

## DataHub: focusing on sample metadata

Rafael Andrade Buono, PhD FAIRDOM User Meeting, 17th October 2022







DATAHI



A (meta)data management platform

- Primary users: Researchers
- Secondary users: Support staff and service facilities
- Tertiary users: Repositories and schema creators



Flora D'Anna (Data Steward)

Rafael Andrade Buono (Data Steward)

Vahid Kiani (Software Developer)



## FAIR-by-design approach: start with the end in mind

Local archive

ENA

European Nucleotide Archive



Multi-omics data submission/brokering

**BioStudies.** BioSamples

**PRIDE** Archive

ArrayExpress

Repositories' metadata requirements

## FAIR-by-design approach: start with the end in mind



## DataHub method: FAIRDOM-SEEK and ISA model



## Building a full experimental flow in DataHub







## Building a full experimental flow in DataHub





Samples are linked through Registered Sample (multiple) Attributes

## Interacting with a Project in a single page



le Investigations Studies Assays	Investigation - hi-per characterization	n based on this one
Mutant re-sequencing Project  Project items  Project items  Project items  Project items  Project  Project Projec	Characterizing the billing mutant ling	Actions - Oreators and Submitter
	SEEK ID: http://localhost:3000/investigations/4 Projects: Mutant re-sequencing Investigation position:	Creator Creator claire s Submitter claire s
Source matchar Samples 2 Sample collection Study sample Samples 2 Samples 2 Samples 2 Samples 2		Activity Views: 48 Created: 19th Jul 2022 at 09:41
Assay - DNA extraction Assay     Samples 2     DNA extraction     Assay     DNA extraction     Assay     Samples 2     Samples 2		
Library preparation		



Copyright © 2008 - 2022 The University of Manchester and HITS gGmbH

https://datahub.test.elixir-belgium.org/



(v.1.13.0-master)

## Interacting with a experimental flow in DataHub





SEEK Samples reference their input

## Interacting with a experimental flow in DataHub





## Dynamic Table for Sample metadata





- Sample input in a table format
- Headers contain descriptions
- Displays validation errors
- Controlled vocabulary
- Ontology terms choices
- Search box for input values
- Copy/export and paste content

#### Future

- Batch upload/update of Samples
- Being able to add/remove/edit Attributes after the table is defined





## Templates: easing adoption and exchange of standards





- Sets of attributes to provide Sample level (meta)data
  - Helps balancing flexibility and compliance
- DataHub provides Templates containing metadata attributes required by EBI data repositories (e.g., ENA, ArrayExpress, MetaboLights, etc.)
- Attributes can be defined by ontology terms (or URLs)
- Users will be able to make custom Templates
- Templates have sharing permission
- (Some) Future goals:
  - One Template that covers many Repositories
  - Template versioning
  - Being able to reuse Attributes (defined by an IRI)
  - Being able to assemble new Templates by picking & choosing existing Attributes



Template based on ENA's ERC000011 Sample Checklist



## DataHub features: Templates to Sample Types





Attributes based on a Repository's checklist/standard

Attributes fitting the specific experimental step

Sample level metadata based on community standards and researcher's specific needs





