

JULICH

Implementing FAIR through a distributed data infrastructure

Carsten Hoyer-Klick



Research Grant 03EI1005A-G





Challenges I: Finding and Accessing Data

- > Many data bases exist, each in its own flavor
 - > Data access
 - > Data format
 - > Data licenses (if at all)
 - > Sometimes hard to find
- > Data collection is a labor intensive task
- Data cleaning, aggregation, etc.
 is repeated by many researchers with different results
- > Data quality is often unknown

Imagine a library without a catalog and systematic numbers on the shelfs



Research Grant 03EI1005A-G



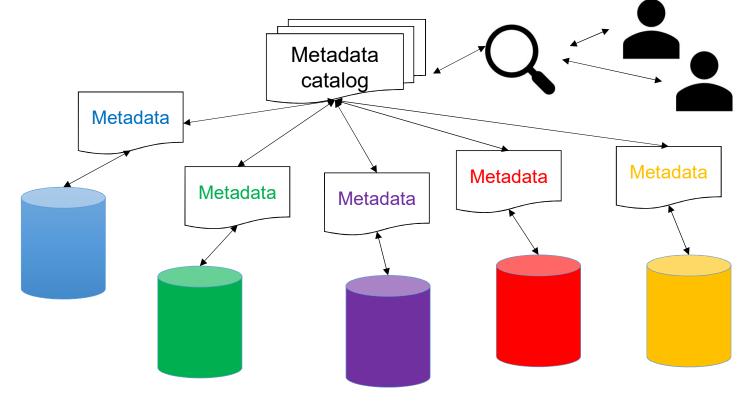
Universität Stuttgart IER Institut für Energiewirtschaft und Rationelle Energieanwendung





Solution I: A Metadata Catalog

- A metadata catalog harvests the (rich) metadata from the available data sources
- The catalog can be used to discover data
- The metadata contains a URI to the actual data
- In case of data bases possibly also an API/Interface description



The domain uses the OEP Metadata string: https://openenergy-platform.org/tutorials/jupyter/OEMetadata/



Research Grant 03EI1005A-G



Universität Stuttgart IER Institut für Energiewirtschaft und Rationelle Energieenwendung

ant Porspan Institute For sowendung P I K





FAIR Data Principles

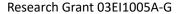
> FAIR stands for

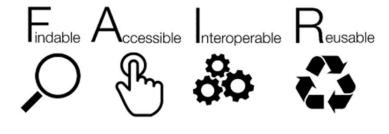
- Findable: Data and Metadata should be easy to find for both humans and computers. Machine readable metadata are essential for automatic discovery of data sets and services
- Accessible: Once the user finds the required data, she/he needs to now how they can be accessed
- Interoperable: The data need to be integrated with other data. In addition the need to interoperate with applications or workflows for analysis, storage and processing
- Reusable: The ultimate goal of FAIR is to optimize the reuse of data. To achieve this, metadata and data should be welldescribed so that they can be replicated and/or combined in different settings.

Deutsches Zentrum für Luft- und Raumfahrt

German Aerospace Center







Universität Stuttgart

IER Institut für Energiewirtschaft

und Rationelle Energieanwendung

InfAl

Institute for Applied Information



Source: go-fair.org



Why Rich Metadata?





Metadata contains vital information

DE	_BadenWu	e DE_	_Bayern D	DE_Berlin	DE	_Brandenburg
processHeat500	t0001	0.000087	0.00007	6 0.000	07	0.00007
processHeat500	t0002	0.000085	0.00007	6 0.000	068	0.000068
processHeat500	t0003	0.000086	0.00007	6 0.000	069	0.000069
processHeat500	t0004	0.000087	0.00007	7 0.000	07	0.00007
processHeat500	t0005	0.000121	0.00011	6 0.000	112	0.000112
processHeat500	t0006	0.000167	0.00016	8 0.000	17	0.00017
processHeat500	t0007	0.000195	0.00019	9 0.000	203	0.000203
processHeat500	t0008	0.000204	0.00020	8 0.000	212	0.000212
processHeat500	t0009	0.000205	0.00020	9 0.000	212	0.000212
processHeat500	t0010	0.000188	0.00019	0.000	191	0.000191
processHeat500	t0011	0.000187	0.00018	8 0.000	189	0.000189
processHeat500	t0012	0.000184	0.00018	7 0.000	187	0.000187
processHeat500	t0013	0.000182	0.00018	4 0.000	186	0.000186
processHeat500	t0014	0.000181	0.00018	3 0.000	185	0.000185

What do the collumns stand for?



Without a label, how would you know which one to open?



