# MDR Development Project for materials science

- National Institute for Materials Science, Japan
- Cottage Labs, UK
- AntLeaf, UK
- iGroup, Taiwan



The MDR team: developers, publishers, researchers - at NIMS Library
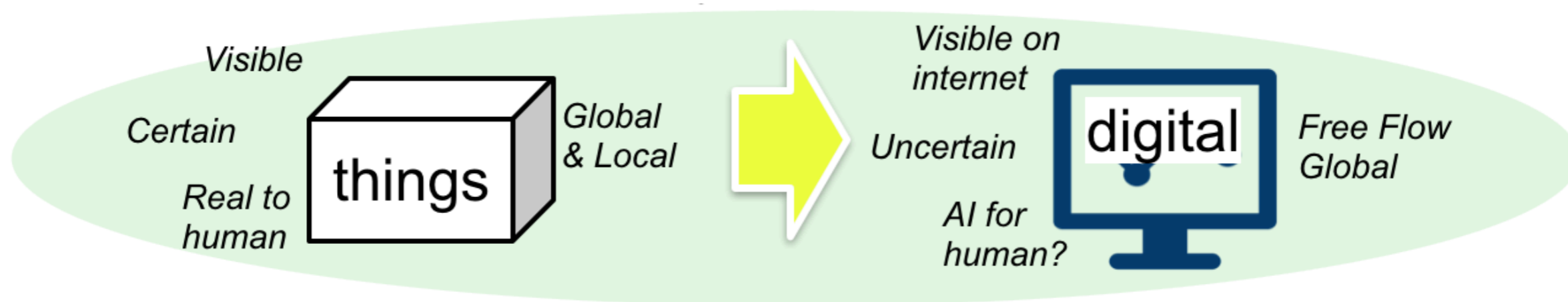
Researchers

Publishers    Developers

Engineers

# 1. Context: NIMS & the MDR

Mikiko Tanifuji

# A landscape of research data – G20 Digital Economy

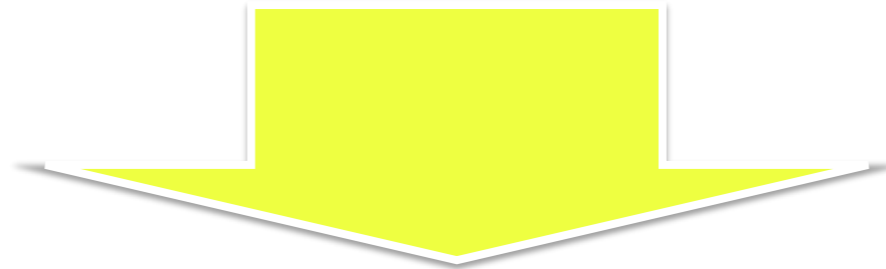- G20 - Trade and Digital Economy, June 8, 2019
  - Human Centric Future Society



- "Data Free Flow with Trust" (DFFT concept)
  - Accumulate data for human society
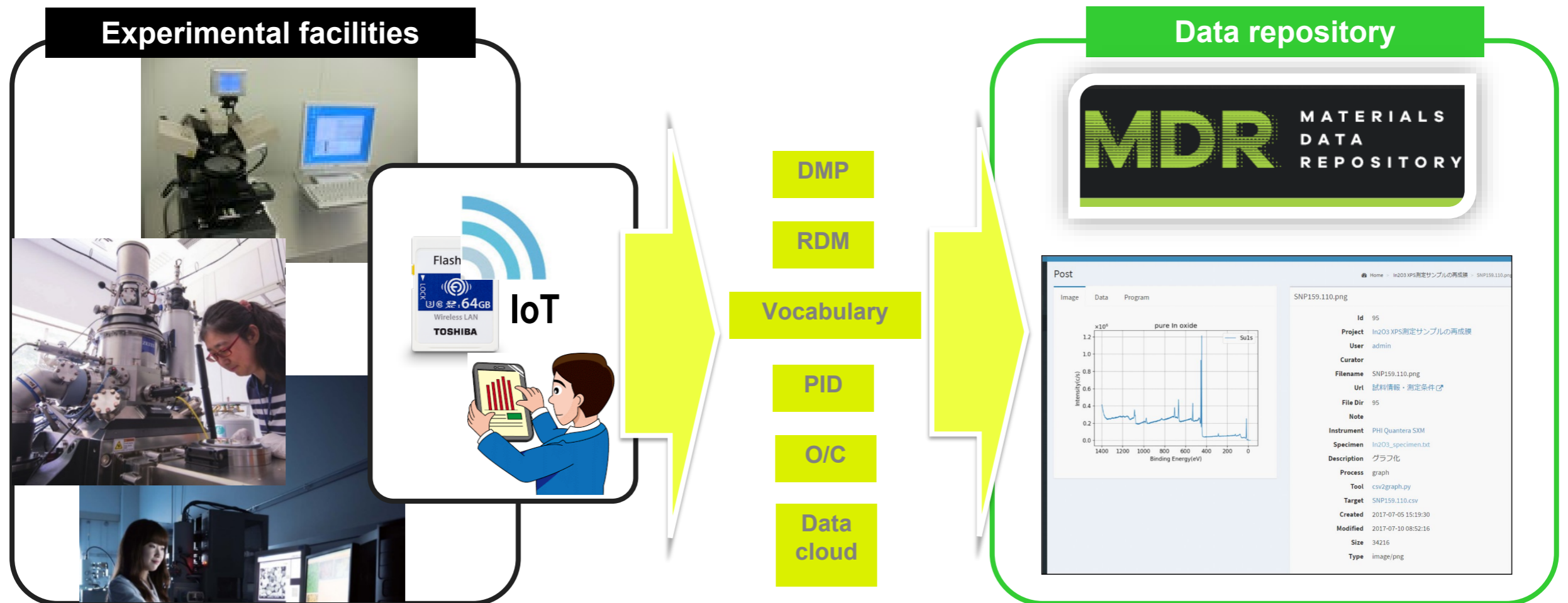  - Appropriate data management and global consensus for how-to-use

# MDR Development Project – Why?

1. A new trend "Data-driven science" >> data science/scientists
2. Not just "machine-readable", move to machine-actionable  >> really FAIR
3. Incentives of "machine-learning" >> must WebAPI, with metadata
4. Not just a database  >> semantic-aware database
5. Not just an archive >> metadata, machine-readable formats, analytics tools

1. Next Generation Repository (NGR) must have machine-actionable data
2. NGR must have researchers' trust-based quality data
3. NGR should/could be repository-tenant concept  Example: res project repository

