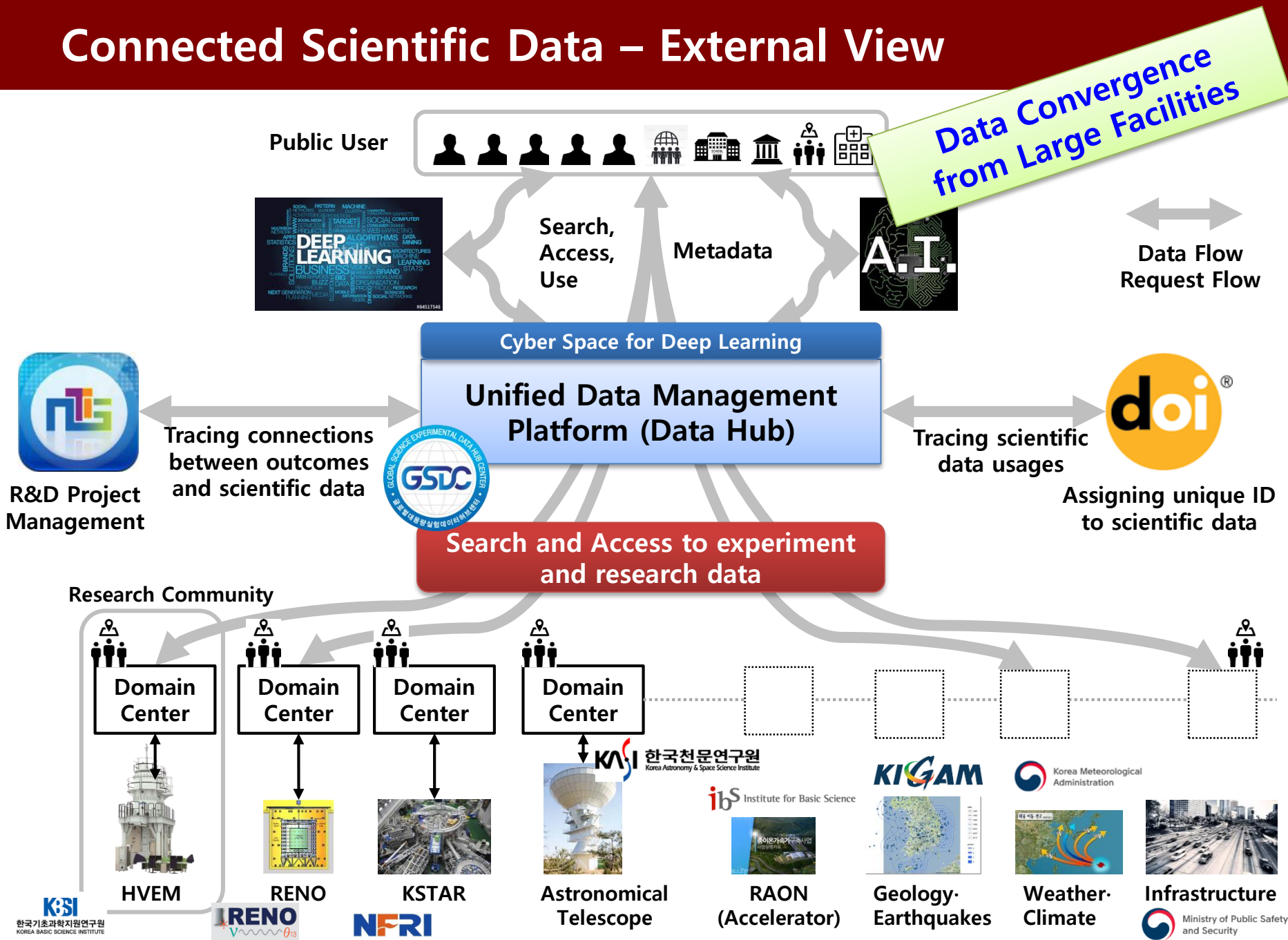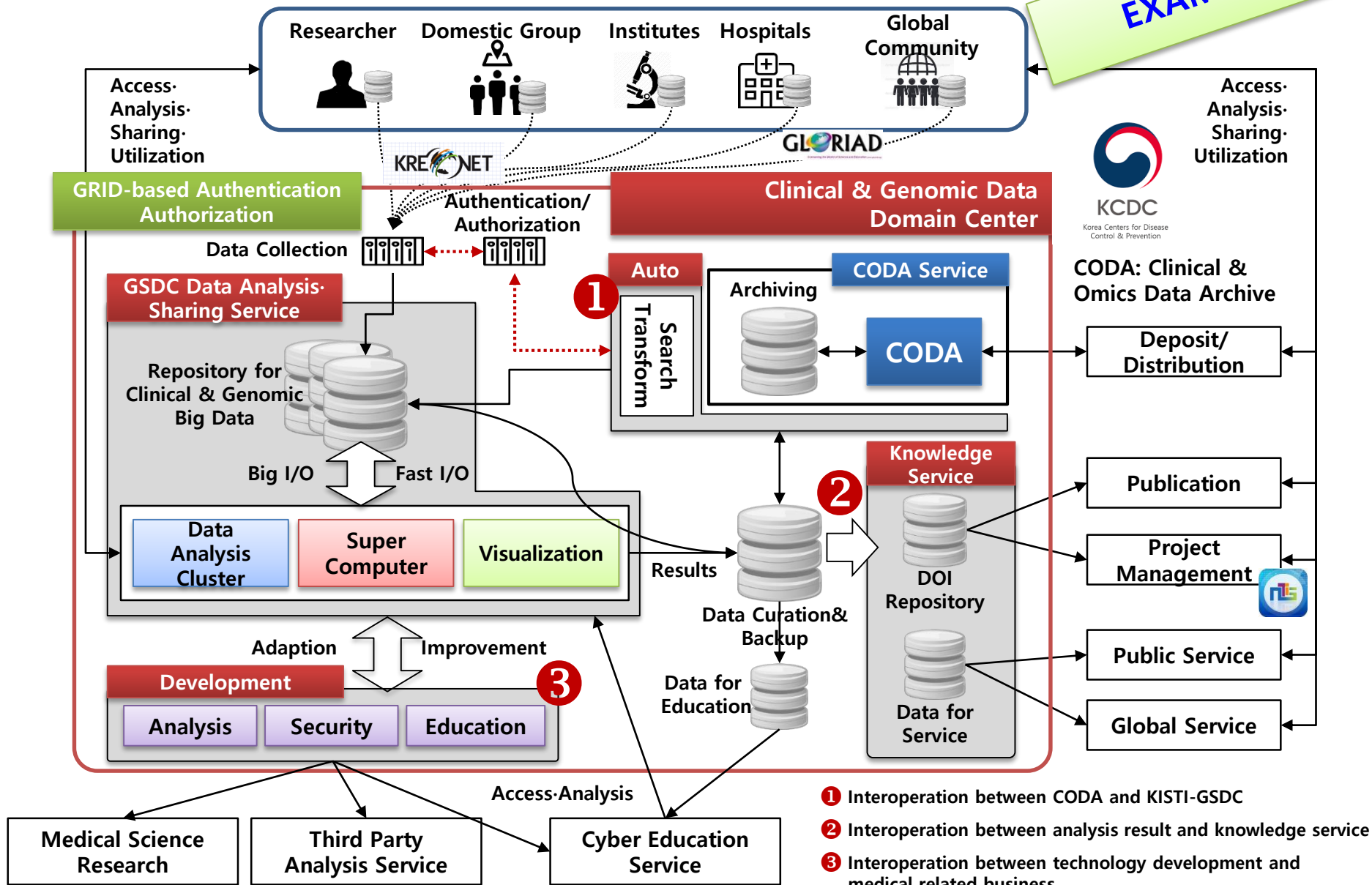3. **Data Convergence from Large Facilities (KISTI-GSDC)**
4. **AI/Big Data**

EXAMPLE

Researcher　Domestic Group　Institutes　Hospitals　Global Community

Access·Analysis·Sharing·Utilization

Access·Analysis·Sharing·Utilization

GLORIAD

KRE●NET

**GRID-based Authentication Authorization**

Authentication/Authorization

**Clinical & Genomic Data Domain Center**

KCDC
Korea Centers for Disease Control & Prevention

**CODA: Clinical & Omics Data Archive**

Data Collection

**GSDC Data Analysis·Sharing Service**

**Auto**

❶

Search Transform

**CODA Service**

Archiving

**CODA**

**Repository for Clinical & Genomic Big Data**

Big I/O　　Fast I/O

**Deposit/ Distribution**

| Data Analysis Cluster | Super Computer | Visualization |

Results

**Knowledge Service**

❷

DOI Repository

**Publication**

**Project Management**

Data Curation& Backup

Adaption　Improvement

**Development**

❸

Data for Education

Data for Service

**Public Service**

**Global Service**

| Analysis | Security | Education |

Access·Analysis

**Medical Science Research**

**Third Party Analysis Service**

**Cyber Education Service**

❶ Interoperation between CODA and KISTI-GSDC

❷ Interoperation between analysis result and knowledge service

❸ Interoperation between technology development and medical related business

⊃ **Centralized data repository** is used for data collection and access

⊃ Central data repository will be **connected with Open Data Platform**, which is also being developed by KISTI.

⊃ Expect each site has **own repositories interoperable with the Open Data Platform** in future



KBSI TEM*

Pohang Accelerator

Extra Equipment

Large-scale Facility

Experiment

Data Generation

Data Convergence

**Integrative Structural Biology 4 Data Convergence Programs**

**(Structural Biology Community)**

**Data Collection**

**(KISTI)**

**Data Analysis**

**(KISTI)**

Education

**Data Challenge School (Data Convergence)**

**(KISTI, KBSI, Structural Biology Community)**

***Transmission Electron Microscope**

# Summary:
## *Response to Requests*

# Summary

**Open Science**

➲ Highlighted openness of access, collaborations and data

➲ Believed to give benefits to scientific community:
Connected Scientific Data → Connected Science

**Data-driven R&D**

➲ Research paradigm is being shifted to data-driven scientific discovery

➲ Data and infrastructure are the key in scientific discovery

**Open Data Platform**

➲ Importance of Open Data is recognized in government level and the TF has launched four pilot projects this year (2018)

➲ Data repositories, setup by pilot projects, will be connected to the platform which makes them interoperable as a final goal.

**Implementation of Open Data Platform is not an easy task…
The pilot projects are not big scale…**

**But, such a trial is a Big Step moving toward to
making Open Science a reality in Korea**

# Responses to Requests

**Identify your initiative as a stakeholder ...**

➲ Policymakers(Government), researchers, ICT specialists are stakeholders.

➲ Policymakers wants to make data, generated from public funding, <u>traceable through Open Data Platform</u>.

➲ Researchers wants to use <u>well managed ICT infrastructure </u>for data analysis and data sharing, in order to accelerate R&D process.

➲ Role of ICT specialized institute like KISTI, expected from government, is <u>well aligned with its mission </u>– promoting Science using ICT technologies.

**Analyze how your initiative addresses a multiplicity of different cultural contexts ...**

➲ Respect differences in R&D domains and encourage scientific communities to make <u>own standard data management plan </u>agreed in community members.

➲ Structural biology community, for example, is developing a guideline including data naming convention, sharing, accessing and management policy.

**Assess the proactive role of your initiative with respect to interoperability ...**

➲ Rather than making data repositories interoperable directly(1-to-1), <u>indirect interoperability</u> is being considered through the Open Data Platform.

➲ Such an approach helps <u>to reduce the burden of standization </u>for all R&D. Interesting data can be accessed through OpenAPI provided by the platform. Data convergence and all R&D activities with data are left to scientist.

Thank you.