

# Accelerating Discovery: Status and Trajectory of Supporting Tools and Systems

Carole Goble

The University of Manchester, UK

European Life Science Research Infrastructures ELIXIR & IBISBA

European Open Science Cloud Life Cluster

(co-author of the infamous FAIR paper 😊)

[carole.goble@manchester.ac.uk](mailto:carole.goble@manchester.ac.uk)

# Summary of Background 20 years of scientific workflow activity

## Development

### Workflow Management Systems

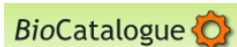


### Support for the Workflow Ecosystem



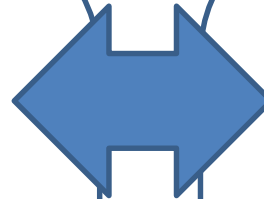
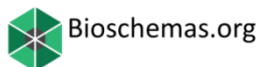
Principles of wf sharing, lifecycles, preservation, packaging

WfMS agnostic repository / registries



Registries of services

Metadata & identifier schemes for data, tools & workflows Standards



## Deployment / Adoption

### Workflows in scientific domains



biodiversity



astrophysics



life science 'omics



IBISBA<sup>1.0</sup>

sys and syn biology  
Industrial Biotech



SYNTHESYS+  
Synthesis of Systematic Resources a DISSCo project

natural history specimens

digital libraries



VPH



biomolecular modelling at exascale



social sciences

Accelerating information exchange

# Summary of Background

## Distributed EU Research Infrastructures for the Life Sciences & the European Open Science Cloud.



23 nations, 230+ orgs  
incl. EMBL-EBI



9 nations, 16 orgs



Everyone

**“Data for Life”** Data Infrastructure for Life Sciences: coordinate, integrate and sustains resources across member states; enable users to access services.

**“Together innovating for sustainable biotechnology”** Accelerate end-to-end bioprocess development, linking distributed R&D facilities: computer-assisted design of biocomponents, construction of enzymes and microbial strains, fermentation, downstream processing.

**“A trusted virtual environment”** to store, share & re-use research information to serve Europe's 1.7 million researchers and 70 million science and tech professionals and to foster interdisciplinary working.

# Example of a Galaxy pipeline. Human in the Loop!!

<https://www.ibisba.eu/>



Delepine et al. *Met. Eng.* 2018

Carbonell et al. *Bioinformatics* 2018  
Mellor et al. *ACS Synth. Biol.* 2016

Carbonell et al. *NAR* 2014

Delepine et al. *Met. Eng.* 2018

Duigou et al. *NAR.* 2018

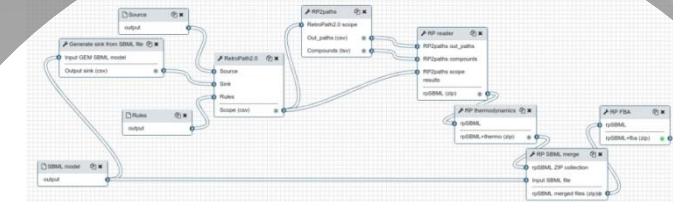
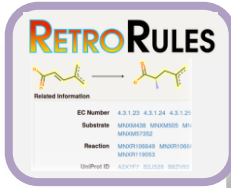
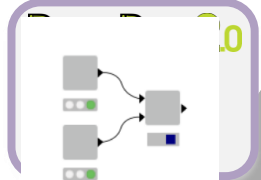
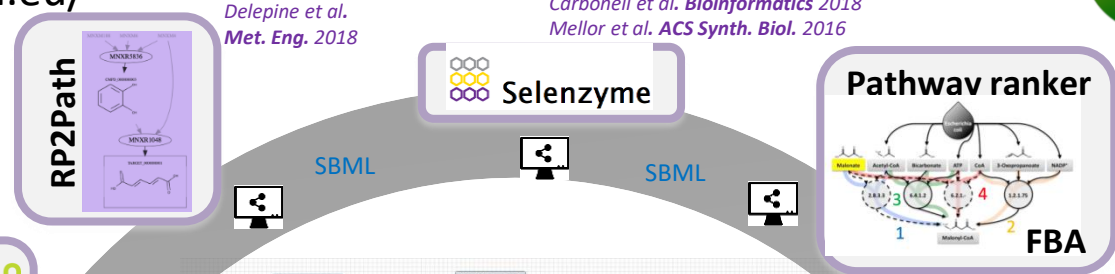
Carbonell et al. *IWBA* 2019

Swainston et al. *Bioinformatics* 2018

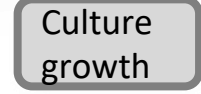
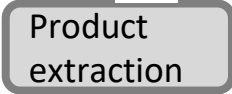
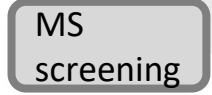
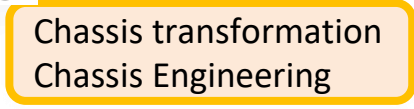
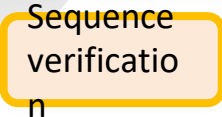
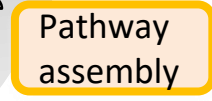
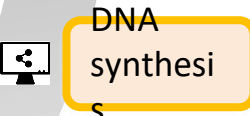
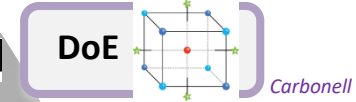
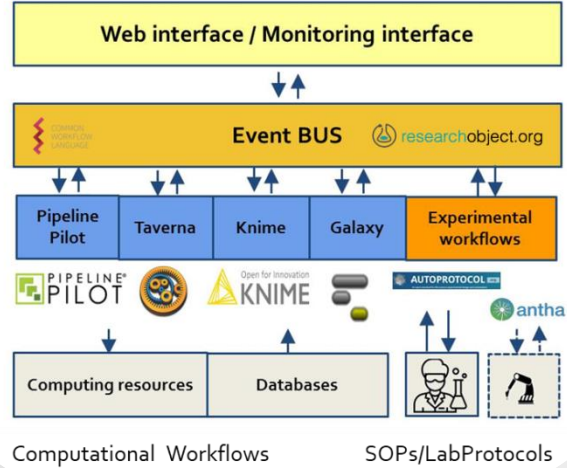
Carbonell et al. *Commun Biol.* 2018