# MDR Development Project - What?

# MDR - a FAIR system of Materials Data Platform



**2019 -**

**NIMS service**
## DCS
Data Curation System

**NIMS service**
## IoT Data
IoT Data Transferring System

**NIMS service**
## LabNote
Online Lab Notebooks

**NIMS service**
## Analytics
High performance computer system

**Public service**

**Public service**

**NIMS service**

**2020 -**

**Public service**
## VocWiki
Vocabulary for Data Management

**NIMS service**
## RDM
Research Data Management

**NIMS service**
## Single Sign-on
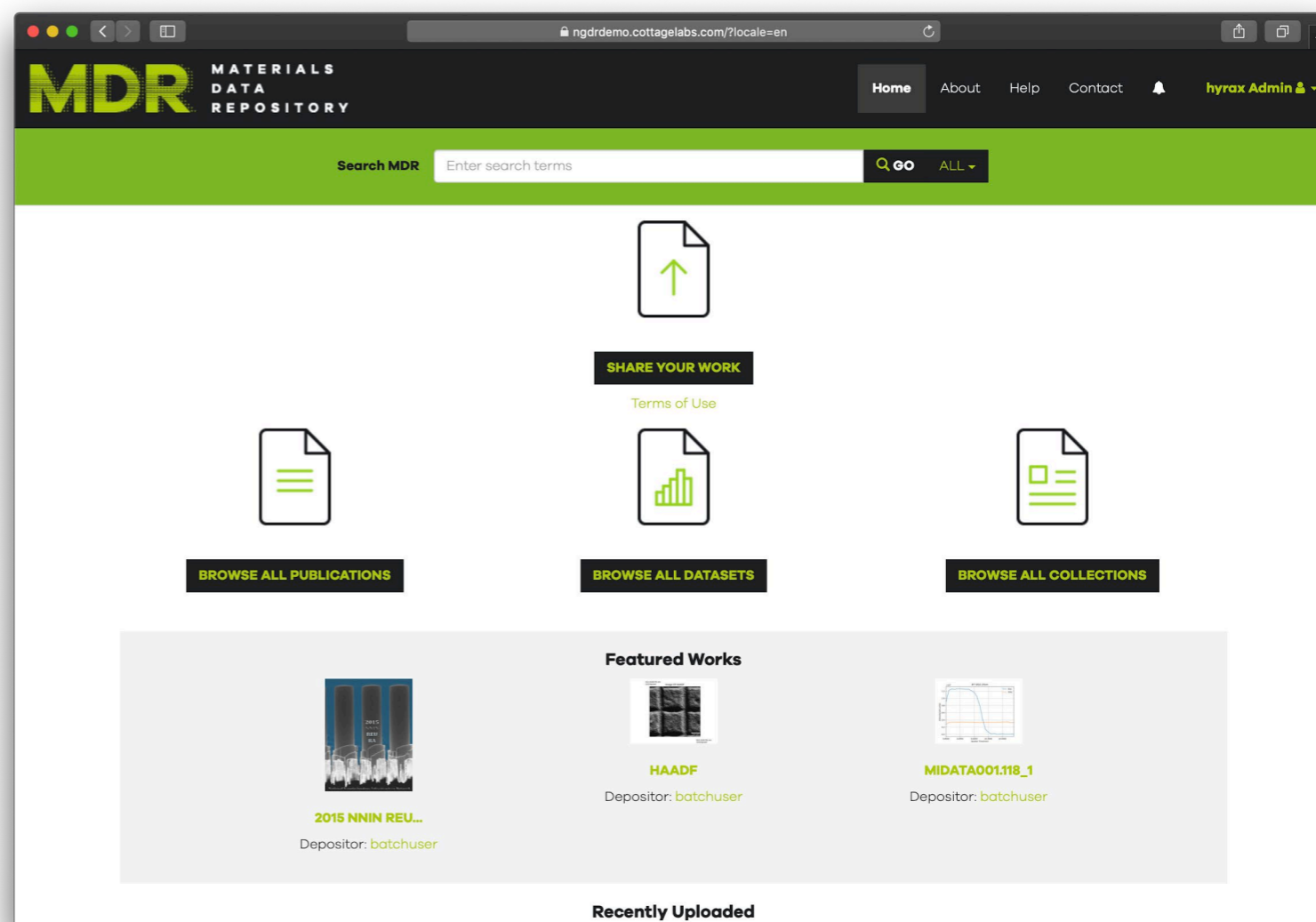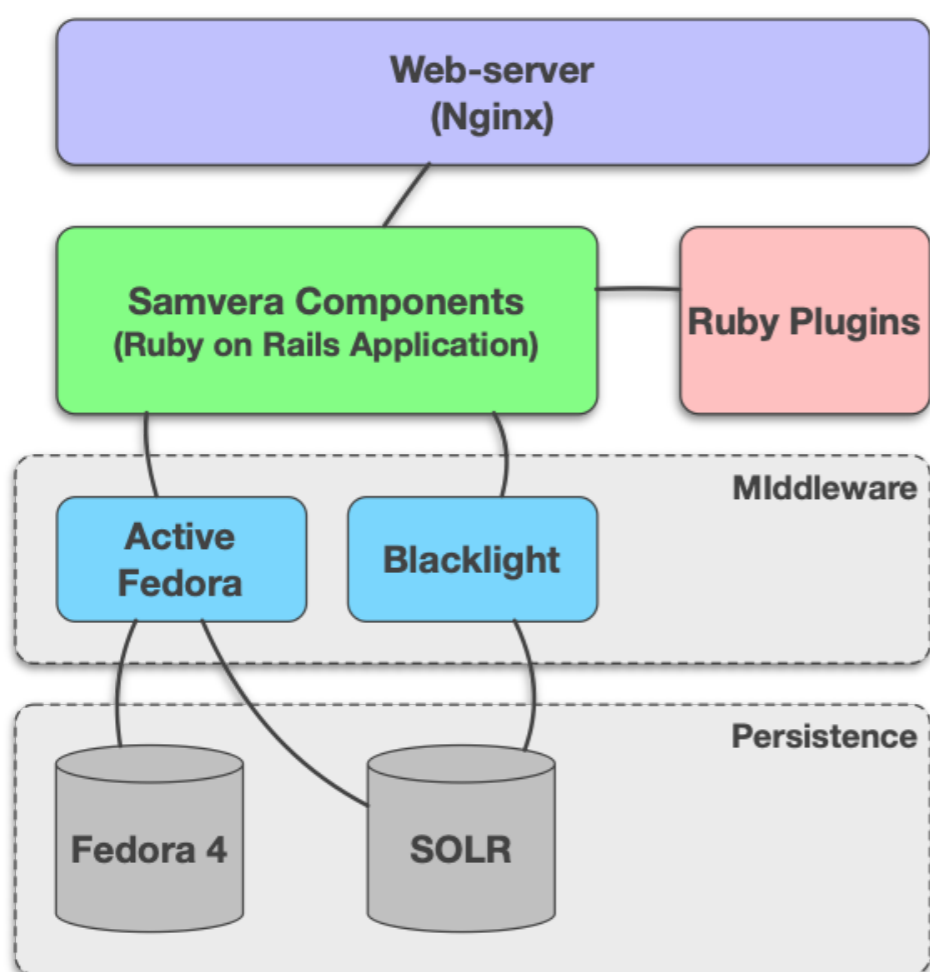A gateway to all data services

