

Reproducibility; Research Objects (RO-Crate) and Common Workflow Language (CWL)



Stian Soiland-Reyes

eScience lab, The University of Manchester

INDElab, University of Amsterdam

@soilandreyes

<https://orcid.org/0000-0001-9842-9718>

<https://slides.com/soilandreyes/>

@soilandreyes

H2020-INFRAEOSC-2018-2 [824087](#)

H2020-INFRAEDI-2018-1 [823830](#)

H2020-INFRAIA-2017-1 [730976](#)

H2020-INFRADEV-2019-2 [871118](#)

H2020-INFRAIA-2018-1 [823827](#)



This work is licensed under a

[Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

WoSSS21: Workshop on
Sustainable Software Sustainability
2021-10-07



Rein in the four horsemen of irreproducibility

Dorothy Bishop *describes how threats to reproducibility, recognized but unaddressed for decades, might finally be brought under control.*

More than four decades into my scientific career, I find myself an outlier among academics of similar age and seniority: I strongly identify with the movement to make the practice of science more robust. It's not that my contemporaries are unconcerned about doing science well; it's just that many of them don't seem to recognize that there are serious problems with current practices. By contrast, I think that, in two decades, we will look back on the past 60 years — particularly in biomedical science — and marvel at how much time and money has been wasted on flawed research.

How can that be? We know how to formulate and test hypotheses in controlled experiments. We can account for unwanted variation with statistical techniques. We appreciate the need to replicate observations.

Yet many researchers persist in working in a way almost guaran-

**MANY RESEARCHERS
PERSIST IN WORKING
IN A WAY ALMOST
GUARANTEED
NOT
TO DELIVER
MEANINGFUL
RESULTS.**

“They ride with what I refer to as the four horsemen of the reproducibility apocalypse:

1. Publication bias
2. Low statistical power
3. P-value hacking
4. HARKing (hypothesizing after results are known)

Assignment of MAGs to reference databases

Four reference databases were used to classify the set of MAGs recovered from the human gut assemblies: HR, RefSeq, GenBank and a collection of MAGs from public datasets. HR comprised a total of 2,468 high-quality genomes (>90% completeness, <5% contamination) retrieved from both the HMP catalogue (<https://www.hmpdacc.org/catalog/>) and the HGG⁸. From the RefSeq database, we used all the complete bacterial genomes available ($n = 8,778$) as of January 2018. In the case of GenBank, a total of 153,359 bacterial and 4,053 eukaryotic genomes (3,456 fungal and 597 protozoan genomes) deposited as of August 2018 were considered. Lastly, we surveyed 18,227 MAGs from the largest datasets publicly available as of August 2018^{13,16,17,18,19}, including those deposited in the Integrated Microbial Genomes and Microbiomes (IMG/M) database⁵². For each database, the function 'mash sketch' from Mash v.2.0⁵³ was used to convert the reference genomes into a MinHash sketch with default k -mer and sketch sizes. Then, the Mash distance between each MAG and the set of references was calculated with 'mash dist' to find the best match (that is, the reference genome with the lowest Mash distance). Subsequently, each MAG and its closest relative were aligned with dnadiff v.1.3 from MUMmer 3.23⁵⁴ to compare each pair of genomes with regard to the fraction of the MAG aligned (aligned query, AQ) and ANI.

Technical reproducibility and cluster quality

A random subset of 1,000 metagenomes (Supplementary Table 1) was tested with two additional approaches to assess the reproducibility of the MAGs generated here. With one of the methods, metagenomes were assembled with MEGAHIT v.1.1.3²⁴ and subsequently binned with MetaBAT 2, MetaBAT 1 and MaxBin v.2.2.4⁶². A refinement step was then performed using the bin_refinement module from MetaWRAP v.1.0²⁵ to combine and improve the results generated by the three binners. The second method involved a modified co-assembly approach, in which individual assemblies from the same study were first merged and dereplicated with CD-HIT v.4.7⁶³ (cd-hit-est with option -c 0.99 defining a sequence identity threshold of 99%). Metagenomic datasets were then mapped to their merged, non-redundant assembly with BWA-MEM to obtain co-abundance information for binning with MetaBAT 2 (with option --minContig 2000). The resulting MAGs with a QS >50 obtained with each method were compared to the MAGs recovered with our main pipeline (individual assembly with SPAdes, plus binning with MetaBAT 2) for the same 1,000 datasets, using the combined Mash and MUMmer workflow described above.

<https://doi.org/10.1038/s41586-019-0965-1>

Reproducibility?

Read classification of the 13,133 human gut metagenomic datasets was performed with sourmash v.2.0.0a4⁶⁸ against the HR, RefSeq and UMGS genome collections. Signature files were generated for both the reference (FASTA) and query (FASTQ) files, with 'sourmash compute --scaled 1000 -k 31 --track-abundance'. For each set of references, a lowest common ancestor database was created ('sourmash lca index --scaled 1000 -k 31'), with each genome representing a unique species lineage. Raw reads were then compared with 'sourmash lca gather' against each database. Species prevalence and abundance was determined with BWA-MEM, where species presence was inferred by assessing the level of genome coverage, mean read depth and depth evenness. First, we calculated depth and variation penalty scores corresponding to the missing coverage ($100\% - \text{genome coverage}$) multiplied by either the $\log(\text{mean depth})$ or the depth coefficient of variation (defined as the standard deviation of read depth divided by the mean), respectively. These metrics allowed us to gauge both coverage and depth simultaneously, as genomes that have a high mean depth (or high depth variation) but are not well covered are less likely to be present in the sample than those that have the same level of coverage with lower read depth. Thresholds for determining genome presence were set at a minimum coverage of at least 60%, and both depth and variation penalty scores at a maximum of the 99th percentile







De novo assembly and binning

Raw reads from each run were first assembled with SPAdes v.3.10.0²⁰ with option --meta²¹. Thereafter, MetaBAT 2¹⁵ (v.2.12.1) was used to bin the assemblies using a minimum contig length threshold of 2,000 bp (option --minContig 2000) and default parameters. Depth of coverage required for the binning was inferred by mapping the raw reads back to their assemblies with BWA-MEM v.0.7.16⁴⁵ and then calculating the corresponding read depths of each individual contig with samtools v.1.5⁴⁶ ('samtools view -Sbu' followed by 'samtools sort') together with the jgi_summarize_bam_contig_depths function from MetaBAT 2. The QS of each metagenome-assembled genome (MAG) was estimated with CheckM v.1.0.7²² using the lineage_wf workflow and calculated as: $\text{level of completeness} - 5 \times \text{contamination}$. Ribosomal RNAs (rRNAs) were detected with the cmsearch function from INFERNAL v.1.1.2⁴⁷ (options -Z 1000 --hmmonly --cut_ga) using the Rfam⁴⁸ covariance models of the bacterial 5S, 16S and 23S rRNAs. Total alignment length was inferred by the sum of all non-overlapping hits. Each gene was considered present if more than 80% of the expected sequence length was contained in the MAG. Transfer RNAs (tRNAs) were identified with tRNAscan-s.e. v.2.0⁴⁹ using the bacterial tRNA model (option -B) and default parameters. Classification into high- and medium-quality

master 2 branches 0 tags

Go to file

Code

 alexmsalmeida Update LICENSE	649913e on 12 Nov	🕒 122 commits
 R	Update funcs_phy-assoc_fig5b.R	13 months ago
 pipelines	Update map2ref.sh	2 years ago
 scripts	Update parse_checkm.py	2 years ago
 LICENSE	Update LICENSE	last month
 README.md	Update README.md	13 months ago

README.md

Analysis of Metagenomic Species (MGS)

Scripts used for characterizing metagenome-assembled genomes (MAGs) used in the following publication:

A Almeida, AL Mitchell, M Boland, SC Forster, GB Gloor, A Tarkowska, TD Lawley and RD Finn (2019) [A new genomic blueprint of the human gut microbiota](#). *Nature* **568**, 499–504

Associated data can also be found in our [FTP server](#).

mashdiff.sh

Compare a set of genomes against a reference database. Selects best representative for whole-genome alignment.

Requirements:

- Mash (tested v2.0)
- MUMmer (tested v3.23)
- Python 2.7

Coded for running within LSF cluster environments.

Usage:

```
mashdiff.sh -i genome_folder/ -r reference.msh -s db_name -p output_prefix
```

About

Analysing Metagenomic Species (MGS)

[Readme](#)[View license](#)

Releases

No releases published

Packages

No packages published

Languages

R 70.7% Python 18.5% Shell 10.8%

Automation

- Automate computational aspects
- Repetitive pipelines, sweep campaigns

Scaling—compute cycles

- Make use of computational infrastructure
- Handle large data

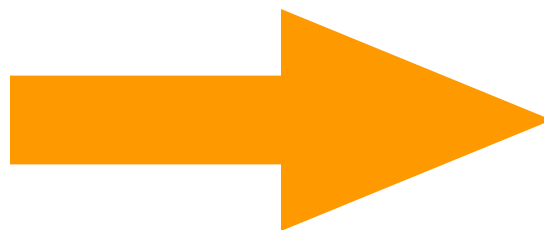
Abstraction—people cycles

- Shield complexity and incompatibilities
- Report, re-use, evolve, share, compare
- Repeat—Tweak—Repeat
- First-class commodities

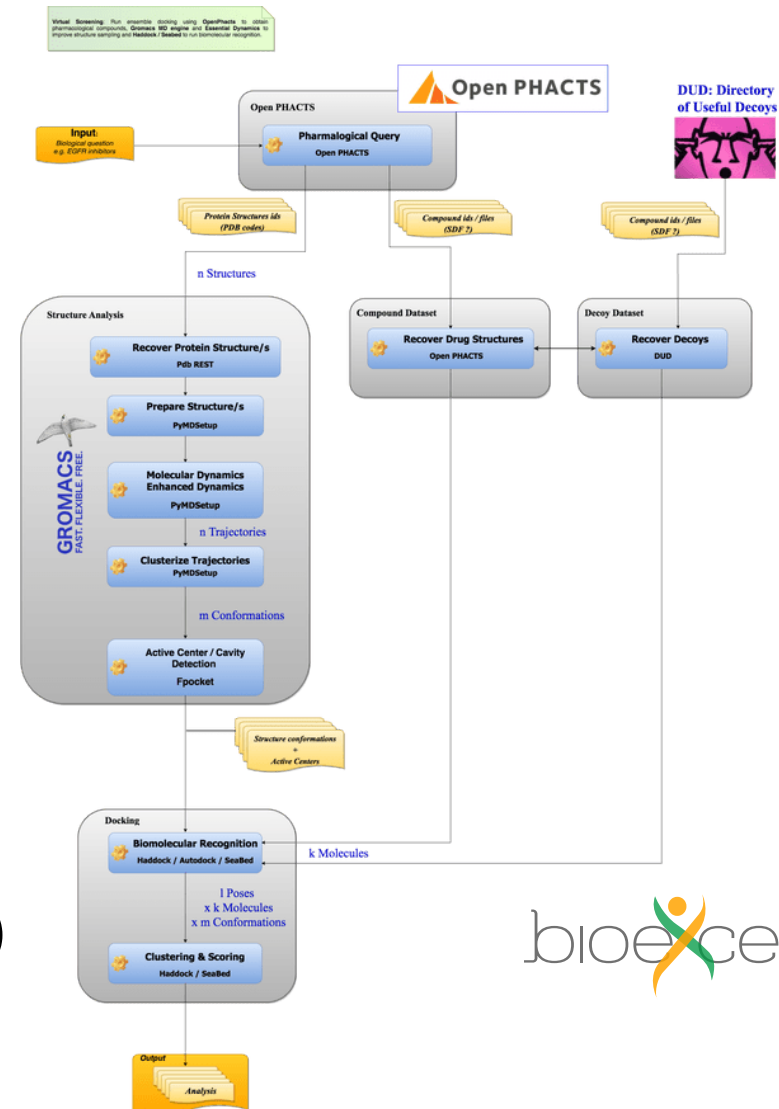
Provenance—reporting

- Capture, report and utilize log and data lineage
- Auto-documentation
- Tracable evolution, audit, transparency
- Reproducible science

Why use workflows?