Borkowski et al. *BioRxiv*, 2019

Koch et al. *J Cheminform.* 2017

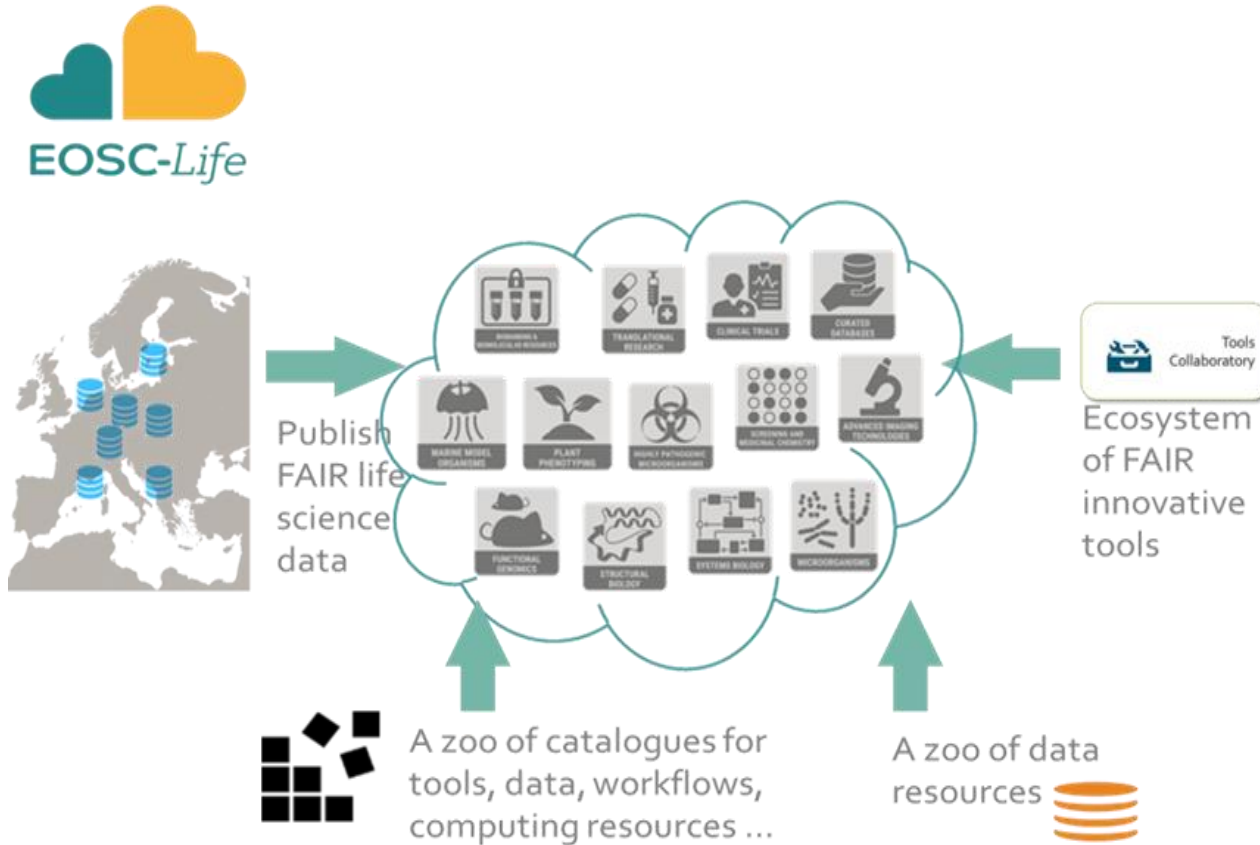


The stages may be computational or lab-based (SOPS, Lab Protocols).