## DataHub features: Templates





## Querying Samples across experimental steps



Select Project(s)					
× Environmental stimuli on yeast					
Select a Template*		Select an Attribute		Enter value	
nucleic acid sequencing		sequencing instrument -		Illumina MiSeq	
Advanced filtering -					
Input sample(s)					
Select a Template		Select an Attribute		Enter value	
source_plant	-	Source Name	-	Yeast culture 1	
Output sample(s)					
Select a Template		Select an Attribute		Enter value	
Select a Template Not selected Query sample(s) visible to you	*	Select an Attribute Not selected	•	Enter value	
Select a Template Not selected Query sample(s) visible to you now 10 v entries	*	Select an Attribute Not selected	•	Enter value	Export table
Select a Template Not selected Query sample(s) visible to you ow 10 entries put	- nucleic acid sequencing	Select an Attribute Not selected sequencing instrument	- Assay Name	Enter value	Export table Raw Data File
Select a Template          Not selected         Query         sample(s) visible to you         now       10         now       10         apput       11	• nucleic acid sequencing	Select an Attribute Not selected sequencing instrument	- Assay Name	Enter value Technology type	Export table Raw Data File
Select a Template          Not selected         Query         sample(s) visible to you         now       10         now       10         orary 1-1         orary 1-1	nucleic acid sequencing     If     RNA sequencing protocol     RNA sequencing protocol	Select an Attribute Not selected sequencing instrument Illumina MiSeq Illumina MiSeq	Assay Name     If yeast sequencing assay 1 yeast sequencing assay 1	Enter value	Export table Raw Data File
Select a Template          Not selected         Query         sample(s) visible to you         now       10         now       10         orary 1-1         orary 1-1         orary 1-2	nucleic acid sequencing       It         RNA sequencing protocol       It         RNA sequencing protocol       It         RNA sequencing protocol       It	Select an Attribute Not selected  sequencing instrument Illumina MiSeq Illumina MiSeq Illumina MiSeq Illumina MiSeq	Assay Name If yeast sequencing assay 1 yeast sequencing assay 1 yeast sequencing assay 1	Enter value	Export table Raw Data File 1-1_R1.fastq.gz 1-1_R2.fastq.gz 1-2_R1.fastq.gz
Select a Template Not selected Query sample(s) visible to you now 10 entries uput 11 orary 1-1 orary 1-2 orary 1-2	nucleic acid sequencing       I1         RNA sequencing protocol       I1         RNA sequencing protocol	Select an Attribute Not selected Sequencing instrument Illumina MiSeq Illumina MiSeq Illumina MiSeq Illumina MiSeq Illumina MiSeq	Assay Name IT yeast sequencing assay 1 yeast sequencing assay 1 yeast sequencing assay 1 yeast sequencing assay 1 yeast sequencing assay 1	Enter value	Export table Raw Data File 1-1_R1.fastq.gz 1-1_R2.fastq.gz 1-2_R1.fastq.gz 1-2_R2.fastq.gz
Select a Template Not selected Query sample(s) visible to you now 10 entries orary 1-1 orary 1-1 orary 1-2 orary 1-2 orary 1-3	nucleic acid sequencing       It         RNA sequencing protocol       It         RNA sequencing protocol	Select an Attribute Not selected Sequencing instrument Illumina MiSeq	Assay Name I1 yeast sequencing assay 1 yeast sequencing assay 1	Enter value	Export table         Raw Data File         1-1_R1.fastq.gz         1-1_R2.fastq.gz         1-2_R1.fastq.gz         1-2_R2.fastq.gz         1-3_R1.fastq.gz
Select a Template Not selected Cuery sample(s) visible to you now 10 entries put 11 orary 1-1 orary 1-2 orary 1-2 orary 1-3	nucleic acid sequencing       I1         RNA sequencing protocol       I1         RNA sequencing protocol	Select an Attribute Not selected Sequencing instrument Illumina MiSeq	Assay Name It yeast sequencing assay 1 yeast sequencing assay 1	Enter value	Export table         Raw Data File       I         1-1_R1.fastq.gz       I         1-1_R2.fastq.gz       I         1-2_R1.fastq.gz       I         1-2_R2.fastq.gz       I         1-3_R1.fastq.gz       I         1-3_R2.fastq.gz       I

"In a Project, show me the sequencing files originated from a Illumina MiSeq instrument and that come from Yeast culture 1"



## Metadata rich submissions: ISA-JSON and RO-Crate





## DataHub features: view for the future





- Improvements and streamlining the GUI
- Integration with external storage systems (e.g., iRODS)
- Integration with useGalaxy.be, WorkflowHub
- Features for data submission/brokering to ELIXIR Deposition Databases
- Packaging of Data and Metadata in RO-Crates
- Integration with external long term archival solutions
- Explore possible link with DMP tools (e.g. DMPonline)



https://datahub.elixir-belgium.org/

## Acknowledgements



**laanderen** DataHub is hosted at the Flemish Supercomputer Center



VLAAMS SUPERCOMPUTER

> DataHub development receives support from the Research Foundation – Flanders (FWO) under project No I002819N



Data brokering efforts have received support from European Union's Horizon 2020 programme under grant agreement No 871075 (ELIXIR-CONVERGE).



This work is licensed under the Creative Commons Attribution 4.0 International License, except where otherwise noted. To view a copy of this license, visit http://creativecommons.org/licenses/by/4.0/.



# Get in touch with us.



