# ELIXIR Belgium
# Building an ecosystem of services in Life Sciences

Frederik Coppens

# The Belgian Node : a short history

# Towards a Belgian Node

- December 2013 - Memorandum of Understanding, Belgium joins ELIXIR

- November 2015 - ELIXIR Consortium Agreement, full ELIXIR Member

- December 2017 - ELIXIR Collaboration Agreement, VIB lead, 7 academic partners
  - VIB legal entity for ELIXIR Belgium
  - All Flemish & Brussels universities on board
  - 1 Walloon university

- December 2018 - Sciensano becomes partner

# A people infrastructure

# ELIXIR Belgium funding

- Belgian Science policy office covers the annual fee
- 2016 – funding to hire a coordinator from Flemish government ((d)HoN in kind by VIB)
  - 1 FTE
- 2017-2018 – infrastructural funding from Flemish government
  - Development of new programme
  - Bridge-the-gap-funding
  - ~5 FTE / year
- 2019-2022 – International Research Infrastructure (IRI) call by FWO
  - Dedicated call for ESFRIs
  - 2-yearly call, funding for 4 years, when funded skip a call
  - ~ 15 FTE / year
  - ~50/50 central node / partners

elixir
BELGIUM

# ELIXIR Belgium funding

- 2023-2026 – renewal of IRI funding
    - Evaluation ongoing
    - ~ 20 FTE / year requested
- 2022-2026 – Belgian Genome Biobank project
    - Lead technical infrastructure – mostly in kind
    - ¼ FTE / year

# Embedded in the (inter)national ecosystem

# International projects

- ELIXIR-Excelerate (done)
- ELIXIR-CONVERGE
  - WP lead RDMkit
- EOSC-Life
  - Task lead tools collaboratory & co-lead WorklfowHub
- BY-COVID
  - WP lead Infectious Diseases Toolkit & (federated) data analysis
- EuroScienceGateway
  - Task leads
- Genomic Data Infrastructure
  - Task lead
  - Lead for technical infrastructure in Belgium
- AgroServ
  - Co-lead of WP on data management

# Lots of different roles & mandates

- Chair of the Belgian 1+ Million Genomes mirror group ICT

- Member of the 1+ Million Genomes working group ICT

- Member of the Belgian 1+ Million Genomes Steering Committee


- Member of the Flemish Supercomputer Center (VSC) User Council

- Member of the Flemish Open Science Board (FOSB), representing ESFRIs


- Belgian representative in the Strategy Working Group for Data, Computing and Digital Research Infrastructures (SWG DIGIT) of the ESFRI Forum

- Member of the EOSC TF Technical interoperability of Data and Services


- Member of the Galaxy Executive Board

- Co-chair of ELIXIR Galaxy Community

# Delivering Services for Flemish Researchers

# ELIXIR Belgium services

# Providing free-to-use Services-on-Top



Services

Hardware infrastructure

# Providing infrastructure services

# Enabling FAIR data by design

# Scalable data management services

**How can we help researchers, data stewards and project managers navigate and contribute to this FAIR data repository landscape?**

# RDMkit https://rdmkit.elixir-europe.org

RDM support throughout the entire **life cycle of projects** as outlined in **DMPs**

**Online focal point** for guidance, information, best practice, examples

Context and **signpost** for FAIR data resources as a Hub for a RDM Knowledge Commons

