# Cross-searching of datasets by linking repository and vocabulary management in materials data platform

Kosuke Tanabe, Asahiko Matsuda, Masashi Ishii

iD https://orcid.org/0000-0002-9986-7223

National Institute for Materials Science (NIMS), Japan

October 23rd, 2024
DCMI 2024 Toronto Best practices session

# NIMS Materials Data Platform "DICE"

- Research data platform for materials science

- Developed and operated by National Institute for Materials Science (NIMS), Japan

- Consists of several research software and services
  - Materials databases
  - Research data sharing service
  - User authentication / authorization

# Selected in this presentation

- Materials Data Repository (MDR)
- MDR XAFS DB
- MatVoc / MatVoc Explorer

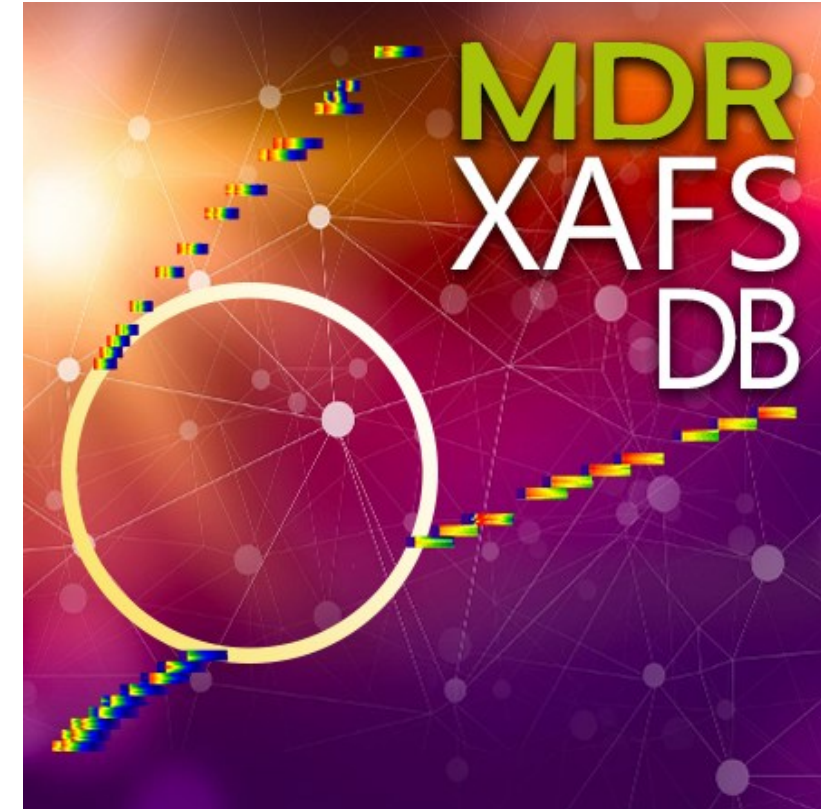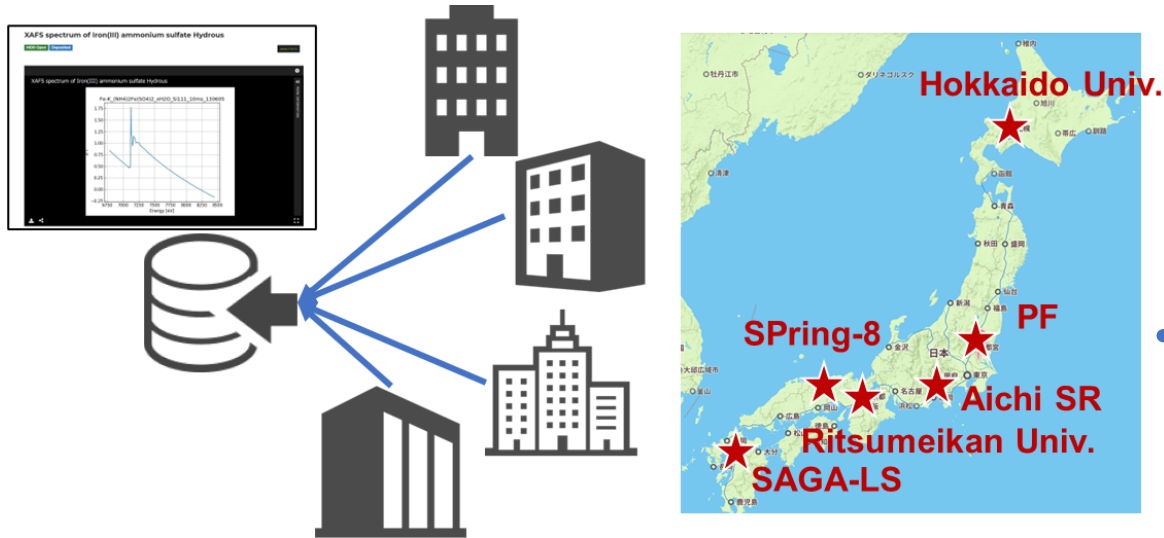# NIMS Materials Data Repository (MDR)