Findable
Accessible
Interoperable
Reusable
(Reproducible)



Adapted from Bertram Ludäscher (2015)

<https://www.slideshare.net/ludaesch/works-2015provenancemileage>

@soilandreyes

<https://doi.org/10.1007/s13222-012-0100-z>

Existing Workflow systems

Ryan edited this page 3 days ago · 299 revisions

Permalink: <https://s.apache.org/existing-workflow-systems>

Cite as (update dates):

Peter Amstutz, Maxim Mikheev, Michael R. Crusoe, Nebojša Tijanić, Samuel Lampa, et al. (2021): **Existing Workflow systems**.

Common Workflow Language wiki, GitHub. <https://s.apache.org/existing-workflow-systems> updated 2021-06-29, accessed 2021-09-03.

Computational Data Analysis Workflow Systems

An incomplete list

Please add new entries at the bottom.

See also: <https://github.com/pditommaso/awesome-pipeline>

1. Arvados - CWL-based distributed computing platform for data analysis on massive data se <https://github.com/arvados/arvados>
2. Apache Taverna <http://www.taverna.org.uk/> <https://taverna.incubator.apache.org/>
3. Galaxy <http://galaxyproject.org/>
4. SHIWA <https://www.shiwa-workflow.eu/>
5. Apache Oozie <https://oozie.apache.org/>
6. DNANexus <https://wiki.dnanexus.com/API-Specification-v1.0.0/IO-and-Run-Specifications-Specification-v1.0.0/Workflows-and-Analyses>
7. ~~BioDT~~ <http://www.biodatomics.com/> (dead link)
8. Agave <http://agaveapi.co/live-docs/>
9. DiscoveryEnvironment <http://www.iplantcollaborative.org/ci/discovery-environment>
10. Wings <http://www.wings-workflows.org/>
11. Knime <https://www.knime.org/>

291. BRANE Framework <https://onnovalkering.github.io/brane/>

292. ApolloWF <https://apollowf.github.io/>

293. IS-EPOS Platform <https://ieeexplore.ieee.org/document/9308147> <https://tcs.>

294. pyinvoke <http://www.pyinvoke.org/>

295. targets R package <https://cran.r-project.org/package=targets>

296. Compi <https://doi.org/10.7717/peerj-cs.593> <https://github.com/sing-group/>

297. TriggerFlow: Event-based Orchestration of Serverless Workflows (<https://gith>

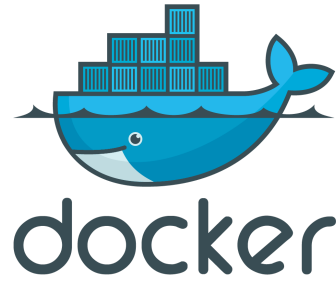
298. Google Cloud Workflows: Orchestrate and automate Google Cloud and HTTP workflows. (<https://cloud.google.com/workflows/docs>)

299. PanDA Workflow Management System: <https://doi.org/10.1051/epjconf/201610801003> <https://github.com/PanDAWMS>

300. Harvester: <https://github.com/HSF/harvester>

<https://s.apache.org/existing-workflow-systems>

COMMON WORKFLOW LANGUAGE



Implementations

Software	Description
cwltool	Reference implementation of CWL
Arvados	Distributed computing platform for data analysis on massive data sets. Using CWL on Arvados
Toil	https://www.commonwl.org/Implementations Toil is a workflow engine entirely written in Python
Apache Taverna	Domain-independent Workflow Management System
Galaxy	Web-based platform for data intensive biomedical research.
AWE	Workflow and resource management system for bioinformatics data analysis.
Funnel	Use Google Genomics Pipeline API with CWL
Rabix Bunny	Reproducible Analyses for Bioinformatics

```
cwlVersion: v1.0
class: Workflow
inputs:
  inp: File
  ex: string

outputs:
  classout:
    type: File
    outputSource: compile/classfile

steps:
  untar:
    run: tar-param.cwl
    in:
      tarfile: inp
      extractfile: ex
    out: [example_out]

  compile:
    run: arguments.cwl
    in:
      src: untar/example_out
    out: [classfile]
```

<http://www.commonwl.org/>

TOOLBOX • 02 SEPTEMBER 2019

Workflow systems turn raw data into scientific knowledge

How workflow tools can make your computational methods portable, maintainable, reproducible and shareable.

BY JEFFREY M. PERKEL

Reinventing the wheel is pointless, but for computational biologists it's sometimes unavoidable. So when Rob Finn and Folker Meyer realized how much their work overlapped, they decided to try something different.

Finn is head of the sequence-families team at the European Bioinformatics Institute (EBI) in Hinxton, UK; Meyer is a computer scientist at Argonne National Laboratory in Lemont, Illinois. Both run facilities that let researchers perform a computationally intensive process called metagenomic analysis, which allows microbial communities to be reconstructed

from shards of DNA. It would be helpful, they realized, if they could try each other's code. The problem was that their analytical 'pipelines' — the carefully choreographed computational steps required to turn raw data into scientific knowledge — were written in different languages. Meyer's team was using an in-house system called AWE, whereas Finn was working with nearly 9,500 lines of Python code.

"It was a horrible Python code base," says Finn — complicated, and difficult to maintain. "Bits had been bolted on in an ad hoc fashion over seven years by at least four different developers." And it was "heavily tied to the compute infrastructure", he says, meaning it was written for specific computational resources and

a particular way of organizing files, essentially unusable outside the EBI. The EBI wasn't using AWE, the reverse was true. Then Finn and Meyer learned about Common Workflow Language (CWL).

CWL is a way of describing analyses, lines and computational tools — more than 250 systems now available, including popular options as Snakemake, Nextflow and Galaxy. Although they speak different languages and support different features, they all have a common aim: to make computational methods reproducible, portable, maintainable and shareable. CWL is essentially a common language that researchers can use to share pipelines for whichever system. For Finn, that

► language brought sanity to his codebase, reducing it by around 73%. Importantly, it has made it easier to test, execute and share new methods, and to run them on the cloud.

There is a learning curve to adopting workflow languages. But, says Brian Naughton, data lead and co-founder of the drug-discovery firm Hexagon Bio in Menlo Park, California, "the energy that you expend learning is more than made up for by the energy you save in having your code be reproducible."

Nature **573**, 149-150 (2019)

@soilandreyes

<https://doi.org/10.1038/d41586-019-02619-z>

<https://github.com/NCBI->

Hackathons/Epigenomics_CWL/blob/master/cwl/tools/preprocess.cwl

```
1 cwlVersion: v1.0
2 class: Workflow
3
4 inputs:
5   toConvert: File
6
7
8 outputs:
9   converted:
10    type: File
11    outputSource: convertMethylation/converted
12   combined:
13    type: File
14    outputSource: mergeSymmetric/combined
15
16
17 steps:
18   convertMethylation:
19     run: interconverter.cwl
20     in:
21       toConvert: toConvert
22     out: [converted]
23   mergeSymmetric:
24     run: symmetriccpgs.cwl
25     in:
26       toCombine: convertMethylation/converted
27     out: [combined]
```

```
cwlVersion: v1.0
class: CommandLineTool
inputs:
  toConvert:
    type: File
    inputBinding:
      prefix: -i
outputs:
  converted:
    type: File
    outputBinding:
      glob: "*.meth"
baseCommand: interconv
arguments: ["-d", $(runtime.outdir)]
hints:
  - class: DockerRequirement
    dockerPull: "quay.io/ncbi/interconverter"
```

```
1 cwlVersion: v1.0
2 class: CommandLineTool
3 inputs:
4   toCombine:
5     type: File
6     inputBinding:
7       prefix: -i
8 outputs:
9   combined:
10    type: File
11    outputBinding:
12      glob: "*.sym"
13
14 baseCommand: symmetriccpgs.sh
15 arguments: ["-d", $(runtime.outdir)]
16
17 hints:
18   - class: DockerRequirement
19     dockerPull: "quay.io/ncbi/symmetriccpgs"
```

@soilandreyes

<https://github.com/NCBI->

Hackathons/Epigenomics_CWL/blob/master/cwl/tools/preprocess.cwl

```
1 cwlVersion: v1.0
2 class: Workflow
3
4 inputs:
5   toConvert: File
6
7
8 outputs:
9   converted:
10    type: File
11    outputSource: convertMethylation/converted
12   combined:
13    type: File
14    outputSource: mergeSymmetric/combined
15
16
17 steps:
18   convertMethylation:
19     run: interconverter.cwl
20     in:
21       toConvert: toConvert
22     out: [converted]
23   mergeSymmetric:
24     run: symmetriccpgs.cwl
25     in:
26       toCombine: convertMethylation/converted
27     out: [combined]
```

```
cwlVersion: v1.0
class: CommandLineTool
inputs:
  toConvert:
    type: File
    inputBinding:
      prefix: -i
outputs:
  converted:
    type: File
    outputBinding:
      glob: "*.meth"
baseCommand: interconv
arguments: ["-d", $(runtime.outdir)]
hints:
  - class: DockerRequirement
    dockerPull: "quay.io/ncbi/interconverter"
```

```
1 cwlVersion: v1.0
2 class: CommandLineTool
3 inputs:
4   toCombine:
5     type: File
6     inputBinding:
7       prefix: -i
8 outputs:
9   combined:
10    type: File
11    outputBinding:
12      glob: "*.sym"
13
14 baseCommand: symmetriccpgs.sh
15 arguments: ["-d", $(runtime.outdir)]
16
17 hints:
18   - class: DockerRequirement
19     dockerPull: "quay.io/ncbi/symmetriccpgs"
```

@soilandreyes

<https://github.com/NCBI->

Hackathons/Epigenomics_CWL/blob/master/cwl/tools/preprocess.cwl

```
1 cwlVersion: v1.0
2 class: Workflow
3
4 inputs:
5   toConvert: File
6
7
8 outputs:
9   converted:
10    type: File
11    outputSource: convertMethylation/converted
12   combined:
13    type: File
14    outputSource: mergeSymmetric/combined
15
16
17 steps:
18   convertMethylation:
19     run: interconverter.cwl
20     in:
21       toConvert: toConvert
22     out: [converted]
23   mergeSymmetric:
24     run: symmetriccpgs.cwl
25     in:
26       toCombine: convertMethylation/converted
27     out: [combined]
```

```
cwlVersion: v1.0
class: CommandLineTool
inputs:
  toConvert:
    type: File
    inputBinding:
      prefix: -i
outputs:
  converted:
    type: File
    outputBinding:
      glob: "*.meth"
baseCommand: interconv
arguments: ["-d", $(runtime.outdir)]
hints:
  - class: DockerRequirement
    dockerPull: "quay.io/ncbi/interconverter"
```

```
1 cwlVersion: v1.0
2 class: CommandLineTool
3 inputs:
4   toCombine:
5     type: File
6     inputBinding:
7       prefix: -i
8 outputs:
9   combined:
10    type: File
11    outputBinding:
12      glob: "*.sym"
13
14 baseCommand: symmetriccpgs.sh
15 arguments: ["-d", $(runtime.outdir)]
16
17 hints:
18   - class: DockerRequirement
19     dockerPull: "quay.io/ncbi/symmetriccpgs"
```

@soilandreyes

<https://github.com/NCBI->

Hackathons/Epigenomics_CWL/blob/master/cwl/tools/preprocess.cwl

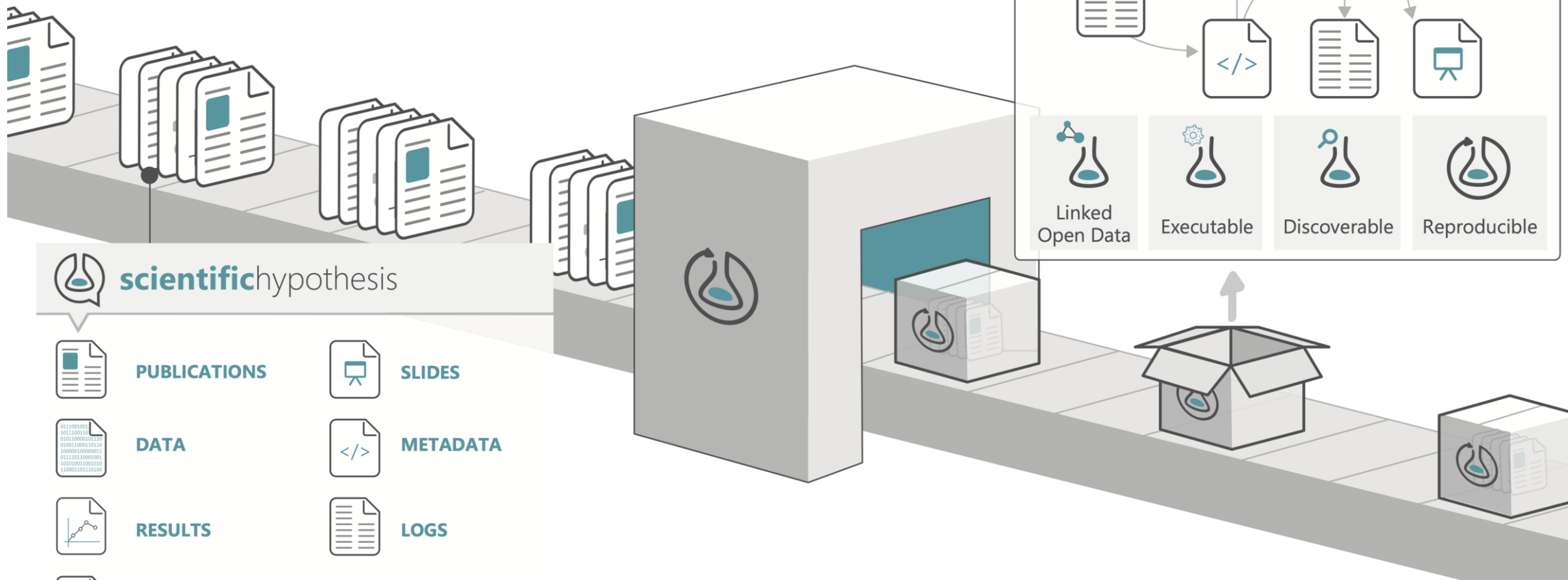
```
1 cwlVersion: v1.0
2 class: Workflow
3
4 inputs:
5   toConvert: File
6
7
8 outputs:
9   converted:
10    type: File
11    outputSource: convertMethylation/converted
12   combined:
13    type: File
14    outputSource: mergeSymmetric/combined
15
16
17 steps:
18   convertMethylation:
19     run: interconverter.cwl
20     in:
21       toConvert: toConvert
22     out: [converted]
23   mergeSymmetric:
24     run: symmetriccpgs.cwl
25     in:
26       toCombine: convertMethylation/converted
27     out: [combined]
```

```
cwlVersion: v1.0
class: CommandLineTool
inputs:
  toConvert:
    type: File
    inputBinding:
      prefix: -i
outputs:
  converted:
    type: File
    outputBinding:
      glob: "*.meth"
baseCommand: interconv
arguments: ["-d", $(runtime.outdir)]
hints:
  - class: DockerRequirement
    dockerPull: "quay.io/ncbi/interconverter"
```

```
1 cwlVersion: v1.0
2 class: CommandLineTool
3 inputs:
4   toCombine:
5     type: File
6     inputBinding:
7       prefix: -i
8 outputs:
9   combined:
10    type: File
11    outputBinding:
12      glob: "*.sym"
13
14 baseCommand: symmetriccpgs.sh
15 arguments: ["-d", $(runtime.outdir)]
16
17 hints:
18   - class: DockerRequirement
19     dockerPull: "quay.io/ncbi/symmetriccpgs"
```

@soilandreyes

 Enabling **reproducible**, transparent research.




researchobject.org

What is RO-Crate?

Addressable
resources



Local Data

 <https://orcid.org/0000-0001-2345-6789>



ID? Title? Description?

  Who created this data?

 What parts does it have?

 When?

 What is it about?

 How can it be reused?

 As part of which project?

 Who funded it?

 How was it made?

 https://en.wikipedia.org/wiki/Scanning_electron_microscope



RO-Crate is method for self-decribed *datasets*
as a *digital object* using a single Linked Data
metadata document

Credit: Peter Sefton

Adapted from <https://arkisto-platform.github.io/standards/ro-crate/>



The dataset may contain *any kind* of resource, about anything, in *any format* as a file, URL or PID

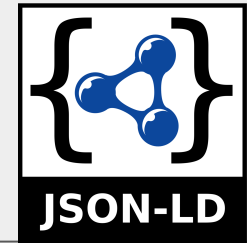
Credit: Peter Sefton

Adapted from <https://arkisto-platform.github.io/standards/ro-crate/>


```
|-- Folder1/  
|         |-- file1.this  
|         |-- file2.that  
|-- Folder2/  
|         -- file1.this  
|         |-- file2.that  
|-2021-04-08 07.58.17.jpg
```



```
{  
  "@id": "2021-04-08 07.58.17.jpg",  
  "@type": "File",  
  "contentSize": 3271409,  
  "dateModified": "2021-04-08T07:58:17+10:00",  
  "description": "",  
  "encodingFormat": [  
    {"@id":  
      "https://www.nationalarchives.gov.uk/PRONOM/x-fmt/391"},  
    "image/jpeg"  
  ],  
  "name": "Cute puppy"  
},
```



Each resource have a **machine readable**
description in JSON-LD format

```
|-- Folder1/
|           |-- file1.this
|           |-- file2.that
|-- Folder2/
|           |-- file1.this
|           |-- file2.that
|-2021-04-08 07.58.17.jpg
```

A screenshot of a web browser window. The address bar shows the file path "/Users/pt/Do...". Below the address bar, there is a button that says "Download all the metadata for in JSON-LD format". Below that, there is a "Download: Cute puppy" button with a downward arrow icon. To the right of this button is a preview image of a black and white puppy looking up. Below the image is a table of metadata:

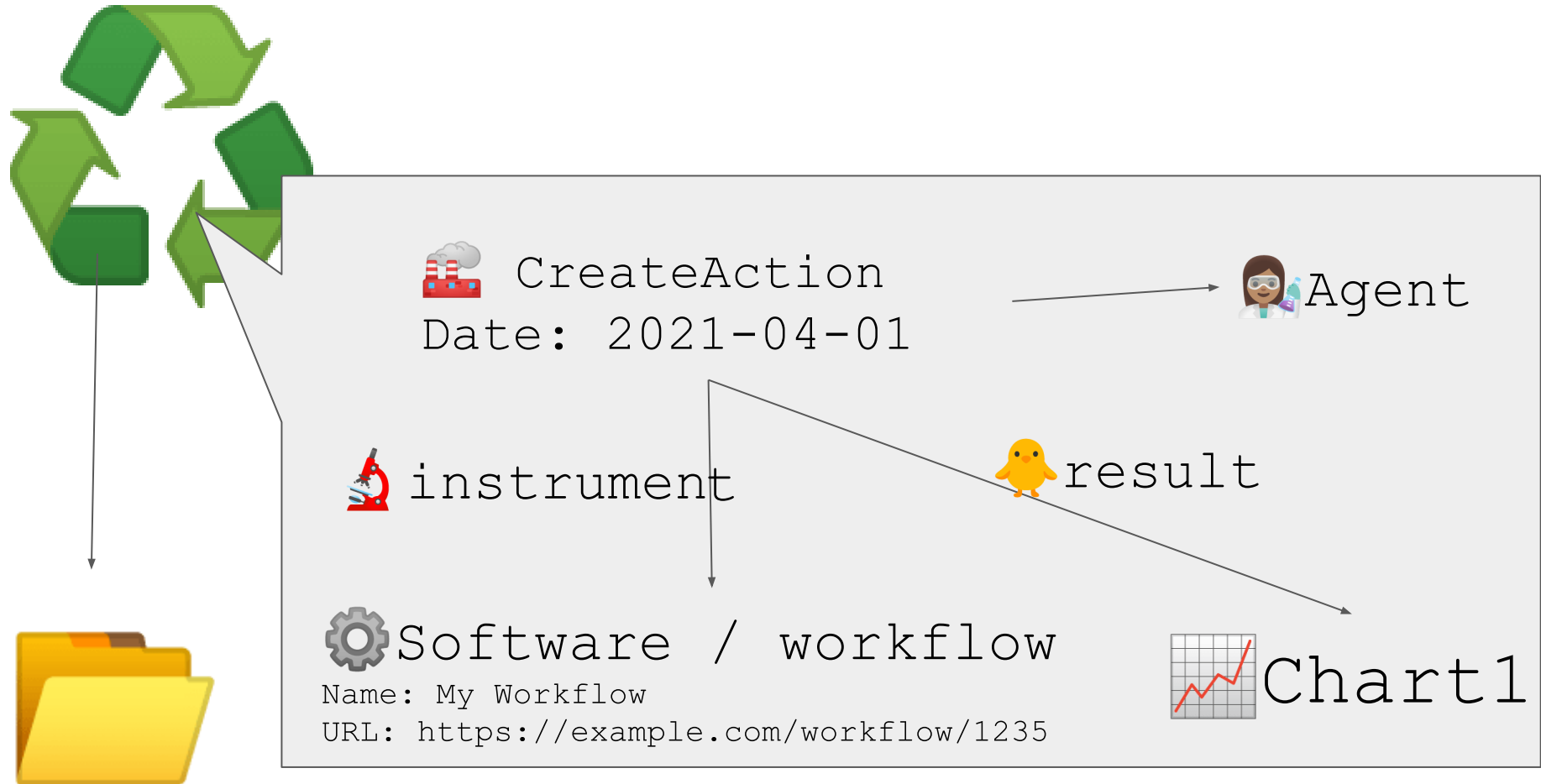
@id	2021-04-08 07.58.17.jpg
name [?]	Cute puppy
@type	File

At the bottom of the browser window, there is a user profile icon and the name "webinar-886835...ics" with an upward arrow, and a "Show all" button with a close icon.

A **human-readable** description/preview in an HTML file that *lives alongside the metadata*

Credit: Peter Sefton

Adapted from <https://arkisto-platform.github.io/standards/ro-crate/>



Provenance and workflow information can be included
– to assist in *re-use* of data and research processes

Credit: Peter Sefton

Adapted from <https://arkisto-platform.github.io/standards/ro-crate/>



RO-Crate Digital Objects may be *packaged* for distribution eg via Zip, Bagit and OCFL
– or simply be *published on the Web*

Credit: Peter Sefton

Adapted from <https://arkisto-platform.github.io/standards/ro-crate/>

Techie deep-dive!

Warning: JSON ahead

RO-Crate Structure

The structure an *RO-Crate* MUST follow is:

```
<RO-Crate root directory>/
|  ro-crate-metadata.json    # RO-Crate Metadata
|  ro-crate-preview.html     # RO-Crate Website
|  ro-crate-preview_files/   # MAY be present
|  | [other RO-Crate Website files]
|  [payload files and directories] # 0 or more
```

Minimal example of RO-Crate

The following *RO-Crate Metadata File* represents a minimal description of an RO-Crate **8. Contextual Entities**

```
{ "@context": "https://w3id.org/ro/crate/1.1/context",
  "@graph": [
    {
      "@type": "CreativeWork",
      "@id": "ro-crate-metadata.json",
      "conformsTo": {"@id": "https://w3id.org/ro/crate/1.1"},
      "about": {"@id": "."}
    },
    {
      "@id": "./",
      "identifier": "https://doi.org/10.4225/59/59672c09f4a4b",
      "@type": "Dataset",
      "datePublished": "2017",
      "name": "Data files associated with the manuscript:Effects of facili",
      "description": "Palliative care planning for nursing home residents",
      "license": {"@id": "https://creativecommons.org/licenses/by-nc-sa/3.0/au/"}
    },
    {
      "@id": "https://creativecommons.org/licenses/by-nc-sa/3.0/au/",
      "@type": "CreativeWork",
      "description": "This work is licensed under the Creative Commons Attr",
      "identifier": "https://creativecommons.org/licenses/by-nc-sa/3.0/au/",
      "name": "Attribution-NonCommercial-ShareAlike 3.0 Australia (CC BY-NC"
    }
  ]
}
```



RO-Crate

<https://w3id.org/ro/crate/1.1>

RO-CRATE 1.1

1. About this document
2. Introduction
3. Terminology
4. RO-Crate Structure
5. Metadata of the RO-Crate
6. Root Data Entity
7. Data Entities

8. Contextual Entities

- Contextual vs Data entities
- Identifiers for contextual entities
- People
- Organizations as values
- Contact information
- Publications via citation property
- Publisher
- Funding and grants
- Licensing, Access control and copyright
- Extra metadata such as Exif
- Places
- Subjects & keywords
- Time
- Thumbnails

9. Provenance of entities
10. Workflows and scripts

Appendix

JSON for Linking Data

Data is messy and disconnected. JSON-LD organizes and connects it, creating a better Web.

<https://json-ld.org/>

⇔ Linked Data

Linked Data empowers people that publish and use information on the Web. It is a way to create a network of standards-based, machine-readable data across Web sites. It allows an application to start at one piece of Linked Data, and follow embedded links to other pieces of Linked Data that are hosted on different sites across the Web.

A Simple Example

```
{
  "@context": "http://json-ld.org/contexts/person.jsonld",
  "@id": "http://dbpedia.org/resource/John_Lennon",
  "name": "John Lennon",
  "born": "1940-10-09",
  "spouse": "http://dbpedia.org/resource/Cynthia_Lennon"
}
```

{ } JSON-LD

JSON-LD is a lightweight Linked Data format. It is easy for humans to read and write. It is based on the already successful JSON format and provides a way to help JSON data interoperate at Web-scale. JSON-LD is an ideal data format for programming environments, REST Web services, and unstructured databases such as CouchDB and MongoDB.

<https://www.w3.org/TR/json-ld/>

<https://schema.org/>

Dataset

[Thing](#) > [CreativeWork](#) > [Dataset](#)

A body of structured information describing some topic(s) of interest.

Usage: Between 100 and 1000 domains



Property	Expected Type	Description
Properties from Dataset		
distribution	DataDownload	A downloadable form of this dataset, at a specific location, in a specific format.
includedInDataCatalog	DataCatalog	A data catalog which contains this dataset. Supersedes catalog , includedDataCatalog . Inverse property: dataset .
Properties from CreativeWork		
about	Thing	The subject matter of the content.
accessibilityAPI	Text	Indicates that the resource is compatible with the referenced accessibility API (WebSchemas).
accessibilityControl	Text	Identifies input methods that are sufficient to fully control the described resource (WebSchemas values).
accessibilityFeature	Text	Content features of the resource, such as accessible media, alternatives and supported enhancements (WebSchemas wiki lists possible values).

Metadata

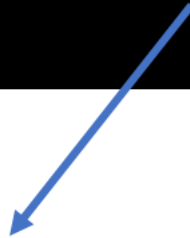


```
{
  "@id": "figure.png",
  "@type": ["File", "ImageObject"],
  "name": "XXL-CT-scan of an XXL Tyrannosaurus rex skull",
  "identifier": "https://doi.org/10.5281/zenodo.3479743",
  "author": {"@id": "https://orcid.org/0000-0002-8367-6908"},
  "encodingFormat": "image/png"
}
```



```
{
  "@id": "https://orcid.org/0000-0002-8367-6908",
  "@type": "Person",
  "affiliation": { "@id": "https://ror.org/03f0f6041" },
  "name": "J. Xuan"
}
```

ORCID



```
{
  "@id": "https://ror.org/03f0f6041",
  "@type": "Organization",
  "name": "University of Technology Sydney",
  "url": "https://www.uts.edu.au/"
}
```

ROR

Data and **Contextual** entities
described *within* RO-Crate Metadata File

Base vocabulary & types: **schema.org**

Cross-references to further contextual entities

RO-Crate **principle**:

Reuse existing PIDs and URLs

..but always **describe entities** which lack a
human-readable resolution

Tooling!

How can I use it?

While we're mostly focusing on the specification, some tools already exist for working with RO-Crates:

- [Describo](#) interactive **desktop application** to create, update and export RO-Crates for different profiles. (~ *beta*)
- [CalcyteJS](#) is a **command-line** tool to help create RO-Crates and HTML-readable rendering (~ *beta*)
- [ro-crate](#) - **JavaScript/NodeJS** library for RO-Crate rendering as HTML. (~ *beta*)
- [ro-crate-js](#) - utility to render HTML from RO-Crate (~ *alpha*)
- [ro-crate-ruby](#) **Ruby** library to consume/produce RO-Crates (~ *alpha*)
- [ro-crate-py](#) **Python library** to consume/produce RO-Crates (~ *planning*)

These applications use or expose RO-Crates:

- [Workflow Hub](#) imports and exports [Workflow RO-Crates](#)
- [OCFL-indexer](#) NodeJS application that walks the [Oxford Common File Layout](#) on the file system, validate RO-Crate Metadata Files and parse into objects registered in Elasticsearch. (~ *alpha*)
- [ONI indexer](#)
- [ocfl-tools](#)
- [ocfl-viewer](#)
- [Research Object Composer](#) is a REST API for gradually building and depositing Research Objects according to a pre-defined profile. (RO-Crate support *alpha*)



Edit: Dataset - data

name

data

description

+ TextArea

<https://uts-eresearch.github.io/describo/>

license

+

CreativeWork

hasPart

+

Dataset

+ File

Dataset

another level

File

DT1-214-A.mp3

File

NT1-20003-002.jpg

File

NT5-TokelauOf-CAT-PDSC_ADMIN.xml

File

NT5-TokelauOf-vid.mp4

author

+

Organization

+ Person

publisher

+

Organization

+ Person

funder

+

Organization

+ Person

Show all available properties

save

This item is connected to:

@type: RootDataset → hasPart

name: my crate



UNIVERSITY OF TECHNOLOGY, SYDNEY

RO-Crate for data scientists

Selected Resource: onedrive:/Sefton sample crate

Build the Collection

Manage Collection Data Files

Browse Collection Entities

Load Root Dataset

</> Add Property

name

Sample dataset for RO-Crate v1.0

citation



+ CreativeWork

+ Text



ScholarlyArticle: This is an example publication with a dodgy DOI



contactPoint



ContactPoint: Contact Peter Sefton



datePublished



+ DateTime

+ Date



June 29, 2017



description



+ Text

This is a simple dataset for demonstration purposes it contains just one image and a directory full of useless text files.



funder



+ Person

+ Organization



Organization: Australian National Data Service



Organization: University of Technology Sydney



Organization: DataCrate Project



Making your own RO-Crate with Describo

<https://arkisto-platform.github.io/describo/>
<https://arkisto-platform.github.io/tools/description/describo->

Credit: Marco La Rosa, Peter Sefton

Dataset: Survey of Victoria Arch, Wombeyan Caves NSW

[Download this Dataset](#)

[Download all the metadata for Survey of Victoria Arch, Wombeyan Caves NSW in JSON-LD format](#)

[Check this crate](#)

<https://doi.org/10.4225/59/5a4d9b76d79f4>

https://data.research.uts.edu.au/examples/ro-crate/0.2/Victoria_Arch_pub/

FAIR is not just machine-readable!

@id <https://dx.doi.org/10.4225/59/5a4d9b76d79f4>

name Survey of Victoria Arch, Wombeyan Caves NSW

@type Dataset

description This data is part of a project by Michael Lake and supported by the Australian Speleological Federation. Data was acquired at Wombeyan Caves by Robert Zlot in January 2014 using the Zebedee 3D Mapping System developed by CSIRO.

https://data.research.uts.edu.au/examples/ro-crate/0.2/Victoria_Arch_pub/ro-crate-preview.html

datePublished 2017-10-01

creator

- [Robert Zlot](#)
- [Mike Lake](#)
- [Lukas Kaul](#)

path /

contactPoint [Contact Mike Lake](#)

contentLocation Victoria Arch

hasPart [README.txt--to--wcc06_archcentre1.laz](#)
[wcc04_archentrance.laz--to--README.html](#)

identifier

- <https://dx.doi.org/10.4225/59/5a4d9b76d79f4>
- dx.doi.org/10.4225/59/5a4d9b76d79f4

license [Copyright CSIRO 2014: Available for research and academic purposes](#)

RO-Crate for digital humanities

Recordings in South Efate

Open (subject to agreeing to PDSC access conditions)

Item: NT1 / 98007

Collection: NT1

Contributors:

- Nick Thieberger (collector)
- Corin Bone (operator)
- Nick Thieberger (recorder)
- Kalsarap Namaf (speaker)
- Ikopeth (speaker)
- John Maklen (speaker)
- Waia Tenene (speaker)

NT1-98007. Text #047 (speaker is John Maklen. Text title: History of villages before Erakor); Text #048 (speaker is John Maklen. Text title: Mantu the flying fox and Erromango); Text #049. Text title: Asaraf (speaker is John Maklen); Text #050. Text title: Mumu and Kotkot (speaker is John Maklen); Text #051. Text title: Natopu ni Erakor—the spirit who lives at Erakor (speaker is John Maklen); Text #038. Text title: The need for respect (speaker is Ikopeth) Stories can be seen at NT8-TEXT. There are time-aligned transcripts of this item and handwritten transcripts by Manuel Wayne scanned as jpg files. [show less](#)

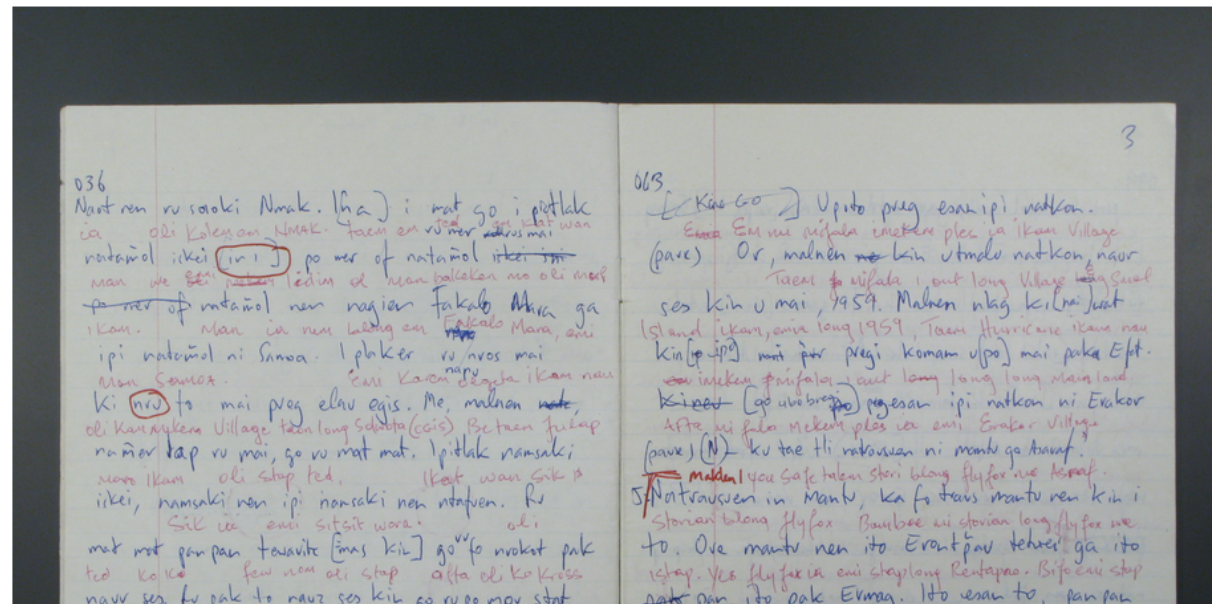
show item metadata

Content Languages



<https://arkisto-platform.github.io/case->

Capturing cultural heritage records as RO-Crates



Images Audio XML Files

NT1-98007-98007A

0:00 / 45:29

NT1-98007-98007A.ixt NT1-98007-98007A.eaf

PLAY (09:22)

Natopu nen kin ito esan ga, ipi natopu nigmam.
The natopu (local spirit) that lives here, it is our natopu.

Natopu	nen	kin	ito	esan	ga,	ipi	natopu	nigmam.
natopu	nen	kin	i= to	esan	ga	i= pi	natopu	nigmam
spirit	that	COMP	3S.RS= stay	place	3S.	3S.RS= be	spirit	1P.POS

PLAY (09:26)

Ser natamöl ni Erakor runomser mtaki natopu nen kin rusoso ki maarik.
Everyone from Erakor is scared of the natopu which they call 'maarik' (mister)

<https://mod.paradisec.org.au/view/NT1/98007?version=v1>

<http://mod.paradisec.org.au/view/NT1/98007>

@soilandreyes

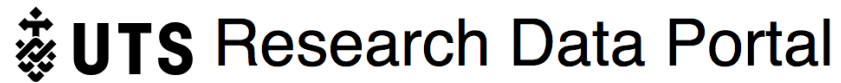
RO-Crate for repositories

RO-Crate as a metadata archival format

Metadata held alongside heterogeneous data

Exchange mechanism (import/export)

Avoid vendor lock-in

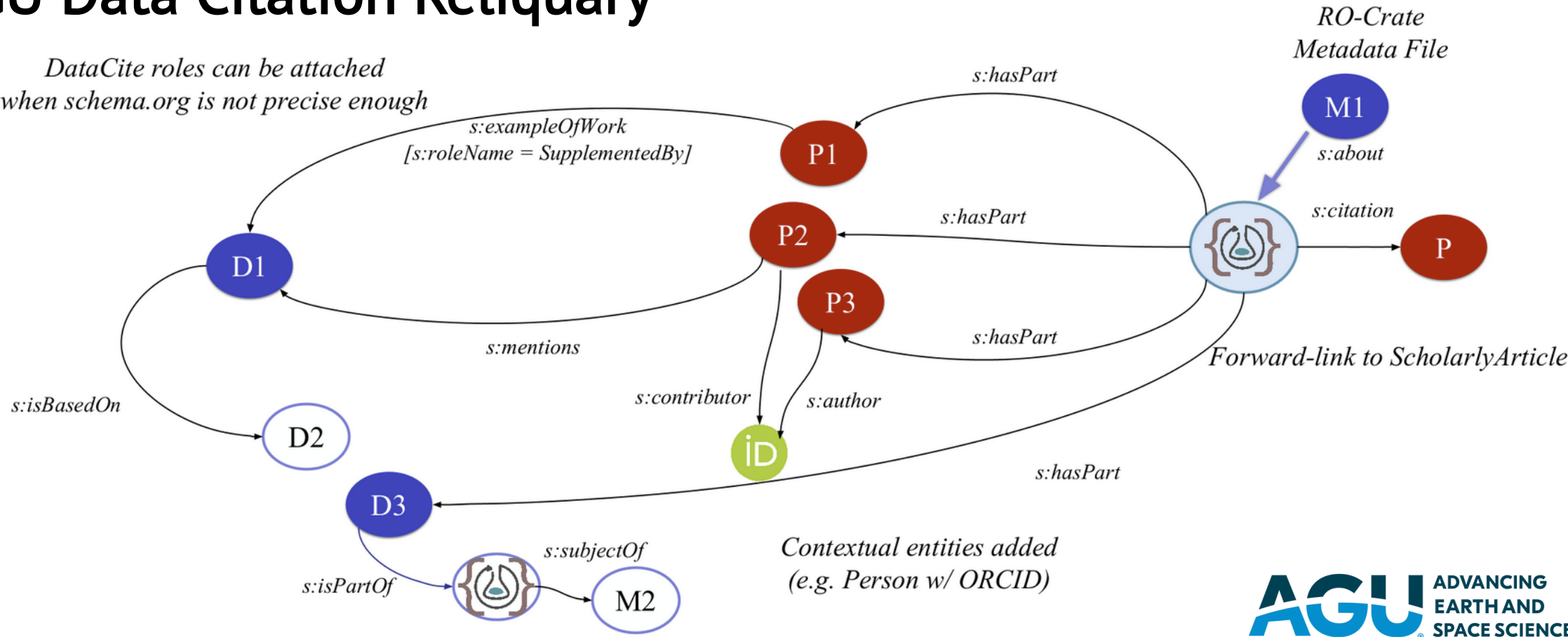


Oxford Common File Layout

Enabling large number of data citations: AGU Data Citation Reliquary

View is centric from particular RO-Crate
(Relations go mainly in one direction)
Uses schema.org vocabulary s:

DataCite roles can be attached when schema.org is not precise enough



Contextual entities added (e.g. Person w/ ORCID)

D3 is also part of another RO-Crate, but is detailed in its own Metadata file



Credit: Paolo Manghi

AGU Data Citation Workshop



@soilandreyes

<https://doi.org/10.5281/zenodo.4916734>



RO-Crate

RO-Crate in Earth Science

Earth Science RO-Crate profile (ongoing work)

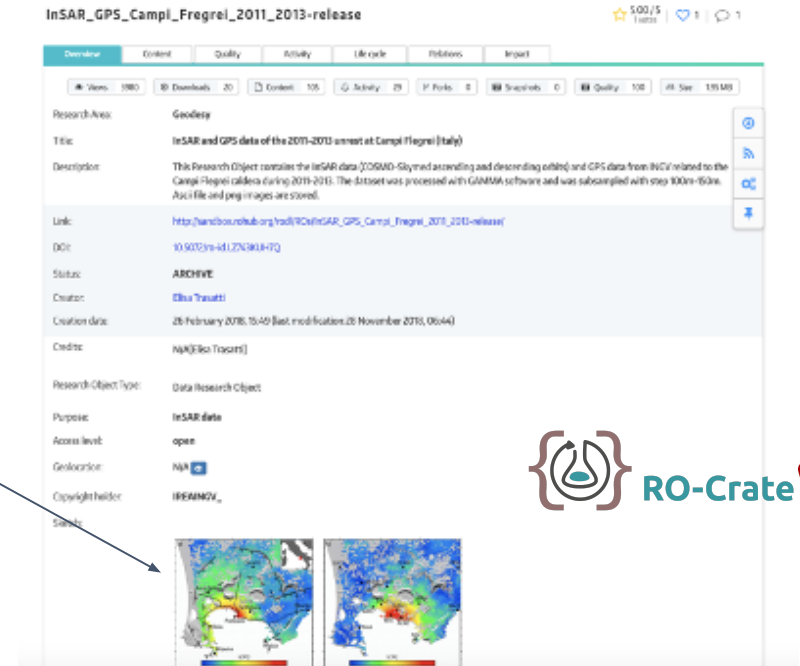
Data Cubes (**Collections and Products**) are very common in Earth Science research

We create the concept of **Data Cube Data Entity**

Relevant **metadata** for a Data Cube

- Simple **geolocation** (sch:contentLocation) and **temporal coverage** (sch:temporalCoverage)
- **Spatial coverage** (extent, resolution, crs)
- **Temporal coverage** (startTime, endTime, resolution-step)
- **Vertical coverage** (highestElevation, lowestElevation, resolution-step)
- Other (processing level, processor, instrument, platform, etc.)

http://www.rohub.org/rodetails/InSAR_GPS_Camp1_Fregrei_2011_2013-release/



- A valid RELIANCE *RO-Crate JSON-LD* graph MUST describe:
1. The [RO-Crate Metadata File Descriptor](#)
 2. The [Root Data Entity](#)
 3. One or more **Data Cube Data Entities**
 4. Zero or more [Data Entities](#)
 5. Zero or more [Contextual Entities](#)

Credit: Oscar Corcho, Carole Goble

<https://doi.org/10.5281/zenodo.4913285>

@soilandreyes

RO-Crate for workflow descriptions



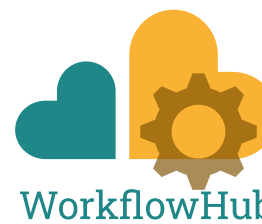
MC_COVID19like_Assembly_Reads

Version 1

Work-in-progress

This WF is based on the official Covid19-Galaxy assembly workflow as available from <https://covid19.galaxyproject.org/genomics/2-assembly/> . It has been adapted to suit the needs of the analysis of metagenomics sequencing data. Prior to be submitted to INDSC databases, these data need to be cleaned from contaminant reads, including reads of possible human origin.

The assembly of the SARS-CoV-2 genome is performed using both the Unicycler and the SPAdes assemblers, similar to the original WV.



[Visit source](#)

[Download RO Crate](#)

Creators and Submitter

Creators

Not specified

Submitter

Matteo Chiara

License

Apache Software License 2.0

Activity

Views: 78 Downloads: 2

Created: 4th Nov 2020 at 18:35

Last updated: 5th Nov 2020 at 07:42

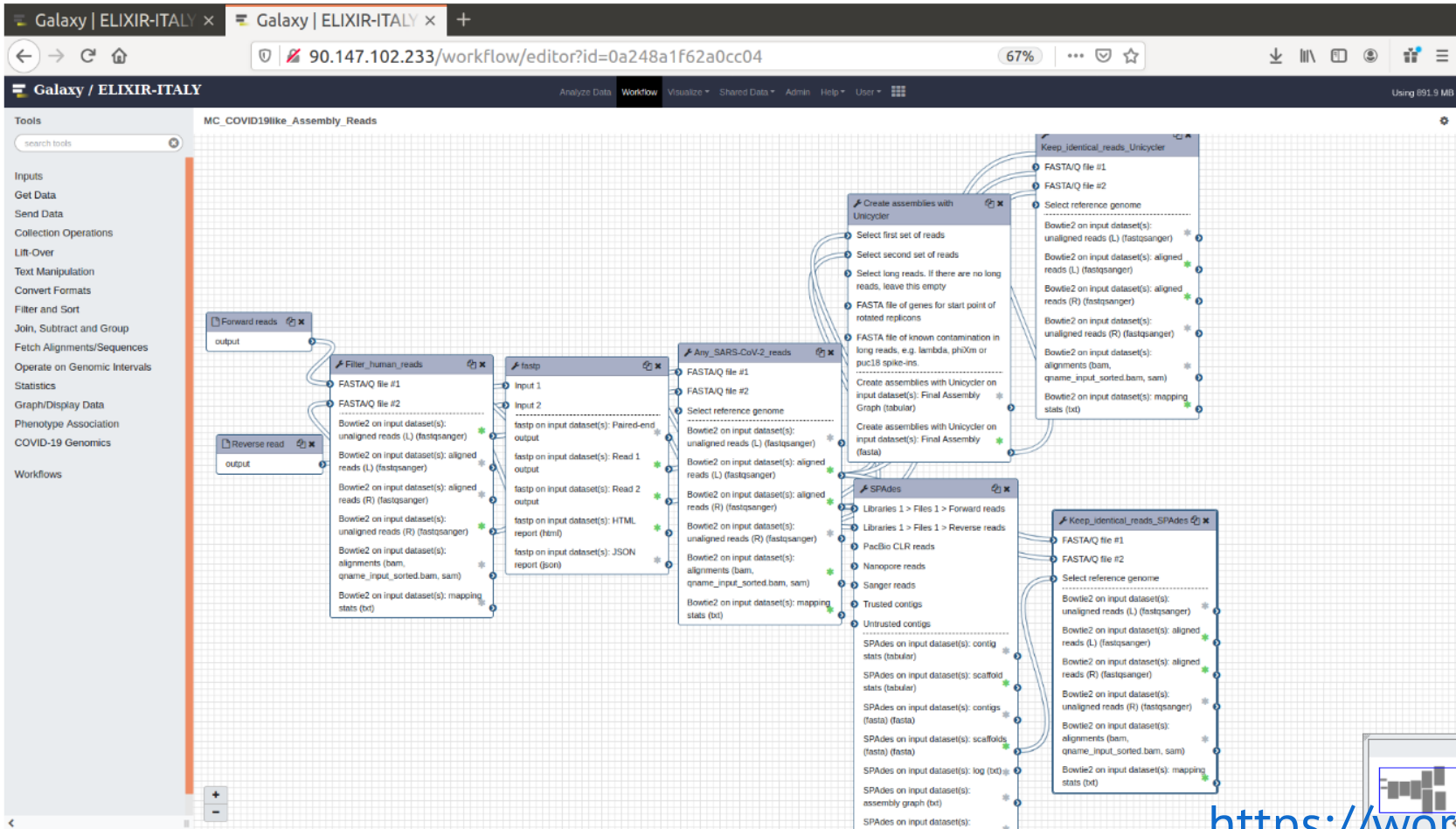
Last used: 2nd Dec 2020 at 06:56

Tags

covid-19

Attributions

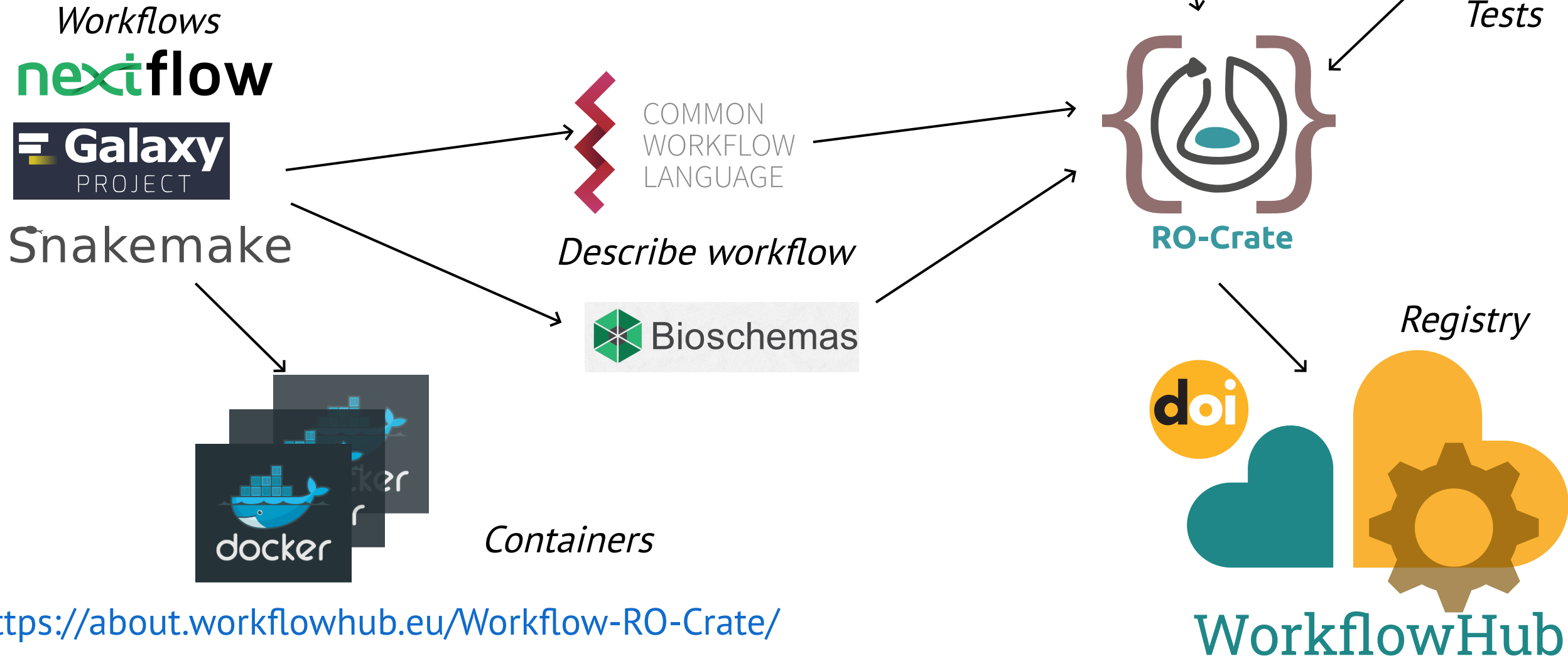
None



SEEK ID: <https://workflowhub.eu/workflows/68?version=1>

<https://workflowhub.eu/workflows/68>

Describing workflows with RO-Crate



<https://about.workflowhub.eu/Workflow-RO-Crate/>

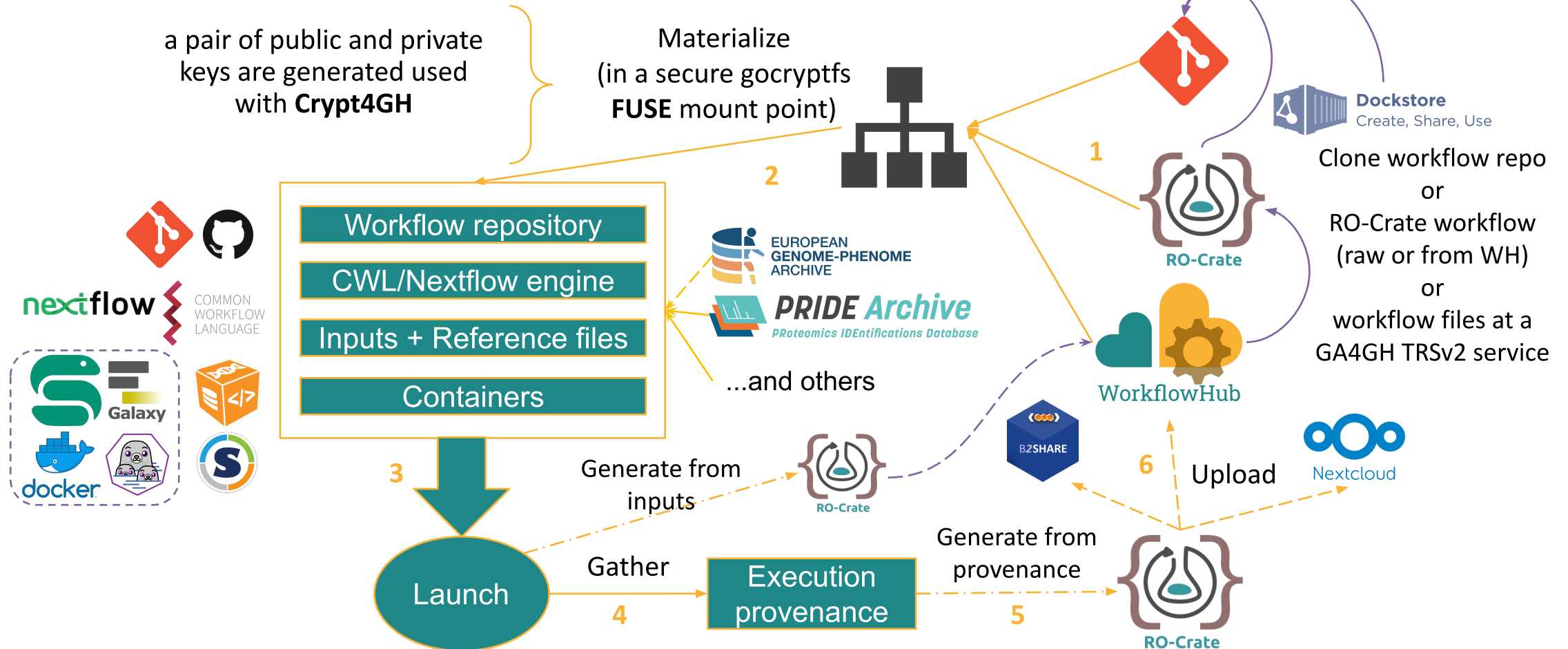
<https://pypi.org/project/galaxy2cwl/>

@soilandreyes

RO-Crate for workflow run provenance

Analysis life-cycle using WfExS

(dashed lines represent ongoing or future flows)



Preserving a Workflow Run as RO-Crate?

- Workflow **language** & **version**
- Workflow **engine** & version (e.g. Toil)
- Workflow **definition**
- **Input** data (or pointers to such)
- **Parameters**? What can be implicit and explicit?
- Tool **Dependencies** to install (mostly implied by workflow?)
- **Container** platform requirement [e.g. Docker, Conda]
- **Operating system** requirement
- **Hardware** requirements (memory, CPU, GPU)
 - Equivalent of AWS **cloud instance** type sufficient?
- **Where to run/submit** (e.g. usegalaxy.eu)
- Explicit/resolved **container IDs**
- **Archive containers** from Docker Hub (protect against image expiration)
- ...

Join discussion in the
Workflow Hub Club community!

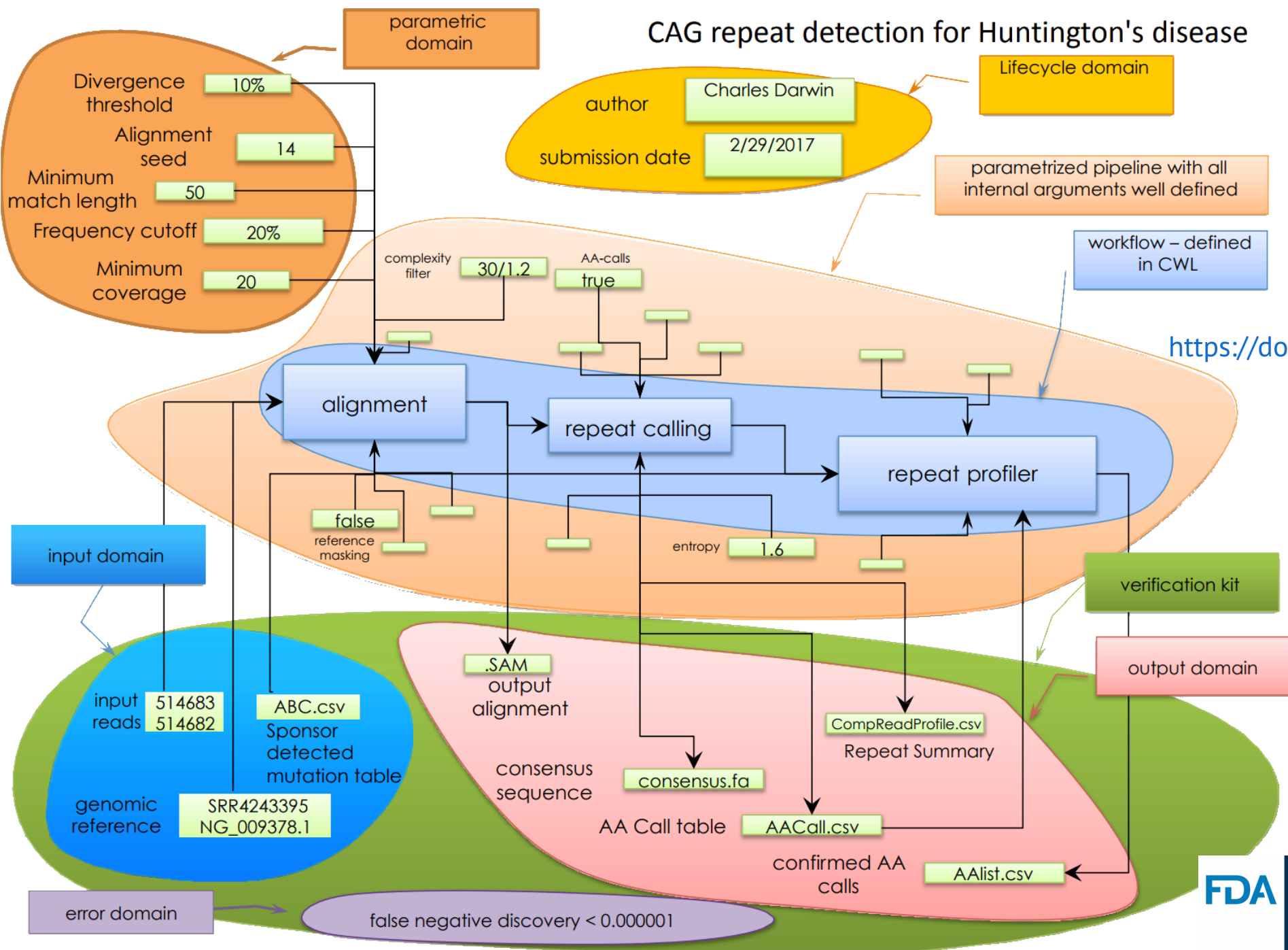
<https://about.workflowhub.eu/>

@soilandreyes



RO-Crate for regulatory sciences

CAG repeat detection for Huntington's disease



<https://www.biocomputeobject.org/>
<https://doi.org/10.1371/journal.pbio.3000099>

**THE GEORGE
WASHINGTON
UNIVERSITY**
WASHINGTON, DC



**U.S. FOOD & DRUG
ADMINISTRATION**

object_id : https://beta.portal.aws.biochemistry.gwu.edu/bco/BCO_00016916
 spec_version : <https://w3id.org/ieeee/ieeee-2791-schema/>
 etag : fea7e938e6bdf9a2cfcba7fa02f5a5fc3973dccb0b03a64319e1ee29966a5b6b

provenance_domain :

embargo :
 created : 2020-08-04T23:50:56.016Z
 modified : 2020-08-04T23:50:56.016Z
 name : Human Healthy Bulk RNA-seq Expression (Bgee)
 version : v-1.0
 obsolete_after : 2020-04-22T23:57:00.000Z
 contributors :
 contribution :
 createdBy
 name : Amanda Bell
 email : amandab2140@gwu.edu
 affiliation : GW HIVE-Lab
 orcid : <http://orcid.org/0000-0002-9920-565X>
 license : Attribution 4.0 International CC BY 4.0

Provenance Domain
description_domain :

keywords :
 Gene Expression
 Gene Expression Regulation
 Tissue specificity
 xref :
 namespace : ensembl
 name : Ensembl Genome Browser
 ids :
 Ensembl gene ID
 access_time : 2020-04-22T14:03:00.000Z
 platform :
 OncoMX
 pipeline_steps :
 step_number : 1
 name : oncomx server
 prerequisite :
 uri :
 description : Process data
 input_list :

Description Domain
error_domain :

empirical_error :
 D168Y: percentage: 0.56, calls: 0.5615, STDEV.P: 0.00075
 algorithmic_error :
 SCORE_threshold: 0.5, QUALITY: 25, COVERAGE: 5000

Error Domain
parametric_domain :

param : grep
 value : -r
 step : 1

Parametric Domain
execution_domain :

environment_variables :
 key : EDITOR
 value : vim
 key : HOSTTYPE
 value : x86_64-linux
 external_data_endpoints :
 url : <https://data.oncomx.org/ONCOMXDS000012>
 name : Human Healthy Bulk RNA-seq Expression (Bgee)
 script :
 uri :
 filename : make-dataset.py
 uri : <http://data.oncomx.org/ln2wwwdata/software/pipeline/integrator/make-dataset.py>
 access_time : 2020-04-22T14:28:00.000Z
 software_prerequisites :
 uri :
 filename : shell
 uri : <https://www.python.org/download/releases/2.7.5>
 access_time : 2020-04-22T14:30:00.000Z
 name : Python
 version : 2.7.5
 script_driver : Python

Execution Domain
io_domain :


input_subdomain :
 uri :
 filename : Homo_sapiens_UBERON:0000066
 uri :
 http://data.oncomx.org/ln2wwwdata/downloads/bgee/current/Homo_sapiens_UBERON:0000066_AFFYMETRIX_RNA_SEQ.tsv
 access_time : 2020-04-22T20:44:00.000Z
 output_subdomain :
 uri :
 filename : human_normal_expression.csv
 uri : <https://data.oncomx.org/ONCOMXDS000012>
 access_time : 2020-04-22T20:50:00.000Z
 mediatype : TEXT/CSV

IO Domain
extension_domain :


dataset_categories :
 category_value : Homo sapiens
 category_name : species
 category_value : normal
 category_name : disease_status
 extension_schema : <https://data.oncomx.org/ONCOMXDS000012>

Extension Domain
usability_domain :


List of human taxid:9606 genes with healthy RNA-Seq and Affymetrix expression data in Bgee; additional documentation available at (https://github.com/BgeeDB/bgee_pipeline/tree/develop/pipeline/collaboration/oncoMX#information-about-the-files-generated-for-oncomx) Only the subset of RNA-Seq data are used to generate the expression profiles for healthy individuals for human used by OncoMX.

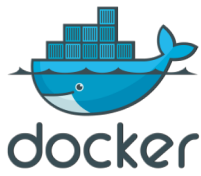


IEEE2791-2020



8 Top Level Domains	
Provenance Domain: Metadata describing the BCO	
Usability Domain: Free text field for researcher to explain the analysis and relevant details.	
Extension Domain: User-defined fields	
Description Domain: Steps of the analysis, external resources needed for the steps, and the relationship of I/O objects	
Execution Domain: Information about the environment in which the analysis was run	
Parametric Domain: Records any parameters that were changed from default values	
Input and Output Domain: A list of global input and output files	
Error Domain: Used for describing errors. Can include the limits of detectability, false positives, false negatives, statistical confidence of outcomes, and description of errors	
Required	Optional





COMMON WORKFLOW LANGUAGE



identifiers.org



```

{
  "id": "0001300",
  "name": "Detection of BRCA (hgnc:3234) gene mutations in human (taxonomy:9606) non-small cell lung carcinoma (doid:1300) patients.",
  "version": "1.0",
  "creator": "peter@genome.gov.edu",
  "digital_signature": "0x1234567890123456",
  "verification_status": "unverified",
  "publication_status": "draft",
  "activity_domain": {
    "name": "R1000",
    "description": "R1000 is a deletion (location:747 to glutamic acid-749) mutations and the single-point substitution mutation 5458 (c.25370G) in exon 21 are the most frequent in BRCA and are termed classical mutations. The point mutation 5708 (c.2840C) accounts for 90% of BRCA and exhibits resistance in about one half of the cases.",
    "url": "https://www.mycancergenome.org/content/diseases/lung-cancer/egfr/1",
    "author": "Liam Voss",
    "url": "https://orcid.org/0000-0002-9311-1508"
  },
  "description_domain": {
    "name": "R1000",
    "description": "Non-small cell lung carcinoma",
    "taxonomy": "9606",
    "doid": "1300",
    "pubmed": "19480293",
    "pubmed_url": "https://pubmed.ncbi.nlm.nih.gov/19480293/"
  },
  "execution_domain": {
    "platform": "HIVE",
    "driver": "https://hive.biochemistry.gwu.edu/hive-python-driver",
    "script": "hive/workflows/BRCA_mutation_detection_hive.py",
    "preconditions": {
      "name": "HIVE-Mezagon",
      "version": "1.4.3",
      "name": "HIVE-Mezagon",
      "version": "1.4.3",
      "env_parameters": { "HIVEV4.3" }
    },
    "parametric_domain": {
      "mezagon_all_align_length": "45",
      "mezagon_align_percent_allowed": "15",
      "mezagon_align_window_size": "10",
      "mezagon_align_coverage_allowed": "30"
    },
    "io_domain": {
      "input_url_list": [
        "https://hive.biochemistry.gwu.edu/hive-read/557031",
        "https://hive.biochemistry.gwu.edu/557365/allcount-aligned.csv"
      ],
      "output_url_list": [
        "https://hive.biochemistry.gwu.edu/557367/BCOProfile.csv",
        "error_domain": { "empirical_error": { "false_negative_discovery": "0.00001" } }
      ]
    }
  }
}

```

Alternate metadata views

Domain-specific explanation: BCO

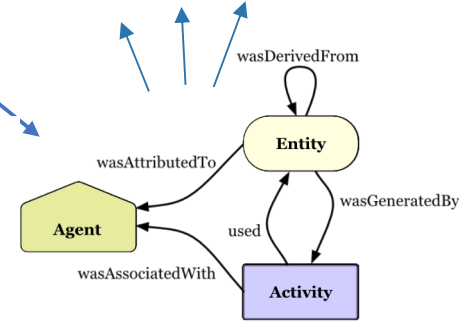
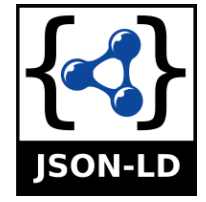
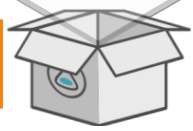
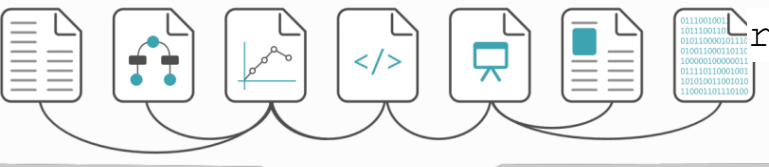
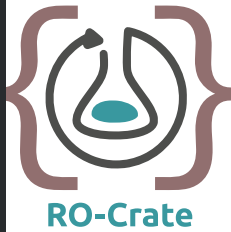
General index: RO-Crate



```

1 {
2   "@context": [
3     "https://w3id.org/ro/crate/1.0/context",
4     {
5       "@vocab": "https://schema.org/"
6     }
7   ],
8   "@graph": [
9     {
10      "@id": "ro-crate-metadata.json",
11      "@type": "CreativeWork",
12      "about": {
13        "@id": "."
14      },
15      "identifier": "ro-crate-metadata.json",
16      "conformsTo": {
17        "@id": "https://w3id.org/ro/crate/1.0"
18      },
19      "license": {
20        "@id": "https://creativecommons.org/licenses/by/4.0/"
21      },
22      "description": "Made with Describo: https://github.com/ro-crate/describo"
23    },
24    {
25      "@type": "Dataset",
26      "author": {
27        "@id": "https://orcid.org/0000-0001-9842-9714"
28      }
29    }
30  ]
31 }

```



<https://biocompute-objects.github.io/bco-ro-crate/>

RO-Crate for computational tools

Software used to create files

To specify which software was used to create or update a file, the software application SHOULD be represented with an entity of type `SoftwareApplication`, with a `version` property, e.g. from `tool --version`.

