# Implementation Challenges and Standards Opportunity for FAIR Principles

*Wo Chang*

**Digital Data Advisor**
**ISO/IEC JTC 1/SC42/WG2 Big Data, Convenor**
**IEEE Big Data Governance and Metadata**
**Management WG, Chair**
**wchang@nist.gov**
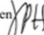
**October 27, 2020**

# Open Data Initiative (2013)



EXECUTIVE OFFICE OF THE PRESIDENT
OFFICE OF SCIENCE AND TECHNOLOGY POLICY
WASHINGTON, D.C. 20502

February 22, 2013

MEMORANDUM FOR THE HEADS OF EXECUTIVE DEPARTMENTS AND AGENCIES

FROM: John P. Holdren
Director

SUBJECT: Increasing Access to the Results of Federally Funded Scientific Research

1. **Policy Principles**

The Administration is committed to ensuring that, to the greatest extent and with the fewest constraints possible and consistent with law and the objectives set out below, the direct results of federally funded scientific research are made available to and useful for the public, industry, and the scientific community. Such results include peer-reviewed publications and digital data.

Scientific research supported by the Federal Government catalyzes innovative breakthroughs that drive our economy. The results of that research become the grist for new insights and are assets for progress in areas such as health, energy, the environment, agriculture, and national security.

Access to digital data sets resulting from federally funded research allows companies to focus resources and efforts on understanding and exploiting discoveries. For example, open weather data underpins the forecasting industry, and making genome sequences publicly available has spawned many biotechnology innovations. In addition, wider availability of peer-reviewed publications and scientific data in digital formats will create innovative economic markets for services related to curation, preservation, analysis, and visualization. Policies that mobilize these publications and data for re-use through preservation and broader public access also maximize the impact and accountability of the Federal research investment. These policies will accelerate scientific breakthroughs and innovation, promote entrepreneurship, and enhance economic growth and job creation.

The Administration also recognizes that publishers provide valuable services, including the coordination of peer review, that are essential for ensuring the high quality and integrity of many scholarly publications. It is critical that these services continue to be made available. It is also important that Federal policy not adversely affect opportunities for researchers who are not funded by the Federal Government to disseminate any analysis or results of their research.

To achieve the Administration's commitment to increase access to federally funded published research and digital scientific data, Federal agencies investing in research and development must have clear and coordinated policies for increasing such access.

**The White House**
Office of the Press Secretary

For Immediate Release                    May 09, 2013

# Executive Order -- Making Open and Machine Readable the New Default for Government Information

EXECUTIVE ORDER

- - - - - - -

MAKING OPEN AND MACHINE READABLE THE NEW DEFAULT
FOR GOVERNMENT INFORMATION

By the authority vested in me as President by the Constitution and the laws of the United States of America, it is hereby ordered as follows:

Section 1. General Principles. Openness in government strengthens our democracy, promotes the delivery of efficient and effective services to the public, and contributes to economic growth. As one vital benefit of open government, making information resources easy to find, accessible, and usable can fuel entrepreneurship, innovation, and scientific discovery that improves Americans' lives and contributes significantly to job creation.

Decades ago, the U.S. Government made both weather data and the Global Positioning System freely available. Since that time, American entrepreneurs and innovators have utilized these resources to create navigation systems, weather newscasts and warning systems, location-based applications, precision farming tools, and much more, improving Americans' lives in countless ways and leading to economic growth and job creation. In recent years, thousands of Government data resources across fields such as health and medicine, education, energy, public safety, global development, and finance have been posted in machine-readable form for free public use on Data.gov. Entrepreneurs and innovators have continued to develop a vast range of useful new products and businesses using these public information resources, creating good jobs in the process.

# Interagency Technical Advisory Group (Dec. 2013)

To provide a forum for Federal agency coordination on operational requirements and insights on how to maximize access to scientific and technical data. Members are Federal employees participating in their individual capacity as subject matter experts and providing their own perspectives from a range of agency and entity settings including:

- NIST (Chair)
- Census
- NIH/NCI
- Smithsonian
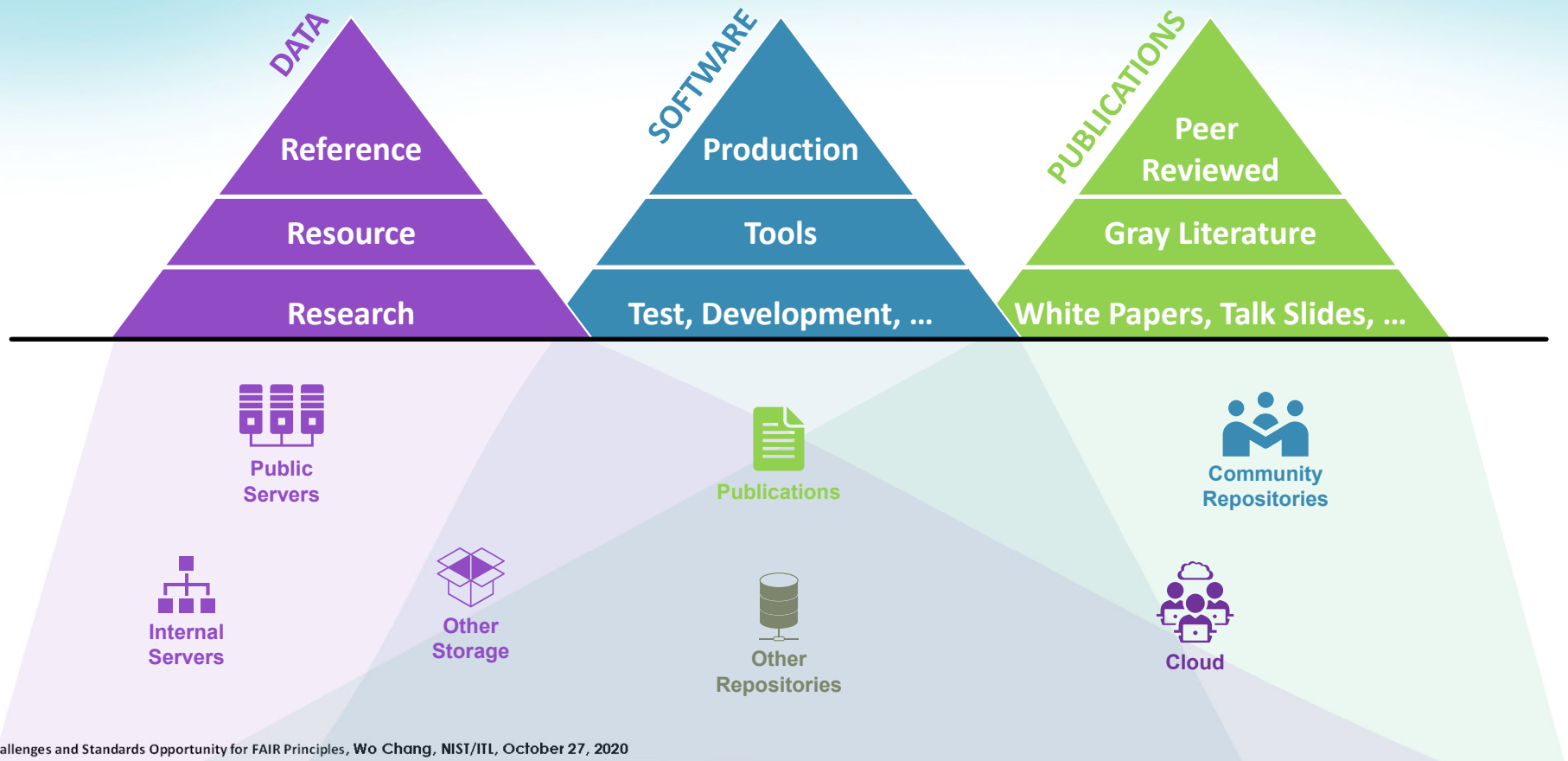- DOE
- NARA
- Treasury
- USDA

**Engage with:**

RDA — RESEARCH DATA ALLIANCE
- PID Information Type WG
- Data Type Registry WG
- Data Fabric IG