- Data repository in the DICE
  - https://mdr.nims.go.jp

- released in June 2020

- About 1,000 articles and 12,800 datasets as of October 2024

- provides facet search UI and REST API

- Built on Samvera Hyrax (OSS)

- Project leader: Kosuke Tanabe

# MDR XAFS DB

- Collection of X-ray absorption fine structure (XAFS) datasets
  - One of the most popular collections in MDR

- Includes approximately 2600 XAFS spectra datasets collected by six research institutions in Japan

- Project leader: Masashi Ishii



- https://doi.org/10.48505/nims.1447

# MatVoc / MatVoc Explorer

- MatVoc: https://matvoc.nims.go.jp
  - Vocabulary service for materials science
  - Built on Wikibase
  - Includes more than 2400 entities
  - Provides SPARQL endpoint

- MatVoc Explorer: https://matvoc.nims.go.jp/explore/en/home
  - Provides browse and search UI
  - Provides a link to MDR and other external databases associated with each vocabulary

- Project leader: Asahiko Matsuda


MatVoc Explorer — Materials Vocabulary

NIMS XAFS DB Project Materials Dictionary > Chemicals > Inorganic materials > Oxide > Zinc oxide

# Q1800: Zinc oxide

Vocabulary ID  http://matvoc.nims.go.jp/entity/Q1800

| Language | Label | Description | Alias |
|---|---|---|---|
| English | Zinc oxide | BENTEN-registered chemicals, ZnO | ZnO, Zinc white, Zinc oxide (ZnO), Zinc oxide, philosopher's wool, flowers of zinc, Chinese white, calamine, 1314-13-2 |
| Japanese | 酸化亜鉛 | BENTEN登録済み化学物質、ZnO | ZnO |

## 📖 Semantic Relatives

▶ Parents

## 🔗 DICE Links

AtomWork Zn+O ,  MDR Zinc+oxide

# Metadata schemas around MDR and MDR XAFS DB

# Metadata (schemas) in MDR and MDR XAFS DB

Scientific metadata

◦ **XAFS Metadata**

◦ **DICE common messaging format**
◦ **MDR Schema**

◦ DataCite Metadata Schema
◦ JPCOAR Schema

Bibliographic metadata

# XAFS Metadata

- https://github.com/xafs-db/xafs-schema

- Metadata schema developed by the Japan XAFS Society
  - https://www.jxafs.org/xafs-database

- Includes domain-specific properties (e.g., instruments, measurement)

```yaml
measurement:
  detectors:

# 1. IC => 電流アンプ - V/F - カウンタ の場合
  - name: I1
    arrangement: ...//Current Amp.//V/F Converter//Counter
    processing_lines:
    - conversion_factor: 1e-14
      conversion_factor_unit: A
      processors:
      - processor: Current Amp.
        type: Average
        manufacturer: NF
        model_number: CA5350
        gain: 1e6
        gain_unit: V/A
        time_constant: 1e-3
        time_constant_unit: s
```

# DICE common messaging format

- https://doi.org/10.48505/nims.3240

- Designed for entire datasets and research areas on the data platform DICE

- Supports many scientific entities and properties including instruments, specimens and experimental methods

- Consists of many deeply-nested properties

- *"We have defined a common metadata schema for research dataset distribution/exchange/storage, with features for not only system-to-system communication but also description of the datasets. The metadata were initially defined by a XML Schema, and has been converted into JSON Schema. The schema has a complex structure, which includes the dataset's bibliographic metadata as well as description of the dataset from the materials science viewpoint."*

