

# NFDI4BIOIMAGE

Carsten Fortmann-Grote<sup>1</sup>

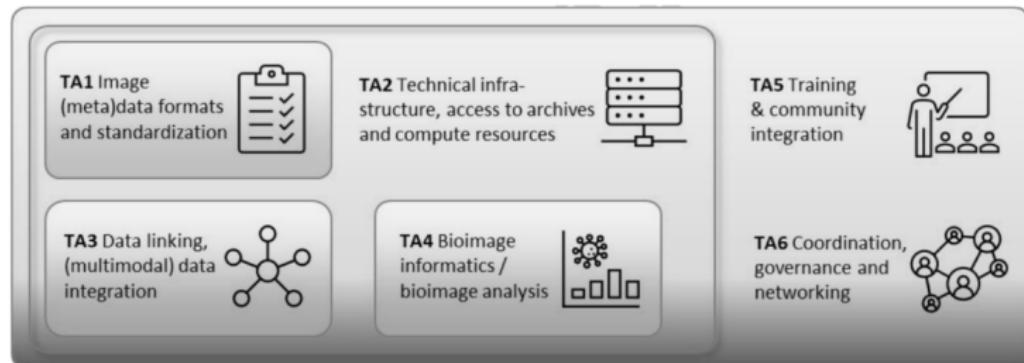
April 18 2024



This presentation is released to the public domain under Creative Commons Attribution License CC-BY-4.0 International.



# The NFDI4BIOIMAGE Consortium



- Harmonization of bioimage data and metadata formats and standards
- FAIR Image Objects
- Multimodal metadata vocabularies
- Standard Operating Procedures for FAIR image processing
- Proliferation of best practices in image data management and processing

# All Hands Meeting Oct. 2023



# No Milestones due yet

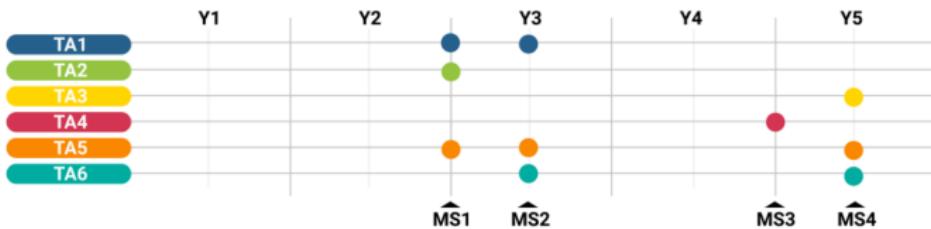


Figure 13: Overview of the project milestones, consisting of the indicated milestone parts from the Task Areas.

- Milestone 1: White paper on the FAIR-IO concept published (MS1-TA1), production services for NFDI4BIOIMAGE available (MS1-TA2) and search index for bioimage RDM training materials implemented (MS1-TA5)
- Milestone 2: Release of RFC process document (MS2-TA1), first NFDI4BIOIMAGE conference (MS2-TA6), and publication of survey results and search index statistics (MS2-TA5)
- Milestone 3: SOPs for FAIR usage of image analysis workflows available (MS3-TA4)
- Milestone 4: Second NFDI4BIOIMAGE conference (MS4-TA6), preview of over-arching database for multimodal metadata (MS4-TA3), and publication of the search index statistics and survey results (MS4-TA5)

# Outline

- FAIRification of image data: 5 star scheme
- Teaching & Training material
- BIDS, ARC, OMERO interfaces
- Data integration at MPI-EvolBio