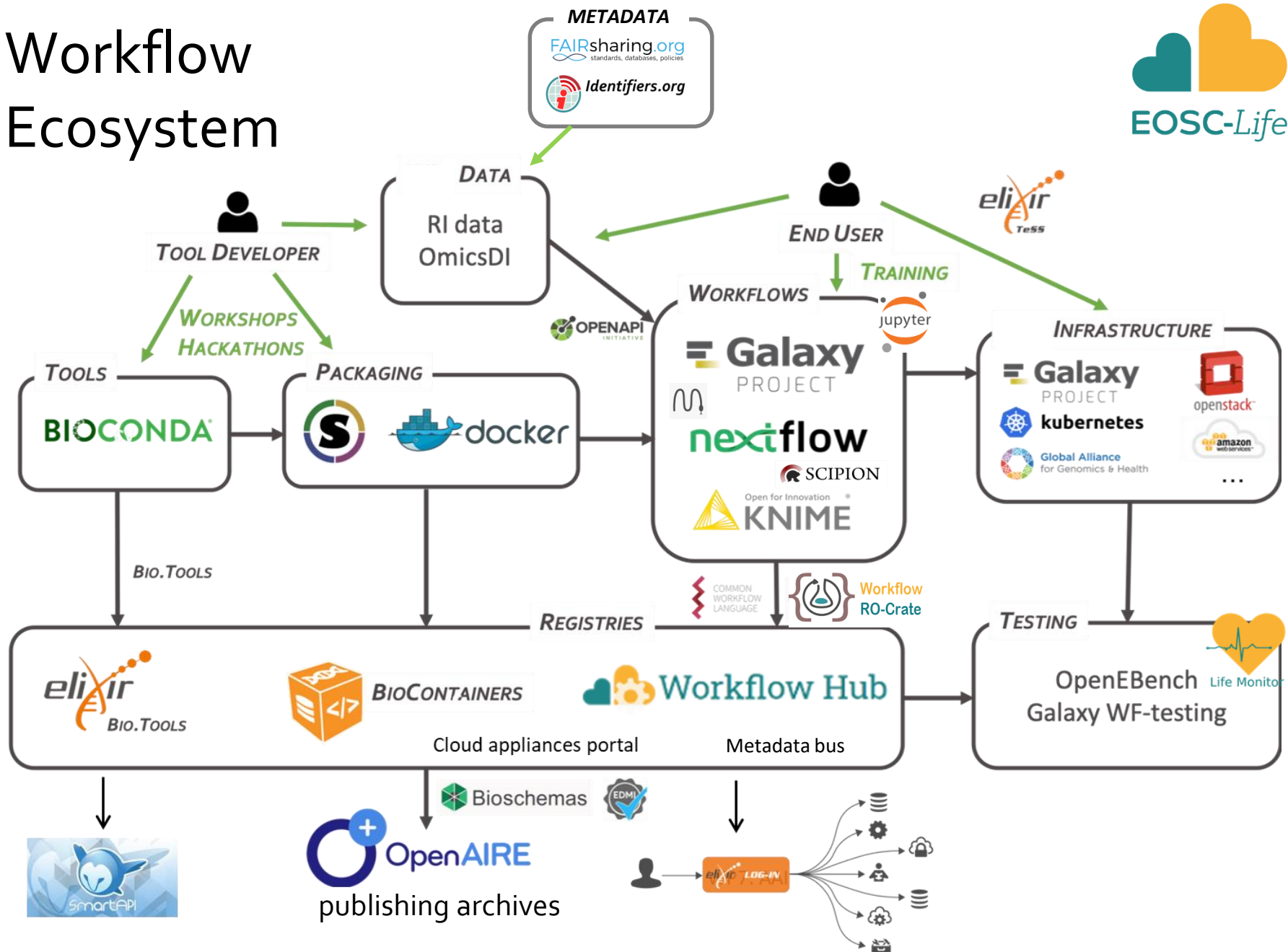
# EOSC Life FAIR Workflow Collaboratory

The Workflow Ecosystem for 13+1 European Research Infrastructures in the Life Sciences



- Democratising the use of workflows
- Accelerating the adoption, exchange and use of workflow methodology
- Accelerating the creation of workflows.
- Supporting the scientific workflow for different communities in the wild, at the reality coal face
- Building workflow capacity

# Workflow Ecosystem



Challenge: Interoperating and Sustaining

# We finally have a tech stack



Registries, Repositories,  
Portals  
Find & launch



Description for  
finding, running,  
reporting



**MIRCW**  
Minimum Information for  
Registering a  
Computational Workflow

Interoperability



Portable packaging



Algorithm Design &  
Execution using WfMS

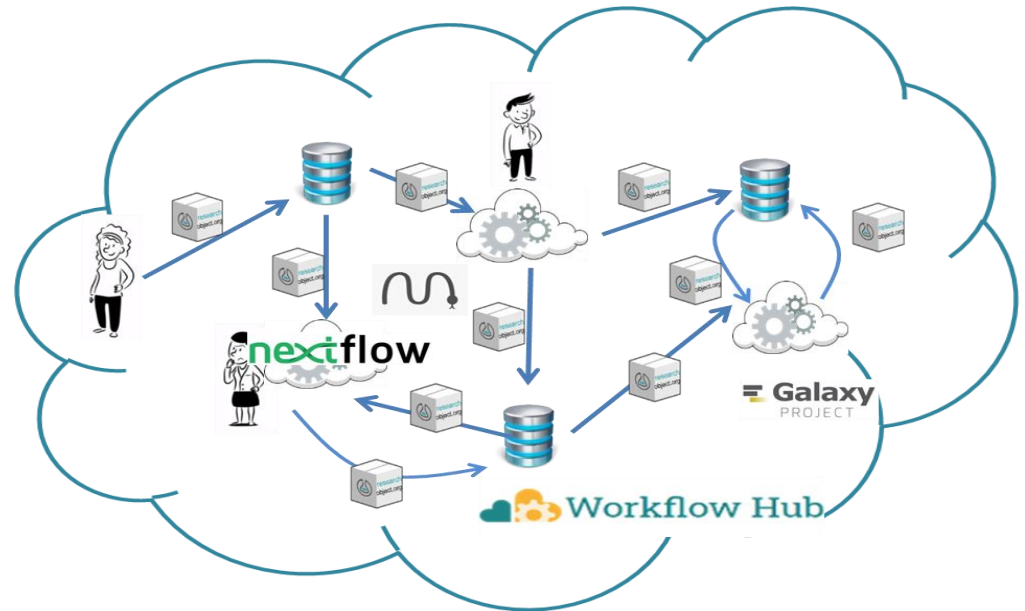


Versioning  
Development

# FAIR Commons using Research Objects as the Currency of Science, supporting interoperability, a metadata bus and Knowledge Graph of workflows and data

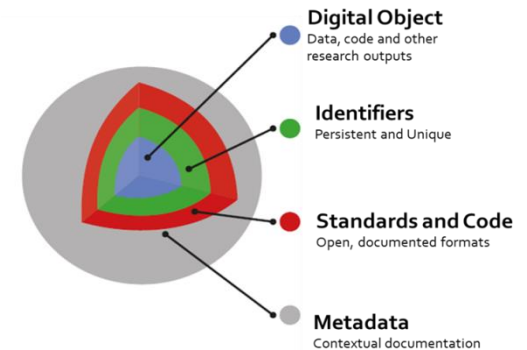


**Describes** workflows to be **portable, scalable & interoperable** with different workflow systems and containerised tools (GA4GH Genomic Data Toolkit)



**Bundles** descriptions, references, files  
**Adds** context, provenance, examples, data  
**Relates** data collections, SOPs, lab protocols  
**Links** descriptions with native workflows  
**Metadata packaging framework**

FAIR Digital Object Framework



EC "Turning FAIR into Reality, 2018"

# What workflows tools / systems have you worked on or been involved with?

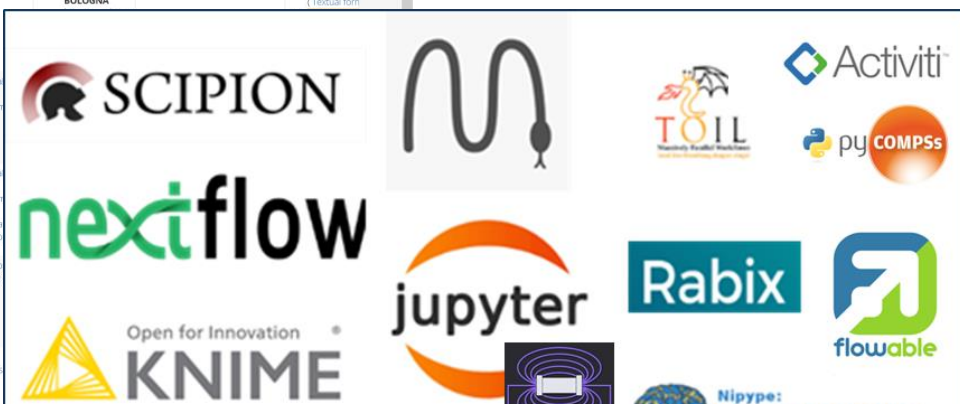
Everything: 'Protein chain' x Search bio.tools 17 tools About Menu - carolegoble -

Sort by Score Updated Added Name Citation Count Publication Date Display as Compact Detailed

Name	Description	Homepage	Version	Topic	Credits & Support	Operation	Input
WS-SNPs and GO	A web server for predicting disease associated variations from protein sequence and structure.	Link +	-	<ul style="list-style-type: none"> <li>Protein properties</li> <li>Pathology</li> </ul>	<ul style="list-style-type: none"> <li>ELIXIR-ITA-BOLOGNA</li> </ul>	<ul style="list-style-type: none"> <li>Variant classification</li> </ul>	<ul style="list-style-type: none"> <li>PDB ID (Textual form)</li> </ul>
I-MUTANT	Neural Network based Predictor of Protein stability Changes upon Single Point Mutation from the Protein Structure.	Link +	2.0	<ul style="list-style-type: none"> <li>Protein folding, sta and design</li> <li>DNA polymorphosm</li> </ul>			
I-MUTANT Suite	Predictor of effects of single point protein mutation on its stability from protein sequences or structures.	Link +	3.0	<ul style="list-style-type: none"> <li>Protein folding, sta and design</li> <li>DNA polymorphosm</li> <li>Genetic variation</li> <li>Protein structure a</li> <li>Structure predictio</li> <li>Proteins</li> <li>Structure predictio</li> <li>Protein secondary structure</li> </ul>			
DisLocate	Prediction of cysteine connectivity patterns in a protein chain.	Link +	1.0	<ul style="list-style-type: none"> <li>Protein folding, sta and design</li> <li>Structural domains</li> <li>Bioinformatics</li> <li>Mathematics</li> </ul>			
Knot-ID	Tool to study the effects of protein chains using the concept of	Link +	-	<ul style="list-style-type: none"> <li>Protein folds and structural domains</li> <li>Bioinformatics</li> <li>Mathematics</li> </ul>			