**Data deposit | Data deposit via IoT | Data search | Data download | Data visualizations | Data analytics & Informatics**
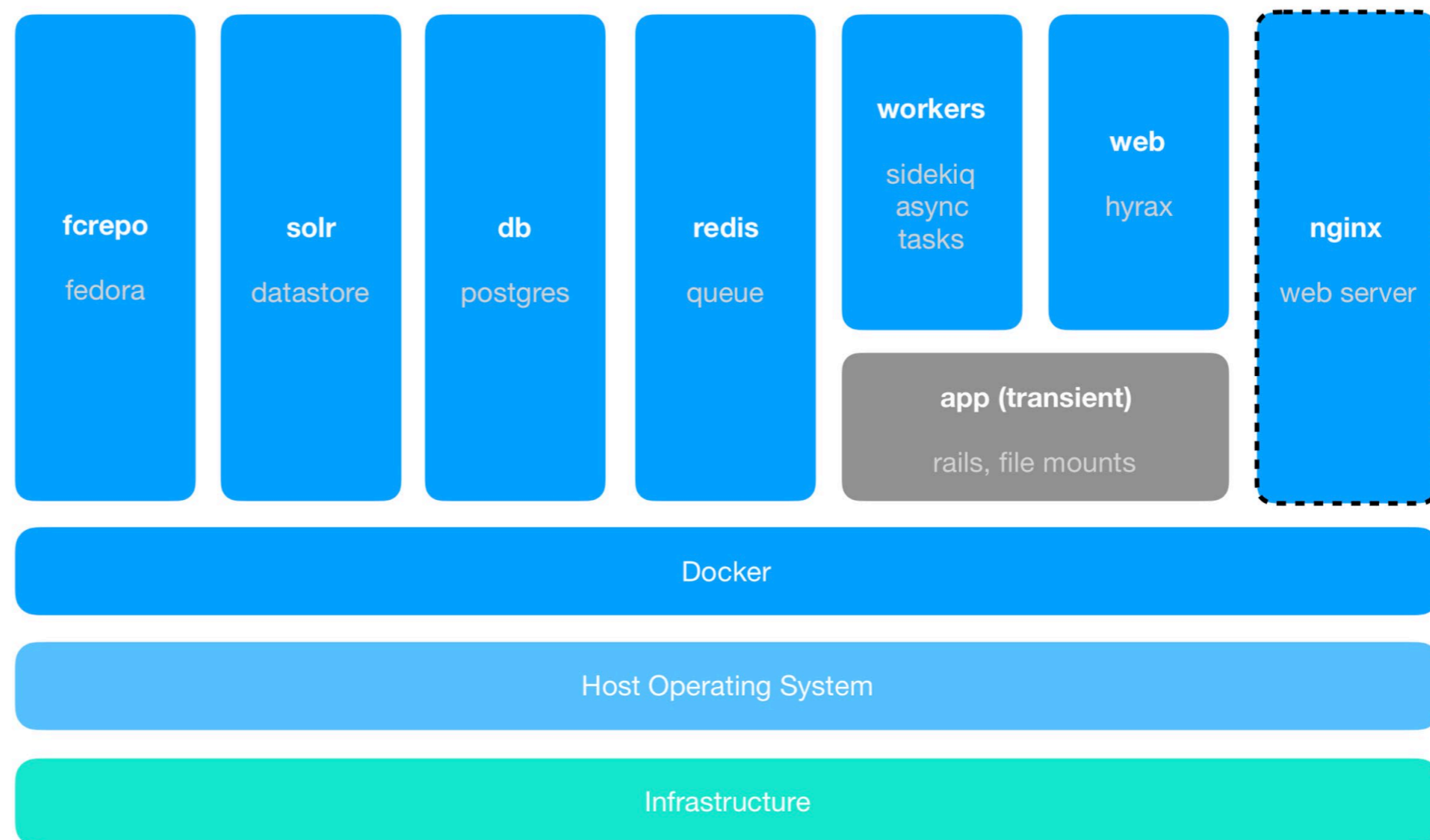
# 2. The MDR system

Steven Eardley
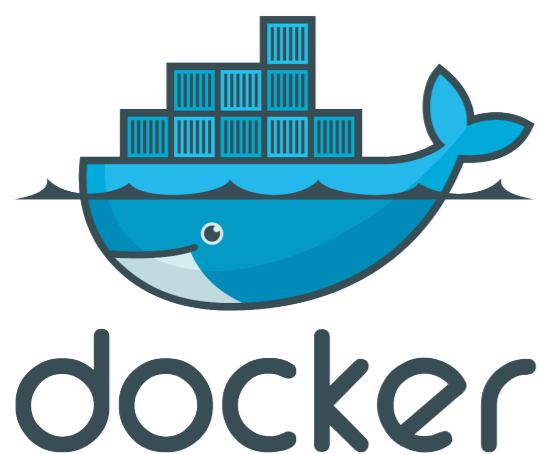
# About the Materials Data Repository (MDR)

- Hyrax (Samvera)

# Nested View

## Add New Dataset

Descriptions   Experimental method   **Instrument**   Specimen type   Files

Relationships   Sharing

**Instrument**

**Title**

**Alternative title**

**Date**

Processed ⇅

**Description**

**Identifier**

choose type ⇅

---

| Description | > |

| Experimental method | ˅ |

**Data origin**
experiments

**Specimen set**
NA

| Instrument | ˅ |

**Instrument**

| **Title** | PHI Quantera SXM |
| **Processed** | 07/06/2017 |
| **Identifier - Local** | 21354144 |
| **Manufacturer** | |
| **Organization** | ULVAC-PHI |
| **Role** | Manufacturer |
| **Operator** | |
| **Name** | ▮▮▮▮▮▮▮▮▮▮@nims.go.jp |
| **NIMS Person ID** | ▮▮▮ |
| **Role** | operator/データ測定者・計算者 |
| **Managing organization** | |
| **Organization** | NIMS |
| **Sub organization** | 材料分析ステーション |
| **Role** | Managing organization |

| Specimen type | > |

---

**Managing organization**

**Organization**

**Sub organization**

✖ Remove

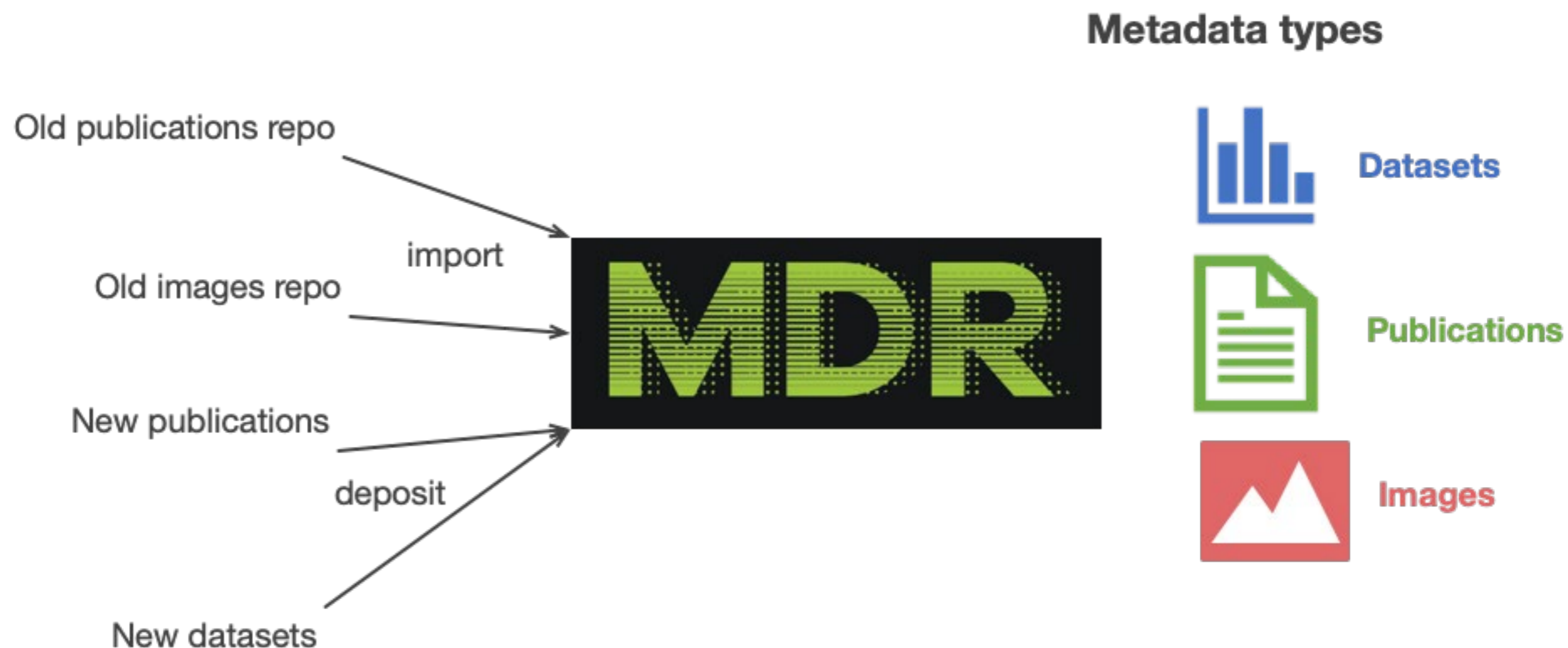➕ Add another instrument

---

# Containerised Development and Deployment

# 3. A focus on metadata

Asahiko Matsuda

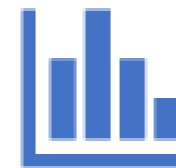# Datasets, publications, & images coexisting in MDR