# FAIRification of Image (meta)data

Quantify the FAIRness of containerized Bioimage research objects

# FAIRification of Image (meta)data

5 stars linked open data



# FAIRification of Image (meta)data

## 5 star bioimage containers

### 5 Stars



Package your data for sharing  
(with permission)



including descriptive labels



that are machine readable



in consistent containers



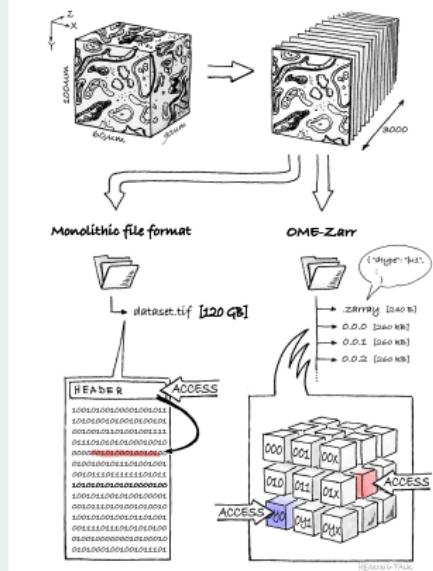
to computationally link resources.

*Images generated by Andra Waagmeester & Josh Moore through AI generation using Midjourney (Sep. 2023)*

# Zarr is the technical backbone of next generation (bio)image file formats

## What is Zarr?

- hdf5 like hierarchical layers
- chunks stored in separate "files" (blobs), linked by json file (vs. monolithic binary blob in hdf5)
- suitable for object storage

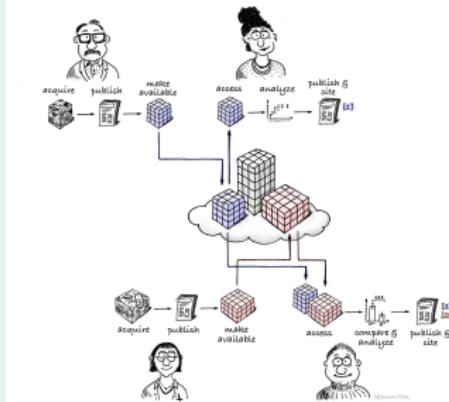


10.5281/zenodo.7037679

Zarr is the technical backbone of next generation (bio)image file formats

## What is Zarr?

- IO latency issues
- NFDI4BI leads community efforts for specification and implementation
- BioImage Archive (BIA), NASA, Open Geospatial Consortium (OGC) have committed to adopting Zarr as community standard

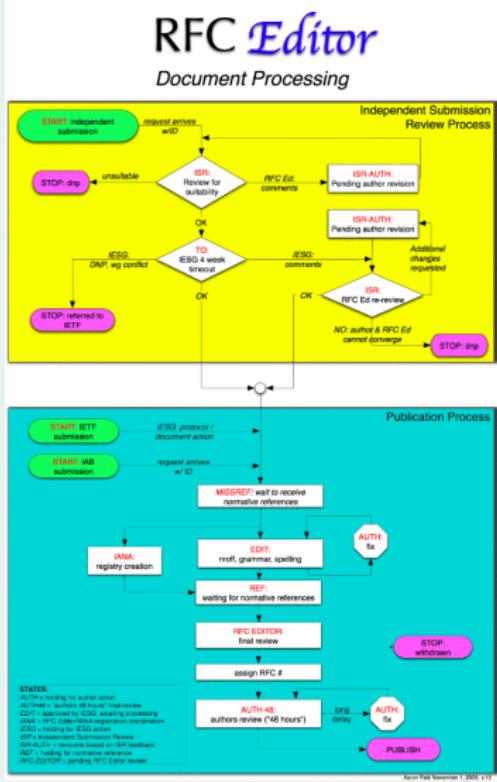


10.5281/zenodo.7037679

# The RFC (Request For Comments) process

- Established Internet Engineering Task Force (IETF) procedure for
  - protocols (HTTP, FTP)
  - best practices (RFC)
  - standards (RDF)
  - file formats (XML, HTML)
  - anything that defines the internet

## The RFC process



## Example: RFC2119

BEST CURRENT PRACTICE  
Errata Exist  
S. Bradner  
Harvard University  
March 1997

## **Key words for use in RFCs to Indicate Requirement Levels**

**Status of this Memo**

This document specifies an Internet Best Current Practices for the Internet Community, and requests discussion and suggestions for improvements. Distribution of this memo is unlimited.