Research Grant 03EI1005A-G



Universität Stuttgart IER Institut für Energiewirtschaft und Rationelle Energiesnwendung





The OEO Metadata String

- Developed within the "The Open Energy Familiy" as part of open science activities within the Domain of Energy Systems Analysis
- > Implemented as a JSON-String -> can still be read and edited by humans
- > Contains all main elements of metadata
 - > **Descriptive**: Descriptions of the content, identification, make the data findable
 - > Structural: References to other (meta-)data, information about the structure of the data
 - > Administrative: Management of the data
 - > Use: How to use the data, protocols, API
 - > Technical: functions of a systems, digitalization, hardware and software used
 - > Preservation: how the data is stored
 - > Rights: Access rights and licenses
 - > Provenance: Source of the data

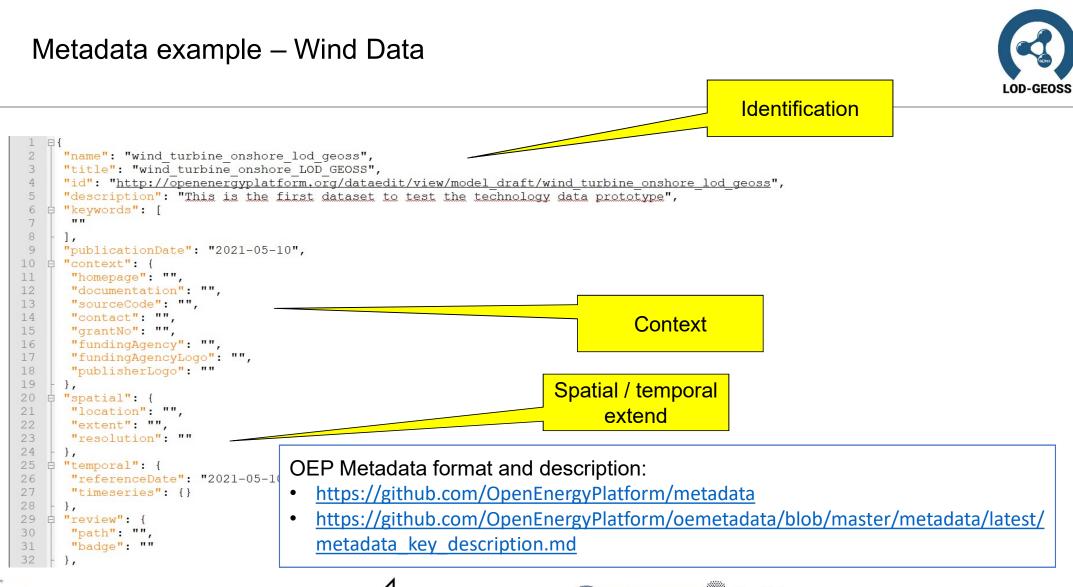


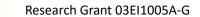




Universität Stuttgart IER Institut für Energiewirtschaft und Rationelle Energieanwendung







für Luft- und Raumfahrt German Aerospace Center InfAl

Institute for Applied Informatics

Deutsches Zentrum

DLR

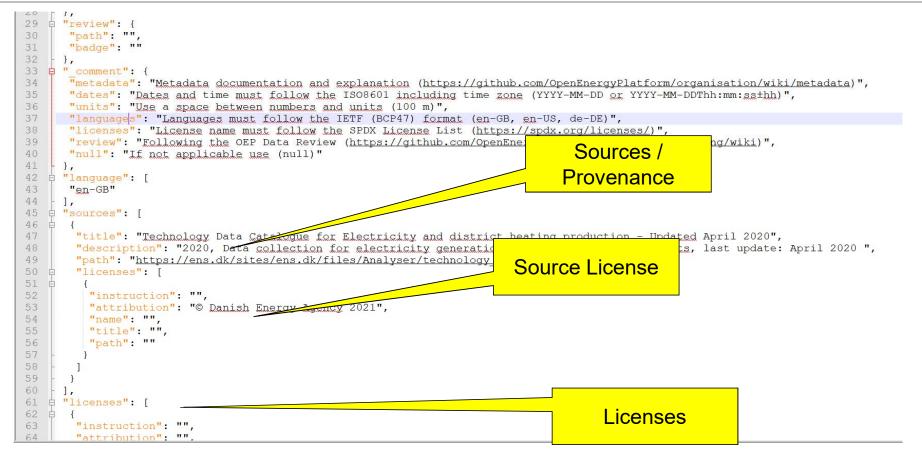
Universität Stuttgart IER Institut für Energiewirtschaft und Rationelle Energieanwendung