ISO IEC
- SC 32 / WG2 Metadata

# Challenging Problem: Concept Model

# Challenging Problem: Logical Model



**Extended Metadata**

**Specialized Metadata**

**Minimum Metadata**

**Primary Digital Objects**

Data.
Extended

Software.
Extended

Pub.
Extended

Software.
SpecialMeta

Data.
SpecialMeta

Pub.
SpecialMeta

Data.
Metadata

Software.
Metadata

Pub.
Metadata

Data

Software

Pub

# NIST Common Access Platform (CAP, 2014)

**Goal:**

Develop an interoperable data infrastructure that is scalable to enable automatic data mashups between heterogeneous datasets from various domains without worrying about the data source and structure.
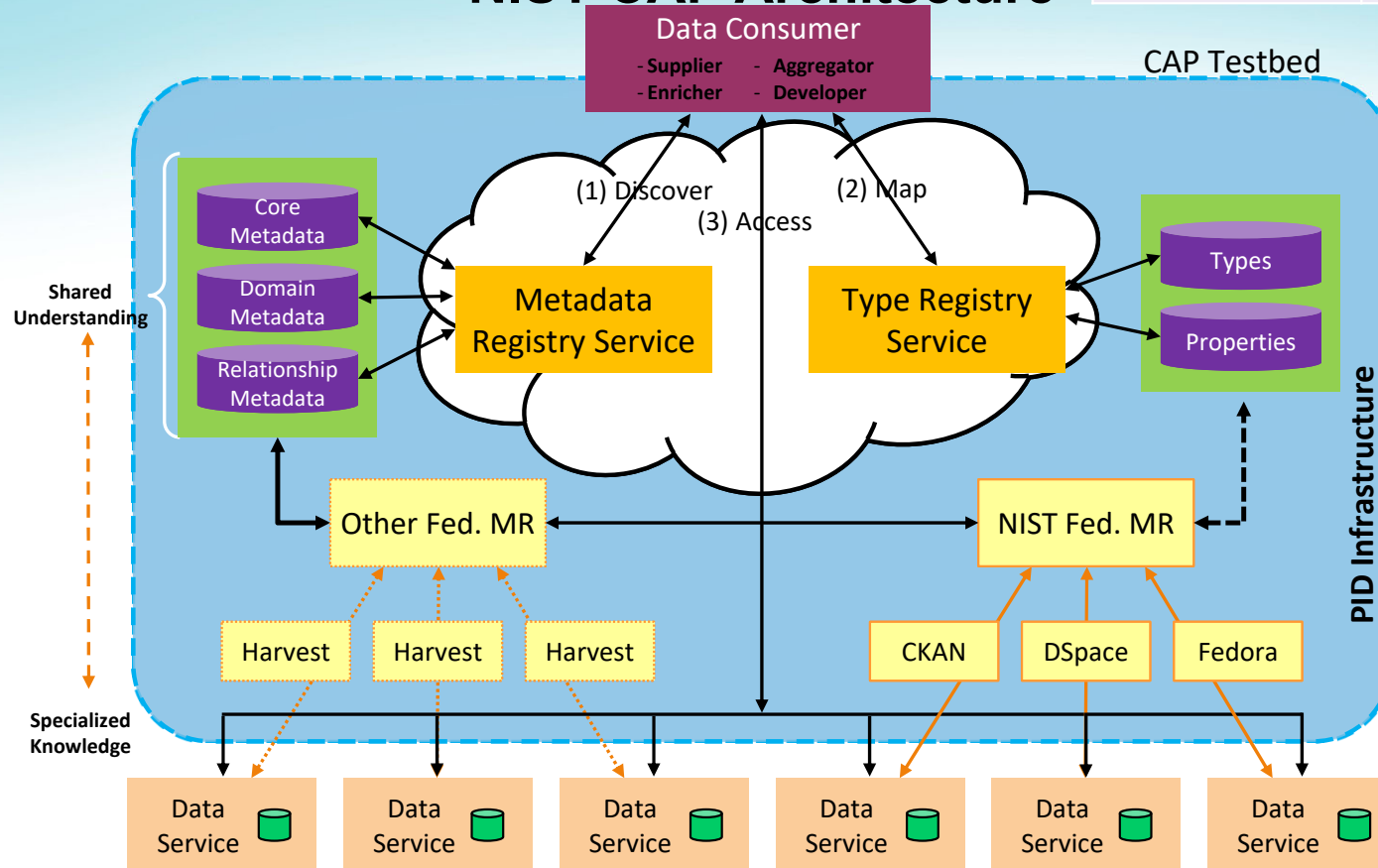
**Approach:**

Provide basic data infrastructure using persistent identifiers to enable:
* Standard **metadata registry** for data <u>discovery</u> using a **machine-readable format**
* Standard **data type registry** that enables data <u>consumption</u> using a **machine actionable format**

**(without standard data type registry, the data is not easily interoperable and re-usable)**

# NIST CAP Architecture

| FAIR | CAP |
|------|-----|
| Findability | Discover |
| Accessibility | Discover |
| Interoperability | Map |
| Reusability | Access |

**Data Consumer**
- Supplier    - Aggregator
- Enricher    - Developer

CAP Testbed

(1) Discover    (2) Map
(3) Access

Shared Understanding

Core Metadata
Domain Metadata
Relationship Metadata

Metadata Registry Service

Type Registry Service

Types
Properties

PID Infrastructure

Other Fed. MR

NIST Fed. MR

Harvest    Harvest    Harvest

CKAN    DSpace    Fedora

Specialized Knowledge

Data Service    Data Service    Data Service    Data Service    Data Service    Data Service

