

COVID-19 PubSeq: Public SARS-CoV-2 Sequence Resource

Andrea Guarracino⁵, Peter Amstutz², Thomas Liener³, Michael R. Crusoe⁴, Adam Novak⁶, Erik Garrison⁶, Tazro Ohta⁷, Bonface Munyoki¹, Danielle Welter⁸, Sarah Zaranek², Alexander (Sasha) Wait Zaranek², Pjotr Prins¹

¹Department of Genetics, Genomics and Informatics, The University of Tennessee Health Science Center, Memphis, TN, USA., ²Curii Corporation Boston, MA, ³Independent, ⁴Department of Computer Science, Faculty of Sciences, Vrije Universiteit Amsterdam, The Netherlands, ⁵Centre for Molecular Bioinformatics, Department of Biology, University Of Rome Tor Vergata, Rome, Italy, ⁶UC Santa Cruz Genomics Institute, University of California, Santa Cruz, CA, USA, ⁷Database Center for Life Sciences, Tokyo, Japan, ⁸Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Luxembourg

As part of the COVID-19 Virtual Biohackathon 2020 we formed a working group to create COVID-19 PubSeq, a Public Sequence Resource for SARS-CoV-2 virus sequences. Our goal was to create a repository that had a **low barrier to entry** for uploading and analyzing sequence data, without imposing **any restriction** on their utilization. We followed FAIR data principles: data are published with public domain (CC0) or creative commons 4.0 (CC-BY-4.0) license, structured metadata is validated against standard ontologies, and, importantly, reproducible workflows are executed after the upload in order to provide up-to-date results rapidly and in standardized data formats.

Data and analysis tools together

Existing data repositories don't enforce strict quality control on the submitted **data** and its **metadata**, and don't add value in terms of running **additional analysis**. In addition, some databases have licenses that place restrictions on the data utilization.



- Data is validated for being in a supported data format, and for **not being duplicated** entries in the resource.
- Structured metadata is strictly validated against **standard ontologies**.

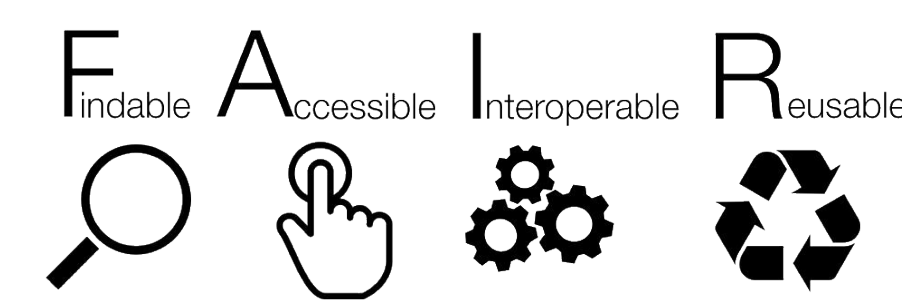
Reproducible and **scalable** CWL workflows are executed on the cloud **after the upload** in order to provide **up-to-date results** rapidly and in **standardized** data formats.

<http://covid19.genenetwork.org>

On COVID-19 PubSeq the data, metadata, and analysis tools live together, **publicly** and **freerly**.

State-of-the-art standards

COVID-19 PubSeq leverages state-of-the-art standards and technologies.



[The FAIR Guiding Principles¹](https://www.fair4life.org/)



<http://www.ontobee.org>



<http://commonwl.org>



<https://creativecommons.org>



[Ontology Lookup Service](https://www.ebi.ac.uk/ols/)



<https://github.com>

COVID-19 PubSeq is running on Arvados, a cloud open source platform for managing, processing, and sharing scientific and biomedical data.