## Abstract

In many standards track documents several words are used to signify the requirements in the specification. These words are often capitalized. This document defines these words as they should be interpreted in IETF documents. Authors who follow these guidelines should incorporate this phrase near the beginning of their document:

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119.

Note that the force of these words is modified by the requirement level of the document in which they are used.

1. **MUST** This word, or the terms "REQUIRED" or "SHALL", mean that the definition is an absolute requirement of the specification.
  2. **MUST NOT** This phrase, or the phrase "SHALL NOT", mean that the definition is an absolute prohibition of the specification.
  3. **SHOULD** This word, or the adjective "RECOMMENDED", mean that there may exist valid reasons in particular circumstances to ignore a particular item, but the full implications must be understood and carefully weighed before choosing a different course.
  4. **SHOULD NOT** This phrase, or the phrase "NOT RECOMMENDED" mean that there may exist valid reasons in particular circumstances when the particular behavior is acceptable or even useful, but the full implications should be understood and the case carefully weighed before implementing any behavior described with this label.

# Training (TA5)

## I3D-bio customizable slide decks

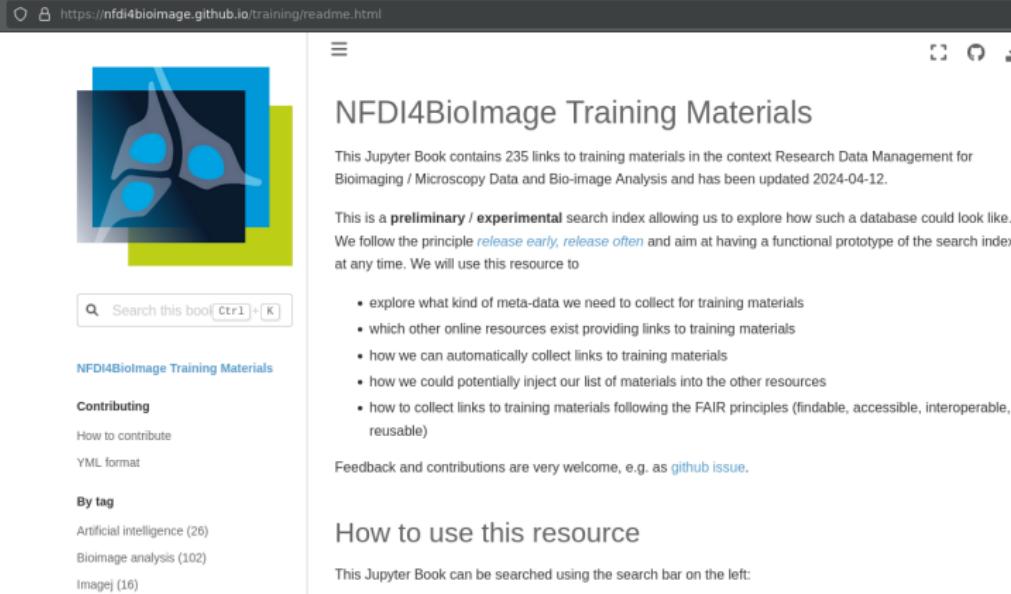
Training material by the I3D:bio team

I3D:bio's OMERO training material: Re-usable, adjustable, multi-purpose slides for local user training

Description	Links	License & Citation
<ul style="list-style-type: none"><li>Record on zenodo with all material for download: <a href="https://zenodo.org/records/8323588">https://zenodo.org/records/8323588</a></li><li>Youtube playlist of the video tutorials: <a href="https://www.youtube.com/playlist?list=PL2k-L-zWPoR7SHjG1HhDlwLZj0MB_stlU">https://www.youtube.com/playlist?list=PL2k-L-zWPoR7SHjG1HhDlwLZj0MB_stlU</a></li></ul>		<p>Research Data Management for Bioimage Data at the <a href="#">ADD INSTITUTE HERE</a></p> <p>What is the image data management platform OMERO?</p> <p><a href="#">ADD AUTHOR / RESPONSIBLE PERSON FROM YOUR INSTITUTE</a></p>