Implementation Challenges and Standards Opportunity for FAIR Principles, Wo Chang, NIST/ITL, October 27, 2020
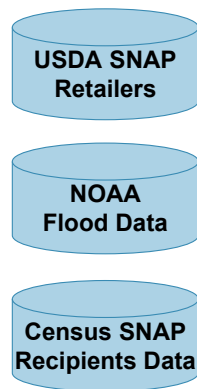
# Implementation Challenges

# FindIT ConnectIT Prototype – Based on NIST CAP Arch.

**Use Case:**
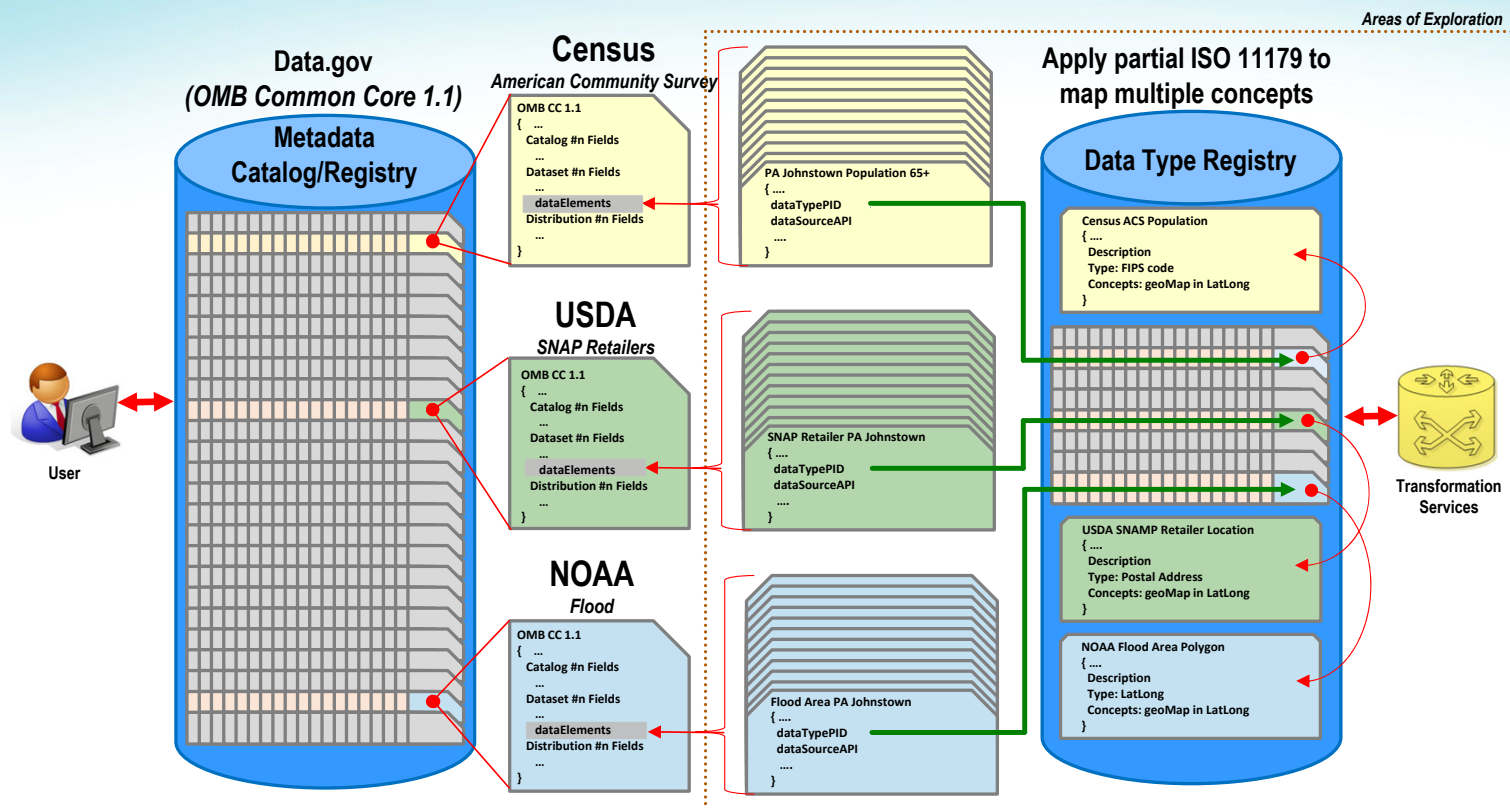**FEMA Emergency Preparedness**

In anticipation of a flood FEMA seeks to understand whether or not a State will require assistance, and what kind of resources will be needed.