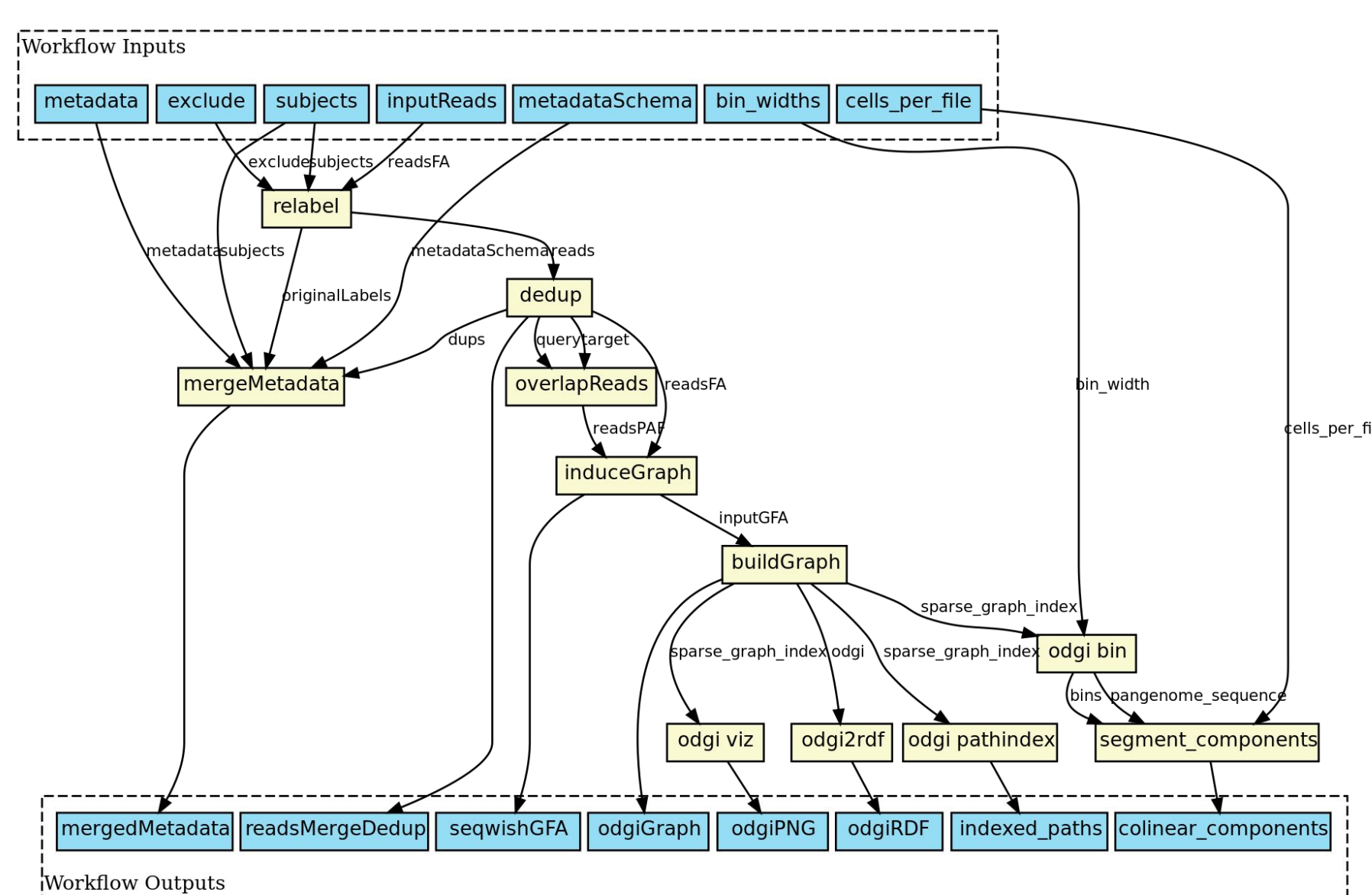


<http://arvados.org>

Pangenome generation workflow

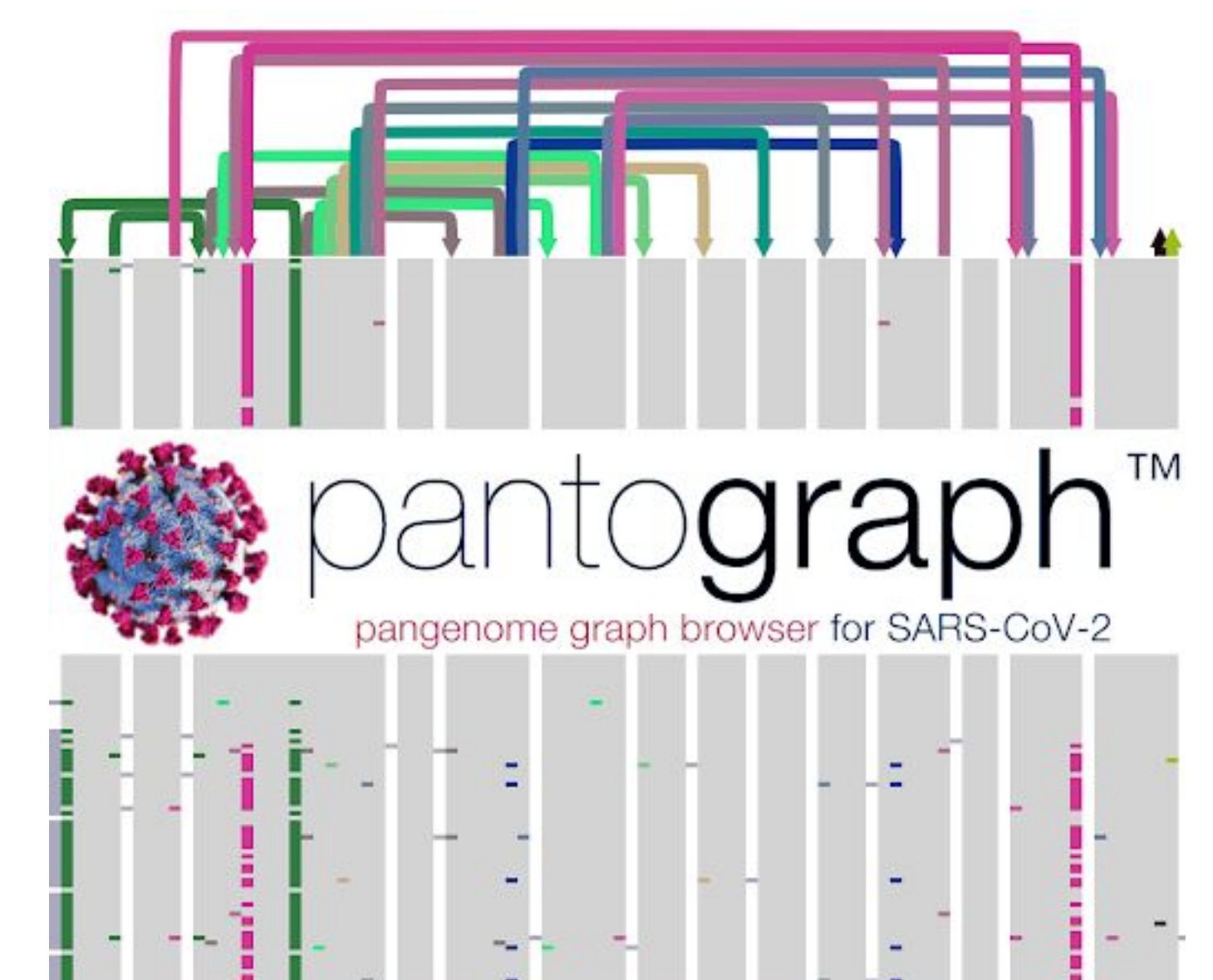
Each time someone uploads a valid sample, it is immediately combined with all the already uploaded SARS-CoV-2 genomes in order to generate an up-to-date SARS-CoV-2 **pangenome** as input for **Pantograph**, an interactive visualization of pangenomes.

