



Standardized Computational Workflows at EBRAINS

Sofia Karvounari
Athena Research Center - Technical Coordination team

Standardization

Define in a common standard way:

- Computational workflows: series of **EBRAINS tools** that can be linked/combined to create Directed Acyclic Graphs (DAGs), loops or branches for accomplishing a scientific objective.
- **EBRAINS tools**: scientific simulations or data analysis (data manipulation) pieces of software **packaged** together with libraries, dependencies, binaries.

Workflows and EBRAINS tools will be **executed** via workflow management systems that can lead to **monitoring**, automation of some procedures and recovery of failing steps.



Two worlds..

SCIENTISTS

TECHNICAL EXPERTS

.. linked together!

SCIENTISTS

TECHNICAL EXPERTS

What standardization offer

FOR SCIENTISTS

Automation

Scalability

Accessibility

Portability

FAIRness is everywhere...

F A I R

Findable

Accessible

Interoperable

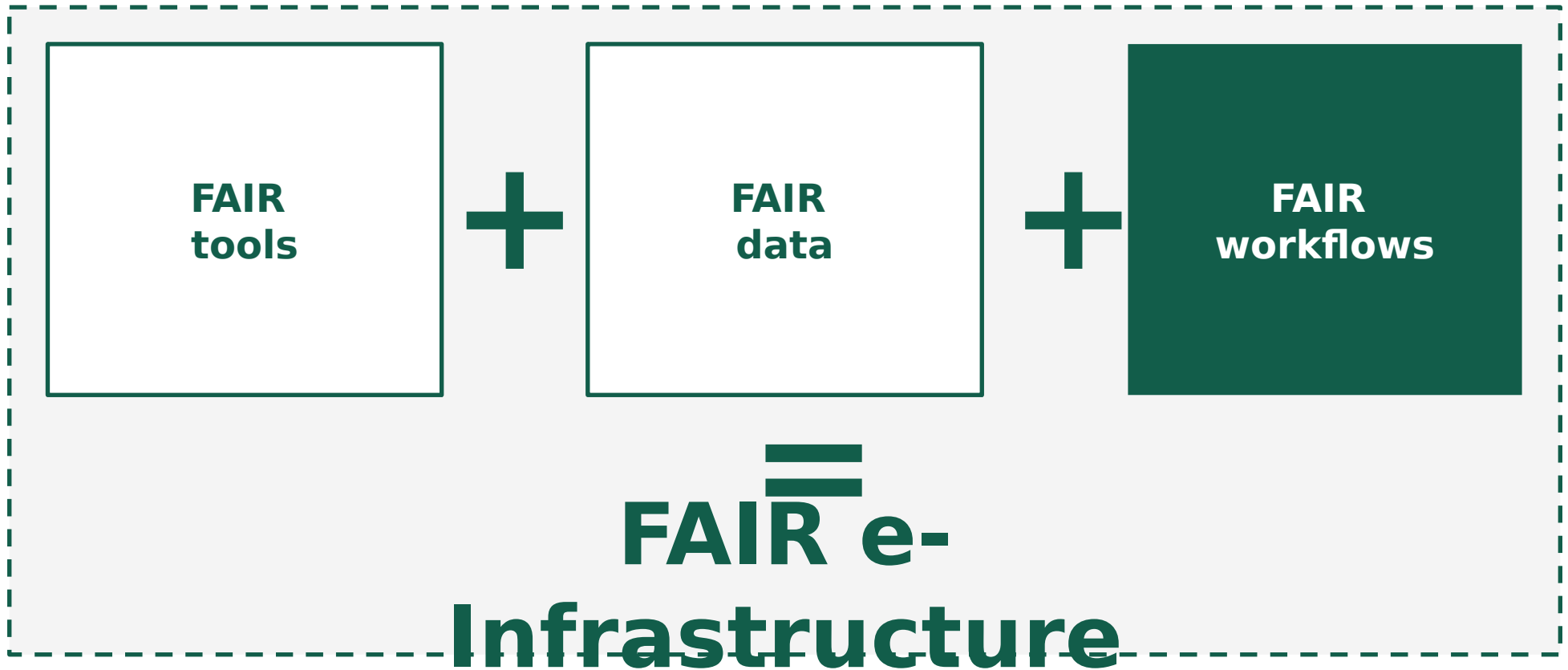
Reusable

[1] [The FAIR Guiding Principles for scientific data management and stewardship](#)

[2] [Sharing Fair Research Objects](#)

...Why standard workflows are important?

FOR EBRAINS



Let's take a step back

"EBRAINS is a new digital research infrastructure, created by the EU-funded Human Brain Project, that gathers an extensive range of data and tools for brain-related research."

Service categories at EBRAINS



[2] https://ebrains.eu/#news_anchor

Current limitations

WORKFLOWS

- No straight forward way to combine tools together
- No way to monitor large jobs that take time
- No way to pause, restart failed tasks, provide logs
- No way to be easily found, accessed and replicated/reproduced
- Configuration adjustments in order to run in different infrastructures
- No graphs, dependencies, parallelism provided by Jupyter Notebooks

NON - INTERACTIVE TOOLS

- Not defined in an interoperable, standard and easy way to (re-)use

(Standard) Computational Workflows



Definition



**Command Line
tools**



Execution



Monitoring



(Standard) Computational Workflows



Definition



Command Line
tools



Execution



Monitoring



Dive in defining



A. Common Workflow Language



B. Specific Use Case written in CWL



C. Rabix composer



A. Common Workflow Language



WHAT IT IS

Emerging open standard/common way to define analysis workflows and tools



ITS VALUES (ASAP)

- Offers **Scalability & Portability**
- Easy **Accessibility**
- Open standard, hence defining workflows and tools become **Automated**



A. Common Workflow Language



WHY USE

- Common declarative format for **tools** and **workflows**
- Community based standards effort. Extensible
- Support containers (Docker, Singularity)
- Data locality
 - input and output files are modeled as rich objects with URI/IRI + metadata
 - platforms that understand CWL can
 - use URI/IRI to send compute near location of data
 - fetch data from a remote location



[1] [Portable workflow creation and deployment with the CWL standards](#)

A. Common Workflow Language



WHY USE

- Other fields (Bioinformatics, Astrophysics) already see the value
- Compatible with plethora of workflow management systems
- Decouples description from execution of a workflow



[1] [Portable workflow creation and deployment with the CWL standards](#)

B. Use Case: Power Spectral Density



WORKFLOW VIA CWL

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: Workflow
5
6 inputs:
7   bucket_id: string
8   input_file: string
9   channels: int[]
10  psd_output_file_name: string
11  output_file_name: string
12  token: string
13
14 outputs:
15   final_output:
16     type: File
17     outputSource: visualization/plot
18
19 steps:
20   fetching_data:
21     run: step1/fetching_data_tool.cwl
22     in:
23       bucket_id: bucket_id
24       object_name: input_file
25       token: token
26     out: [fetched_file]
27
28   analysis:
29     run: step2/analysis_tool.cwl
30     in:
31       input_file: fetching_data/fetched_file
32       output_file_name: psd_output_file_name
33       channels: channels
34     out: [output_file]
35
36   visualization:
37     run: step3/visualization_tool.cwl
38     in:
39       input_file: analysis/output_file
40       output_file_name: output_file_name
41       channels: channels
42     out: [plot]
```

Thanks to Arnau Manasanch
(WP2)

B. Use Case: Power Spectral Density Fetching Tool

WORKFLOW STEP

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: Workflow
5
6 inputs:
7   bucket_id: string
8   input_file: string
9   channels: int[]
10  psd_output_file_name: string
11  output_file_name: string
12  token: string
13
14 outputs:
15   final_output:
16     type: File
17     outputSource: visualization/plot
18
19 steps:
20   fetching_data:
21     run: step1/fetching_data_tool.cwl
22     in:
23       bucket_id: bucket_id
24       object_name: input_file
25       token: token
26     out: [fetched_file]
27
28   analysis:
29     run: step2/analysis_tool.cwl
30     in:
31       input_file: fetching_data/fetched_file
32       output_file_name: psd_output_file_name
33       channels: channels
34     out: [output_file]
35
36   visualization:
37     run: step3/visualization_tool.cwl
38     in:
39       input_file: analysis/output_file
40       output_file_name: output_file_name
41       channels: channels
42     out: [plot]
```

TOOL DEFINITION

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: CommandLineTool
5 baseCommand: fetching_data.py
6 hints:
7   DockerRequirement:
8     dockerPull: docker-registry.ebrains.eu/tc/cwl-workflows/psd_workflow_fetching_data:latest
9   ResourceRequirement:
10     ramMin: 2048
11     outdirMin: 4096
12 inputs:
13   bucket_id:
14     type: string
15     inputBinding:
16       position: 1
17   object_name:
18     type: string
19     inputBinding:
20       position: 2
21   token:
22     type: string
23     inputBinding:
24       position: 3
25 outputs:
26   fetched_file:
27     type: File
28     outputBinding:
29       glob: ${inputs.object_name}
```



B. Use Case: Power Spectral Density Analysis Tool

WORKFLOW STEP

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: Workflow
5
6 inputs:
7   bucket_id: string
8   input_file: string
9   channels: int[]
10  psd_output_file_name: string
11  output_file_name: string
12  token: string
13
14 outputs:
15   final_output:
16     type: File
17     outputSource: visualization/plot
18
19 steps:
20   fetching_data:
21     run: step1/fetching_data_tool.cwl
22     in:
23       bucket_id: bucket_id
24       object_name: input_file
25       token: token
26     out: [fetched_file]
27
28   analysis:
29     run: step2/analysis_tool.cwl
30     in:
31       input_file: fetching_data/fetched_file
32       output_file_name: psd_output_file_name
33       channels: channels
34     out: [output_file]
35
36   visualization:
37     run: step3/visualization_tool.cwl
38     in:
39       input_file: analysis/output_file
40       output_file_name: output_file_name
41       channels: channels
42     out: [plot]
```

TOOL DEFINITION

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: CommandLineTool
5 baseCommand: analysis.py
6 hints:
7   DockerRequirement:
8     dockerPull: docker-registry.ebrains.eu/tc/cwl-workflows/psd_workflow_analysis:latest
9   ResourceRequirement:
10    ramMin: 2048
11    outdirMin: 4096
12 inputs:
13   input_file:
14     type: File
15     inputBinding:
16       position: 1
17   output_file_name:
18     type: string
19     inputBinding:
20       prefix: --output_file
21       position: 2
22   channels:
23     type: int[]
24     inputBinding:
25       prefix: --channels
26       position: 3
27 outputs:
28   output_file:
29     type: File
30     outputBinding:
31       glob: ${inputs.output_file_name}
```



B. Use Case: Power Spectral Density Visualisation Tool

WORKFLOW STEP

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: Workflow
5
6 inputs:
7   bucket_id: string
8   input_file: string
9   channels: int[]
10  psd_output_file_name: string
11  output_file_name: string
12  token: string
13
14 outputs:
15   final_output:
16     type: File
17     outputSource: visualization/plot
18
19 steps:
20   fetching_data:
21     run: step1/fetching_data_tool.cwl
22     in:
23       bucket_id: bucket_id
24       object_name: input_file
25       token: token
26     out: [fetched_file]
27
28   analysis:
29     run: step2/analysis_tool.cwl
30     in:
31       input_file: fetching_data/fetched_file
32       output_file_name: psd_output_file_name
33       channels: channels
34     out: [output_file]
35
36   visualization:
37     run: step3/visualization_tool.cwl
38     in:
39       input_file: analysis/output_file
40       output_file_name: output_file_name
41       channels: channels
42     out: [plot]
```

TOOL DEFINITION

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: CommandLineTool
5 baseCommand: analysis.py
6 hints:
7   DockerRequirement:
8     dockerPull: docker-registry.ebrains.eu/tc/cwl-workflows/psd_workflow_analysis:latest
9   ResourceRequirement:
10    ramMin: 2048
11    outdirMin: 4096
12 inputs:
13   input_file:
14     type: File
15     inputBinding:
16       position: 1
17   output_file_name:
18     type: string
19     inputBinding:
20       prefix: --output_file
21       position: 2
22   channels:
23     type: int[]
24     inputBinding:
25       prefix: --channels
26       position: 3
27 outputs:
28   output_file:
29     type: File
30     outputBinding:
31       glob: ${inputs.output_file_name}
```



B. Use Case: Power Spectral Density

WORKFLOW INPUTS

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: Workflow
5
6 inputs:
7   bucket_id: string
8   input_file: string
9   channels: int[]
10  psd_output_file_name: string
11  output_file_name: string
12  token: string
13
14 outputs:
15   final_output:
16     type: File
17     outputSource: visualization/plot
18
19 steps:
20   fetching_data:
21     run: step1/fetching_data_tool.cwl
22     in:
23       bucket_id: bucket_id
24       object_name: input_file
25       token: token
26     out: [fetched_file]
27
28   analysis:
29     run: step2/analysis_tool.cwl
30     in:
31       input_file: fetching_data/fetched_file
32       output_file_name: psd_output_file_name
33       channels: channels
34     out: [output_file]
35
36   visualization:
37     run: step3/visualization_tool.cwl
38     in:
39       input_file: analysis/output_file
40       output_file_name: output_file_name
41       channels: channels
42     out: [plot]
```

INPUT PARAMETERS

```
1 bucket_id: 'fetching-data-from-bucket-and-kg'
2 input_file: 'data_mat/sub-1_chs-32_hem-RH_ana-ISO000_stim-SPN.mat'
3 channels: [0, 1, 2, 3, 4]
4 psd_output_file_name: 'psd.json'
5 output_file_name: 'output.png'
6 token: 'COPY_TOKEN_HERE'
```

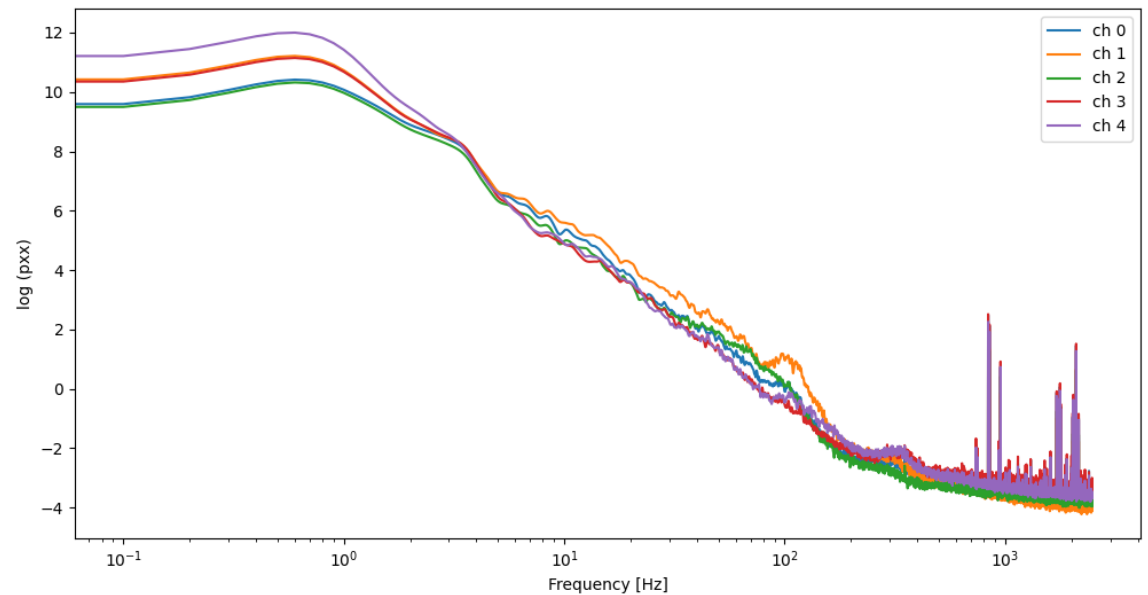


B. Use Case: Power Spectral Density

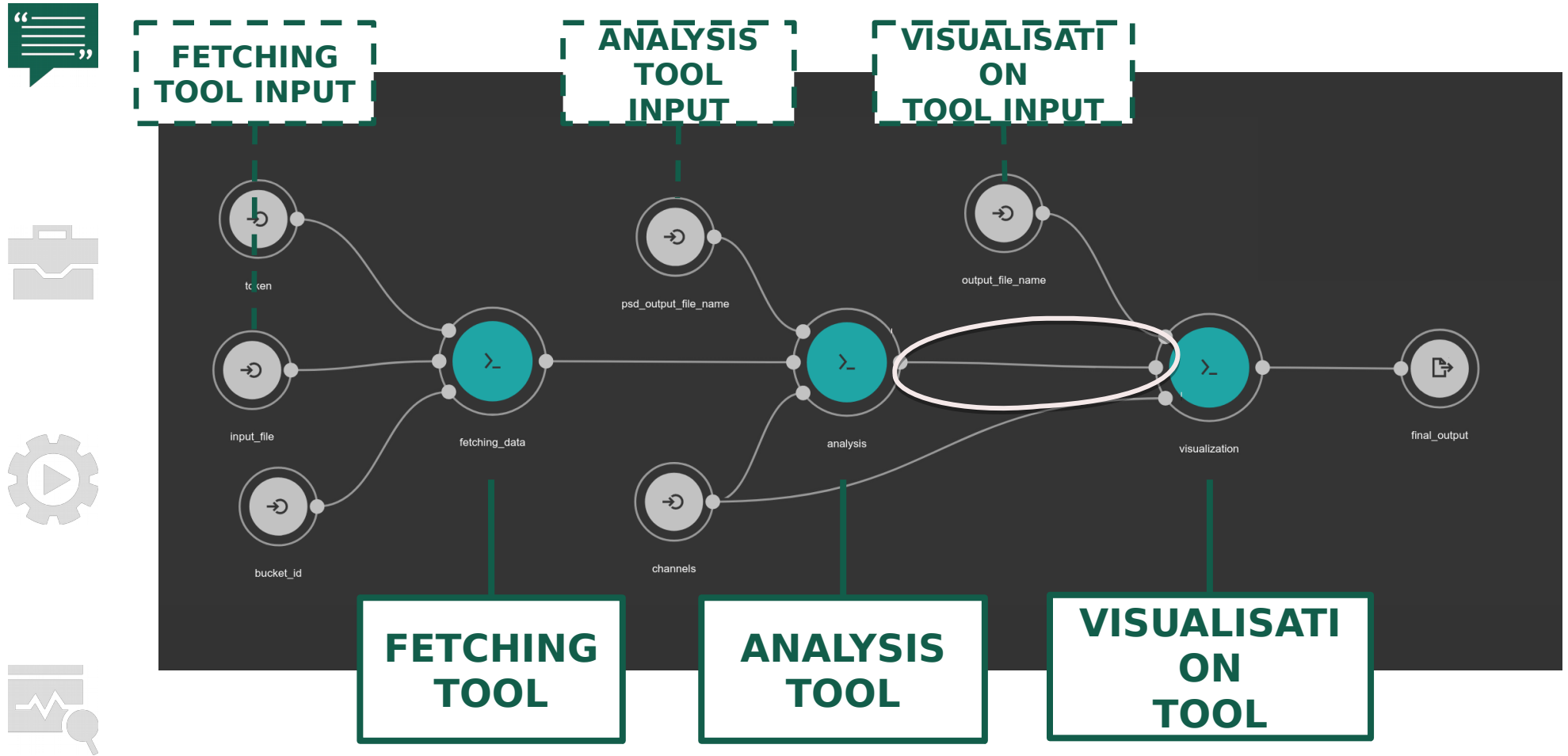
WORKFLOW OUTPUT

```
#!/usr/bin/env cwltool
#
# cwlversion: v1.0
# class: Workflow
#
# inputs:
#   bucket_id: string
#   input_file: string
#   channels: int[]
#   psd_output_file_name: string
#   output_file_name: string
#   token: string
#
# outputs:
#   final_output:
#     type: File
#     outputSource: visualization/plot
#
# steps:
#   fetching_data:
#     run: step1/ fetching_data_tool.cwl
#     in:
#       bucket_id: bucket_id
#       object_name: input_file
#       token: token
#     out: [ fetched_file ]
#
#   analysis:
#     run: step2/analysis_tool.cwl
#     in:
#       input_file: fetching_data/fetched_file
#       output_file_name: psd_output_file_name
#       channels: channels
#     out: [ output_file ]
#
#   visualization:
#     run: step3/visualization_tool.cwl
#     in:
#       input_file: analysis/output_file
#       output_file_name: output_file_name
#       channels: channels
#     out: [ plot ]
```

PLOT



C. Rabix composer



[1] [Rabix](#)

(Standard) Computational Workflows



Definition



**Command Line
tools**



Execution



Monitoring



Dive in Command Line tools



A. Types of tools

B. Descriptions via CWL



C. Packaged vs Unpackage



D. Why Docker?

How a tool can be packaged via Docker?

What else is there?



E. Why Packaged Command Line tools important?

A. Types of tools:



NON-INTERACTIVE

User (possibly) interacts with the program before the runtime

INTERACTIVE

User interacts with the program during the runtime

B. Descriptions via CWL

WORKFLOW

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: Workflow
5
6 inputs:
7   bucket_id: string
8   input_file: string
9   channels: int[]
10  psd_output_file_name: string
11  output_file_name: string
12  token: string
13
14 outputs:
15   final_output:
16     type: File
17     outputSource: visualization/plot
18
19 steps:
20   fetching_data:
21     run: step1/fetching_data_tool.cwl
22     in:
23       bucket_id: bucket_id
24       object_name: input_file
25       token: token
26     out: [fetched_file]
27
28   analysis:
29     run: step2/analysis_tool.cwl
30     in:
31       input_file: fetching_data/fetched_file
32       output_file_name: psd_output_file_name
33       channels: channels
34     out: [output_file]
35
36   visualization:
37     run: step3/visualization_tool.cwl
38     in:
39       input_file: analysis/output_file
40       output_file_name: output_file_name
41       channels: channels
42     out: [plot]
```

COMMAND LINE TOOL

```
1 #!/usr/bin/env cwl-runner
2
3 cwlVersion: v1.0
4 class: CommandLineTool
5 baseCommand: echo
6 inputs:
7   message:
8     type: string
9     inputBinding:
10      position: 1
11   stdout: output.txt
12 outputs:
13   out:
14     type: stdout
```



B. Descriptions via CWL

COMMAND LINE TOOL

```
1 #!/usr/bin/env cwltool
2
3 cwlVersion: v1.0
4 class: CommandLineTool
5 baseCommand: analysis.py
6 hints:
7
8
9 ResourceRequirement:
10   ramMin: 2048
11   outdirMin: 4096
12 inputs:
13   input_file:
14     type: File
15     inputBinding:
16       position: 1
17   output_file_name:
18     type: string
19     inputBinding:
20       prefix: --output_file
21       position: 2
22   channels:
23     type: int[]
24     inputBinding:
25       prefix: --channels
26       position: 3
27 outputs:
28   output_file:
29     type: File
30     outputBinding:
31       glob: $(inputs.output_file_name)
```



C. Packaged vs Unpackaged?



```
1  #!/usr/bin/env cwl-runner
2
3  cwlVersion: v1.0
4  class: CommandLineTool
5  baseCommand: echo
6  inputs:
7    message:
8      type: string
9      inputBinding:
10       position: 1
11  stdout: output.txt
12  outputs:
13    out:
14      type: stdout
```

- Tool needs to be in \$PATH in order for CWL to access it
- All dependencies, libraries must be taken care by the developer manually
- Tool already deployed & ready to be used

C. Packaged vs Unpackaged?



```
6 hints:
7   DockerRequirement:
8     dockerPull: docker-registry.ebrains.eu/tc/cwl-workflows/psd_workflow_analysis:latest
```

D. Why Docker?

```
1  #!/usr/bin/env cwltool
2
3  cwlVersion: v1.0
4  class: CommandLineTool
5  baseCommand: analysis.py
6  hints:
7    DockerRequirement:
8      dockerPull: docker-registry.ebrains.eu/tc/cwl-workflows/psd_workflow_analysis:latest
9    ResourceRequirement:
10     ramMin: 2048
11     outdirMin: 4096
12  inputs:
13    input_file:
14     type: File
15     inputBinding:
16       position: 1
17    output_file_name:
18     type: string
19     inputBinding:
20       prefix: --output_file
21       position: 2
22    channels:
23     type: int[]
24     inputBinding:
25       prefix: --channels
26       position: 3
27  outputs:
28    output_file:
29     type: File
30     outputBinding:
31       glob: ${inputs.output_file_name}
```

- CWL works great with Docker containers, input/output files are taking care
- Dependencies, libraries, binaries, code: all packaged
- Docker is supported by other containerized methods like Singularity (HPC)



D. Why Docker?

How can a tool be packaged via Docker?

```
1 FROM python:3.8-slim
2
3 # install dependencies
4 RUN pip install --no-cache-dir argparse numpy pandas scipy
5
6 # copy python script, make executable and add to path
7 COPY analysis.py /home/tool/analysis.py
8 RUN chmod +x /home/tool/analysis.py
9 ENV PATH="/home/tool:$PATH"
10
11 CMD [ "/bin/bash" ]
12
```



1. Create a simple "Dockerfile" with all the commands to install the tool dependencies.



2. "build" (create) the image:

```
docker build -t docker-registry.ebrains.eu/tc/cwl-workflows/analysis -f Dockerfile
```

or "tag" (rename) if it already exists:

```
docker tag old_image_name docker-registry.ebrains.eu/tc/cwl-workflows/analysis
```



3. "push" (upload) the image:

```
docker login docker-registry.ebrains.eu
docker push docker-registry.ebrains.eu/tc/cwl-workflows/analysis
```

D. Why Docker?

What else is there?



FUTURE USE CASE?

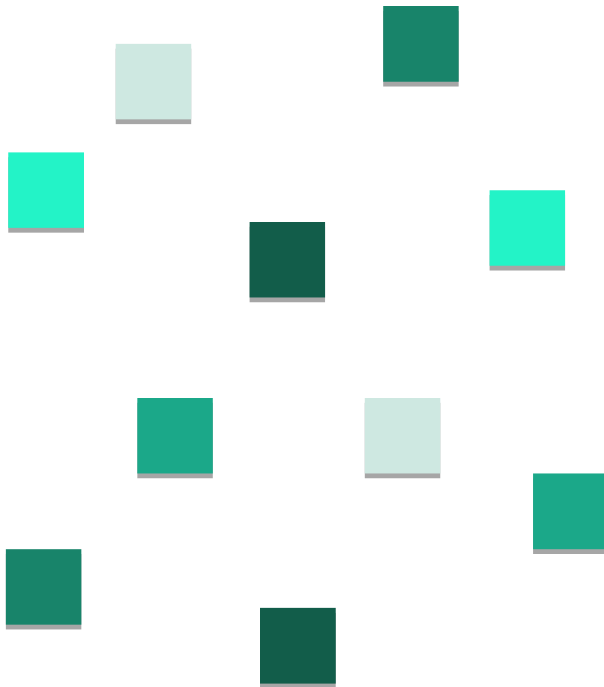
- Singularity
- Modules (Spack)
- Automatic installation of packages via IRIs (beta)



E. Why Packaged Command Line tools important?

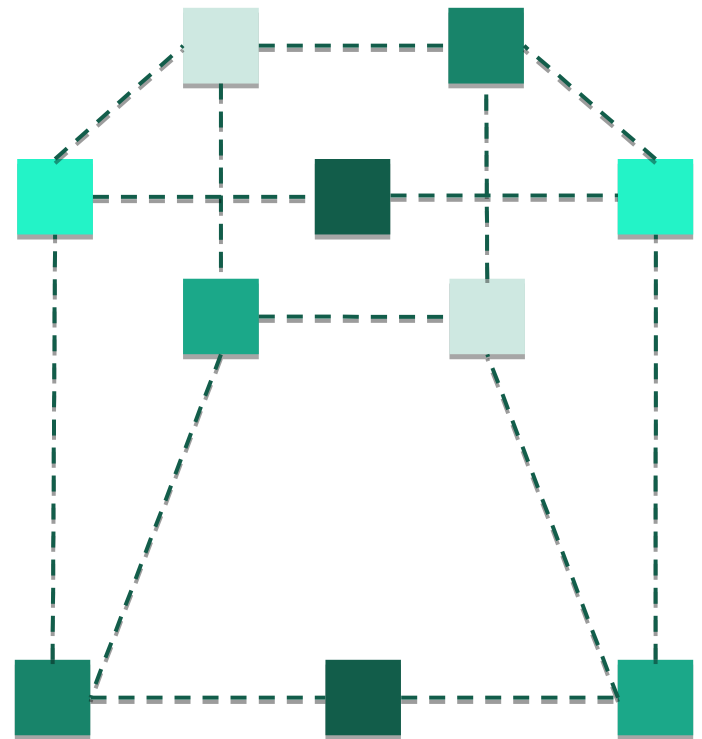


Repository / Catalogue



Command line tools

Workflows



(Standard) Computational Workflows



Definition



Command Line
tools



Execution



Monitoring



Dive in executing



A. HPC



B. OpenShift



A. HPC

(EBRAINS) USERS: SCIENTISTS

SHORT TERM

LONG TERM

PyUnicore
(Jupyter Notebook)

Dashboard

UNICORE

Toil

HPC - 1

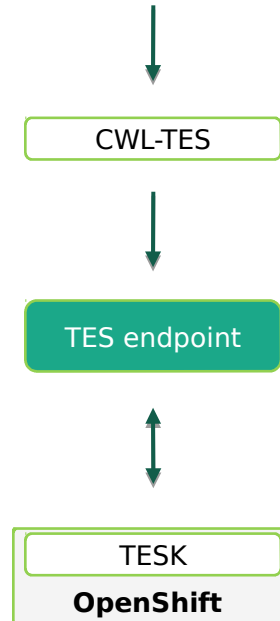
UNICORE

HPC - 1



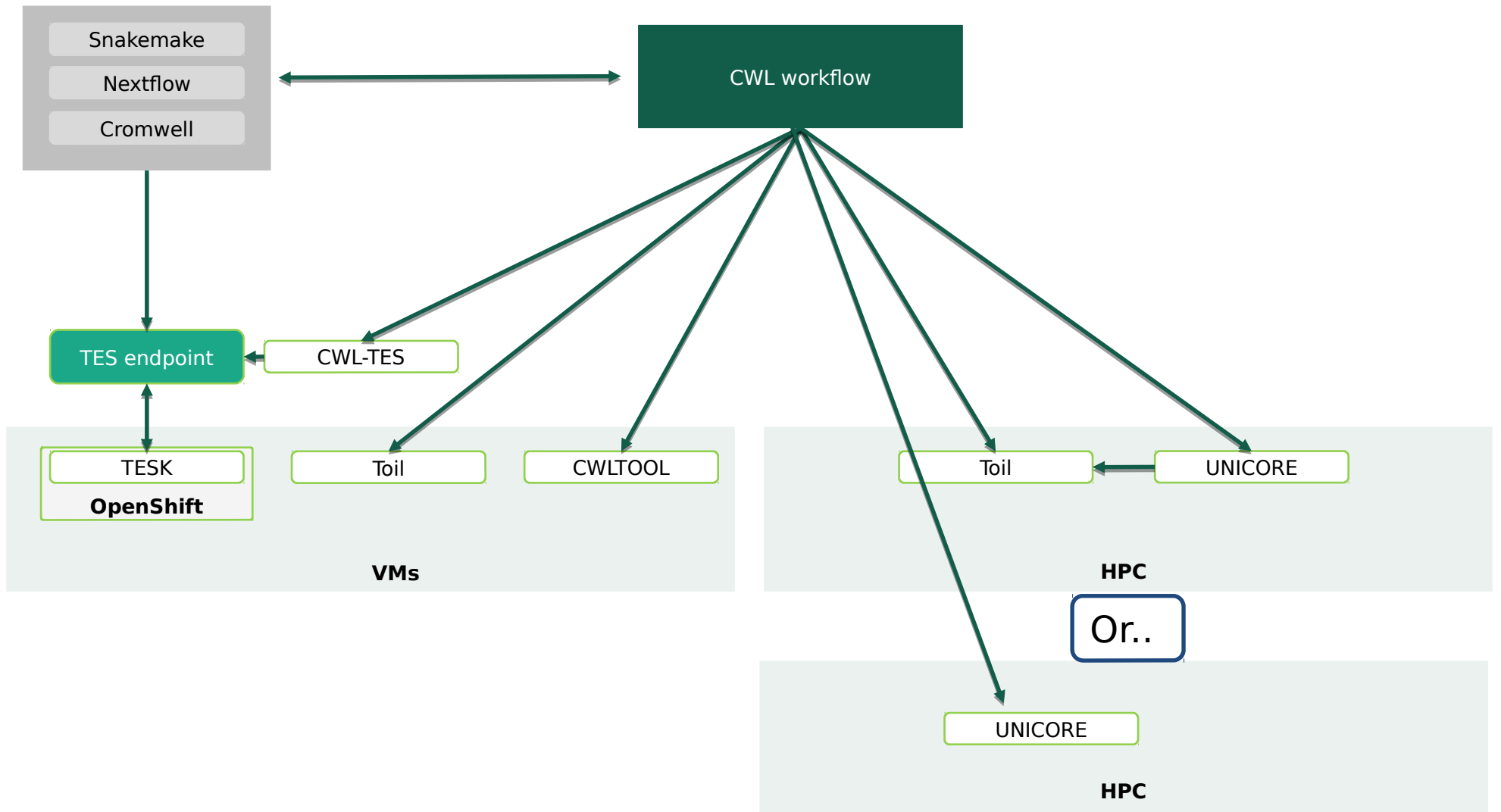
B. OpenShift

(EBRAINS) USERS: SCIENTISTS



In the making...

Another view..



(Standard) Computational Workflows



Definition



Command Line
tools



Execution



Monitoring



Dive in monitoring

A. Dashboard from other RIs



Dive in monitoring

The screenshot displays the SCHEMA SCHEDULER interface. At the top, there is a navigation bar with tabs for Software, Workflows (selected), Data, Job history, Help, Admin Options, and Logout (vergoulis). Below the navigation bar, there is a '+ New workflow' button and a 'Resources from project:' dropdown menu showing 'bufet-diploma-thesis (9776 remaining jobs)' with a '+ New project' button.

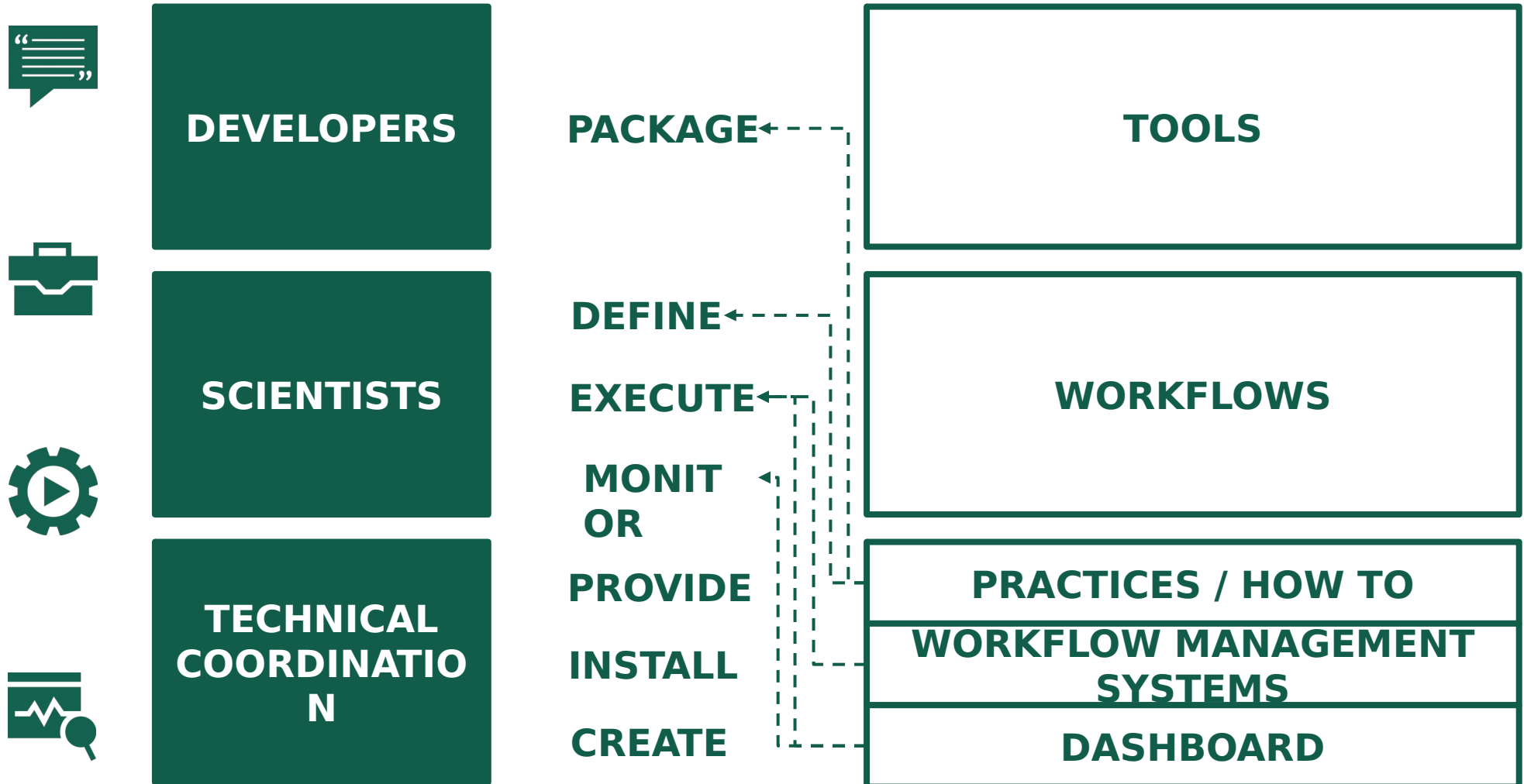
The main content area shows a table of workflows:

Software Name	Version	Uploader	Run	Edit	Delete
GWAS-workflow ?	1.0	kostis_zagganas	▶ Run	✎ Edit	✕ Delete
Hashsplitter ?	1.0	kostis_zagganas	▶ Run	✎ Edit	✕ Delete
MSA ?			▶ Run	✎ Edit	✕ Delete
RNA-seq-workflow ?			▶ Run	✎ Edit	✕ Delete

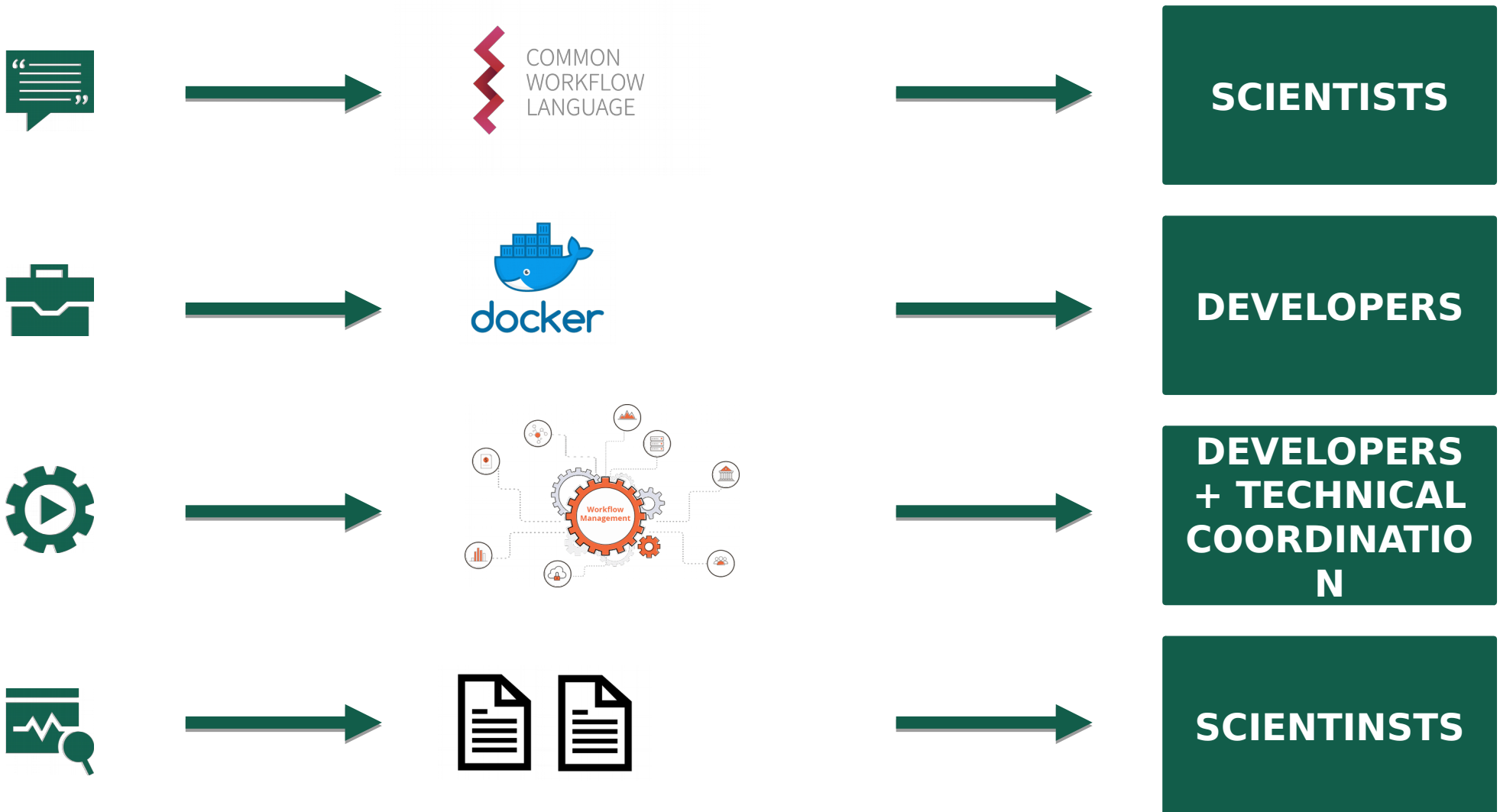
A modal dialog box titled 'Run workflow (GWAS-workflow v.1.0)' is open, showing the 'Output directory' field with 'Select' and 'Clear' buttons, and 'Arguments' for 'metadata' and 'variants' with 'Run' and 'Run example' buttons. A blue arrow points from the 'Run' button in the table to the dialog box.



Our Vision



Our vision(short term)



Our vision (long term)

WorkflowHub

Workflows

Query	
Search here...	Go
Created At	
Any time	▼
Workflow Type	
Common Workflow Language	x
Galaxy	72
Jupyter	14
Nextflow	14
Shell Script	9
Snakemake	3
More...	
Tag	
CWL	9
covid-19	5
INDELS	4
SNPs	4
Alignment	2
Amplicon	2
More...	
Submitter	
Ambarish Kumar	4
Bart Nijse	4
Laura Rodriguez-Navas	3
Jasper Koehorst	3
Irene Sánchez	3
Douglas Lewis	2


25 Workflows matching the given criteria: [\(Clear all filters\)](#)

Workflow type: Common Workflow Language x

⌵ Creation date (Ascending) ▼

← Previous **1** 2 Next →

Default Condensed Table


▲  **Virus genome assembly with Unicycler and Spades.**


Stable Download


Virus genome assembly with Unicycler and Spades. The 2 assemblers works in parallel. The graph visualization is made with Bandage. workflow git repository : <https://github.com/fjmoreews/cwl-workflow-SARS-CoV-2/blob/master/Assembly/workflow/assembly-wf-virus.cwl> Based on https://github.com/galaxyproject/SARS-CoV-2/blob/master/genomics/2-Assembly/as_wf.png