t Potsdam Institute For Cuimare Immach Research





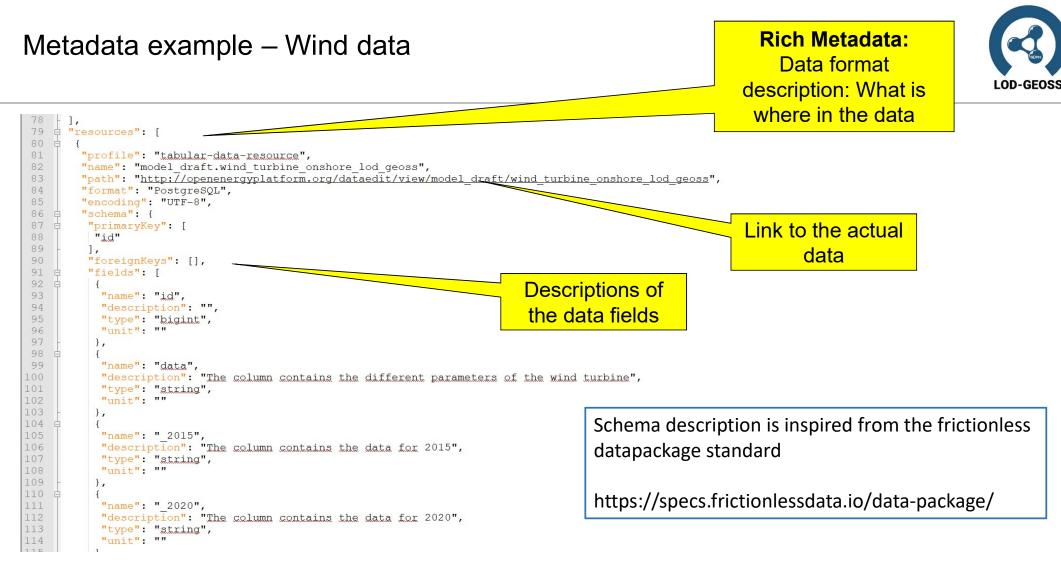
Metadata example – Wind data





Research Grant 03EI1005A-G

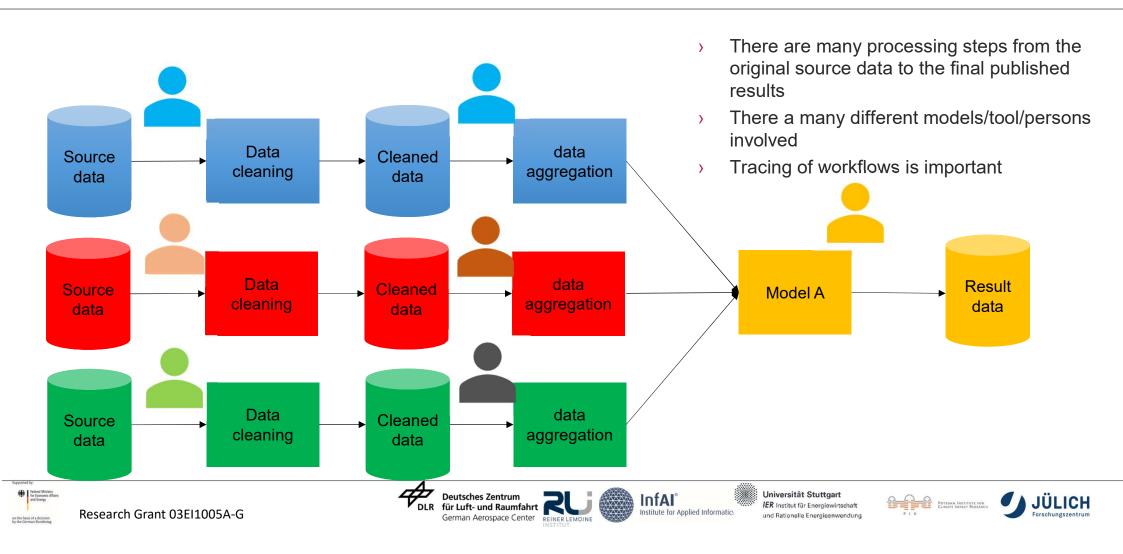






Research Grant 03EI1005A-G

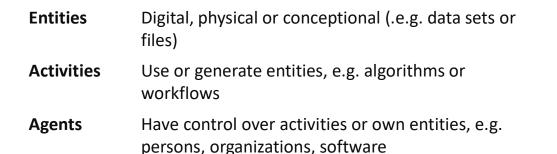




Challenge II: Who has done what to the data?



PROV of the W3 Provenance Working Group



Entities, Activities and Agents are related through relations with each other.