pangenome-generate.cwl

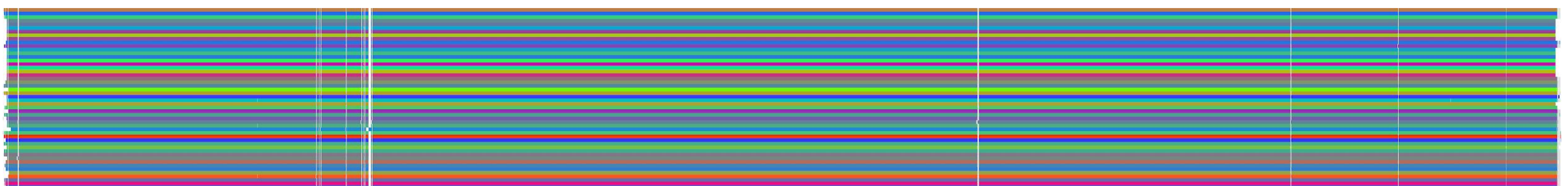


- A graphical pangenome² models the **full set of genomic elements** in a given species or clade.
- Pangenomes naturally express **genome rearrangements**, therefore Pantograph allows researchers to browse the entire genetic diversity in a SARS-CoV-2 population that would otherwise be underestimated.

<https://graph-genome.github.io>



Graphical representation of a SARS-CoV-2 pangenome of 100 genomes realized with [vgteam/odgi](https://github.com/vgteam/odgi).



COVID-19 Integrated knowledge base

Swiss Institute of Bioinformatics

Metadata can be downloaded as **Turtle RDF** which can be loaded into any RDF triple-store. The **Swiss Institute of Bioinformatics** has included the COVID-19 PubSeq data in <https://covid-19-sparql.expasy.org>.

Information on a specific sample

```
PREFIX pubseq: <http://biohackathon.org/bh20-seq-schema#MainSchema/>
PREFIX sio: <http://semanticscience.org/resource/>
select distinct ?predicate ?object
{
  ?sample sio:SIO_000115 "MT326090.1" .
  ?sample ?predicate ?object .
}
```

References

1. Wilkinson et al. (2016). The FAIR guiding principles for scientific data management and stewardship. *Scientific data*, 3, 160018.
2. Eizenga et al. (2020). Pangenome graphs. *Annual Reviews of Genomics and Human Genetics*, 21.

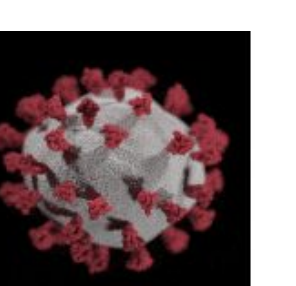
Acknowledgments



COVID-19 PubSeq

COVID-19 PubSeq [home page](#) showing the number of submitted sequences and where they come from.

COVID-19 PubSeq: Public SARS-CoV-2 Sequence Resource
Database contains 6753 public sequences!



PUBSEQ DOWNLOAD UPLOAD STATUS DEMO EXPORT BLOG ABOUT

Make your sequence data **FAIR**. Upload your SARS-CoV-2 sequence (FASTA or FASTQ formats) with metadata (JSONLD) to the **public sequence resource**. The upload will trigger a recompute with all available sequences into a Pangenome available for **download!**

Your uploaded sequence will automatically be processed and incorporated into the public pangenome with metadata using workflows from the High Performance Open Biology Lab defined [here](#). All data is published under a **Creative Commons license**. You can take the published (GFA/RDF/FASTA) data and store it in a triple store for further processing. Clinical data can be stored securely at **REDCap**.

Note that form fields contain web **ontology URI's** for **disambiguation** and machine readable metadata. For examples of use, see the **BLOG**.



Leaflet | © OpenStreetMap contributors, CC BY, Imagery © Mapbox