# A Showcase for ELIXIR's FAIR Services

# DataHub : collection of metadata throughout the data life cycle
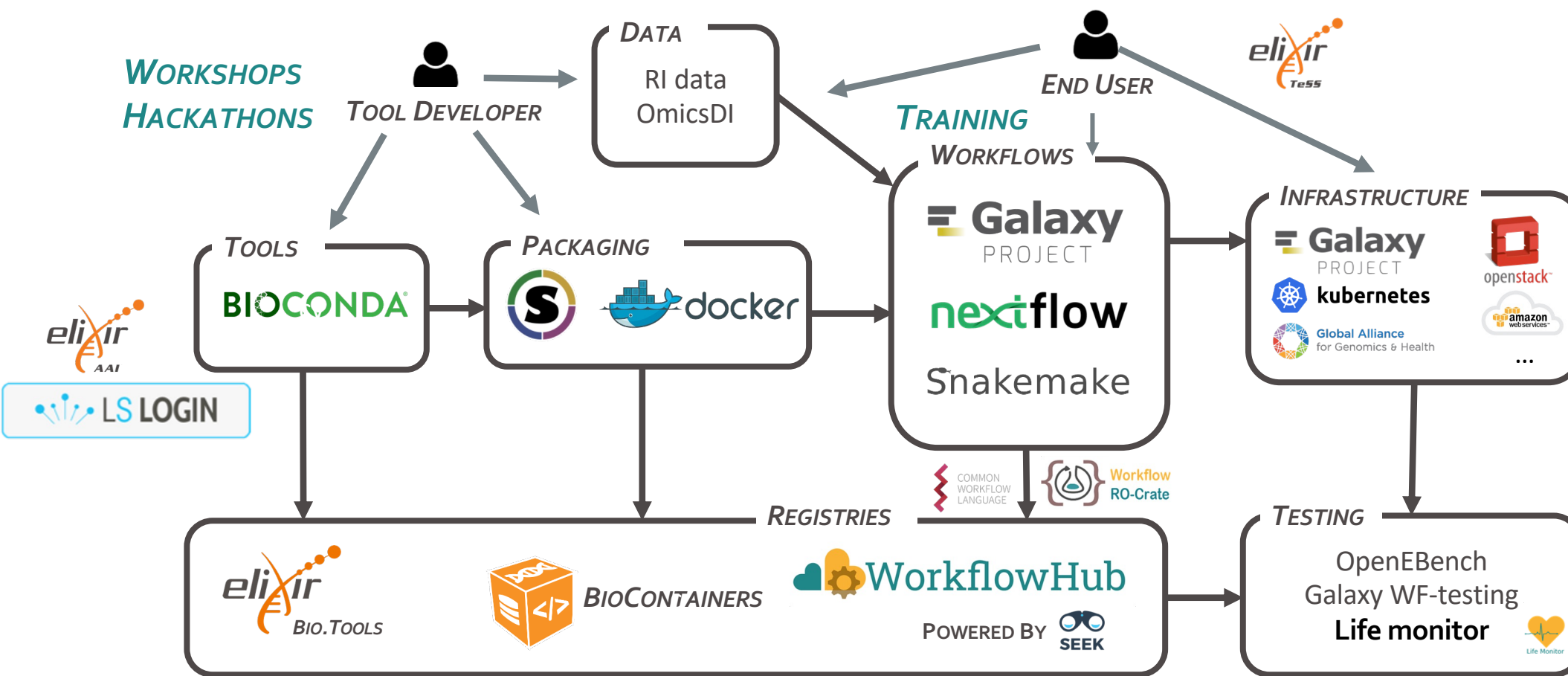
# Reproducible analysis services

# EOSC-Life Tools Collaboratory Roadmap

# EOSC-Life Tools Collaboratory Roadmap

# The Tools Platform ecosystem

# Workflows as entry point for end-users

Galaxy is an **open**, web-based virtual research environment for **accessible**, **reproducible**, and **transparent** computational biomedical research

**Galaxy | Workflow Editor**

usegalaxy.org/workflow/editor?id=7930192474610810

Galaxy

Analyze Data    Workflow    Visualize ▾    Shared Data ▾    Admin    Help ▾    User ▾

Using 4.1 TB

⚠ Certain large memory tools are temporarily running with reduced memory (RNA STAR, SPAdes, Unicycler) or have been temporarily disabled (Trinity).

**Tools**

search tools

Inputs
Get Data
Collection Operations
Expression Tools

GENERAL TEXT TOOLS
Text Manipulation
Filter and Sort
Join, Subtract and Group
Datamash

GENOMIC FILE MANIPULATION
FASTA/FASTQ
FASTQ Quality Control
SAM/BAM
BED
VCF/BCF
Nanopore
Convert Formats
Lift-Over

COMMON GENOMICS TOOLS
Operate on Genomic Intervals
Fetch Sequences/Alignments

GENOMICS ANALYSIS
Assembly
Annotation

**COVID-19: PE Variation**

**Details**

🔧 Convert VCF to VCF_BGZIP (Galaxy Version 1.0.2)

Label

67%

---

**Galaxy @ Belgium**

Analyze Data    Workflow    Visualize ▾    Shared Data ▾    Admin    Help ▾    User ▾

Using 9%

**Tools**

variant

⬆ Upload Data

Show Sections

using plot-vcfstats

**SnpEff** Variant effect and annotation

**Delly call** and genotype structural variants

**Lofreq filter** called variants posteriorly

**medaka variant tool** Probability decoding

**medaka variant pipeline** via neural networks

**TB Variant Filter** M. tuberculosis H37Rv VCF filter

**VCFdistance:** Calculate distance to the nearest variant

**VCFdistance:** Calculate distance to the nearest variant

**FreeBayes** bayesian genetic variant detector

**Delly filter** somatic or germline structural variants

**SnpEff eff:** annotate variants for SARS-CoV-2

**Call variants** with LoFreq

**Call specific mutations in reads:** Looks for reads with mutation at known positions and calculates frequencies and stats.