# The structure of DICE common messaging format



| **Mandatory** | Common metadata | | | | |
|---|---|---|---|---|---|
| | ID, Depositor, Specimen, Instrument, Data origin… | | | | |
| **Domain-specific** | **Characterization metadata**<br>Method, Environment… | **Specimen metadata**<br>Material type, Structural features | **Property metadata**<br>Characteristic properties | **Synthesis/Process metadata**<br>Processed date Temperature | **Calculation metadata**<br>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 |

S. Kikuchi et al., IEICE Tech. Rep. vol. 119, no. 66, SC2019-2, pp. 7-17, 2019, in Japanese
https://ken.ieice.org/ken/paper/20190531k1nc/
(also in Ranganathan et al., 14th Int. Conf. Open Repositories, 2019. https://doi.org/10.5281/zenodo.3553963 )

# An example of the DICE common messaging format

```
C: > Users > kosuke > Downloads > {} meta.json > {} data > {} included > {} additional-attachment-pointer > {} spared-description-for-token
  2        "data": {
 24            "attributes": {
 25                "common-term": {
 32                    "basic-data-description": {
190                        "deposit-execution-identifier": {
191                            "description": "API-FWK",
192                            "identifier-type": "nims-internal",
193                            "nims-identifier": "urn:USER_IDENTIFIER.dpfc.nims.go.jp:2bf8fb6a-44a3-4892-9d33-e0dac3c99f0f"
194                        }
195                    },
196                    "instrument-description": [
197                        {
198                            "name": "BL14B2_XAFS",
199                            "identifier": {
200                                "identifier-type": "local-identifier",
201                                "local-identifier": "RDEtagid-64568318"
202                            },
203                            "manufacturer": {
204                                "description": "",
205                                "identifier-type": "descriptional",
206                                "organization-identifier": "",
207                                "organization-description": {
208                                    "title-sub-organization": "",
209                                    "title-major-organization": "Japan Synchrotron Radiation Institute"
210                                }
211                            },
212                            "process-date": {
213                                "date": "2020-02-06",
214                                "time": "00:00:00+09:00"
215                            },
```

# Initial plan in depositing XAFS datasets to MDR

- All research datasets and applications are expected to use the DICE common messaging format to describe and share its metadata



S. Kikuchi et al., IEICE Tech. Rep. vol. 119, no. 66, SC2019-2, pp. 7-17, 2019, in Japanese
https://ken.ieice.org/ken/paper/2019053 1k1nc/

# Difficulties in applying the huge metadata schema

- Some of the XAFS datasets were deposited to MDR using the schema, but it took much more efforts than expected

- It was extremely difficult for researchers and application engineers to understand its complex specification
  - Encountered implementation incompatibilities in handling the metadata and schema between the application software
  - Encountered many limitations and bugs in the software (customized OSS)

- After some operations, it was officially abandoned in 2022

**DICE common message format schema**

A common metadata schema for distribution, exchange, and storage of research data among DICE systems.

With the discontinuation of the system that used this schema, this schema was also discontinued as of December 2022.

https://dice.nims.go.jp/about.html

# Technical limitations in MDR

- MDR basically supports only keyword search

- MDR doesn't fully support nested metadata structure in the DICE common messaging format
  - e.g., cannot update a value in a nested property

- These came from technical limitations in MDR's software stack

- "…What should we do?"

# MDR Schema

- https://github.com/nims-dpfc/mdr-schema

- Simplified version of the DICE common messaging format

- Solely designed for depositing datasets to MDR
  - Mostly compatible with the existing metadata on MDR

- Eliminates deeply-nested hierarchy structure

- Expects to be written in YAML

# An example of metadata in MDR Schema

```yaml
managing_organization:
  ror: https://ror.org/026v1ze26
  organization: National Institute for Materials Science

instruments:
- name: BL14B2_XAFS
  description: SPring-8 Engineering Science Research Beamline XAFS setup

experimental_methods:
# x-ray photoelectron spectroscopy
- category_vocabulary: https://matvoc.nims.go.jp/entity/Q31

specimens:
- name: HAVAR
  description: Standard Sample

chemical_compositions:
- description: W(CO)6

structural_features:
- description: radial distribution function
```

# Validating MDR Schema YAML

- Using Yamale to create and validate the metadata schema
  - https://github.com/23andMe/Yamale

```
$ yamale metadata-sample.yaml
Validating /home/kosuke/mdr-schema/metadata-sample.yaml...
Validation failed!
Error validating data '/home/kosuke/mdr-schema/metadata-sample.yaml'
        titles.0.title: Required field missing
$ |
```

# Schema definition of MDR Schema

```
159    # 人物
160    person:
161      name: str()
162      orcid: str(required=False)
163      e_rad: str(required=False)
164      organization: str(required=False)
165      department: str(required=False)
166      ror: str(required=False)
167      role: enum('author', 'editor', 'translator', 'depositor',
168
```

```
228    # 試料
229    specimen:
230      name: str()
231      description: str(required=False)
232      identifier: str(required=False)
233      material_type_vocabulary: str(required=False)
234      material_type_description: str(required=False)
235
236    ---
237    # 試料の化学組成
238    chemical_composition:
239      specimen_identifier: str(required=False)
240      identifier: str(required=False)
241      category_vocabulary: str(required=False)
242      category_description: str(required=False)
243      description: str(required=False)
```

# Why YAML?

- human-readable and (almost) writable
  - Easy to create and distribute a metadata template file since it supports single-line and inline comments
  - It is painful for human to read and write a spreadsheet that has a lot of columns

- It can be converted to JSON or other formats easily

# An example of MDR metadata migration

- "instruments"->"function" ->"category" is moved to "function_category"

- "instruments" ->"managing_organization" is moved to "instrument_management_ organization"

- Assigns local identifier "instrument_00001" to the instrument and refers it from the "instrument_management_ organization"

DICE common messaging format (converted to YAML)

```
instruments:
- name: BL14B2_XAFS
  description: SPring-8
  function:
  - category: spectroscopy
    identifier: https://matvoc.nims.go.jp/entry/Q30
  managing_organization:
    organization: JASRI
```

MDR Schema

```
instruments:
- identifier: instrument_00001
  name: BL14B2_XAFS
  description: SPring-8
  function_category:
  - https://matvoc.nims.go.jp/entry/Q30
instrument_managing_organization:
  instrument_identifier: instrument_00001
  organization: JASRI
```

# Comparison of metadata schemas around MDR XAFS DB

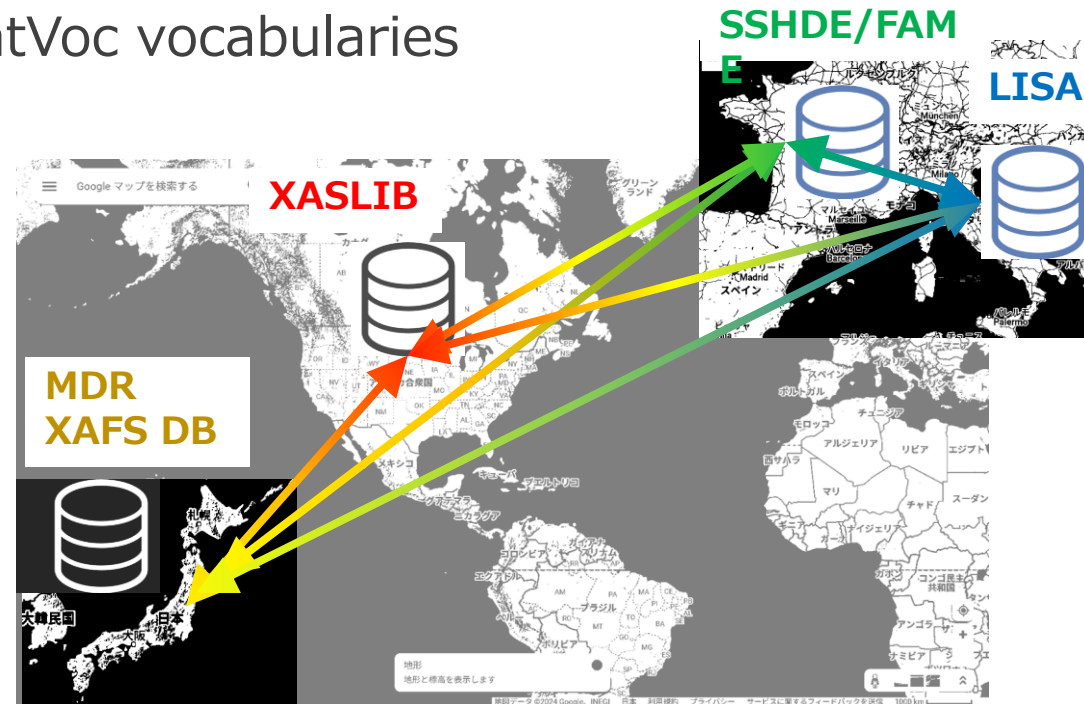| | XAFS Metadata | DICE common messaging format | MDR Schema |
|---|---|---|---|
| Support scientific metadata properties? | Yes | Yes | Partially yes |
| Research domain specific? | Yes | No | No |
| Format | YAML, JSON or TSV | XML, JSON | YAML, JSON |
| Metadata structure | deeply-nested | deeply-nested | single-nested |
| Human-friendly? | ? | No | Yes(?) |