Analysts must rapidly identify and assemble data for presentation in the dashboard

**Multiple Data Sources**

USDA SNAP Retailers

NOAA Flood Data

Census SNAP Recipients Data



Multiple data sets to access

Populating GEO shapes with data

Set of variable & context driven questions

_Flood area shown in Johnstown, PA with total population and SNAP retailers displayed_

At risk population: Under 18; Over 65?

Where is the flood impact area?

How many people in impacted area?

What is the predominant language?

What percentage is below the poverty level?

What are the closest SNAP Retailer NOT in Flood area?

What SNAP Stores are located in the Flood area?

Match Bounds    Census B01003_001E: Total Po
5,058

# NIST Efforts to Develop CAP Arch. For Interoperable Data Structure

**1**

Discover (Findability and Accessibility) is doable but mapping (Interoperability) different information models is very hard (e.g., the way geographic information was encoded and captured by Census vs USDA). Once mapping can be done correctly, access (reusability) is doable. We only hard coded the mapping but we need structure and automatic mapping.

**7**

Others…

**6**

How should the needed metadata be created and by whom?

**Lessons Learned: very hard problem**

1
2
3
4
5
6
7

**2**

Parameterized, summarized, or query-based data sets are more challenging than static ones to work with and describe using data types.

**3**

How to represent complex/composite data types so the semantics and technical details can be queried by automated search engines?

**4**

What level of granularity should be exposed when complex/composite datatypes exist?

**5**

Specificity vs. re-usability of data type

# Standards Opportunity

# IEEE Big Data Governance and Metadata Mgt (BMGMM) WG

https://standards.ieee.org/project/2957.html#Working (Sept. 2020)

**Goal:**

To enable data integration/mashup among heterogeneous datasets from diversified domain repositories to make data discoverable, accessible, and usable through a machine-readable and actionable standard data infrastructure.

Big data provides key characteristics in *Volume*, *Velocity*, *Variety*, and *Variability*, commonly referred to as the *V*s of Big Data. BDGMM is focusing on data from a single source or *Varieties* of data from multiple sources to create an integrated data source for analytics and AI machine learning consumption.

From the new global Internet Big Data economy opportunity in Internet of Things, Smart Cities, and other emerging technical and market trends, it is critical to have a standard reference architecture for Big Data Governance and Metadata Management to support the scalable FAIR (Findability, Accessibility, Interoperability, Reusability) foundation principles.

# IEEE Big Data Governance and Metadata Mgt (BMGMM) WG
## https://standards.ieee.org/project/2957.html#Working (Sept. 2020)

**Approach:**

Apply metadata for scalable and machine actionable to
**FAIR** (Findability, Accessibility, Interoperability, Reusability) principles

**BMGMM**

**S F A I R**

**Scalable** — Utilizes machine-readable and machine actionable formats

**Governance** — Provides authoritative, control, and shared decision making over the management of data assets

**Metadata-based**

**Metadata** — Provides PID-based standard data infrastructure to enable:
* Catalog Registry for data discovery using a machine-readable format
* Data Type Registry for data consumption using a machine actionable format

# IEEE BDGMM Roadmap White Paper (July 2020)

# IEEE BDGMM Roadmap White Paper (July 2020)

https://standards.ieee.org/content/dam/ieee-standards/standards/web/governance/iccom/bdgmm-standards-roadmap-2020.pdf

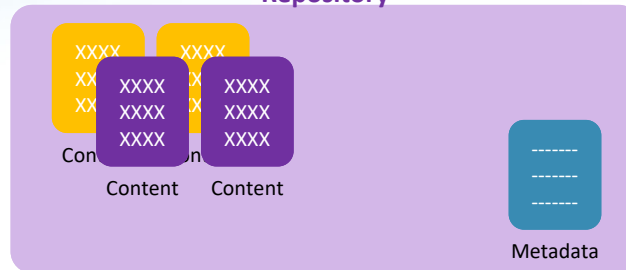# Potential Standard BDGMM Reference Architecture

# Catalog Registry



Content owner adds content to a repository.
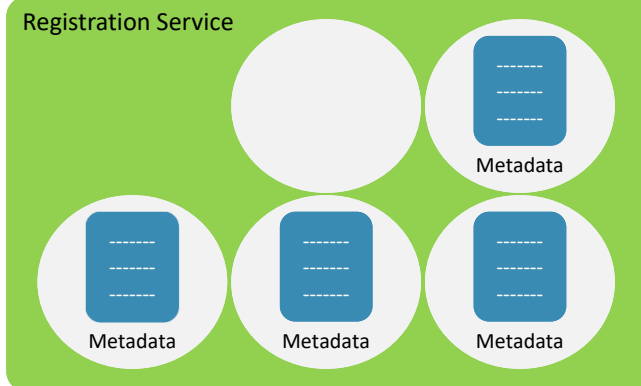
Content

## Repository

Content  Content

Metadata

Metadata for that content is generated in the Repository and pushed, via a Registration Service, into Metadata Registry, creating a digital object.

## Catalog Registry

Registration Service

Metadata

Metadata  Metadata  Metadata

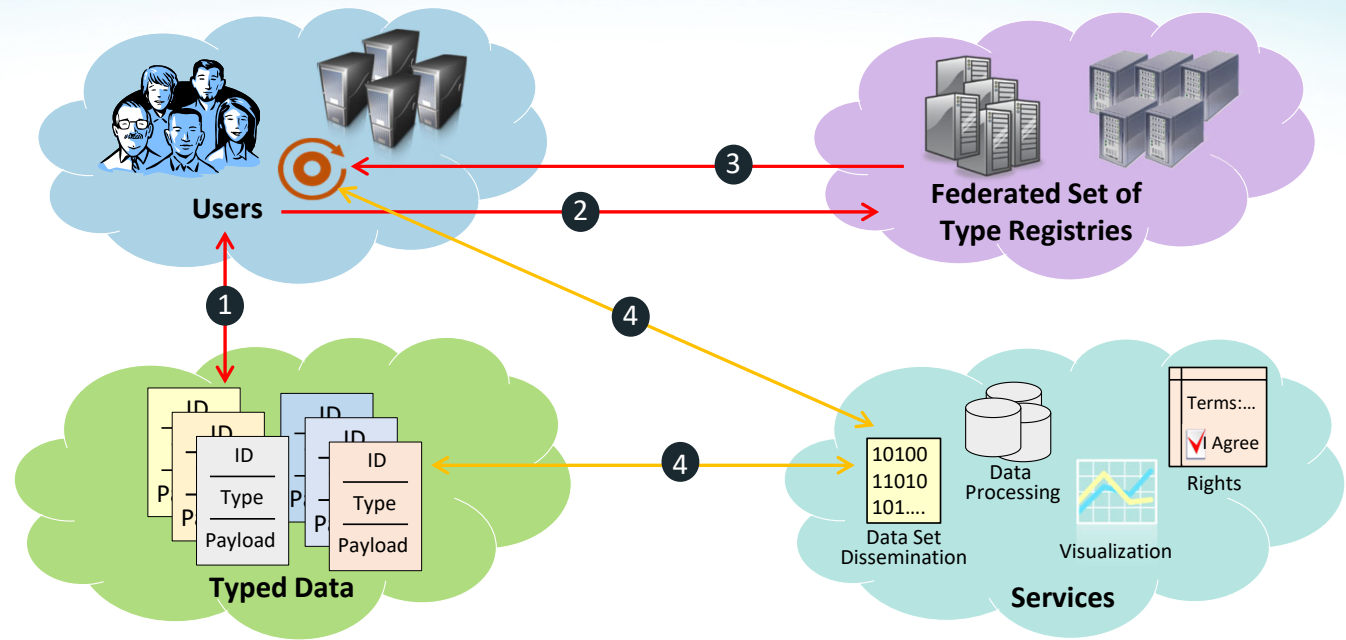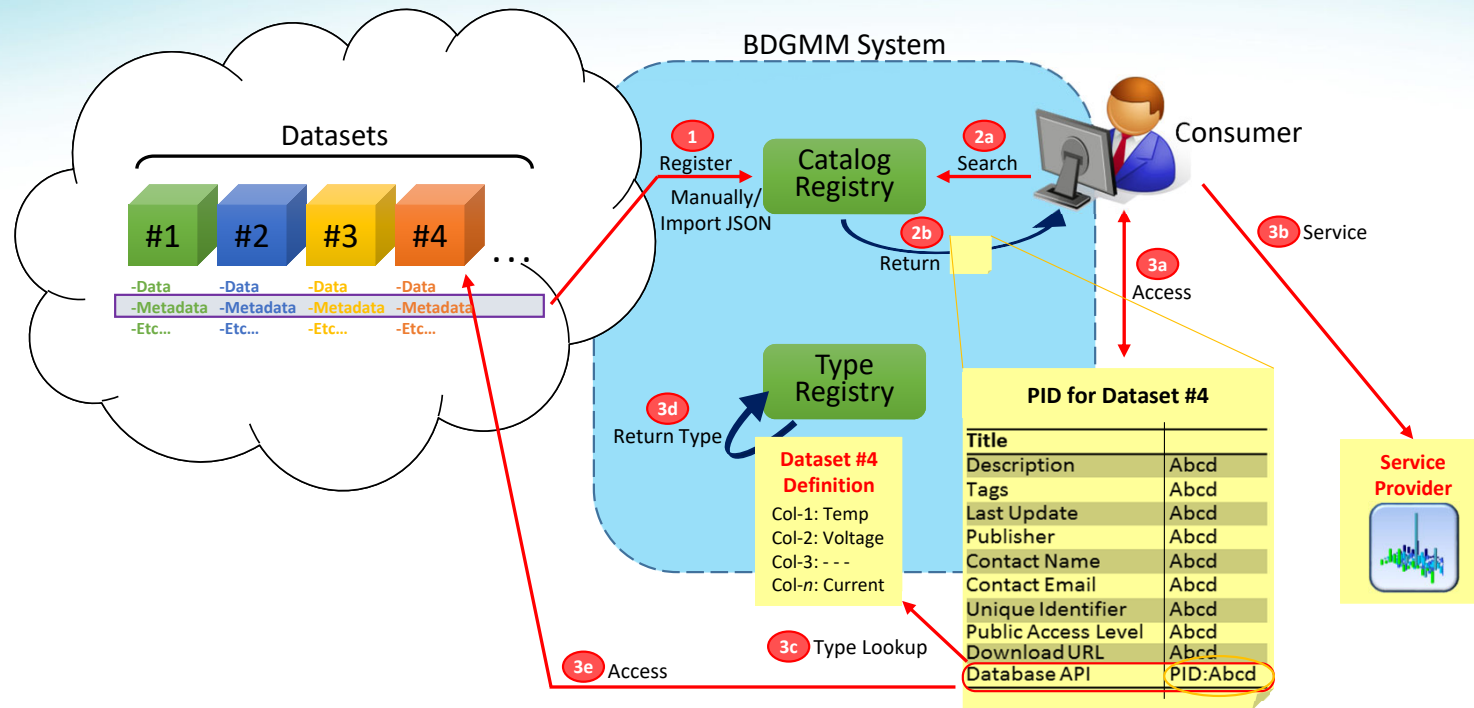Information Management Service