17169 tools

ED AM Ontology



COMMON WORKFLOW LANGUAGE

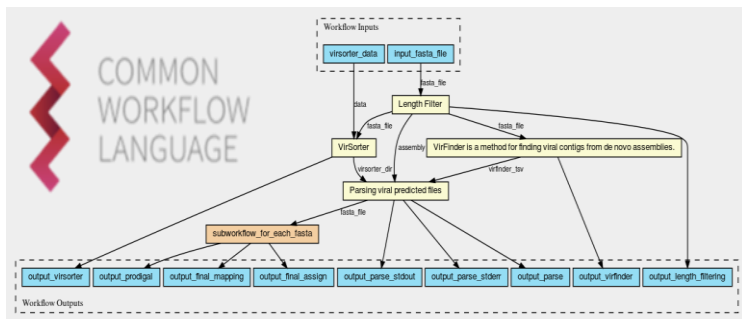
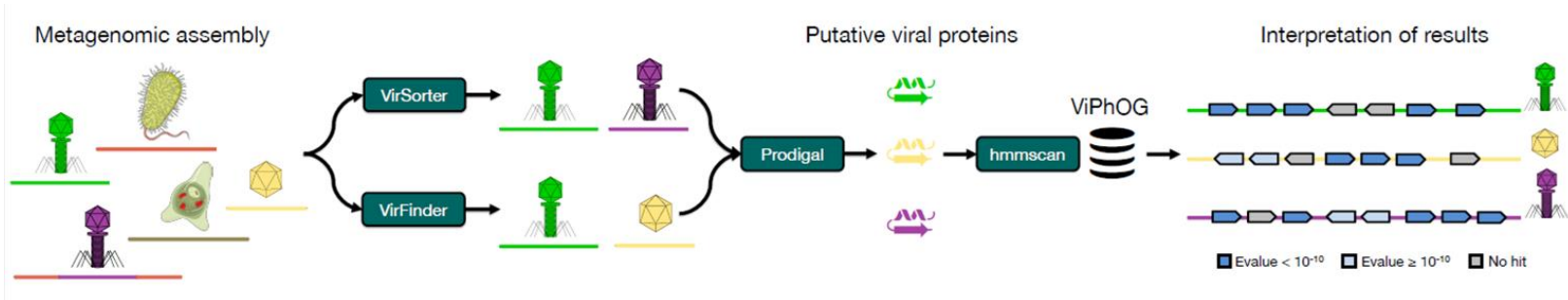



Registries & Repositories

A Zoo!

And don't forget the different data formats, identifiers, repositories & catalogues!

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



<https://github.com/EBI-Metagenomics/emg-virify-pipeline/blob/master/CWL/WorkFlow/pipeline.cwl>



CWL Execution Engine  
CWLTool, CWL Exec

Metagenomic Assemblies  
Thanks to Rob Finn, Mgnify, EBI

Putative viral proteins

Interpretation of results

```

516 lines (425 sloc) | 17.6 KB
1 #!/usr/bin/env nextflow
2 nextflow.preview.dsl-2
3
4 /*
5  * nextflow -- Virus Analysis Pipeline
6  * Author: hoelzer.marting@gmail.com
7  */
8
9 /-----/
10 * Help messages & user inputs & checks
11 /-----/
12
13 /* Comment section
14 First part is a terminal print for additional user information, like
15 second part is file channel input. This allows via --list to all the
16 add csv instead: name,path or name,path01,path02 in case of illumina
17 */
18
19 // terminal prints
20 println ""
21 println "\u001B[32mProfile: $workflow.profile\u001B[0m"
22 println ""
23 println "\u001B[2mCurrent user: $workflow.userName"
24 println "nextflow version: $nextflow.version"
25 println "Starting time: $nextflow.timestamp"
26 println "workdir: $location"
27 println " "
28 println " $workflow.workdir/\u001B[0m"
29
30 if [ $workflow.profile == "standard" ] {
31     println "\u001B[2mCPU to use: $params.cores"
32     println "Output dir name: $params.output\u001B[0m"
33     println " "
34 }
35 println "\u001B[2mDev VIPHOG database: $params.version\u001B[0m"
36 println " "
37
38 if [ !nextflow.version.matches("20.01.*") ] {
39     println "This workflow requires Nextflow version 20.01 or greater -- You are running version $nextflow.version"
40     exit 1
41 }
42
43 if [ params.help ] { exit 0; helpMSG() }
44 if [ params.profile ] {
45     exit 1; "--profile is wrong use -profile"
46 }
47 if [ params.illumina == "" && params.fasta == "" ] {

```



<https://github.com/hoelzer/virify/blob/master/virify.nf>

Designing and building  
reusable libraries



**12,000+**  
users

**2,000+**  
scientific tools

**6+ Million**  
analysis jobs executed

**13+ Million**  
datasets uploaded

**~7 TB**  
of reference data

Integrated with  
ELIXIR AAI, RStudio,  
Jupyter Notebooks

Disruptor: Yes – workflow environments ARE used.  
Sustained, Human in the Loop, Scale out, Professionalisation, RWEs

# If you make workflow environments usable and professionally curated and managed then they are used and make a difference



bioRxiv  
THE PREPRINT SERVER FOR BIOLOGY

HOME | ABOUT | SUBMIT | ALERTS / RSS | CHANNELS

Search    
Advanced Search

bioRxiv is receiving many new papers on coronavirus 2019-nCoV. A reminder: these are preliminary reports that have not been peer-reviewed. They should not be regarded as conclusive, guide clinical practice/health-related behavior, or be reported in news media as established information.

New Results

[Comment on this paper](#)

[Previous](#)

[Next](#)

## No more business as usual: agile and effective responses to emerging pathogen threats require open data and open analytics

Galaxy and HyPhy developments teams, Anton Nekrutenko, Sergei L. Kosakovsky Pond

doi: <https://doi.org/10.1101/2020.02.21.959973>

This article is a preprint and has not been certified by peer review [what does this mean?].