For example:

```
{
  "@id": "https://www.imagemagick.org/",
  "@type": "SoftwareApplication",
  "url": "https://www.imagemagick.org/",
  "name": "ImageMagick",
  "version": "ImageMagick 6.9.7-4 Q16 x86_64 20170114 http://www.imagemagick.org"
}
```

The software SHOULD be associated with the `File` (or other `data entities`) it created as an `instrument` of a `CreateAction`, with the `File` referenced by a `result` property. Any input files SHOULD be referenced by the `object` property.

In the below example, an image with the `@id` of `pics/2017-06-11%2012.56.14.jpg` was transformed into an new image `pics/sepia_fence.jpg` using the *ImageMagick* software application as "instrument". Actions MAY have human-readable names, which MAY be machine generated for use at scale.

```
{
  "@id": "#Photo_Capture_1",
  "@type": "CreateAction",
  "agent": {
    "@id": "https://orcid.org/0000-0002-3545-944X"
  },
  "description": "Photo snapped on a photo walk on a misty day",
  "endTime": "2017-06-11T12:56:14+10:00",
  "instrument": [
    {
      "@id": "#EPL1"
    }
  ]
}
```

<https://www.researchobject.org/rocrate/1.1/provenance.html#software-used-to-create-files>

<https://bioschemas.org/profiles/ComputationalTool/>

Use case:
HPC simulation Workflow
using building blocks



BioExcel Building Blocks has a repository on [Github](#), a code hosting platform for version control and work together on projects from anywhere.

For further information about all the blocks, please go to the [Source & Docs page](#).

[Download and Install BioExcel Building Blocks from Github](#)

INSTALL AND RUN FROM DOCKER



BioExcel Building Blocks is available as a Docker image, a container technology for Linux that allows a application with all of the parts it needs.

For further information about all the blocks, please go to the [Source & Docs page](#).

[Install and run BioExcel Building Blocks from Quay](#)

INSTALL AND RUN FROM BIOCONDA



BioExcel Building Blocks is available on [Bioconda](#), a channel for the conda package manager specialization.

For further information about all the blocks, please go to the [Source & Docs page](#).

[Install BioExcel Building Blocks from Bioconda](#) [Run BioExcel Building Blocks in Bioconda](#)

INSTALL AND RUN FROM SINGULARITY



BioExcel Building Blocks is available as a [Singularity](#) image, that enables users to have full control of the containers can be used to package entire scientific workflows, software and libraries, and even data.

For further information about all the blocks, please go to the [Source & Docs page](#).

[Install and run BioExcel Building Blocks from Singularity](#)

DOWNLOAD AND INSTALL FROM BIOCONTAINERS



BioExcel Building Blocks is available on [Biocontainers](#), an open source and community-driven framework executable environments for bioinformatics software.

[Download and Install BioExcel Building Blocks from Biocontainers](#)

DOWNLOAD BIOBB CLOUD



A [website](#) where the BioExcel Building Blocks workflows are integrated has been packaged in a Virtual [download](#) and run in house.

[Download BioExcel Building Blocks Cloud \(soon\)](#)

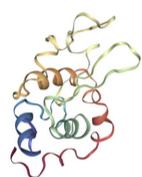
Making Canonical Workflow Building Blocks interoperable across workflow languages



BioExcel Building Blocks workflows

Home • Workflows

PROTEIN MD SETUP



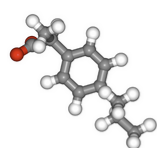
This tutorial aims to illustrate the process of setting up a simulation system containing a protein, step by step, using the BioBlocks library (biobb). The particular example used is the Lysozyme protein (PDB code 1AKI).

[WorkflowHub](#) [Jupyter Notebook](#)

[Execute in binder *](#) [View tutorial](#) [Open Github repository](#) [Open documentation](#)

(*) Binder for biobb is a small installation and to promote fair use of our resources, one user is allowed to run only one notebook server at a time. Launching a new the previous one. Users cannot see the notebooks run by other users, but please avoid entering secret data to the notebooks.

AUTOMATIC LIGAND PARAMETERIZATION



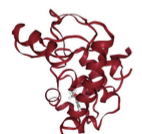
This tutorial aims to illustrate the process of ligand parameterization for a small molecule, step by step, using the Bio library (biobb). The particular example used is the Ibuprofen small compound (3-letter code IBP, Drugbank code DB0 anti-inflammatory drug (NSAID) derived from propionic acid and it is considered the first of the propionics.

[WorkflowHub](#) [Jupyter Notebook](#)

[Execute in binder *](#) [View tutorial](#) [Open Github repository](#) [Open documentation](#)

(*) Binder for biobb is a small installation and to promote fair use of our resources, one user is allowed to run only one notebook server at a time. Launching a new the previous one. Users cannot see the notebooks run by other users, but please avoid entering secret data to the notebooks.

PROTEIN-COMPLEX MD SETUP



This tutorial aims to illustrate the process of setting up a simulation system containing a protein in complex with a lig, using the BioExcel Building Blocks library (biobb). The particular example used is the T4 lysozyme L99A/M102Q protei complex with the 2-propylphenol small molecule (3-letter Code JZ4).

[WorkflowHub](#) [Jupyter Notebook](#)

[Execute in binder *](#) [View tutorial](#) [Open Github repository](#) [Open documentation](#)

(*) Binder for biobb is a small installation and to promote fair use of our resources, one user is allowed to run only one notebook server at a time. Launching a new the previous one. Users cannot see the notebooks run by other users, but please avoid entering secret data to the notebooks.



<https://mmb.irbbarcelona.org/biobb/>

@soilandreyes

<https://doi.org/10.5281/zenodo.5542553>

https://raw.githubusercontent.com/bioexcel/biobb_hpc_workflows/condapack/ro-crate-preview.html

<https://stain.github.io/ro-crate-paper/>



Accepted

<https://stain.github.io/ro-crate-paper/>

<https://doi.org/10.5281/zenodo.5146228>

arXiv:2108.06503

@soilandreyes

- Peter Sefton <https://orcid.org/0000-0002-3545-944X> (co-chair)
- Stian Soiland-Reyes <https://orcid.org/0000-0001-9842-9718> (co-chair)
- Eoghan Ó Carragáin <https://orcid.org/0000-0001-8131-2150> (emeritus chair)
- Oscar Corcho <https://orcid.org/0000-0002-9260-0753>
- Daniel Garijo <https://orcid.org/0000-0003-0454-7145>
- Raul Palma <https://orcid.org/0000-0003-4289-4922>
- Frederik Coppens <https://orcid.org/0000-0001-6565-5145>
- Carole Goble <https://orcid.org/0000-0003-1219-2137>
- José María Fernández <https://orcid.org/0000-0002-4806-5140>
- Kyle Chard <https://orcid.org/0000-0002-7370-4805>
- Jose Manuel Gomez-Perez <https://orcid.org/0000-0002-5491-6431>
- Michael R Crusoe <https://orcid.org/0000-0002-2961-9670>
- Ignacio Eguinoa <https://orcid.org/0000-0002-6190-122X>
- Nick Juty <https://orcid.org/0000-0002-2036-8350>
- Kristi Holmes <https://orcid.org/0000-0001-8420-5254>
- Jason A. Clark <https://orcid.org/0000-0002-3588-6257>
- Salvador Capella-Gutierrez <https://orcid.org/0000-0002-0309-604X>
- Alasdair J. G. Gray <https://orcid.org/0000-0002-5711-4872>
- Stuart Owen <https://orcid.org/0000-0003-2130-0865>
- Alan R Williams <https://orcid.org/0000-0003-3156-2105>
- Giacomo Tartari <https://orcid.org/0000-0003-1130-2154>
- Finn Bacall <https://orcid.org/0000-0002-0048-3300>
- Thomas Thelen <https://orcid.org/0000-0002-1756-2128>
- Hervé Ménager <https://orcid.org/0000-0002-7552-1009>
- Laura Rodríguez-Navas <https://orcid.org/0000-0003-4929-1219>
- Paul Walk <https://orcid.org/0000-0003-1541-5631>
- brandon whitehead <https://orcid.org/0000-0002-0337-8610>
- Mark Wilkinson <https://orcid.org/0000-0001-6960-357X>
- Paul Groth <https://orcid.org/0000-0003-0183-6910>
- Erich Bremer <https://orcid.org/0000-0003-0223-1059>
- LJ Garcia Castro <https://orcid.org/0000-0003-3986-0510>
- Karl Sebby <https://orcid.org/0000-0001-6022-9825>
- Alexander Kanitz <https://orcid.org/0000-0002-3468-0652>
- Ana Trisovic <https://orcid.org/0000-0003-1991-0533>
- Gavin Kennedy <https://orcid.org/0000-0003-3910-0474>
- Mark Graves <https://orcid.org/0000-0003-3486-8193>
- Jasper Koehorst <https://orcid.org/0000-0001-8172-8981>
- Simone Leo <https://orcid.org/0000-0001-8271-5429>
- Marc Portier <https://orcid.org/0000-0002-9648-6484>
- Paul Brack <https://orcid.org/0000-0002-5432-2748>
- Milan Ojsteršek <https://orcid.org/0000-0003-1743-8300>
- Bert Droesbeke <https://orcid.org/0000-0003-0522-5674>
- Chenxu Niu <https://orcid.org/0000-0002-2142-1731>
- Kosuke Tanabe <https://orcid.org/0000-0002-9986-7223>
- Tomasz Miksa <https://orcid.org/0000-0002-4929-7875>
- Marco La Rosa <https://orcid.org/0000-0001-5383-6993>
- Cedric Decruw <https://github.com/cedricdcc>
- Andreas Czerniak <https://orcid.org/0000-0003-3883-4169>
- Jeremy Jay <https://orcid.org/0000-0002-5761-7533>
- Sergio Serra <https://orcid.org/0000-0002-0792-8157>
- Ronald Siebes <https://orcid.org/0000-0001-8772-7904>
- Shaun de Witt <https://orcid.org/0000-0003-4196-3658>
- Shady El Damaty <https://orcid.org/0000-0002-2318-4477>
- Douglas Lowe <https://orcid.org/0000-0002-1248-3594>
- Sergio Serra <https://orcid.org/0000-0002-0792-8157>
- Xuanqi Li <https://orcid.org/0000-0003-1498-6205>

RO-Crate Community is open for anyone to [join us!](#)


<https://www.researchobject.org/ro-crate/community>

@soilandreyes

More documentation

- 5 minute intro
- Tutorial / training material (Carpentries-style)
- Showcases of use

- Home
- About
- Background
- Community
- Examples
- Outreach and Publications
- Profiles
- RO-Crate 1.1 (latest)
- RO-Crate In Use
- Specification
- Tools

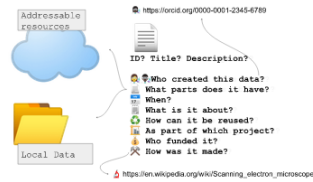


5' intro to RO-Crate

Scroll down through our quick introduction to learn what RO-Crate can do for you.

Describe data

RO-Crate is a method for describing a dataset as a digital object using a **single Linked Data metadata document**.



Add any resource

The dataset may contain any kind of data resource about anything, in any format as a file or URL

Machine-readable

Each resource can have a machine readable description in JSON-LD format.

```
--- Folder1/ [ro:file] (316)
  --- Folder2/ [ro:file] (316)
  --- Folder3/ [ro:file] (316)
  -2021-04-08 07:08:37.399
  {
    "@type": "File",
    "name": "Cute puppy"
  }
```





Future directions

Ongoing:

- CS3mesh4EOSC (Describo Online w/ cloud storage)
- EOSC-Life (Life Science workflows)
- BY-COVID (COVID-19)
- SYNTHESYS+/DiSSCO (FAIR Digital Objects to digitize museum collections)
- GA4GH & ELIXIR Cloud (workflow execution)

Next round (proposals and collaborations):

- Workflows as FAIR Digital Objects
- Digital Twins in Biodiversity
- Minimal metadata standards for Semantic Interoperability
- Cross-repository integrations

@soilandreyes