Discovery Service

The Catalog Registry provides Information Management and Discovery Services for users.

# Data Types Registry

1. Client (process or people) encounter data of an unknown type

2. Resolved the Type to Type Registry

3. Response includes type definitions, relationships, properties, and possibly service pointers. Response can be used locally for processing, or, optionally

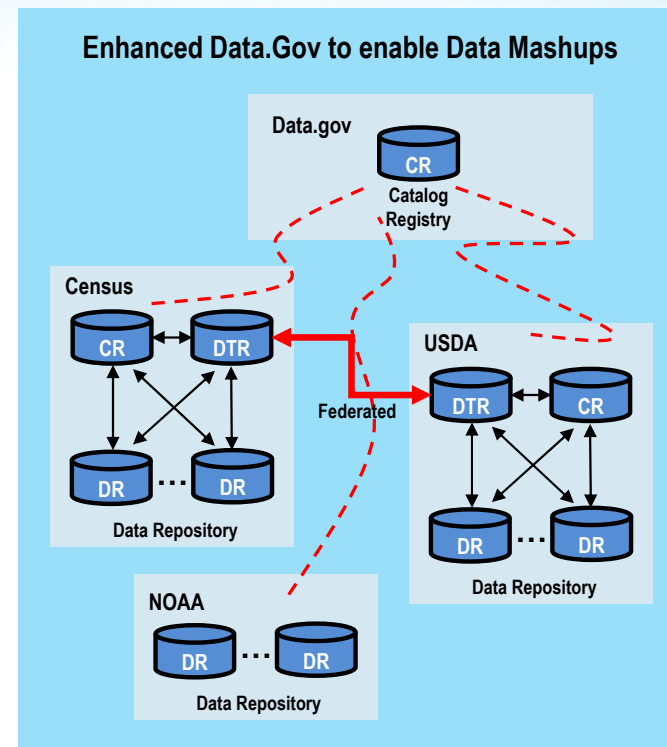4. Typed data or reference to typed data can be sent to service provider

Users

Federated Set of Type Registries

ID
Type
Payload

Typed Data

10100
11010
101....