# Metadata for...

**Publications**

- Title
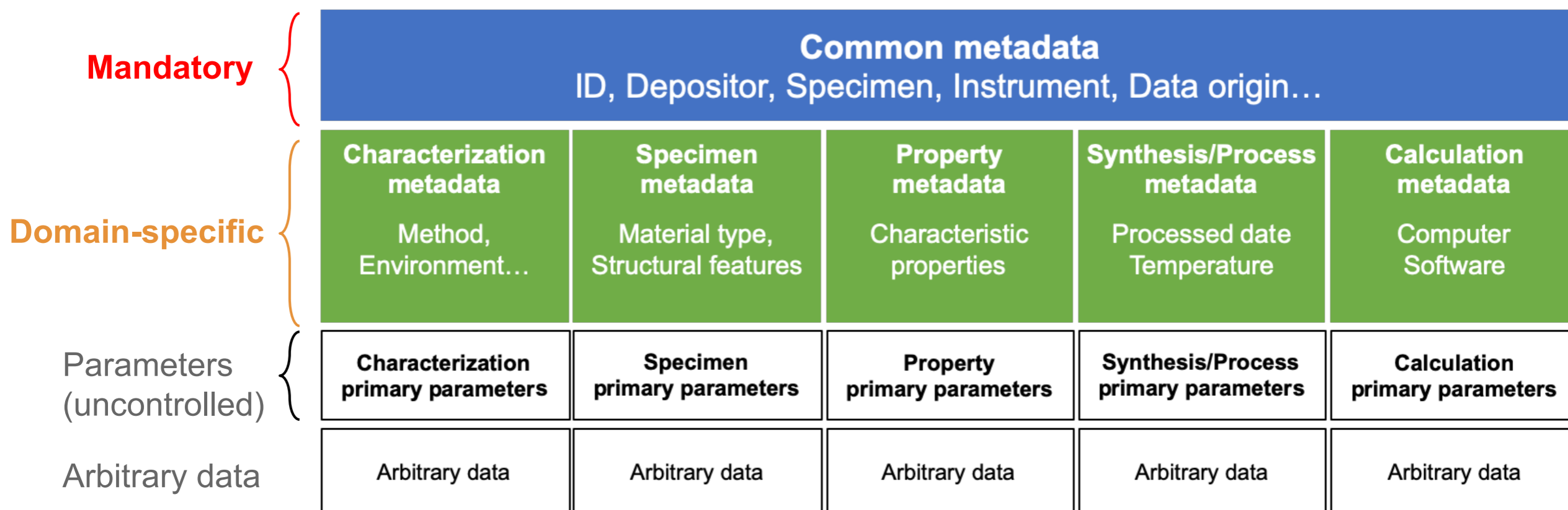- Authors
- Publication
- Issue
- Date
- ...

**Datasets**

- Method
- Specimen
- Facility
- Temperature
- Acceleration energy
- ...

Extremely domain-specific !
How can we model this ?
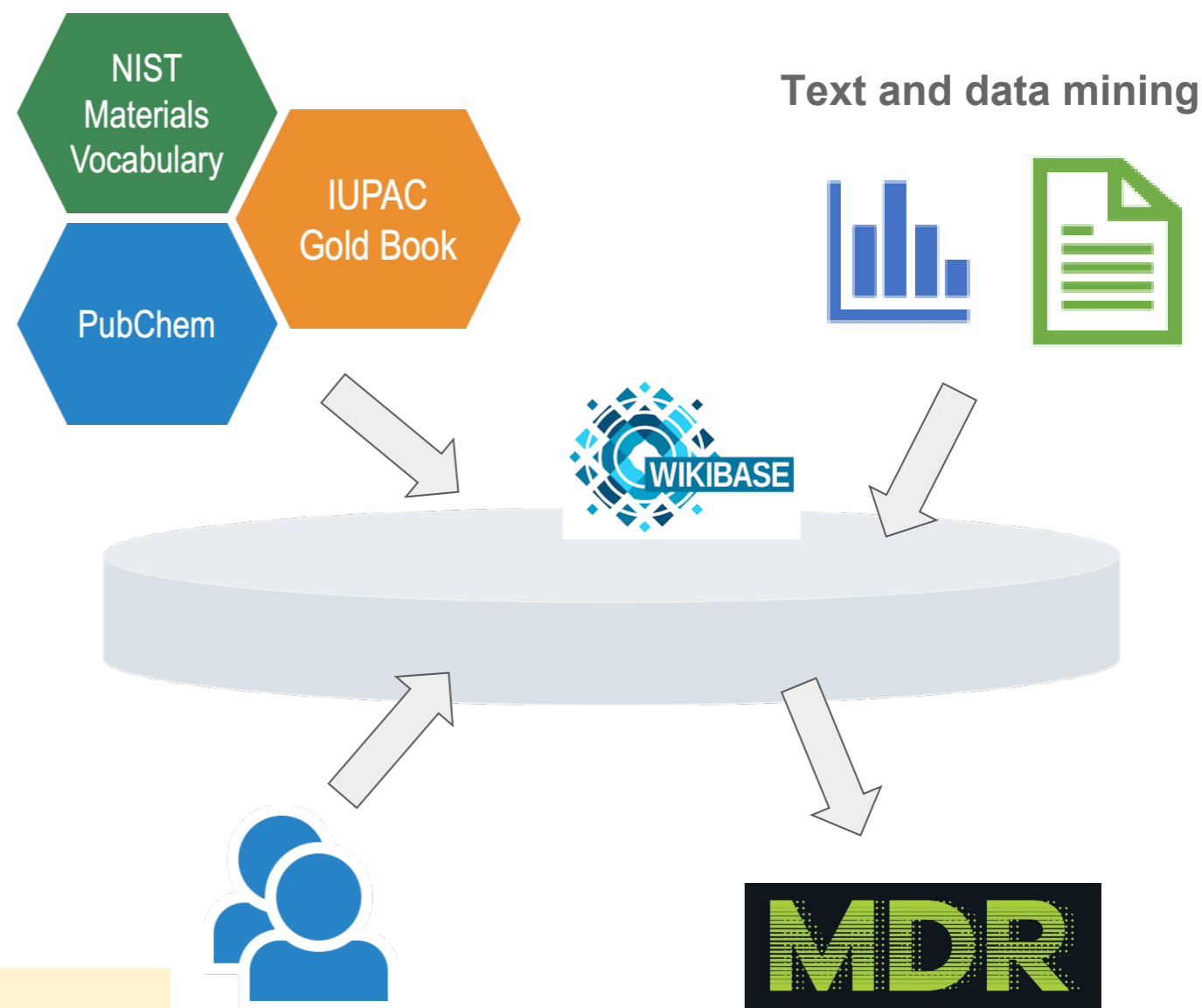
# Tiered and nested metadata model for datasets



**Mandatory**

**Common metadata**
ID, Depositor, Specimen, Instrument, Data origin…

**Domain-specific**

| Characterization metadata | Specimen metadata | Property metadata | Synthesis/Process metadata | Calculation metadata |
|---|---|---|---|---|
| Method, Environment… | Material type, Structural features | Characteristic properties | Processed date Temperature | Computer Software |

Parameters (uncontrolled)

| Characterization primary parameters | Specimen primary parameters | Property primary parameters | Synthesis/Process primary parameters | Calculation primary parameters |
|---|---|---|---|---|
| Arbitrary data | Arbitrary data | Arbitrary data | Arbitrary data | Arbitrary data |