**DCS mutations to SSCS stats:** Extracts all tags from the single stranded

**SnpEff eff: annotate variants for SARS-CoV-2 (Galaxy Version 4.5covid19)**

☆ Favorite    ▾ Options

**Sequence changes (SNPs, MNPs, InDels)**

No vcf or bed dataset available.

**Input format**
VCF

**Select an annotated Coronavirus genome**
NC_045512.2: COVID19 Severe acute respiratory syndrome coronavirus 2 isolate Wuhan-Hu-1

**Output format**
VCF (only if input is VCF)

**Create CSV report, useful for downstream analysis (-csvStats)**
No

**Upstream / Downstream length**
No upstream / downstream intervals (0 bases)
(-ud)

**Annotation options**

☐ Select/Unselect all

☐ Use 'EFF' field compatible with older versions (instead of 'ANN')
☐ Use Classic Effect names and amino acid variant annotations (NON_SYNONYMOUS_CODING vs missense_variant and G180R vs p.Gly180Arg/c.538G>C)
☐ Override classic and use Sequence Ontology terms for effects (missense_variant vs NON_SYNONYMOUS_CODING)
☐ Override classic and use HGVS annotations for amino acid annotations (p.Gly180Arg/c.538G>C vs G180R)
☐ Old notation style notation: E.g. 'c.G123T' instead of 'c.123G>T' and 'X' instead of '*'
☐ Use one letter Amino acid codes in HGVS notation. E.g. p.R47G instead of p.Arg47Gly
☐ Use transcript ID in HGVS notation. E.g. ENST00000252100:c.914C>G instead of c.914C>G
☐ Do not shift variants according to HGVS notation (most 3prime end)
☐ Do not add HGVS annotations
☐ Only use canonical transcripts
☐ Only use protein coding transcripts
☐ Use gene ID instead of gene name (VCF output)
☐ Disable IUB code expansion in input variants
☐ Add OICR tag in VCF file
☐ Do not add LOF and NMD annotations

**History**

search datasets

**covid-19 original data**
33 shown, 2 deleted, 72 hidden
13.28 GB

99: data 97 converted to fastqsanger (READ2)

98: data 97 converted to fastqsanger (READ1)

97: MergeSamFiles on data 96, data 95, and data 94: Merged BAM dataset

93: Filter SAM or BAM, output SAM or BAM on collection 89: bam

89: Map with BWA-MEM on collection 60 (mapped reads in BAM format)
a list with 3 items

88: MultiQC on data 86, data 78, and data 70: Web page

87: MultiQC on data 86, data 78, and data 70: Stats
a list with 3 items

62: fastp on collection 3: JSON Report
a list with 3 items

61: fastp on collection 3: HTML Report

# Applying infrastructure services a pandemic use case

# Submission of viral data

# Submission overview



Raw reads
Consensus sequence
Standardized metadata

# Different ways of usage

## Deploy the container

Laptop
Cloud infrastructure

## Brokering

covid19.usegalaxy.be

*elixir BELGIUM*

## Public Galaxy instance

useGalaxy.eu
useGalaxy.be

```
docker run -p "8080:80" --privileged quay.io/galaxy/ena-upload
```

github.com/ELIXIR-Belgium/ena-upload-container

VLAAMS SUPERCOMPUTER CENTRUM

Vlaanderen is supercomputing

*elixir BELGIUM*

# rdm.elixir-belgium.org

- Data Management in Simple Steps
- Data Management Plan ⌄
- Data Management for Omics Data ⌄
- **Covid-19** ⌃
  - SARS-Cov-2 raw reads submission
  - SARS-Cov-2 assembly submission

# Covid-19 data submission

- The tools
- Overview of the submission process

ELIXIR supports the European Corona action plan and plays an important role in the development of the COVID-19 Data Portal. As the life-science data Research Infrastructure in Europe, ELIXIR is in a unique position to help increase the amount of publicly available Covid-related data and facilitate its processing, publication and reuse.