# advanced case study: International XAFS DB Portal

# International XAFS DB Portal

- Cross-search for MDR XAFS DB and other three XAFS databases in the world
  - SSHADE/FAME (France)
  - XASLIB (US)
  - LISA XAS Database (Italy)

- All datasets are linked with MatVoc vocabularies

- https://ixdb.jxafs.org/

# International XAFS DB Portal

## Material name containing: 銅

Q1426 : Copper
Q1412 : Copper acetate
Q2319 : Copper bis(2,2,6,6-tetramethyl-3,5-heptanedionate)
Q1393 : Copper chromite
Q1409 : Copper fluoride
Q1428 : Copper molybdate
Q889 : Copper nickel
Q1413 : Copper nitrate, hydrous
Q1417 : Copper nitride
Q1410 : Copper phthalocyanine
Q1433 : Copper tungstate
Q1851 : Copper(I) chloride, anhydrous

# International XAFS DB Portal

## Links for Copper(I) sulfide

https://xaslib.xrayabsorption.org/spectrum/86/ (IXAS)
https://mdr.nims.go.jp/concern/datasets/9p290d467 (KEK)
https://mdr.nims.go.jp/concern/datasets/bk128f46k (KEK)
https://mdr.nims.go.jp/concern/datasets/sx61dq71s (KEK)
https://mdr.nims.go.jp/concern/datasets/9306t179v (SPring-8)
https://mdr.nims.go.jp/concern/datasets/pv63g2653 (SPring-8)

## Related crystal structures

F m -3 m
P 1 21/c 1
P 43 21 2
P 63/m m c

Link to MDR XAFS DB and other XAFS databases associated with the entity

# Related COD ID List for Copper(I) sulfide F m -3 m

return

Return To Top

1530508
1532316

Link to Crystallography Open Database

## COD — Crystallography Open Database

**COD Home**
Home
What's new?

**Accessing COD Data**
Browse
Search
Search by structural formula

**Add Your Data**
Deposit your data
Manage depositions
Manage/release prepublications

**Documentation**
COD Wiki
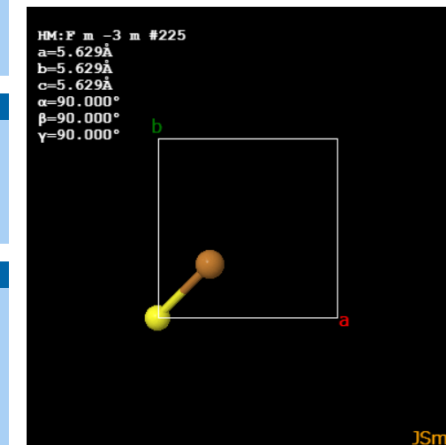Obtaining COD
License
Privacy and GDPR
Querying COD
Citing COD
COD Mirrors
Advice to donators
Useful links