## Jupyter Book + Search Index



The screenshot shows a web browser displaying the NFDI4BioImage Jupyter Book search index at <https://nfdi4bioimage.github.io/training/readme.html>. The page has a dark header with the title "NFDI4BioImage Training Materials". On the left, there's a sidebar with a logo of three stylized brain cells, a search bar, and sections for "Contributing", "How to contribute", "YML format", and "By tag" (with categories like Artificial intelligence, Bioimage analysis, and ImageJ). The main content area contains an introduction about the search index, a bulleted list of goals, and a note about feedback. Below that is a section titled "How to use this resource" with a note about searching.

NFDI4BioImage Training Materials

This Jupyter Book contains 235 links to training materials in the context Research Data Management for Bioimaging / Microscopy Data and Bio-image Analysis and has been updated 2024-04-12.

This is a **preliminary / experimental** search index allowing us to explore how such a database could look like. We follow the principle *release early, release often* and aim at having a functional prototype of the search index at any time. We will use this resource to

- explore what kind of meta-data we need to collect for training materials
- which other online resources exist providing links to training materials
- how we can automatically collect links to training materials
- how we could potentially inject our list of materials into the other resources
- how to collect links to training materials following the FAIR principles (findable, accessible, interoperable, reusable)

Feedback and contributions are very welcome, e.g. as [github issue](#).

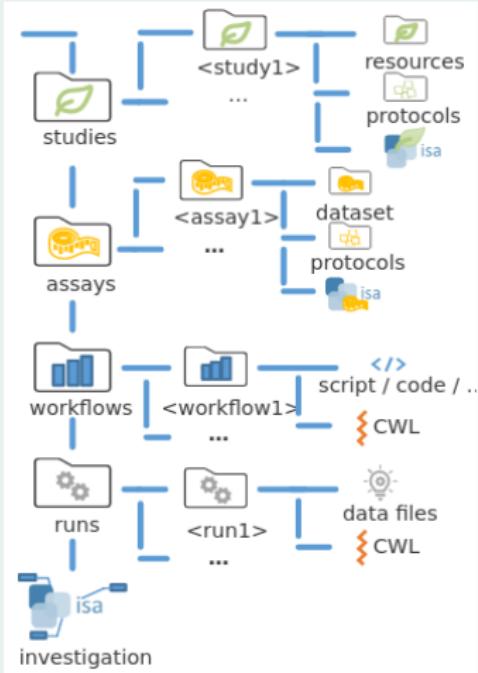
### How to use this resource

This Jupyter Book can be searched using the search bar on the left:



# BIDS-ARC-OMERO (TA1, TA3)

## ARC (Annotated Research Context, NFDI4Plants)



## OMERO

The screenshot shows the OMERO interface with the following details:

**Explore** tab is selected.

**Image ID:** 6752790  
**Owner:** Carsten Fortmann-Grete

**Image Details** panel:

- Acquisition Date:** 2023-12-15
- Import Date:** 2023-12-15
- Dimensions (XY):** 2048x2048
- Pixels Type:** uint8
- Pixels Size (XYZ) (µm):** 1.0x1.0x1.0
- Z-sections/Timepoints:** 1
- Channels:** 1
- ROI Count:** 0

**Tags** panel: 0

**Key-Value Pairs** panel: 1

Key	Value
Phenotype	WS
Genotype	awsX
Type	still
Organism	Pseudomonas
Strain	SBW25
MPB	15456
Parent	0
Investigation	Deep P

**File List:** awsX-01.jpg, awsX-02.jpg, awsX-03.jpg, awsX-04.jpg, awsX-05.jpg, awsX-06.jpg, awsX-07.jpg, awsX-08.jpg, awsX-09.jpg, awsX-10.jpg, awsX-11.jpg, awsX-12.jpg, awsX-13.jpg, awsX-14.jpg, awsX-15.jpg, awsX-16.jpg, awsX-17.jpg, awsX-18.jpg, awsX-19.jpg, awsX-20.jpg, awsX-21.jpg, awsX-22.jpg, awsX-23.jpg, awsX-24.jpg, awsX-25.jpg, awsX-26.ico

# Many open questions

- One-to-one mapping of ARC entities to/from OMERO?
- What about non-image data in ARC?
- OMERO as metadata manager on top of ARC?
- round trip possible?
- large image files (think TB) in git?