Data Set Dissemination

Data Processing

Visualization

Terms:...
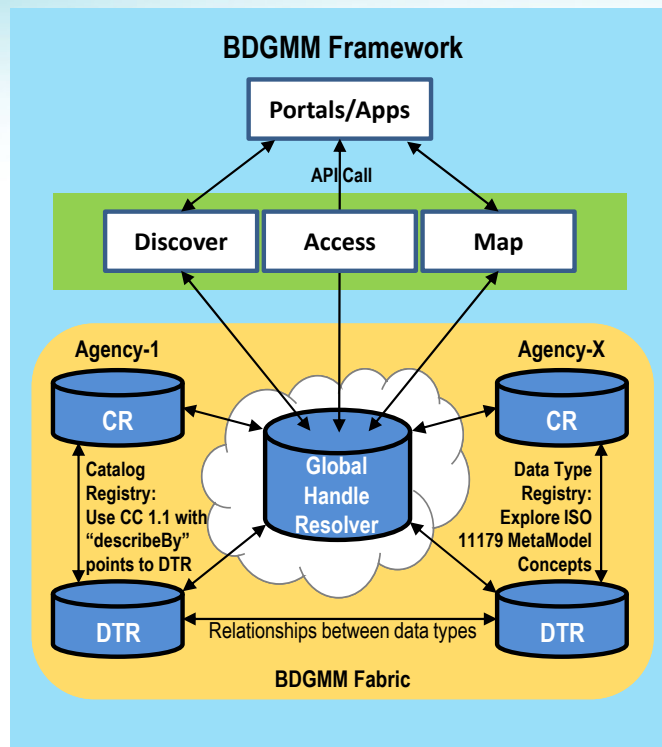☑ I Agree
Rights

Services

# BDGMM Workflow

# Operational Model:
## Interested parties could deploy federated Metadata Registries and Data Type Registries

# Recommendations for Standards Development

**Big Data Governance Management**

Description: Streamlining governance of IT and data are essential for organizations to meet the challenges of the digital era and whoever could govern and manage such resources effectively can reduce the organization's burden and maximize the customers' needs.
Key recommendations may include the following:

- Adopt/develop standard interface for human readable and machine-actionable to access corporate data catalog that provides detail description and linkage to datasets and their usage.
- Utilize best practices standard networking protocols to support open and multi-levels of security for accessing datasets for (a) end-to-end over the net, (b) at repository, (c) at dataset, (d) at data record/element, etc.
- Adopt/develop extensible PID with scalable resolver to handle massive PID resolution.
- Adopt/develop revision control on datasets with backward and forward compatibility.

# Recommendations for Standards Development

**Big Data Metadata Management**

Description: Supporting diversified metadata schemas and models for various datasets would be essential to organizations to meet the ever-growing customers' needs and whoever could manage these metadata cohesively across all datasets can reduce corporate burden. In addition, providing computable object workflow functionality between data elements of various datasets would be a great additional service to customers for monitoring events, trigger certain conditions, etc.
Key recommendations may include the following:

- Utilize best practices standard metadata as much as possible to capture precise description, data types, properties, unit of measurement, characteristics, etc. for given data elements.
- Adopt/develop standard federated metadata registries to support catalogs and types registries.
- Adopt/develop standard interface to support online data element definition.
- Adopt/develop standard computable object workflow functionality to trigger certain conditions including privacy and ethical issues in datasets.

# Recommendations for Standards Development

**Big Data Integration Framework**

Description: Supporting data integration or data mashup among heterogeneous datasets would be critical for analytics to discover new patterns or knowledge and whoever could manage these rich resources effectively would gain much insights into better decision making.
Key recommendations may include the following:

- Adopt/develop standard interface to access data at record level regardless of data at rest or in motion (streaming) from public or secured repositories.
- Adopt/develop standard scalable metadata model to map individual data model across heterogeneous datasets from multiple data sources.

# Recommendations for Standards Development

**Persistent Identifier Framework**

Description: Tagging datasets as persistent identifier (PID) at any level (dataset itself, data record, data element, data type, data property, etc.) would be essential in enabling Findability, Accessibility, Interoperability, and Reusability. Having a standard PID framework would enable interoperability among all heterogeneous datasets across all data repositories.
Key recommendations may include the following:

- Adopt/develop standard PID framework that provides organizational namespace with flexible and extensible structure to meet organizational needs.
- Adopt/develop scalable PID resolver to handle massive PID resolution in a millisecond time interval.

**Questions?**

Please contact: wchang@nist.gov