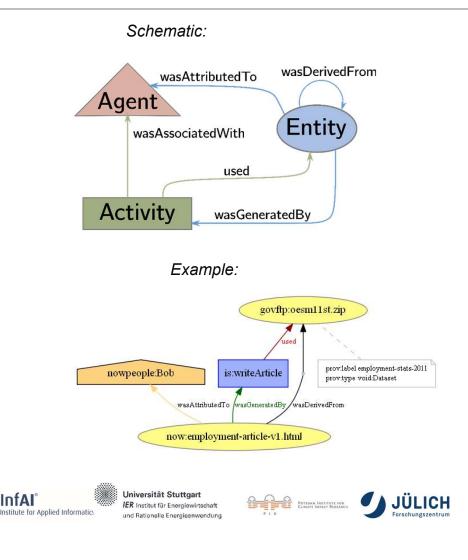
Provenance Information can be collected automatically and visualized in graphs, e.g. with ProvStore

Deutsches Zentrum

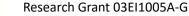
für Luft- und Raumfahrt

German Aerospace Center

A record **of who (Agent)** processed with which **Tools (Activity)** which **data (Entities, input and output)**



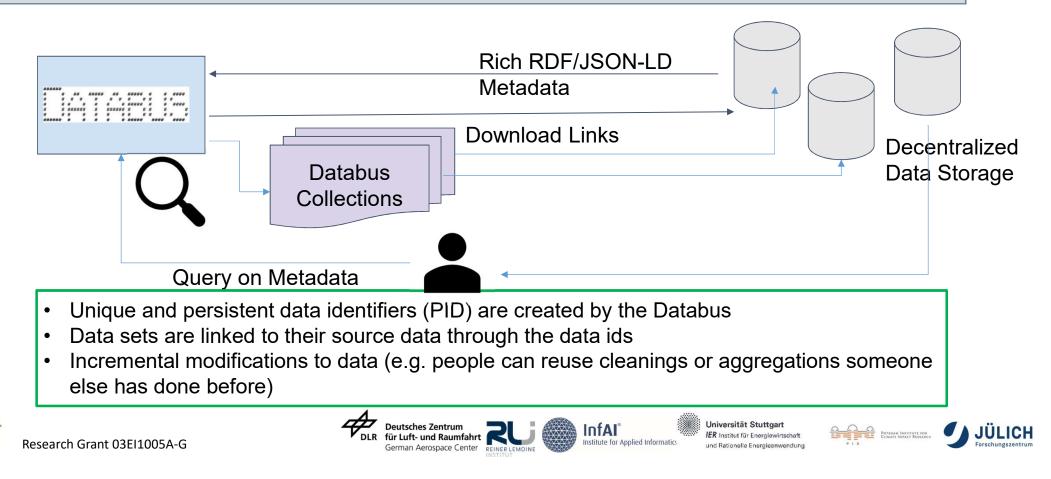




The Databus Platform

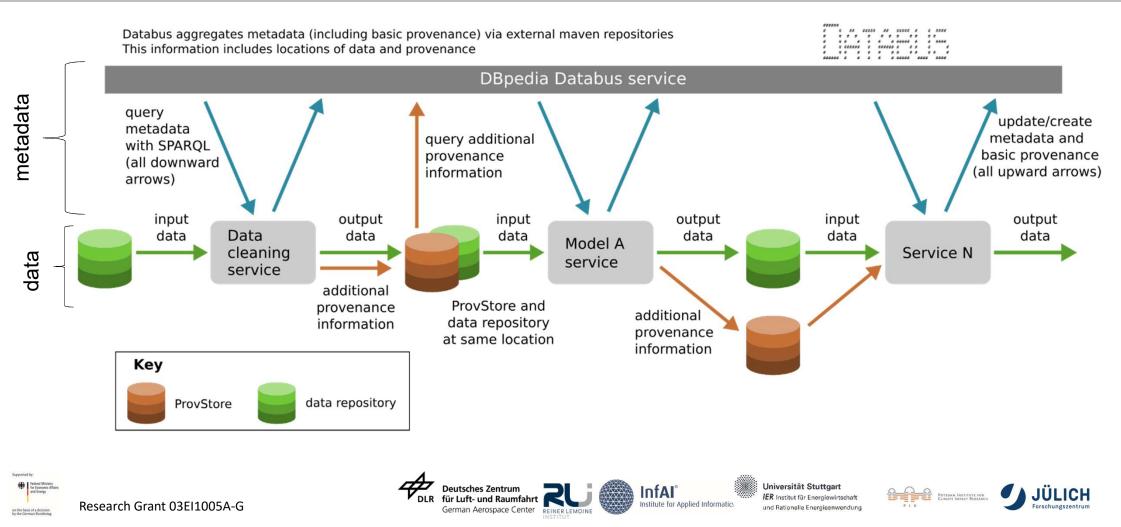


Databus is a virtual bus. It can address files on the web and coordinate dataflows based on DataID metadata. No actual data is uploaded to the bus.