## Mapping (attempts)

## Electronic labbook

Cleaving protein with TCO-PPO-4 and tetra-butanol-transmethylidene peroxisine in 1x PBS (phosphate buffered saline) at 7.41  $\mu$ M pre-wash with 40% PFA (paraformaldehyde) (10 min fixation in 4% PFA at 4°C  
storage of the tissue in 1x PBS until the start of cleavage in 1x PBS -> 30 minutes

Permeabilization solution (2.4)

- 1.2% Goat serum,
- 0.5% Triton X-100,
- 0.1% 3-*methyl-1-butyl-*4-hydroxy phenol,
- 0.05% sodium azide in 1x PBS

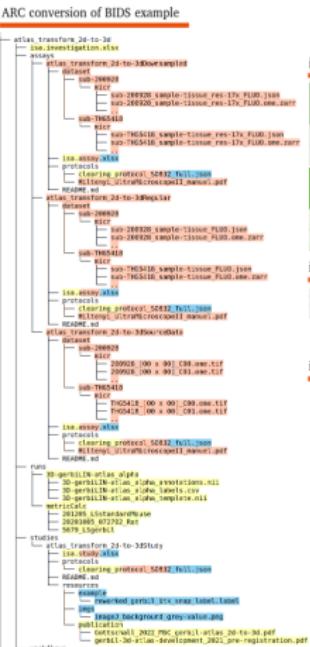
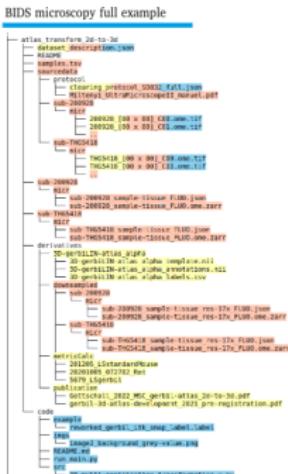
at 30°C, shake rotator 10 ml mixed rim tubes

incubation solution (5 d)  
 - (-) polymerization solution plus TD-PBO-1, 1.1008)  
 at 27°C, slow rotation, dark

- Washing in wash solution (3 x 2 h):
  - 0.5% Goat serum,
  - 0.3% Triton X 100,
  - 0.05% Sodium azide in 1% PBS

## OME-XML Metadata (from ome.tif File)

**LEGEND FIGURE 3 & 4** (blue background) dataset components necessary for the creation, transformation, or analysis of the data, such as code or parameters necessary for the analysis; (light red background) dataset data itself (as well as its descriptive metadata); (yellow background) results, figures created from the data, and components essential for the publication of the results; (purple background) legacy file content meant to be read and reorganized to multiple other files; (part of figures below grey bar) these contents are part of the dataset but not exclusive or mandatory for the IEEE standard; (part of figures below blue bar) these contents are recommended or mandatory structures of the IEEE standard; (part of figures below red bar) valid recommended or mandatory structures of the IEEE standard.



**FIGURE 4** (Left) representation of the directory structure of a valid dataset in the ABC standard manually converted from BIDS directory (compare **FIGURE 3** right); (top-right) small section of *cleaning\_protocol\_MR002\_MR002\_MR002*.json manifest transferred from electronic labbook and COMET-3MR metadata to *iss assay.xls*; (File using the PRIVATE key); (bottom) section of *BIDS samples* and *cleaning\_protocol\_MR002\_MR002\_MR002*.json file transferred to *iss assay.xls*.

Character							Term Source
Id	Parameter	Instrument	Component	Characteristic	RDF	URI	
Name	Type	Model	(Multi)unit	Physical	Physical	http://www.w3.org/2002/07/owl#ObjectProperty	
Jab-200B	Tissue	Light Microscopy	salt solution	salinity	salinity	CHEBI	
Jab-TMS418	Tissue	Light Microscopy	salt solution	permeability	permeability	CHEBI	
Jab-TMS418	Tissue	Light Microscopy	salt solution	resistance	resistance	CHEBI	

Term Accession							Characteristic	Characteristic
Characteristic								
organelle	cell	cell	cell	cell	cell	cell	percentage	percentage
organelle	cell	cell	cell	cell	cell	cell	based unit	based unit
HEB_12397	4 hour unit	percentage	percentage					
organelle	cell	cell	cell	cell	cell	cell	percentage	percentage

## sa.study.xlsx SWATe spreadsheet

sa.investigation.xlsx ARC structure metadata	
<b>ENCRYPTION SOURCE REFERENCE</b>	
Investigation	alias_investigation
Investigation Title	alias_investigation_title
Investigation Subtitle	alias_investigation_subtitle
<b>INVESTIGATION PUBLISHERS</b>	
Investigation Person (Last Name)	alias_investigation_publisher
Investigation Person (First Name)	alias_investigation_publisher
Investigation Person (Middle)	alias_investigation_publisher
Investigation Person (Title)	alias_investigation_publisher
Investigation Person (Affiliation)	alias_investigation_publisher
Investigation Person (Role)	alias_investigation_publisher
Investigation Person (Name_1st)	alias_investigation_publisher
Investigation Person (Name_2nd)	alias_investigation_publisher
<b>STUDY</b>	
Study Title	alias_study_study_title
Study Description	alias_study_study_description
Study Subtitle	alias_study_study_subtitle
Study File Name	alias_study_study_file_name
<b>STUDY ASSAYES</b>	
Study Assay	alias_study_assay
Study Assay Measurement Type	alias_study_assay_measurement_type
Study Assay Measurement Value	alias_study_assay_measurement_value
Study Assay Term	alias_study_assay_term
Study Assay Term Category	alias_study_assay_term_category
Study Assay Term Source Ref	alias_study_assay_term_source_ref
Study Assay File Name	alias_study_assay_file_name

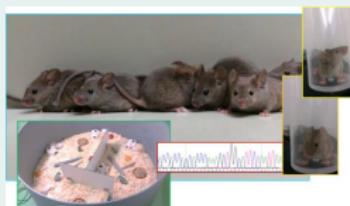
# Conclusions so far

- conversion OMERO → BIDS → ARC seems possible
- only for already specified measurements (BIDS standard)
- requires constant reformatting (e.g. fastq to csv)
- ⇔ Support by format specific processing clients

# MPI-EvolBio's activities

## Evolutionary Genetics (D. Tautz (Emeritus))

- Model organisms: *Mus domesticus*, *mus musculus*
- Behavioural genomics
- Population genetics

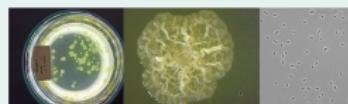


credits:

[https://www.evolbio.mpg.de/3039130/group\\_evolanimalbehpers](https://www.evolbio.mpg.de/3039130/group_evolanimalbehpers)

## Microbial Population Biology (P. Rainey)

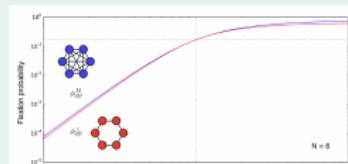
- Model organism: *Pseudomonas fluorescens*, *Bacillus sub.*
- Evolution of communities
- Host-microbe interactions
- Genetics



credits: Theodosiou (left), Schwarz (middle), Grote (right)