ELIXIR Belgium promotes and encourages the publication of all scientific data related to the Covid pandemic and provides the tools, know-how and brokering services for researchers to do so. Our first action is to support the submission of SARS-Cov-2 nucleotide sequences to public repositories.
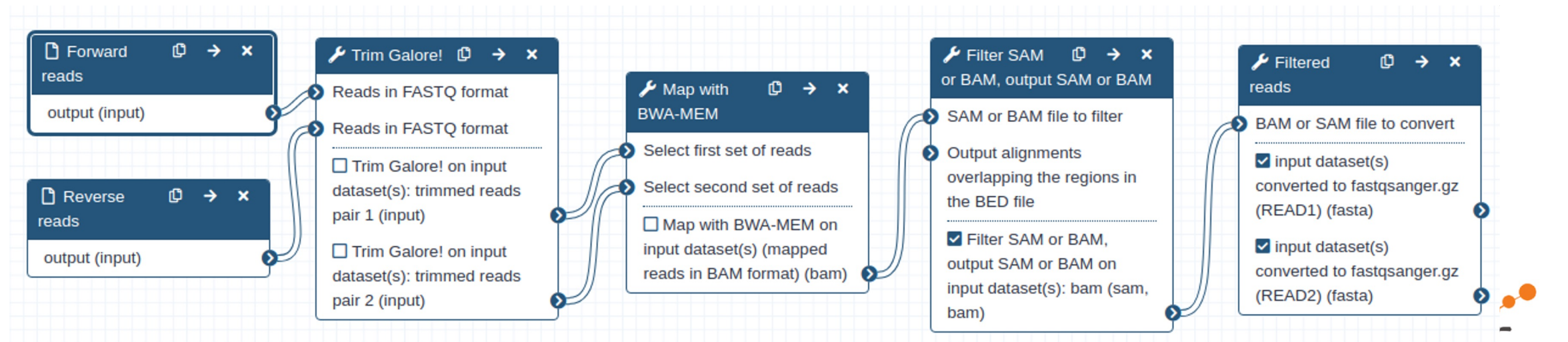
To achieve this, we have collaboratively developed and compiled Galaxy tools and workflows necessary to clean, assemble and submit SARS-CoV-2 sequences to the European Nucleotide Archive (ENA). There are many advantages of using Galaxy including a graphical user interface, access to tools and workflows for pre-processing, downstream

# Enabling data analysis

# Removing Human reads

- Standard preprocessing step in different pipelines.
  - https://covid19.galaxyproject.org/genomics/1-preprocessing/
- Reached out to ENA for a standard method
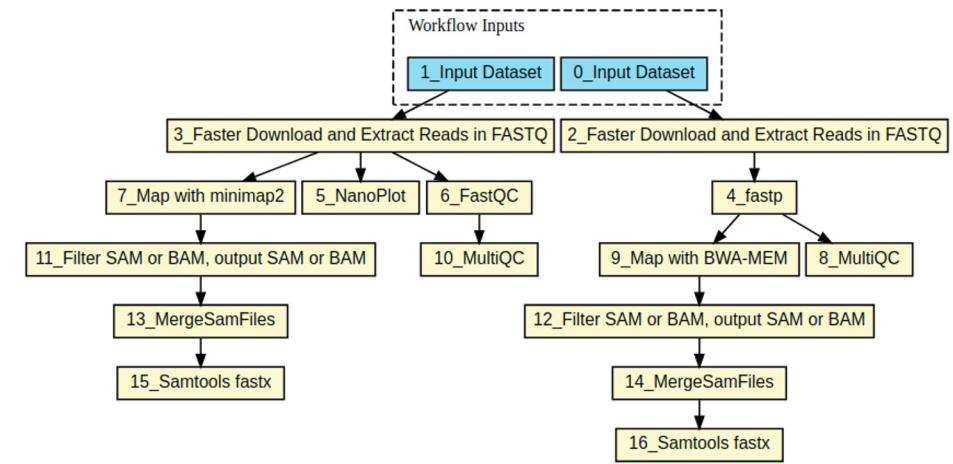  - Repurposed https://github.com/Finn-Lab/Metagen-FastQC

# Workflows

- