

Materials Data Platform overview: metadata, vocabulary, and repository

Mikiko Tanifuji

Materials Data Platform Center

Div. of Materials Data and Integrated System (MaDIS)

National Institute for Materials Science





Who we are



NIMS National Institute for Materials Science

*Established in 2001 by merger of two national institutes:
(metals + inorganic materials) Now covers materials in general*

MaDIS Research & Services Division of
Materials Data and Integrated System

Established in 2017 to focus on materials data and integration

DPFC Materials Data Platform Center

Budget 2017 – 2020, 3 billion yen



Who we are: a Facebook

Materials Development

- 高分子材料設計
- 熱制御材料
- 金属系構造材料
- 半導体
- 多孔質材料

Materials Architectures

- 計算シミュレーション：自動化
- SIP-MI：プロセスから一貫予測
- SMILES X：分子人工設計
- 計測インフォマティクス

Data Science

- スパースモデリング（モデル選択）
- 画像解析～パターン認識、深層学習
- 回帰技術～機械学習、ベイズ推定
- 最適化技術～能動学習等
- 自然言語処理



Data infrastructure

- Data structure and modeling
- Data curation
- Data collection and FAIRable
- Data mining from publications
- Data system technology and development

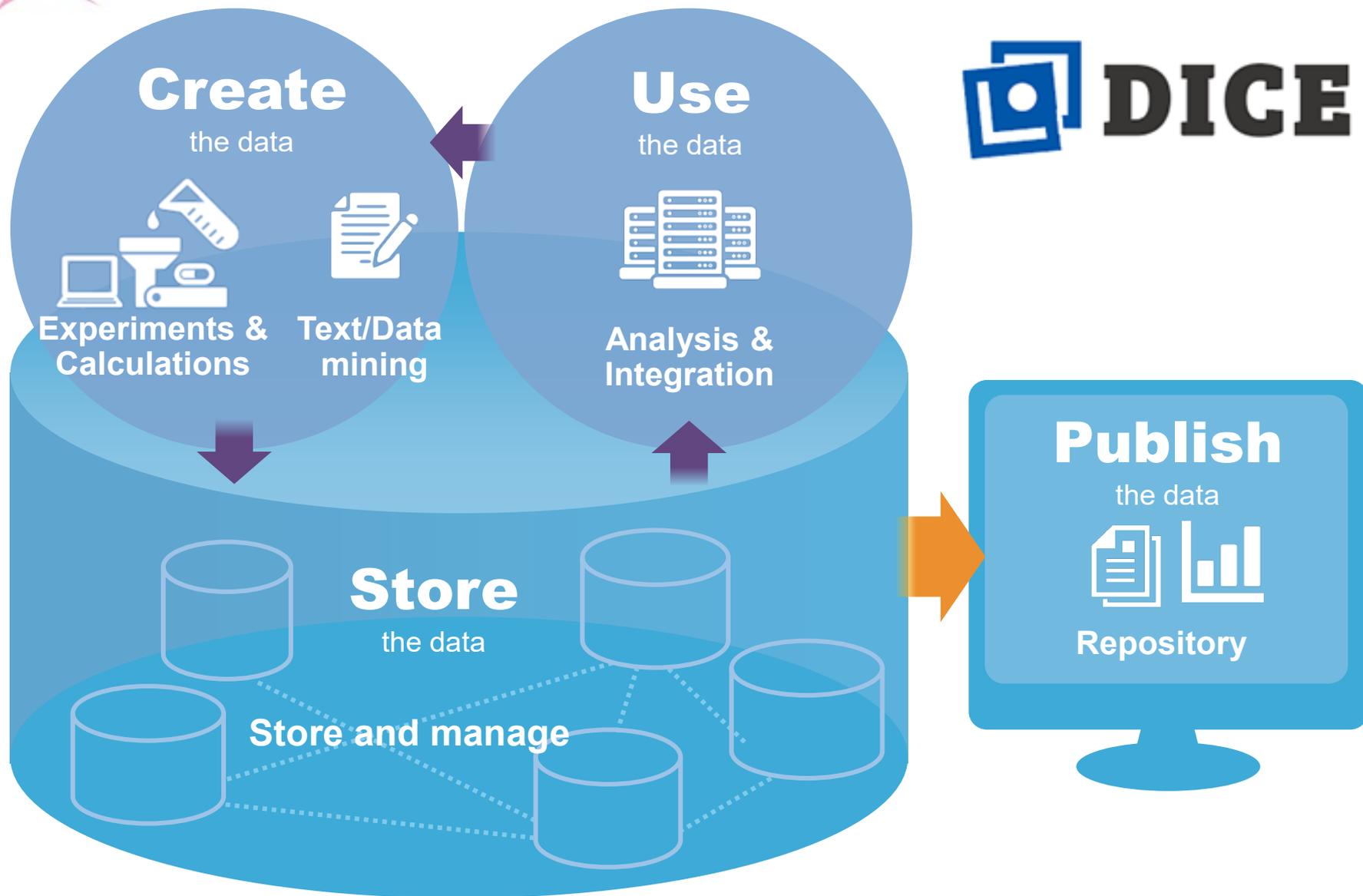
Data-driven people at MaDIS

66 staff at Materials Data Platform Center





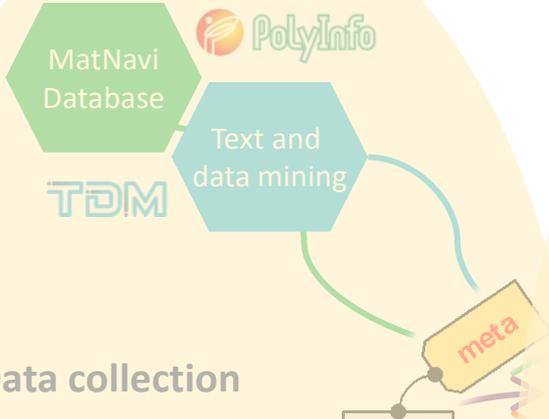
Materials Data Platform at NIMS



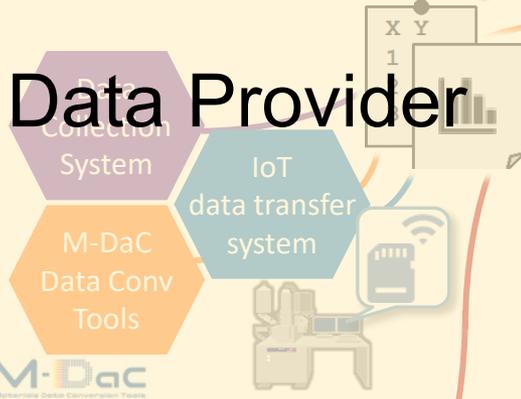


Materials Data Platform overview

Building advanced databases



Data collection

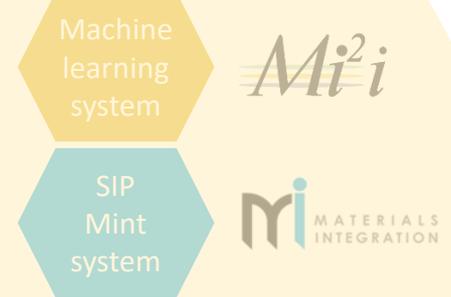


Large-scale facilities



Data Cloud
Materials
data hub
(2021 -)

Analyses and materials integration



Publishing open science

Data Science
Academia-
Private sectors





Four actions mapped to the platform components





DICE Common Message Format

METADATA

Mandatory metadata

Bibliographic metadata

Administrative metadata

Subject material

Common metadata model

Domain-specific metadata

+

Characterization metadata

Method,
Environment...

Specimen metadata

Material type,
Structure...

Property metadata

Physical properties,
Units...

Synthesis/Process metadata

Processed date,
Temperature...

Calculation metadata

Computer software,
Version...

Implemented as data model

DATA

Primary parameters

+

Characterization primary params

Data

Specimen primary params

Data

Property primary params

Data

Synthesis/Process primary params

Data

Calculation primary params

Data

Save as files

After lots of discussions (still ongoing), schema file published at <https://dice.nims.go.jp/>



DICE Common Message Format

Direct MDR upload

Experimental data files



Common metadata model



RDM



Importer



Files



Metadata form



Data

MDR metadata list

- Val CO
- MDR stores and manages materials data using a metadata set implemented by referring to the [DICE common message format schema](#), system of the data platform to another.

Metadata list

- 2020-07-22 [MDR-metadata-20200717-ext.xlsx](#)

Manual

- Dataset metadata
- Publication metadata

dice.nims.go.jp



Ongoing discussion about “primary parameters”

- Highly domain-specific parameters such as
 - Which absorption edge for a XAFS measurement?
 - Which basis set for a Gaussian calculation?

Extremely difficult to agree how to include in the common schema!



At least, save as files?

As a 2 x n key-value CSV

Software	Gaussian09
Calculation	B3LYP
Basis set	cc-pVDZ

The screenshot shows the MDR (Materials Data Repository) interface. At the top, there's a 'Preview: data.csv' section with a search bar and a table of key-value pairs. The table has two columns: 'Key' and 'Value'. The rows are: Software: Gaussian09, Calculation: B3LYP, and Basis set: cc-pVDZ. Below the preview is an 'Items' section with a table listing files. The first item is 'data.csv', uploaded on 18/10/2020, with a size of 56 Bytes and a visibility of 'NIMS'. There is a 'SELECT AN ACTION' button next to it.

As a Schema.org-like JSON-LD file



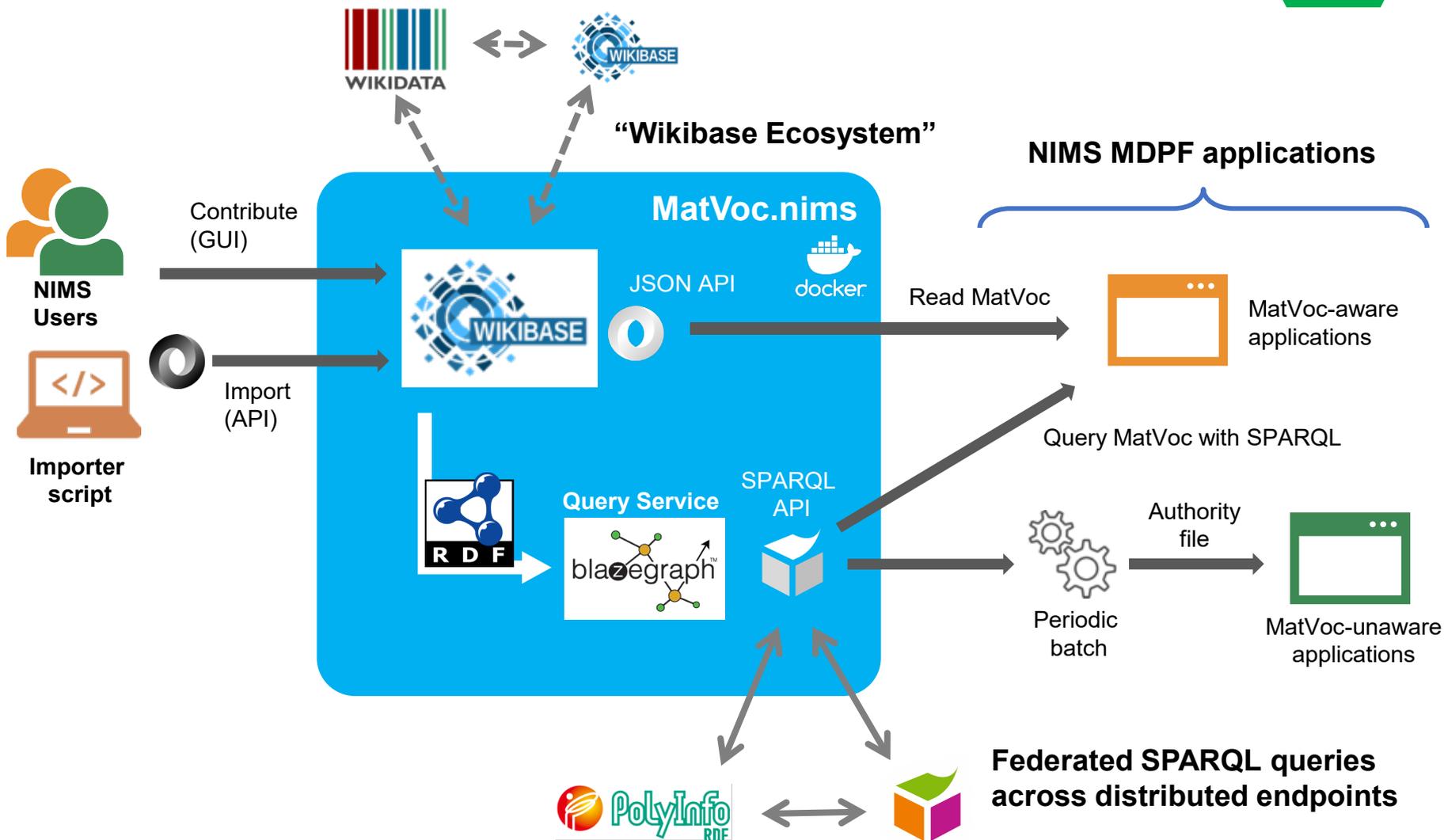
```

{ "@graph": [
  { "@id": "./",
    /* Bibliographic metadata */
    "name": ...,
    "author": ...,
    /* Scientific metadata */
    "variableMeasured": ...,
    "hasPart": { "@id": "data.dat" }
  }, .....
]}