## Theoretical Biology (A. Traulsen)

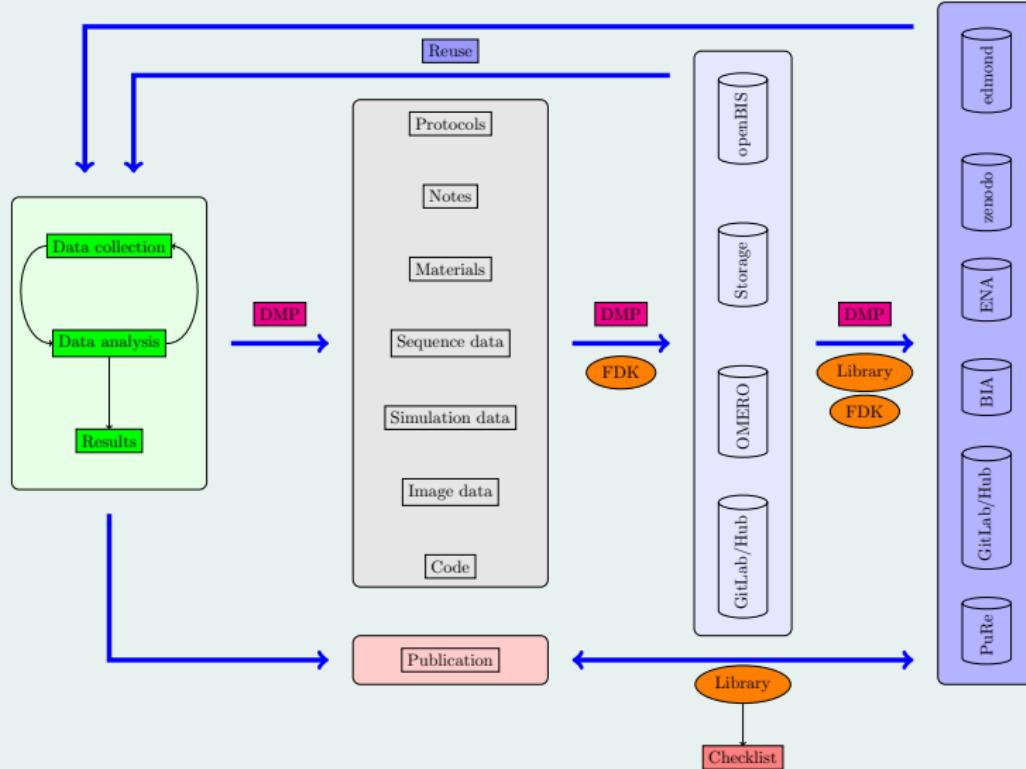
- Population Structure and Game Theory
- Metaorganisms
- Cancer Evolution



Hindersin et al, PLOS Comp. Bio (2019)

10.1371/journal.pcbi.1004437

# Data Management Policy enforces deposition of **all** research data



# Multimodal, multidimensional data in experimental evolution research

- Wildtype clone and/or genetically modified strain(s)
- Timelapse microscopy from growing microbial communities
- Time resolved whole genome / core gene NGS data
- Transcription profiles
- Optical density measurements from plaque assays
- Functional annotations

Need for integrated analysis of multiomics-multimodal data

"Find images, ELN entries for  $\Delta$  mreB strains and annotations for mreB"

# Integrating data the SPARQLing way

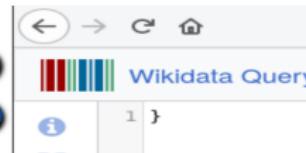
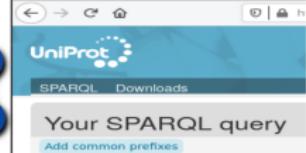
Pseudomonas fluorescens SBW25 genome database (Tripalv3)



JSON API wrapper  
RDF



Remote SPARQL endpoints



Internal DBs & web resources



API ???  
RDF



omero-rdf  
RDF



csv2rdf  
RDF



DBs with or without dedicated API



# Data integration of internal sources

DB	API	LOD ready?	RDF conversion
Tripal Genome DB	JSON-LD	yes	SPARQL SPARQL-anything virtualization <b>JSON to RDF serialization</b>
OMERO	JSON-LD	yes	omero-rdf
OpenBIS	JSON-LD	yes	??? ("semantic annotation")
StrainDB	None	no	csv dump → csv2rdf