Arbitrary data

Metadata view and deposit form also reflect this model

# Metadata used for faceted browsing & searching

# Enriching metadata with vocabularies

- 3 sources of vocabulary terms:
  1. Controlled vocabularies
     - Community governed
  2. Machine-generated
     - Terms extracted by text/data-mining
  3. Crowd-sourced
     - User-generated terms
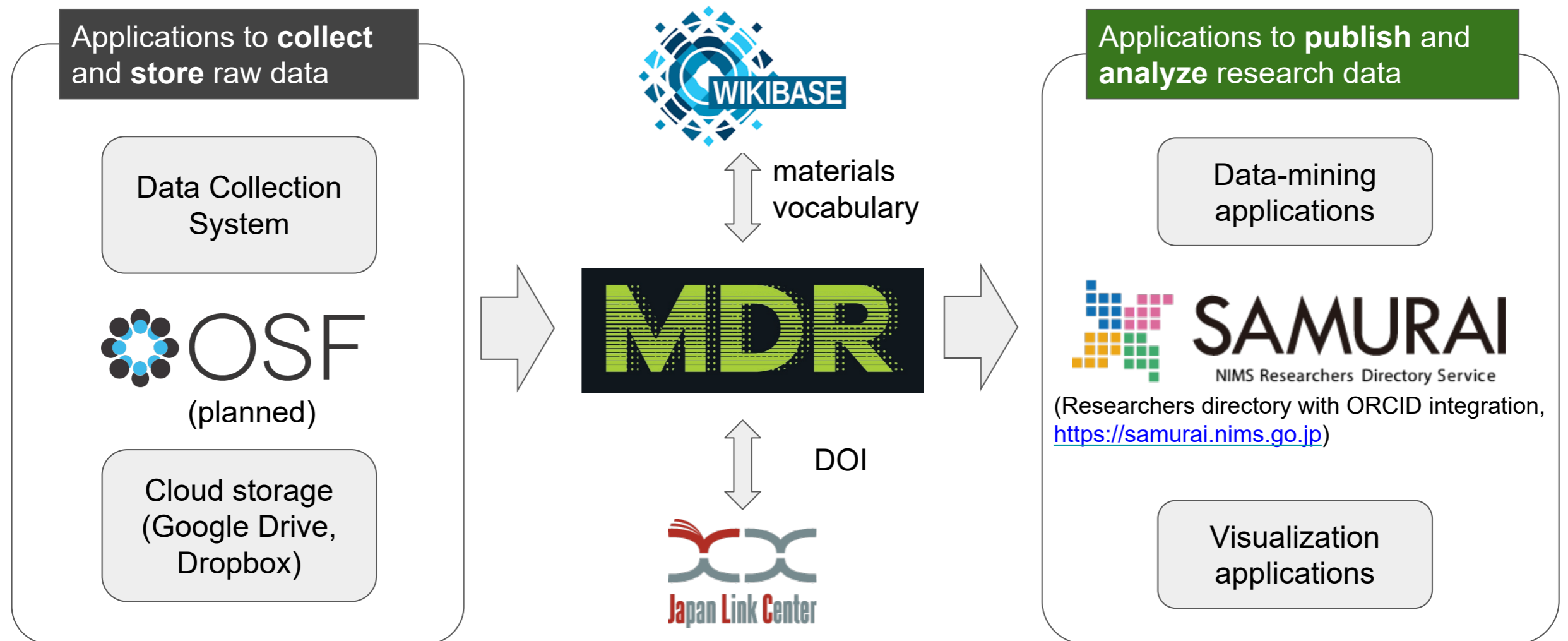     - From NIMS research community
     - "Folksonomy"

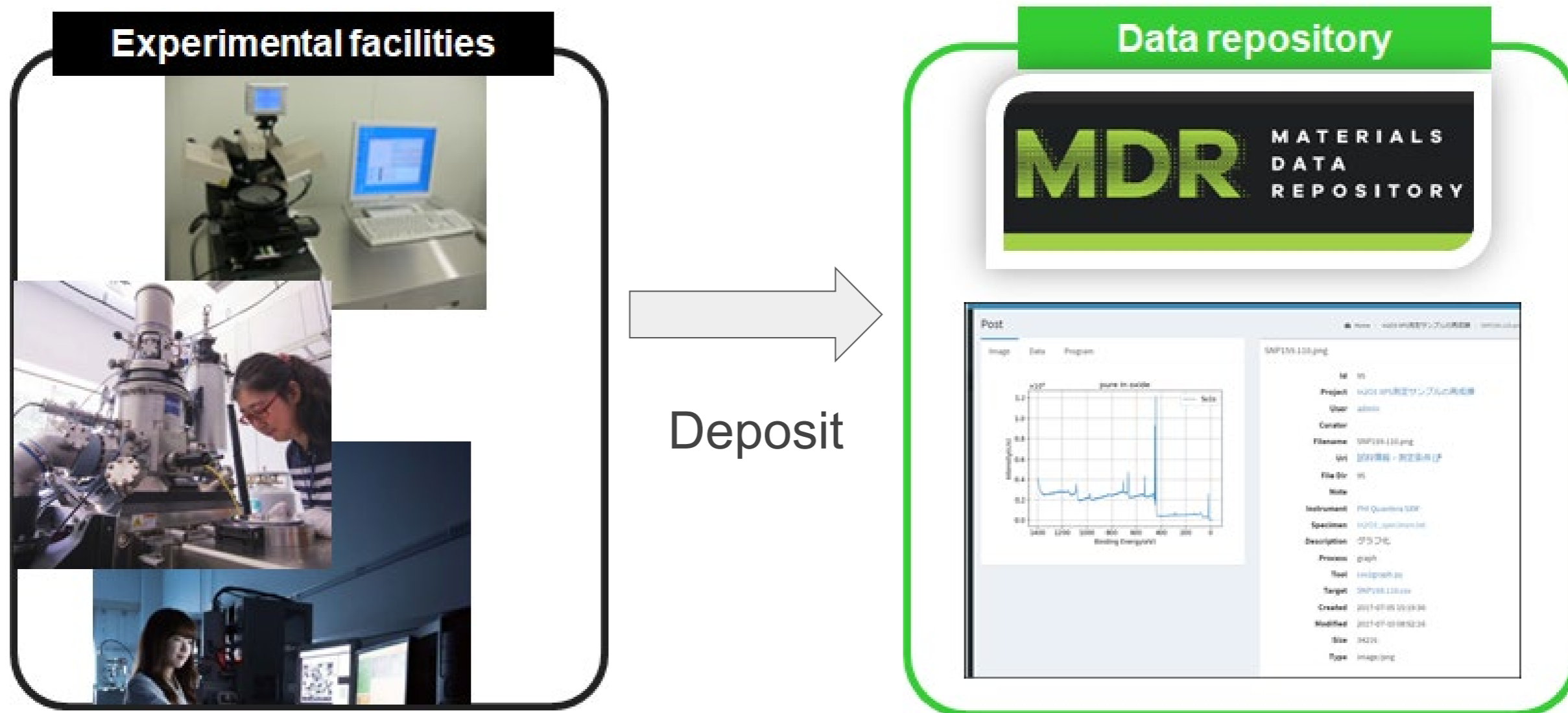We have a separate poster focusing on this.