[Abstract](#) [Full Text](#) [Info/History](#) [Metrics](#) 

### Abstract

The current state of much of the Wuhan pneumonia virus (COVID-19) research shows a regrettable lack of data sharing and considerable analytical obfuscation. This impedes research cooperation, which is essential for tackling public health emergencies, and

Posted February 25, 2020.

## Preprocessing of raw SARS-CoV-2 reads

The raw reads available so far are generated from bronchoalveolar lavage fluid (BALF) and are metagenomic in nature: they contain human reads, reads from potential bacterial co-infections as well as true COVID-19 reads.

### Live Resources

usegalaxy.org	usegalaxy.eu	usegalaxy.org.au	usegalaxy.be
<a href="#">workflow run</a>	<a href="#">workflow run</a>	<a href="#">workflow run</a>	<a href="#">workflow run</a>
<a href="#">history view</a>	<a href="#">history view</a>	<a href="#">history view</a>	<a href="#">history view</a>

### What's the point?

Assess quality of reads, remove adapters and remove reads mapping to human genome.

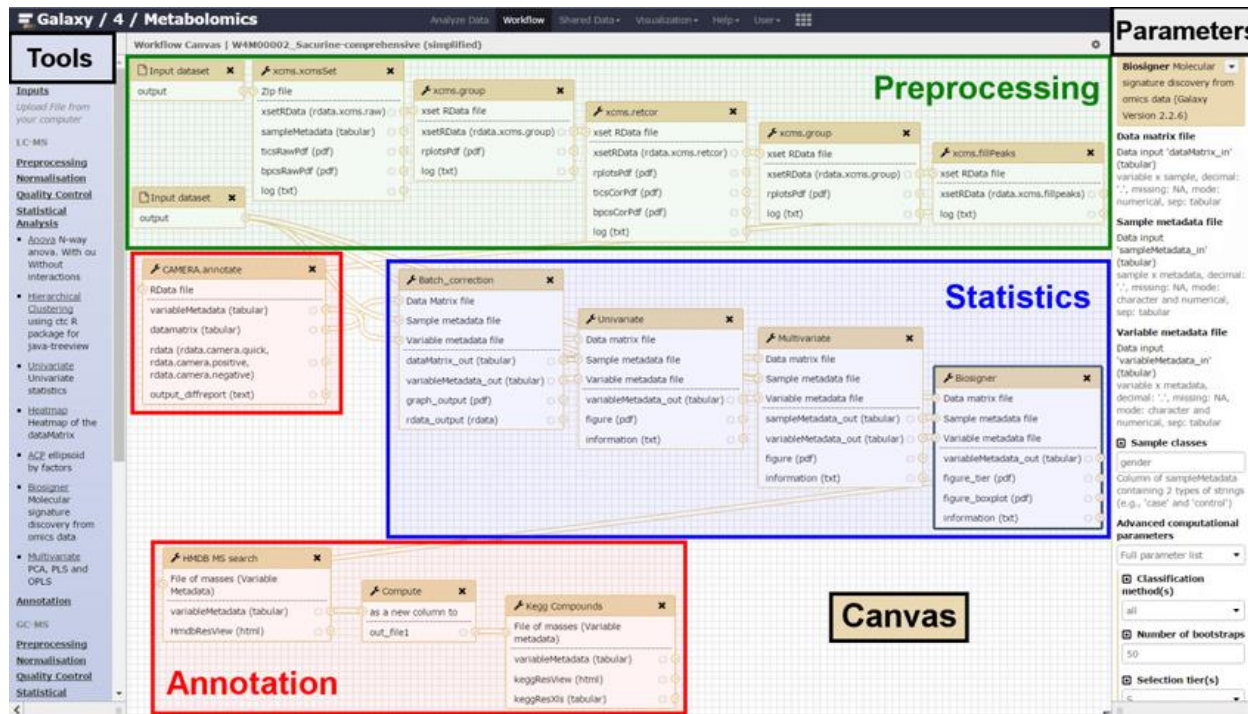
### The outline

Illumina and Oxford nanopore reads are pulled from the NCBI SRA (links to SRA accessions are available [here](#)). They are then processed separately as described in the [workflow section](#).

<https://github.com/galaxyproject/SARS-CoV-2/>

# Clustering around platforms for economies of scale, skill and sustainability

## Workflow collections. Workflow as a Service



The screenshot displays the Galaxy 4 Metabolomics workflow canvas. The interface is divided into several sections:

- Tools:** A sidebar on the left lists various tools categorized into Inputs, Preprocessing, Normalisation, Quality Control, Statistical Analysis, Annotation, and GC-MS.
- Workflow Canvas:** The main area shows a sequence of tools connected by arrows. The workflow is organized into three main stages:
  - Preprocessing:** Includes tools like 'xcms.xcmsSet', 'xcms.group', 'xcms.retcor', and 'xcms.filterpeaks'.
  - Statistics:** Includes tools like 'CAMERA.annotate', 'Batch\_correction', 'Univariate', 'Multivariate', and 'Biosigner'.
  - Annotation:** Includes tools like 'HMDB MS search', 'Compute', and 'Keggs Compounds'.
- Parameters:** A panel on the right shows the configuration for the 'Biosigner' tool, including options for 'Data matrix file', 'Sample metadata file', and 'Variable metadata file'.