```

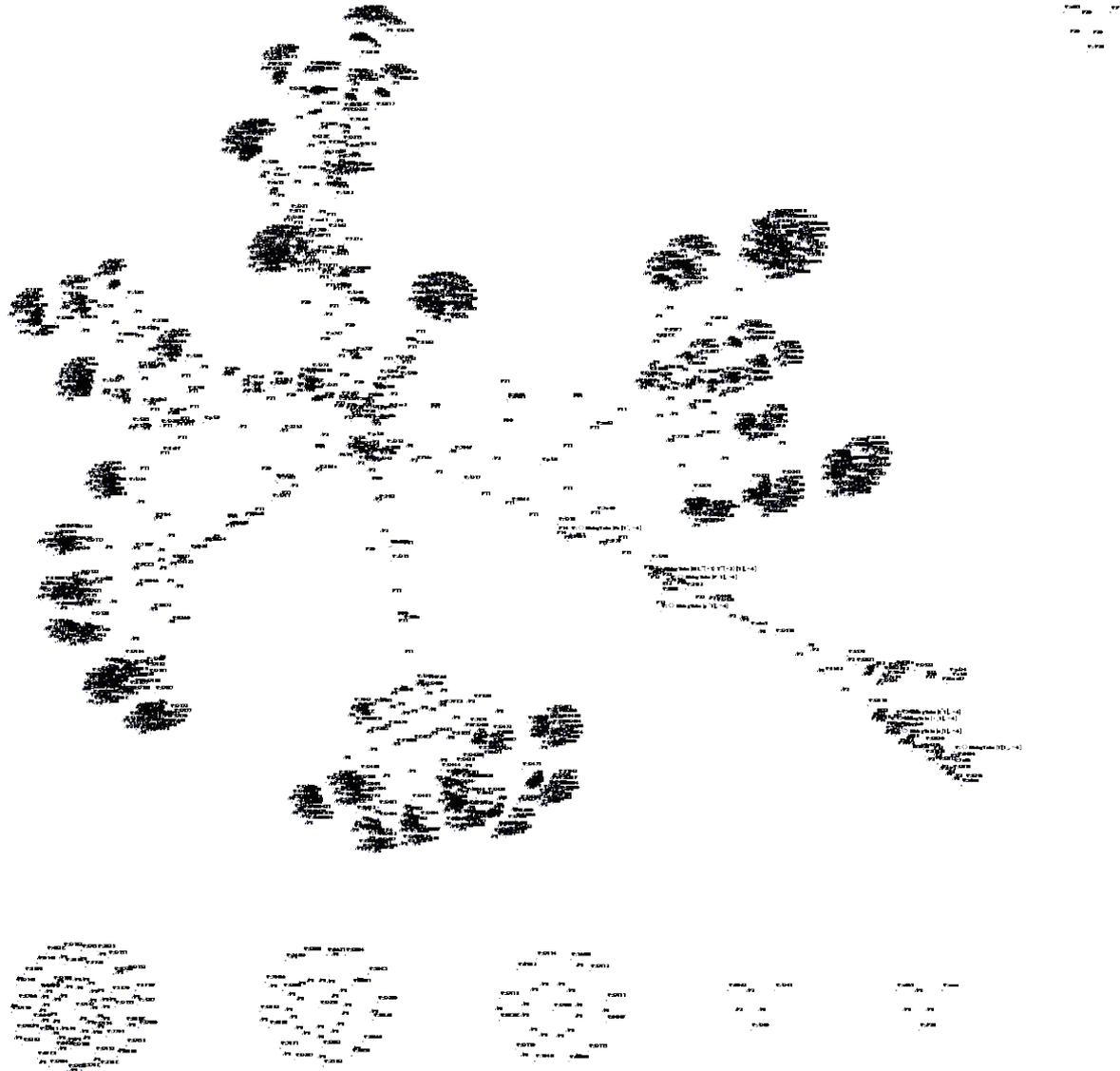


Vocabulary system overview





Bird's eye view of our vocabulary now





Vocabulary-based metadata transfer between systems

Data Collection: Research Data Express

Common message format depositUploadReq.json

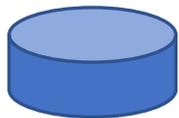
```

{
  "methods-category": [
    {
      "analysis-field-description": "",
      "analysis-field-items": [
        "0688"
      ],
      "energy-level-transition-structure": [
        {
          "measurement-environment": "0542",
          "measurement-environment-description": "",
          "method-category": [
            {
              "method-main-category-code": "030",
              "method-subcategory-code": "0386"
            }
          ],
          "reference-source": [

```

Materials Data Repository

Method	
Characterization methods	
spectroscopy → x-ray absorption spectroscopy	
Instruments	
Instrument	
Title	BL14B2_XAFS
Description	SPring-8 産業利用ビームライン XAFSセットアップ
Instrument function	
Category	spectroscopy
Sub category	x-ray absorption spectroscopy
Manufacturer	
Organization	Japan Synchrotron Radiation Institute
Managing organization	
Organization	Japan Synchrotron Radiation Institute
Specimen details	
Specimen type	
Title	Copper



RDE local dictionary

synced

MatVoc



refer

API query

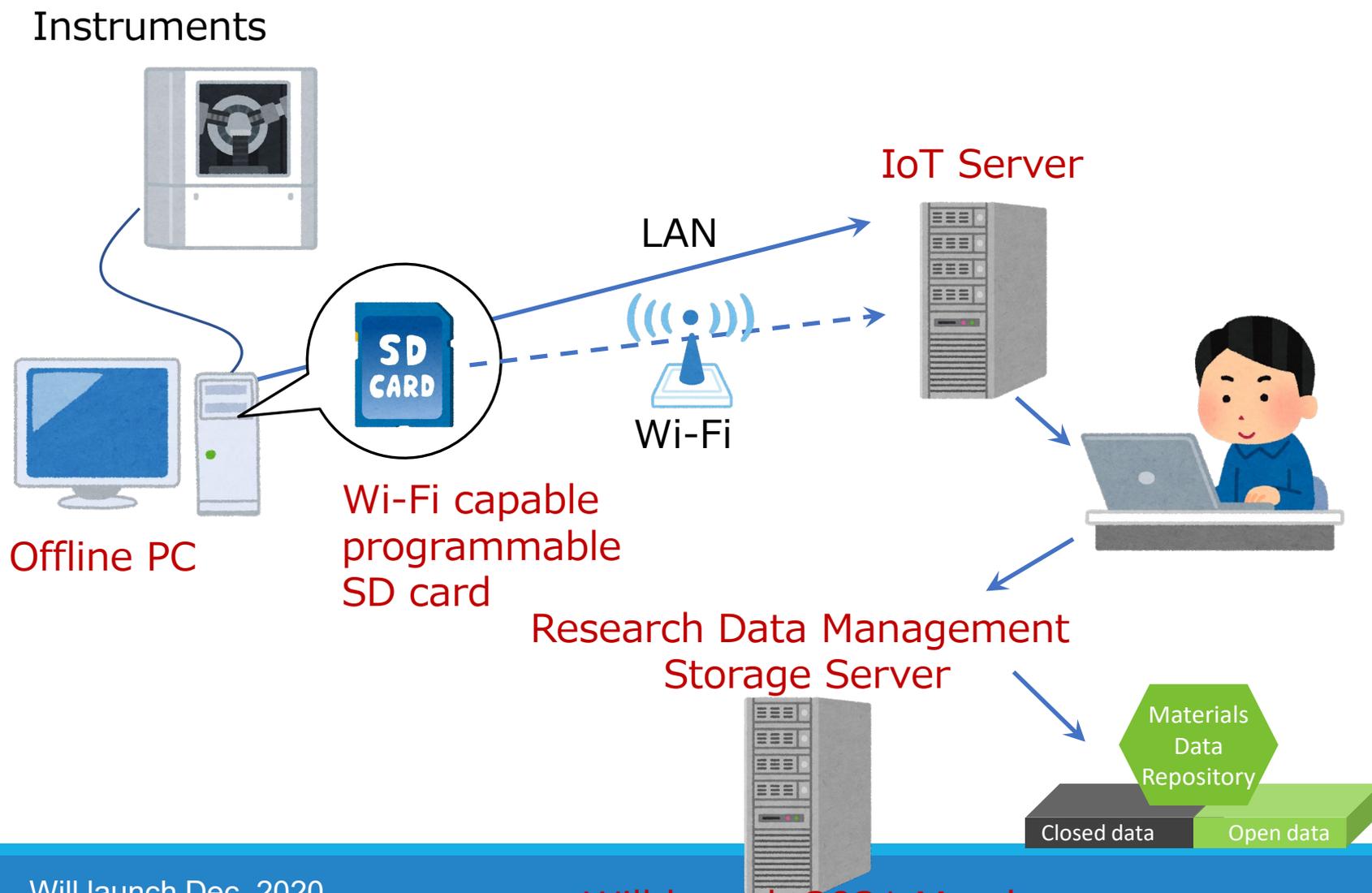
MatVoc ID: Q386

x-ray absorption spectroscopy (Q386)			
No description defined			
XAS			
= In more languages Configure			
Language	Label	Description	Also known as
English	x-ray absorption spectroscopy	No description defined	XAS
日本語	X線吸収分光法	No description defined	XAS
Statements			
has broader	spectroscopy	edit	
		= 1 reference source item	NIST Materials Data Vocabulary

M. Ishii, H. Nagao, A. Matsuda, K. Tanabe, H. Yoshikawa, 23rd XAFS Forum (2020)



Automatic data collection using WiFi-SD cards towards FAIR data

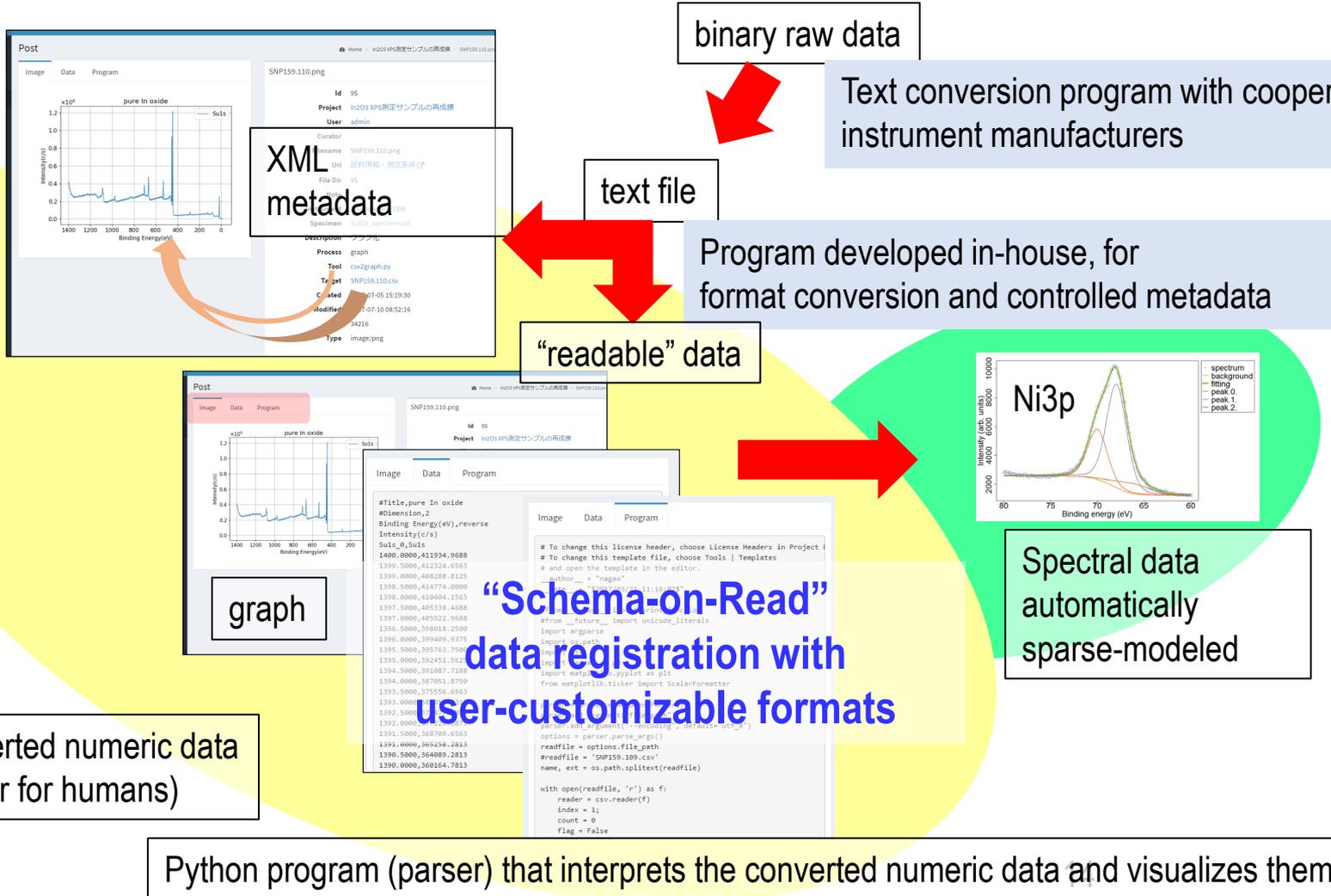


Will launch Dec, 2020

Will launch 2021 March



Data Collection System for efficient measurement data collection and automatic conversion



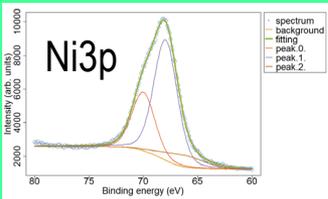
binary raw data

Text conversion program with cooperation of the instrument manufacturers

text file

Program developed in-house, for format conversion and controlled metadata

"readable" data



Spectral data automatically sparse-modeled

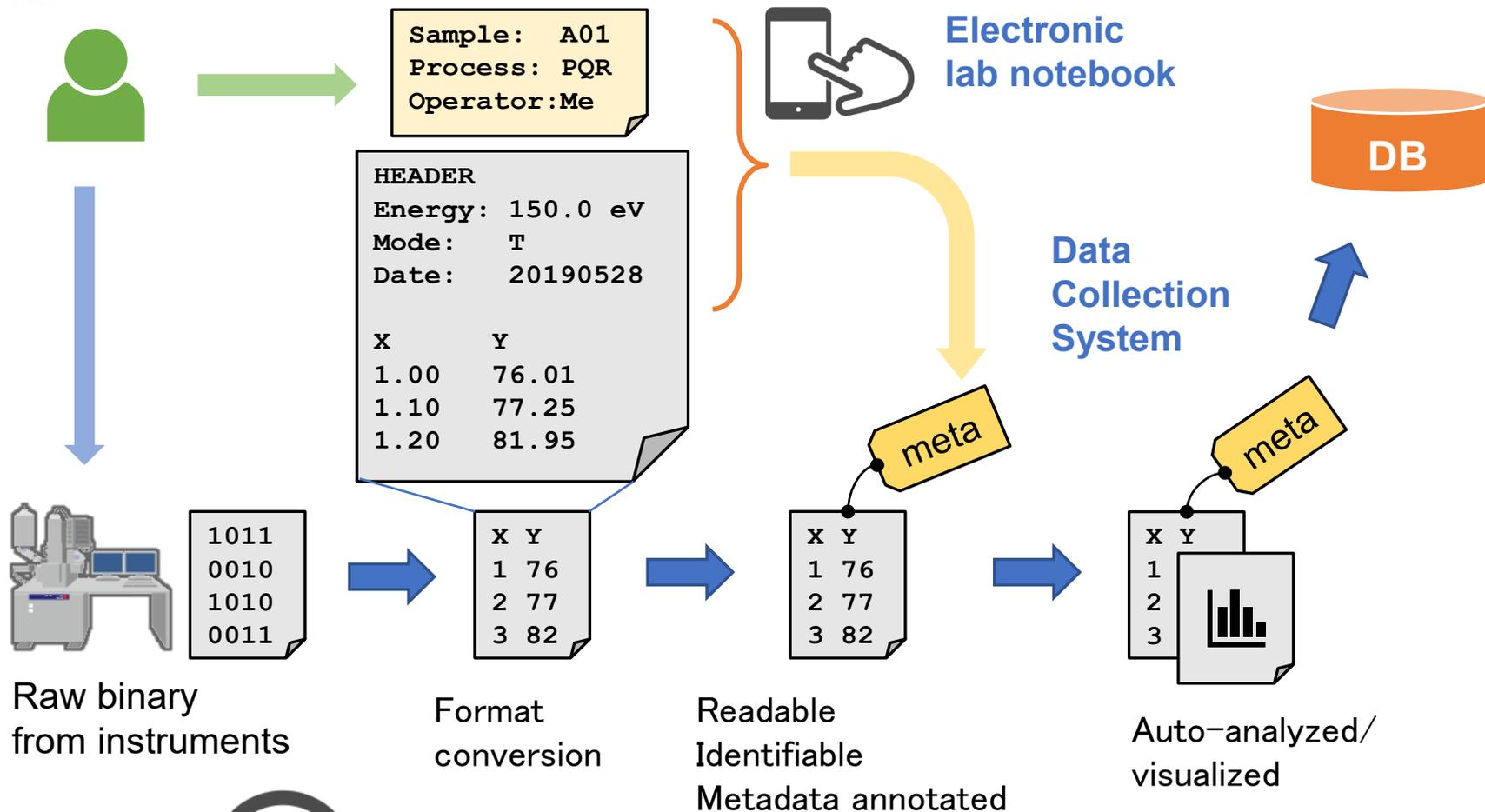
"Schema-on-Read" data registration with user-customizable formats

Converted numeric data (easier for humans)

Python program (parser) that interprets the converted numeric data and visualizes them



Collecting experimental data



Automatic data transfer using IoT



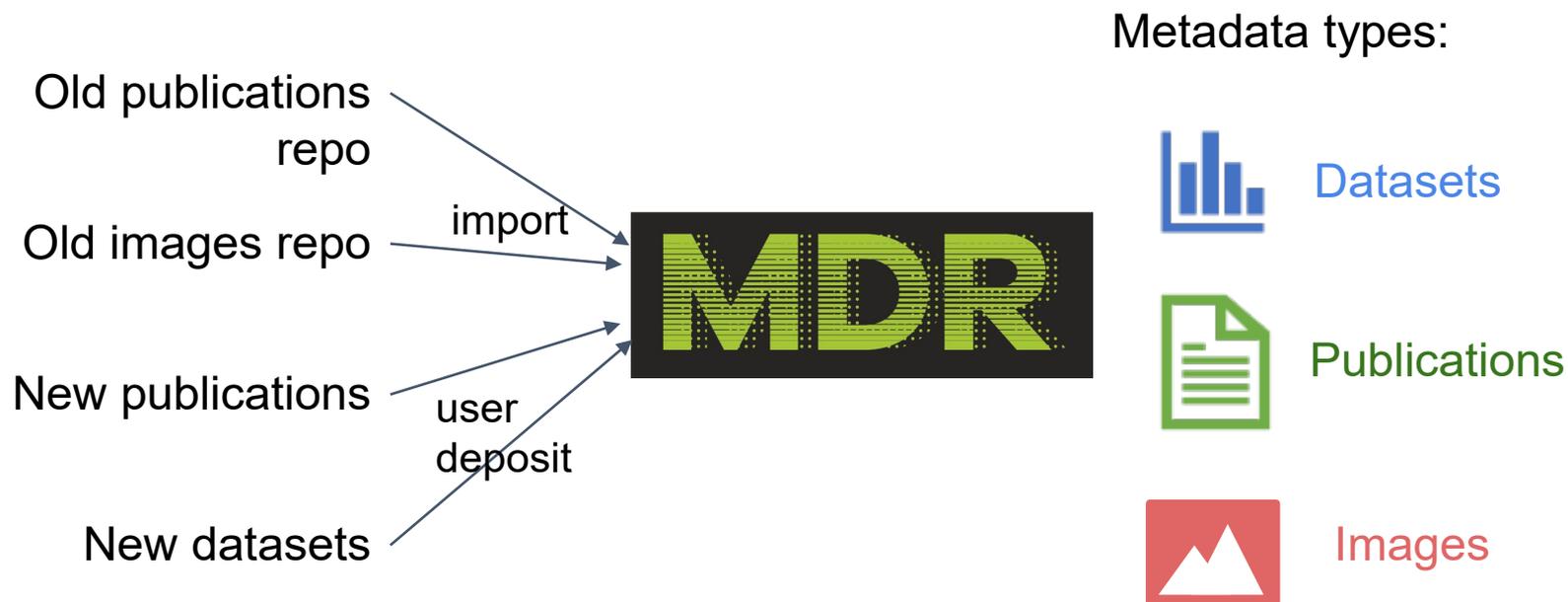
Data Conversion Tools on [GitHub.com/nims-dpfc](https://github.com/nims-dpfc)





Materials Data Repository – transition stage 1

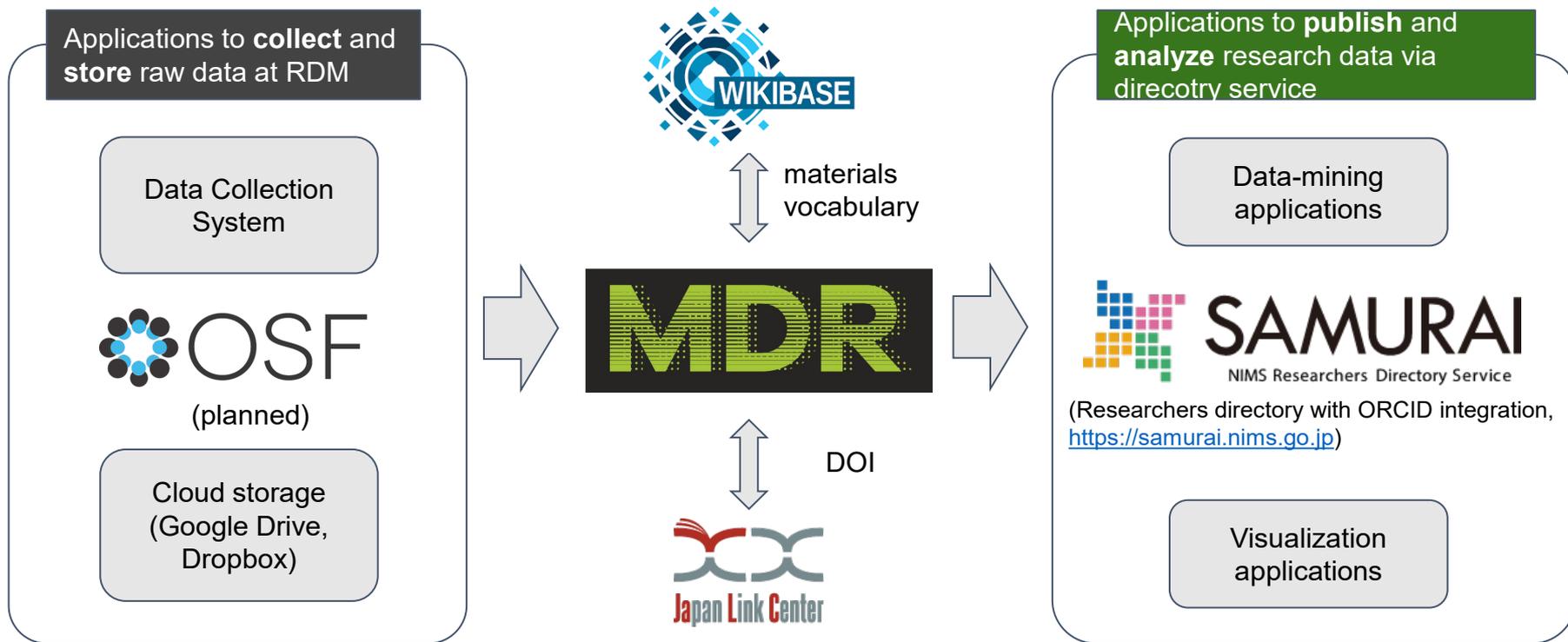
Datasets, publications, and images published on the Materials Data Repository

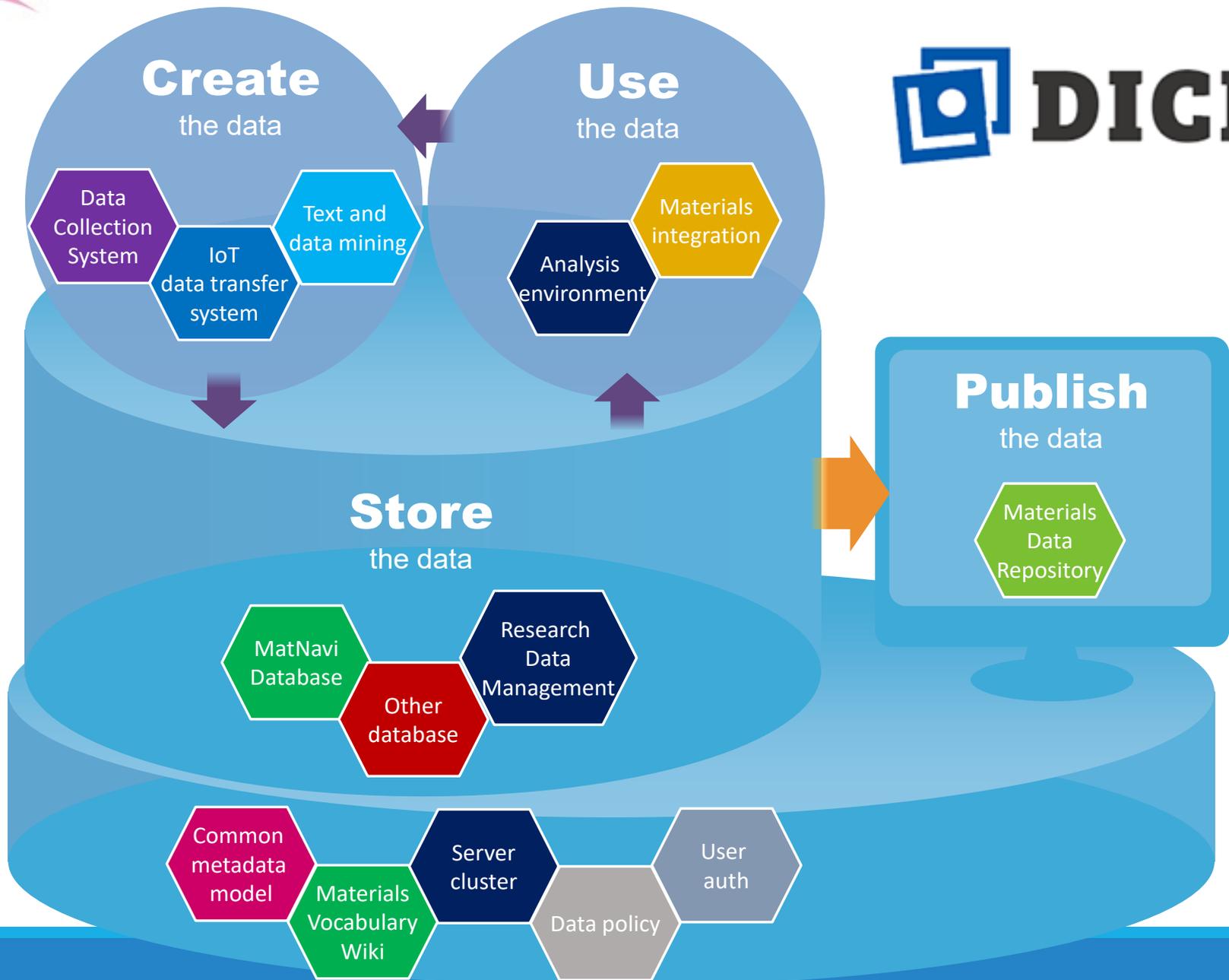




Materials Data Repository – transition stage 2

Integration ↔ FAIRable ↔ Accumulation

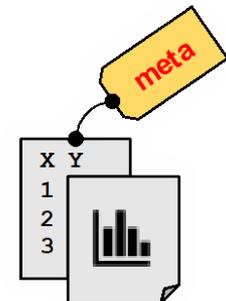






Summary

- Materials Data Platform, DICE as FAIRable platform
- Public and internal services will be launched 2020 - 2021
- DICE as R&D platform for academia and industries
- DICE aims to Japan-wide data hub with:
 - library of data models and meta data schemas for target materials
 - data quality guideline (recommendation) in respond to what data scientists need between what materials scientists can cope with.



*Store Analyze*

DICE is a data platform for all experts offering quality data and applications for materials science

*Publish Discover**Apply Create Use*

"As a researcher, it is tedious to provide a universal, machine-comprehensible metadata to allow my data to be used in serendipitous ways."

"As a materials scientist, there is a particular material I am looking for, but I don't know how to design a model for AI-assisted research."

"As a data scientist, I am looking for materials data that I can use for my data-driven operations."

DICE provides three things: high-quality data, applications, and knowledge -- bridging between the domain experts to aid rapid advancement in materials science.



Dr. Takuya Kadohira

Dr. Kosuke Tanabe

Dr. Asahiko Matsuda



And

Dr. Hideki Yoshikawa