# 4. Integration

Kosuke Tanabe

# Overview of integrations



Applications to **collect** and **store** raw data

Data Collection System

OSF

(planned)

Cloud storage (Google Drive, Dropbox)

WIKIBASE

materials vocabulary

MDR

DOI

Japan Link Center

Applications to **publish** and **analyze** research data

Data-mining applications

SAMURAI
NIMS Researchers Directory Service

(Researchers directory with ORCID integration, https://samurai.nims.go.jp)

Visualization applications

MDR MATERIALS DATA REPOSITORY

NIMS · ANTLEAF · COTTAGE LABS

# Use case for depositing experimental data



Deposit

# Data Collection System (DCS)

- A system to convert raw measurement data, assign metadata, draw a graph, and hand them over to MDR

- NIMS researchers' home-grown application

# Metadata from DCS to MDR

```xml
- <!--
    ``post`` に紐づいた ``work_type=specimen`` の ``work`` の ``Crystalographic_Structure`` を対応させる。
  -->
  <crystallographic-structure>https://komorido.nims.go.jp/wiki/Item:Q35</crystallographic-structure>
  <!-- polycrystal（多結晶）-->
- <chemical-composition>
  - <!--
      ``post`` に紐づいた ``work_type=specimen`` の ``work`` の ``CAS_No_etc`` を対応させる。
    -->
  - <chemical-composition-identifier>
      <identifier-type>CAS</identifier-type>
      <cas-number>7440-22-4</cas-number>
    </chemical-composition-identifier>
  - <!--
      ``post`` に紐づいた ``work_type=specimen`` の ``work`` の ``Chemical_composition`` を対応させる。
    -->
    <description>Ag</description>
  </chemical-composition>
```

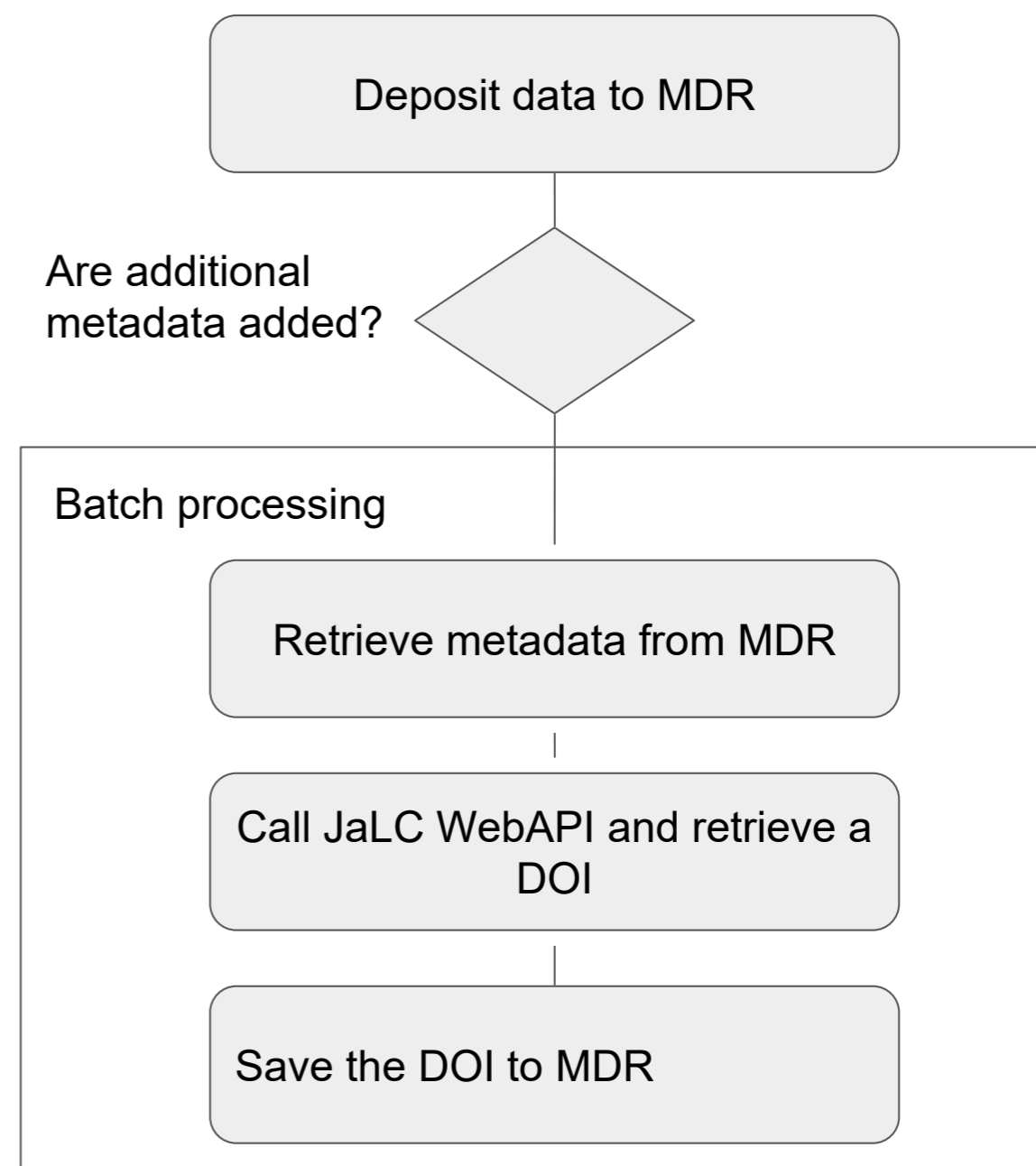URL of a vocabulary term provided by Wikibase

# Dataflow between DCS and MDR

# Integration with DOI Registration System

- MDR supports JaLC DOI



  https://japanlinkcenter.org/
  (DOI RA in Japan)

- Only datasets with both mandatory and domain-specific metadata will be minted DOIs

- The DOI minting is processed by a batch script invoked by MDR

```
Deposit data to MDR
        │
        ◇  Are additional
           metadata added?
        │
┌─────────────────────────────┐
│ Batch processing            │
│                             │
│  Retrieve metadata from MDR │
│           │                 │
│  Call JaLC WebAPI and       │
│  retrieve a DOI             │
│           │                 │
│  Save the DOI to MDR        │
└─────────────────────────────┘
```
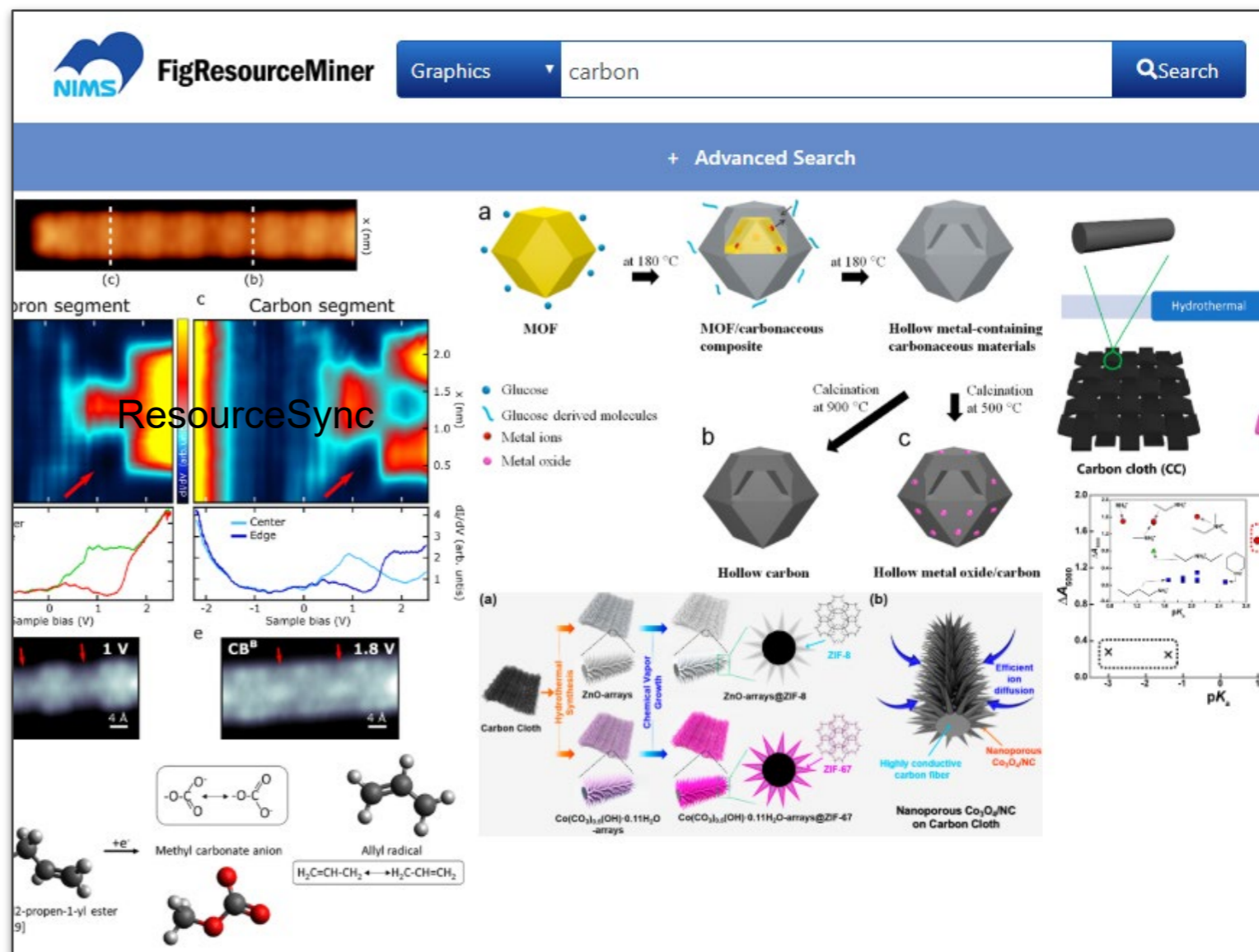
# Application using data on MDR: FigResourceMiner

- Data mining service

- Extract text information from figures and images in articles and datasets

- FigResourceMiner harvests files from MDR

ResourceSync

# Challenge in integration

- Depositing huge data from collaborators outside NIMS network

  - Sometimes over 4TB

  - Collaborators are expected to deposit those data to their local repository, then we can harvest metadata for search

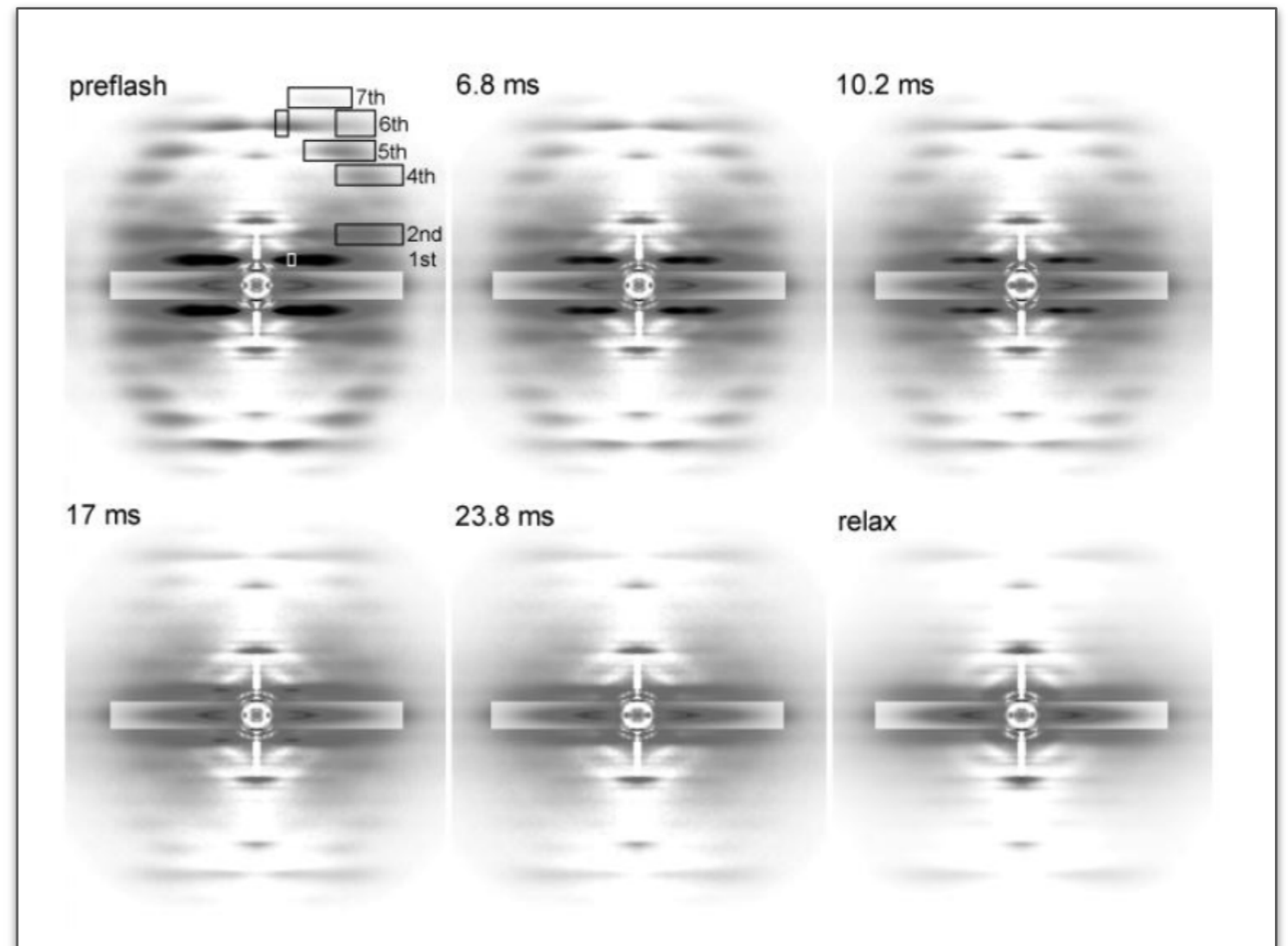  - Don't we need actual data (not just metadata) for data mining?



Image data files generated by the X-ray beamline in SPring-8, located outside NIMS

http://www.spring8.or.jp/wkg/BL40XU/solution/lang/SOL-0000001622

# 5. Supporting discovery

Paul Walk

# COAR and *Next Generation Repositories*

- Defined "behaviours":
  - Exposing Identifiers
  - Declaring Licenses at the Resource Level
  - **Discovery Through Navigation**
  - Interacting with Resources (Annotation, Commentary, and Review)
  - Resource Transfer
  - **Batch Discovery**
  - Collecting and Exposing Activities
  - Identification of Users
  - Authentication of Users
  - Exposing Standardized Usage Metrics
  - Preserving Resources

**NEXT GENERATION**
**REPOSITORIES**

# Discovery Through Navigation (for humans)

- Faceted browsing and searching

- Using vocabulary terms derived from:

  - Controlled vocabularies

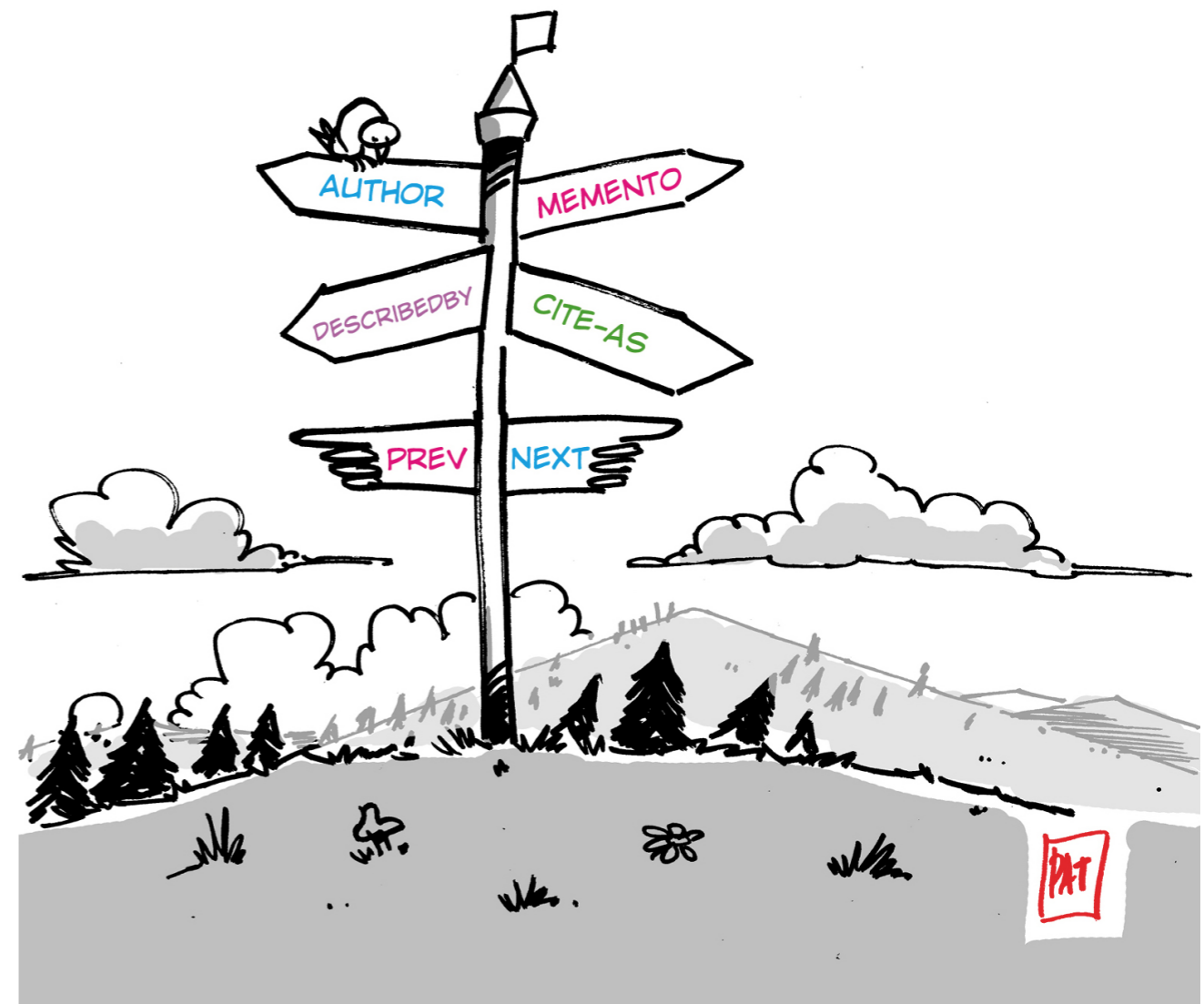  - Terms extracted algorithmically

  - Crowd-sourced keywords

# Discovery Through Navigation (for machines)

- *Signposting* has defined patterns relating to bibliographic resources:
  - Author
  - Bibliographic Metadata
  - Identifier
  - Publication Boundary
  - Resource Type

- It does define a "dataset" resource type…. but...

- **How do we navigate heterogeneous & complex datasets (multiple files)?**
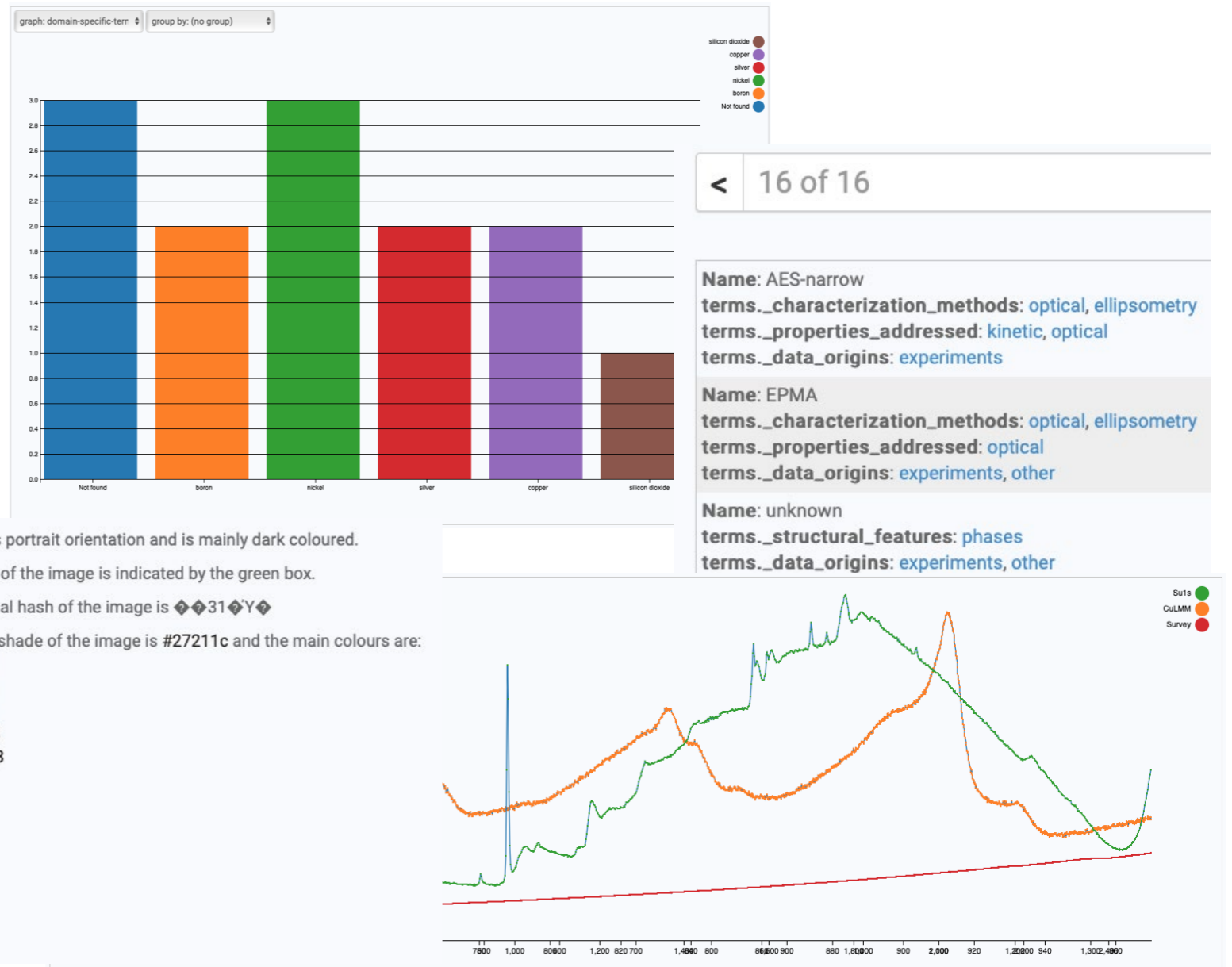
"Signposting the Scholarly Web"

# Batch Discovery (1)

- Aggregation is still an important tactic in the "knowledge commons"

  - mitigates network latency and facilitates processing at scale

- Many conceivable services built on research data will require the data to be harvested and aggregated

- OAI-PMH does not support the harvesting of content

- ResourceSync is an important technology for this

- Implemented in the MDR, about to be tested in collaboration with the Open University Core service

# Batch Discovery (2)

- Once the data is enabled for batch discovery, many new interfaces, tools etc are possible….

# Conclusions

- By September 2019, we will have launched the Materials Data Repository, which:

  - Is a platform to collect and showcase the work of NIMS's researchers

    - Shows some of COAR's Next Generation Repository behaviours

    - Is integrated with a number of other NIMS systems

    - Is playing its part as a significant 'node' in the global knowledge commons

  - **By April, 2020 April, MDR is scheduled to be opened to public**

    - a <u>publicly accessible</u>  platform for R&D of materials

ありがとうございました

Arigatō

Danke schön!

Thank you!