# Flattening a deeply nested JSON-LD graph

```
[12]: Client
Client-73a06e67-d796-11ee-8cd8-e454e06b0546
Connection method: Direct
Dashboard: http://172.16.5.46:42575/status
Launch dashboard in JupyterLab
Scheduler Info

[27]: def get_features(label, serialize=False, distribute=False):
    uri = URLRef('http://pflu.evolbio.mpg.de/web-services/content/v0.1/{}'.format(label))

    logger.info("Getting %s.", uri)
    graph = get_features_graph(uri)
    logger.info("Received %d terms.", len(graph))

    so_id = query_so_id(label)
    so_uri = SO.term(so_id)

    logger.info("%d terms for %s is %s", label, so_id)

    subjects = [s for s in graph.subjects(predicate=rdfs_type,
                                           object=so_uri,
                                           unique=True)]

    logger.info("Preparing to get %d %s features.", len(subjects), label)

    if distribute:
        logger.info("Running distributed jobs...")
        for g in client.gather(client.submit(get_feature_graph, s) for s in subjects):
            graph += g
        logger.info("Distributed jobs ended.")
    else:
        for s in tqdm(subjects):
            graph += get_feature_graph(s)

    logger.info("Distributed job finished. Setting bindings and clean up.")
    set_bindings(graph)
    cleanup(graph)

    if serialize:
        fname = f"({label})_{(today)}.ttl"
        graph.serialize(fname)
        logger.info("Serialized graph to %s.", fname)

    return graph
```

- rdflib for loading json-ld into graph structure
- iteratively parses linked graph uris
- dask for distributed computing
- serialized to turtle format
- -> 2.1 million Triples in ca. 12hrs

# SparNatural for visual query generation

## Explore the *Pseudomonas fluorescens* SBW25 Knowledge Graph

SPARQL Endpoint: <http://microop046:3030/PFLUKG/sparql>

The screenshot shows the SparNatural interface for generating SPARQL queries. The query is built using a visual builder with three main components:

- Top Level:** CDS → Organism → *Pseudomonas fluorescens* +
- Middle Level:** And (CDS → locus tag → Locus Tag → PFLU\_0001) +
- Bottom Level:** And (CDS → GO annotation → GO Annotation → Any)

A large orange button at the bottom right contains a play icon, likely for executing the query. A blue button at the bottom left says "Toggle SPARQL query".

Table		Response	5 results in 0.283 seconds	Page size: 50	...
SO_0000316_1				Locus_Tag_4	Annotation_6
1	<http://pflu.evolbio.mpg.de/web-services/content/v0.1/CDS/11845>	PFLU_0001	<http://purl.obolibrary.org/obo/GO_0003688>		
2	<http://pflu.evolbio.mpg.de/web-services/content/v0.1/CDS/11845>	PFLU_0001	<http://purl.obolibrary.org/obo/GO_0005524>		
3	<http://pflu.evolbio.mpg.de/web-services/content/v0.1/CDS/11845>	PFLU_0001	<http://purl.obolibrary.org/obo/GO_0005737>		
4	<http://pflu.evolbio.mpg.de/web-services/content/v0.1/CDS/11845>	PFLU_0001	<http://purl.obolibrary.org/obo/GO_0006270>		
5	<http://pflu.evolbio.mpg.de/web-services/content/v0.1/CDS/11845>	PFLU_0001	<http://purl.obolibrary.org/obo/GO_0006275>		

Showing 1 to 5 of 5 entries



# Tasks and contributions to NFDI4BI

- P.flu SBW25 Knowledge Graph as a Use Case
- Controlled vocabularies for microbial population biology and evolutionary cell biology
- Contributions to linked open data software, clients, recommendations