### Information card for entry 1530508

1530507 << 1530508 >> 1530509

**Preview**

HM:F m -3 m #225
a=5.629Å
b=5.629Å
c=5.629Å
α=90.000°
β=90.000°
γ=90.000°

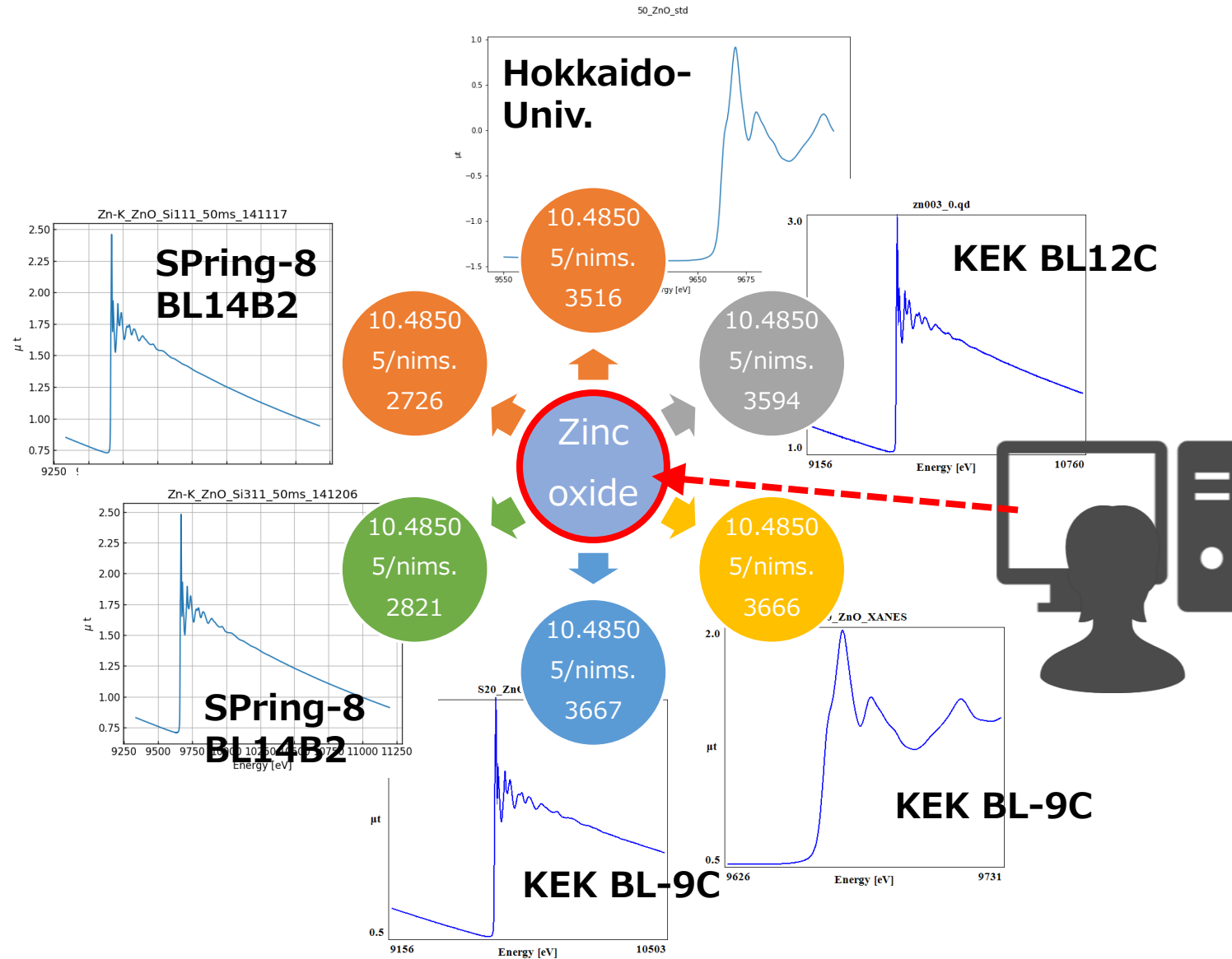JSmol

**Coordinates**    1530508.cif
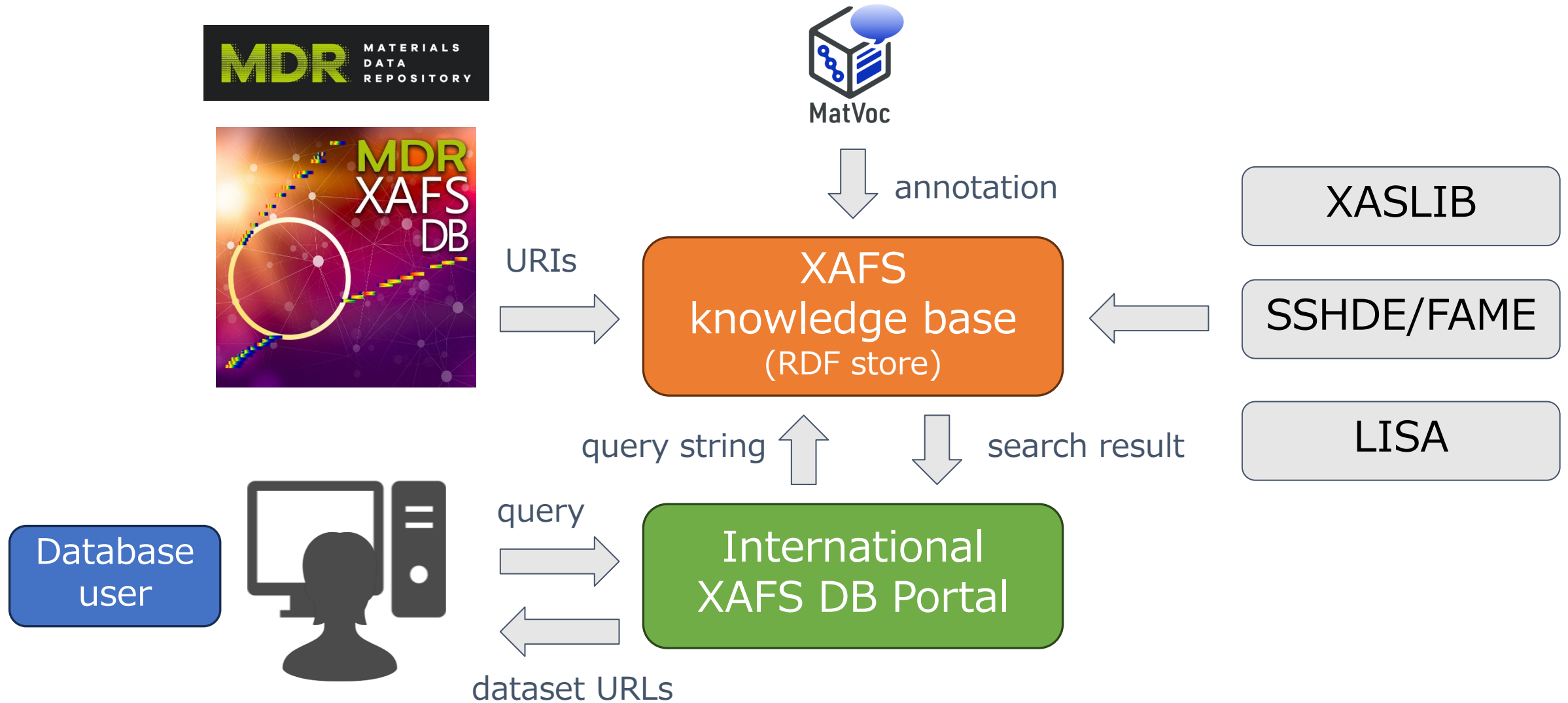
https://www.crystallography.net/cod/1530508.html

# Linking metadata using MatVoc

# XAFS knowledge base

# Lessons learned

- One big metadata schema does not fit all actually

- Everyone says, "Metadata is important", but the meaning depends on the context
  - Researchers, engineers, and librarians have different viewpoints on metadata and its use
  - Especially we should consider researcher-friendly metadata creation environment

- Browse and search features are essential in research data management even if they are basic implementations and have basic metadata
  - It is difficult for even researchers and data curators to grasp the whole picture of their datasets

- Identifier (for datasets and vocabularies) is the key to tame diverse research data and their domain-specific metadata

# **Further plans**

- Support querying MatVoc URIs in MDR REST API

- Add MatVoc URIs to datasets previously deposited on MDR

- Add MatVoc URIs to datasets on other DICE services

- Expand MatVoc vocabulary for a focused domain of materials and enhance findability of those data

- Linking materials to their crystallographic information through MatVoc

- **Cross-linking more XAFS databases around the world**

# TANABE.Kosuke@nims.go.jp