Type: Common Workflow Language
Creators: None
Submitter: francois moreews

Created: 10th Apr 2020 at 11:45, Last updated: 30th Jun 2020 at 09:05

▼  **VIRify**


▼  **var-PE**

▼  **Protein MD Setup tutorial using BioExcel Building Blocks (biobb)**

[1] <https://workflowhub.eu/>

Our vision (long term)

A Workflows and Tools Hub

 Search (e.g. brain or neuroscience)



SEARCH

CATEGORIES

- Project 121
- Dataset**
- Subject
- Sample
- Model
- Software

Tools
2106
Workflows
1245

- Homo sapiens 829
- Mus musculus 182
- Rattus norvegicus 124
- Macaca fascicularis 28
- Macaca mulatta 14
- Mustela putorius furo 4
- Chlorocebus aethiops sabaeus 2
- Mustela putorius 1
- macaca fuscata 1

DATA ACCESSIBILITY

Viewing 1-20 of 1190 results

Sort by Relevance ▾

F

A




I

R






GapMap Frontal to Temporal II (v11.3)

This dataset contains the "GapMap Frontal to Temporal II" in the individual, single subject template of the MNI Colin 27 as well as the MNI ICBM 152 2009c nonlinear asymmetric reference space. In order to provide whol...

Keywords :

-  brain mapping
-  histology
-  imaging

Methods :

-  cytoarchitectonic mapping
-  probability mapping
-  maximum probability mapping
-  magnetic resonance imaging (MRI)
-  silver staining

Our vision

Monitoring + Executing

Dashboard

In the making...



Human Brain Project



EBRAINS

Thank You!

www.humanbrainproject.eu

www.ebrains.eu