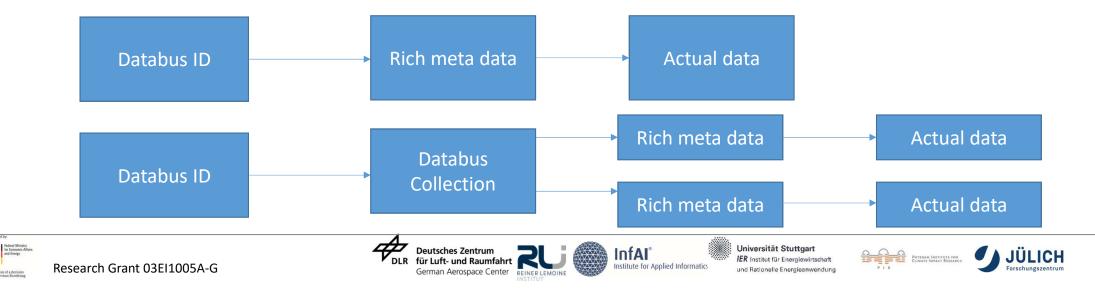
Architecture Concept





Databus as PID / Fair Digital Object (FDO) provider

- > Uploading meta data to the Databus creates a unique ID for the data set which can be used as a persistent data identifier (PID) to cite and share the data
 - > Databus URI (PID): <u>https://databus.dbpedia.org/account_name/group/artifact/version /</u>
 - > Databus collections can be used to group data
 - > The actual data is linked within the meta data
- > Databus ID can be used to access all necessary information
- > Databus ID is a pointer to a (FAIR) digital object
- > Persistent ID to make data citable





Challenge III: What is the data about?

- Each data source comes with its own annotation >
- Example from solar meteorology:)
 - > GHI: Global Horizontal Irradiation
 - > Global: Could also be Global Horizontal Irradiation
 - > Surface downward irradiation: The usual term in climate science for what we usually call GHI
- Taxonomies or ontologies create a data language to annotate data >
- Ontologies can describe relations: *direct radiation* is a part of the global irradiation reaching the surface >
- Ontologies make data interpretable, also by machines and algorithms >
- Good ontologies are created on a consensus building and open development process within the community. >
 - > We use the 'Open Energy Ontology', <u>https://openenergy-platform.org/ontology/</u>, https://doi.org/10.1016/j.egyai.2021.100074
- New JASON-LD Meta data string can link meta data to an ontology)









Implementing the FAIR Principles



>	Findable	Development and setup of the data bus as metadata catalog for data in energy systems analysis
>	Accessible	 Descriptions of the data formats and interfaces to data bases, Best Practice Guides, links to the data sources in the metadata Databus architecture and meta data Citable PIDs
>	Interoperable	 Enhancement of on Open Energy Ontology as a common data language, open data format descriptions The Open Energy Ontology Schema in the metadata
>	Reusable	Improvement of data licenses, here with a special focus on GEOSS data, provenance information with the data













> Global:

- > @id URI that links the resource via the dbpedia Databus / Persistent ID / PID "@id": "https://databus.dbpedia.org/denis/lod-geoss-example/api-example/2021-05-10/apiexample_type=turbineData.json"
- > Subject: an OEO ID for the dataset subject": "OEO:0000448"
- > Schema
 - Is_about: OEO links to the relenvat terms







Metadata Tools for the Open Energy Platform

- Data-Wizzard for uploading data to the Open Energy Platform >
 - > https://openenergy-platform.org/dataedit/wizard/# -> Edit metadata
- Stand-alone tool for creating metadata >
 - > https://meta.rl-institut.de/meta_creator/
- Both tools are not supporting the 1.5 version yet, this still needs to be done. >

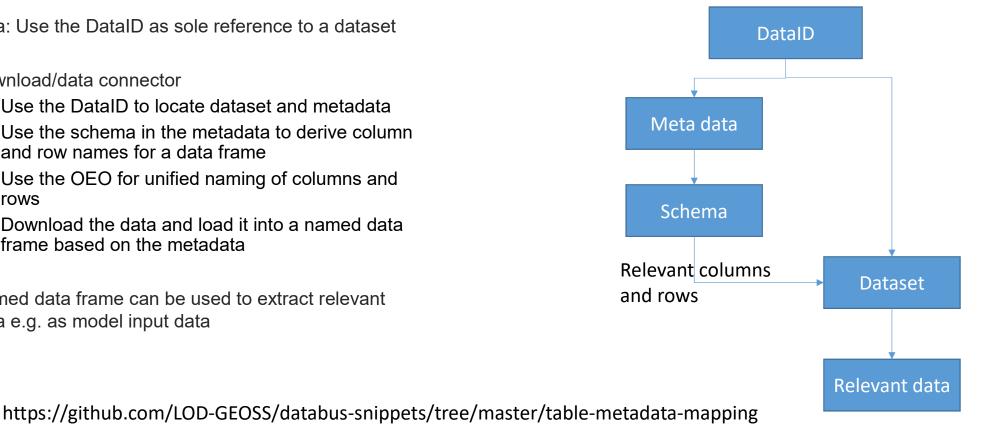






FDO application example

- Idea: Use the DataID as sole reference to a dataset 5
- Download/data connector >
 - > Use the DataID to locate dataset and metadata
 - > Use the schema in the metadata to derive column and row names for a data frame
 - > Use the OEO for unified naming of columns and rows
 - > Download the data and load it into a named data frame based on the metadata
- Named data frame can be used to extract relevant > data e.g. as model input data



Research Grant 03EI1005A-G



Universität Stuttgart IER Institut für Energiewirtschaft und Rationelle Energieanwendung

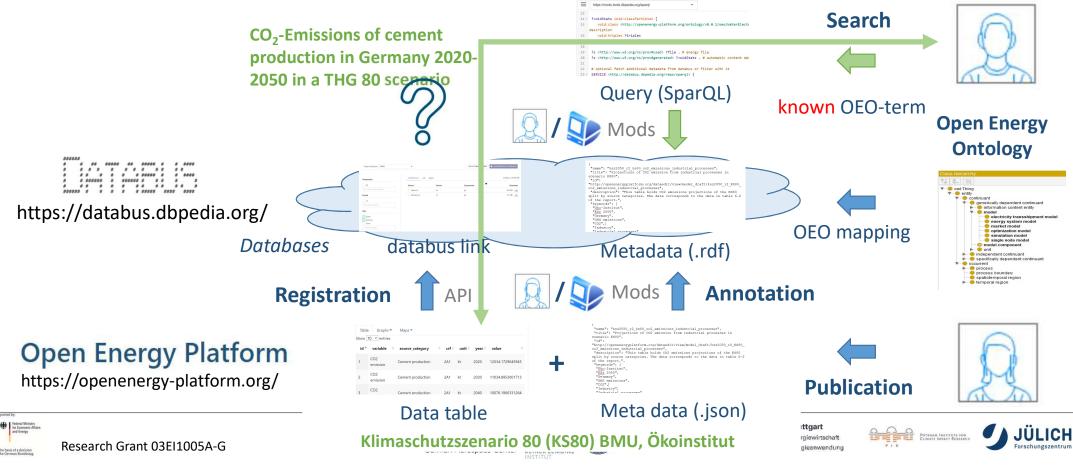




Demonstrator: Publication of a Data Set Using the databus

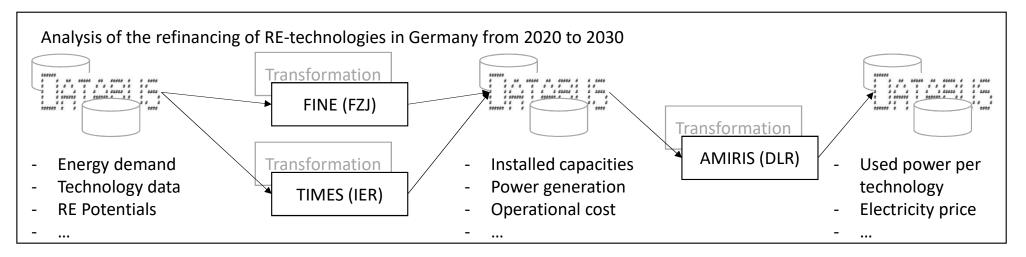


 Goal: Demonstration of the improved visibility and improved discovery of a data set through the registration in the databus





Demonstrator Model Coupling



- > Goal of the demonstrator
 - Provision of a holistic application example of a databus based coupling of heterogeneous energy system models for answering a research question.
 - > Databus and Metadata description as a standard interface for data between models
- > Challenges in the handling of data
 - > Preparation and description of data content (data and metadata standards)
 - > Preparation and description of data formats (open licenses, machine readable formats, ...)
 - > Discovery, interpretation and further processing of the data