- [rna.usegalaxy.eu](http://rna.usegalaxy.eu)
- [clipseq.usegalaxy.eu](http://clipseq.usegalaxy.eu)
- [metagenomics.usegalaxy.eu](http://metagenomics.usegalaxy.eu)
- [hicexplorer.usegalaxy.eu](http://hicexplorer.usegalaxy.eu)
- [cheminformatics.usegalaxy.eu](http://cheminformatics.usegalaxy.eu)
- [proteomics.usegalaxy.eu](http://proteomics.usegalaxy.eu)
- [Imaging.usegalaxy.eu](http://Imaging.usegalaxy.eu)
- [metabolomics.usegalaxy.eu](http://metabolomics.usegalaxy.eu)
- [ecology.usegalaxy.eu](http://ecology.usegalaxy.eu)
- [nanopore.usegalaxy.eu](http://nanopore.usegalaxy.eu)
- [singlecellomics.usegalaxy.eu](http://singlecellomics.usegalaxy.eu)
- [humancellatlas.usegalaxy.eu](http://humancellatlas.usegalaxy.eu)
- [streetscience.usegalaxy.eu](http://streetscience.usegalaxy.eu)

Galaxy interactive workflow collection for Metabolomics

# Workflow making is still largely artisanal

Producing the **FAIR ecosystem** so we can semi-automatically build workflows that can help automate beyond “heroic” one-offs

- AI’s role in building/running the ecosystem
- Recommend workflows, matchmaking....
- Automated discovery of workflows / data
- The capacity and capability of researchers to use and reuse workflows
- Professionalisation of workflow making/mgt

Example: Galaxy Workflow Platform

- Galaxy-ML (<https://github.com/goeckslab/Galaxy-ML>)
- Tool recommender system in Galaxy using deep learning (<https://doi.org/10.1101/838599>)

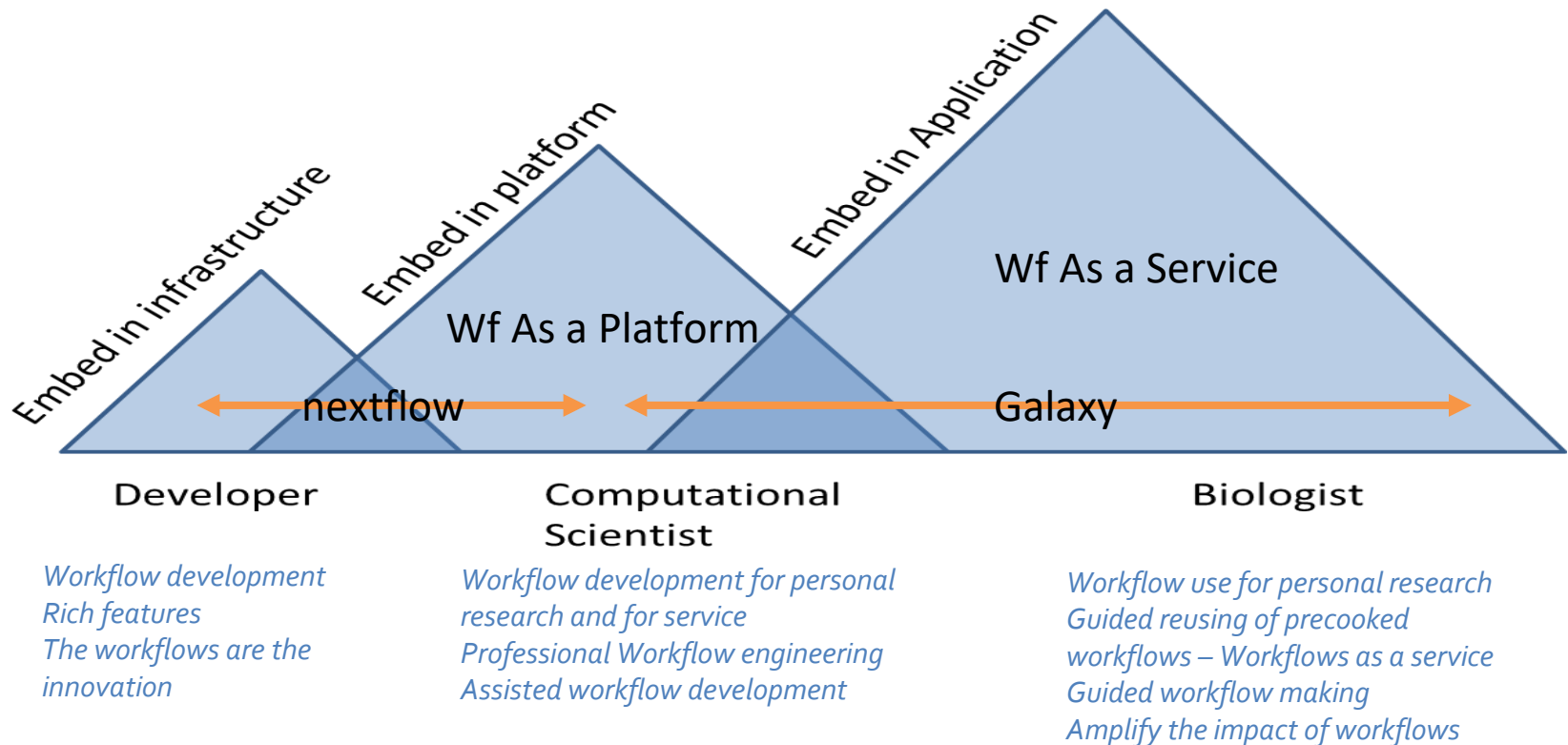


- AI steps in the workflows
- AI enabled wf design
- AI enabled wf set up
- AI enabled wf execution
- AU enabled wf analytics
- AI enabled wf workflow ecosystem

# FAIR Workflows

# Who are the Workflow Users?

Capacity Building, Professionalisation, Economics



different systems for different users, spanning the spectrum, enabling researchers to across the spectrum

# Trajectories & Challenges: FAIR Workflows

- **A slew of workflow systems**
  - Domain specific languages are on the rise, Script-like workflows & problem-specific ad-hoc data models + Adoption of industry standards.
  - Take up dependent on the “plugged-in” availability of data type specific codes, optimised processing and everyone else using it
- **At the same time take-up of professionally managed cloud platforms**
  - Professionalisation of Workflow Engineering
- **Catalogues of workflow components (and data and SOPs)**
  - Designed to work together, Quality Assured, curated and maintained
  - tools / sub workflows, and the data they operate over accessible and licensed.
- **Interoperability: neat and tidy interfaces**
  - Shared (standardised) data formats, identifier schemes and interfaces for tool compatibility, incl. ids for workflows
- **Just enough metadata on everything**
  - Semantic representations of Workflows, Data and Tasks: schema.org, JSON-LD
- **Portable / accessible execution – containers**
  - tools can execute next to each other independent of WfMS or cluster installation

# Acknowledgements

Ian Cottam  
Frederik Coppens  
Björn Grüning  
Mark Wilkinson  
Stian Soiland-Reyes  
Nick Juty  
Rob Finn  
Michel Dumontier  
Michael Crusoe  
Finn Bacall  
Alan Williams  
Katy Wolstencroft  
Pinar Alper  
Sarah Cohen-Boulakia  
Stuart Owen  
Jiten Bhagat  
Don Cruickshank  
David De Roure  
Niall Beard  
Aleks Nenandic  
Paolo Missier  
Marco Roos  
Khalid Belhajjame  
Sarah Cohen-Boulakia  
Yolanda Gil  
Ilkay Altintas  
Daniel Garijo  
Melchior Du-Lac

And all the projects and many more

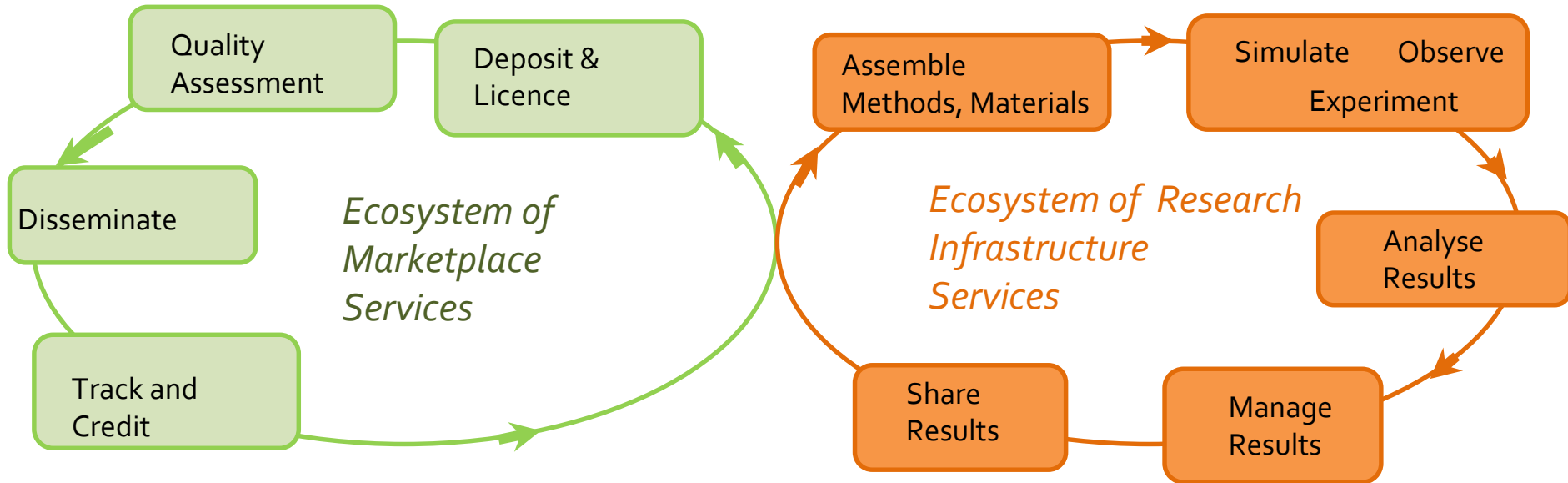


# BONUS SLIDES

# What are we accelerating?

- Accelerating workflow creation
- Accelerating workflow use and reuse, sharing and exchange
- Accelerating the workflows themselves

# Building a FAIR Research Commons



Fixing this to be an 'in flight" feedback loop

# For the Panel

- What are your observations on the impact of systems and tools to automate workflows, with respect to scientific discovery?
- What is the role of artificial intelligence (AI) and automation in scientific discovery? We'd like you to consider both AI-driven workflows as well as workflows that include AI as part of the scientific discovery process.
- What are the current challenges faced by these systems?
- What do you see as disruptors that will affect how workflow tools are being used by the scientific community?
- What is your vision for AI in workflows in the future of scientific discovery? What research is needed to advance this area?