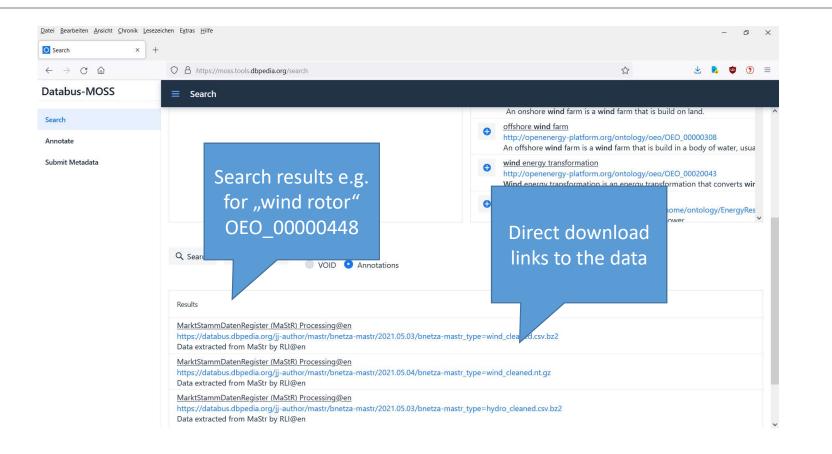


Searching data with the MOSS





Searching with MOSS





Research Grant 03EI1005A-G







- a service to manage and search registered metadata)
- Persistent identifies for tracing data processing and citing data >
- Databus as pointers to digital objects >
- The databus can improve the data economy and scientific efficiency > as we can build better on existing previous work, improve the data instead of repeating work others have done before.
- The databus supports the implementation of FAIR principles in the > **Domain of Energy Systems Analysis**
- The developed architecture in conjunction with the use of the Open > Energy Ontology enables semantic searches for data in the domain of energy systems analysis
- The developed architecture can improve open data exchange, > model coupling, tracing of workflows and collaboration for better scientific results

Deutsches Zentrum

für Luft- und Raumfahrt

German Aerosnace Center





- https://lod-geoss.github.io
- https://databus.dbpedia.org >
- https://openenergyplatform.org/tutorials/jupyter/OEMetadata/
- https://openenergy-platform.org/ontology/
- https://doi.org/10.1016/j.egyai.2021.100074
- https://moss.tools.dbpedia.org/search
- https://github.com/LOD-GEOSS/databussnippets



InfAl









Contact us

German Aerospace Center

Energy Systems Analysis

Carsten Hoyer-Klick

Email: carsten.hoyer-klick@dlr.de

Supported by:

Federal Ministry for Economic Affairs and Energy

on the basis of a decision by the German Bundestag Research Grant 03EI1005A-G

License



© German Aerospace Center

Except where otherwise noted, this work and its content (texts and illustrations) are licensed under the

Attribution 4.0 International (CC BY 4.0)

See license text for further information.

Please cite as:

"Implementing FAIR through a distributed data infrastructure" © Hoyer-Klick, et. Al. | CC BY 4.0

InfAl

Institute for Applied Informatics

Federal Ministry for Economic Atlairs and Energy on the basis of a decision

Research Grant 03EI1005A-G



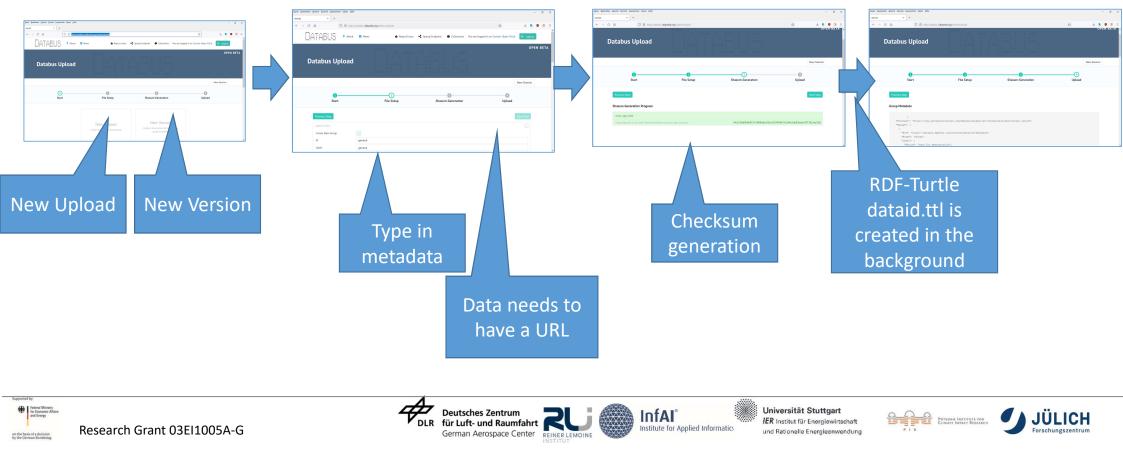
Universität Stuttgart IER Institut für Energiewirtschaft und Rationelle Energiesnwendung



Uploading data to the Databus https://github.com/LOD-GEOSS/databus-snippets



https://databus.dbpedia.org/system/upload





Uploading data to the Databus https://github.com/LOD-GEOSS/databus-snippets



2. Using the API

- 1. Create DataID in RDF-Turtle
 - > Create group-meatadata: Documentation for the group
 - > Create the actual dataID for the data with the URI to the data
- 2. Fetch a Bearer Token from the databus, done by a http GET Requrest TOKEN=\$ (curl -s -d 'client_id=upload-api' -d 'username=XXXXXX' -d 'password=XXXXXXXXX' -d 'grant_type=password' https://databus.dbpedia.org/auth/realms/databus/protocol/openidconnect/token | cut -d'"' -f 4).
- 3. Use bearer token for deploying to the databus with http PUT Request curl -H '{Authorization: Bearer \$TOKEN}' -X PUT <GROUP-OR-VERSION-URI> -d 'DATAID JSONLD CONTENT'
- 3. Using the python3 script use the databus_api_example.py and dev_Databus_api_example.py by modifying the variables below if__NAME=="__main__"







Contents of GroupID and DataID



> GroupID

- > ID
- > Label
- > Abstract
- > Description

- > DataID
 - > ID
 - > Label
 - > Title
 - > Publisher (autogenerated)
 - Artifact (autogenerated)
 - Group (autogenerated)
 - > Abstract
 - > Version
 - > License
 - > License URI
 - > License Documentation
 - > Files (URIs to the data)



Research Grant 03EI1005A-G



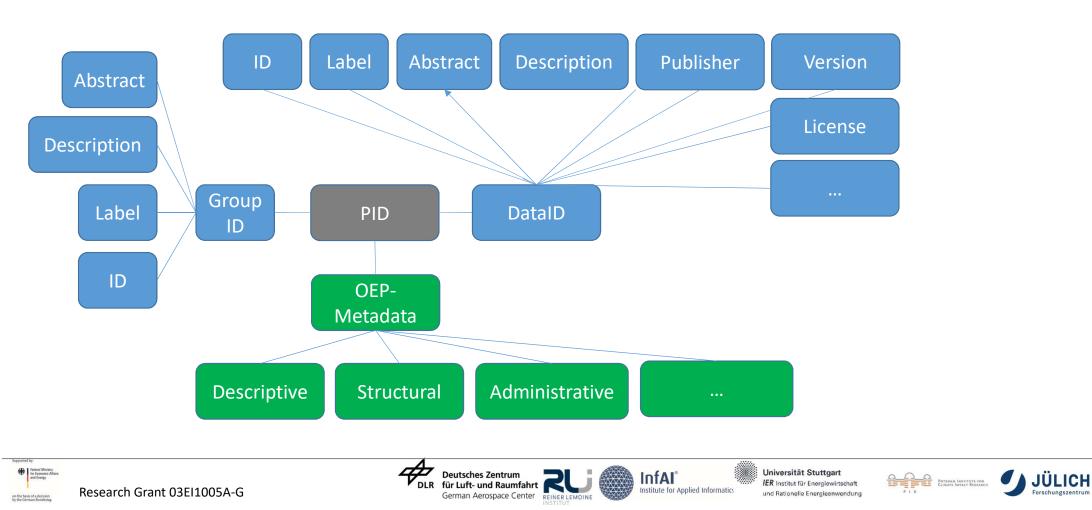






Extending the Graph by Additional Metadata







Annotation of data: MOSS Metadata Overlay Search System

- Only basic metadata (license, publisher, etc) is used for the registration of data on the Databus
- Databus Mods (an additional software layer) are used to add additional information to a data set
- > Additional metadata can be added to any data set following different domain standards via databus mods
- MOSS: Metadata Overlay Search System: moss.tools.dbpedia.org
 - 1. Generate Metadata in JSON-LD (e.g. the updated OEP metadata string)
 - 2. Find the URI of the dataset and enter it into the MOSS submission form at: <u>https://moss.tools.dbpedia.org/submit-data</u>







Universität Stuttgart IER Institut für Energiewirtschaft und Rationelle Energieanwendung





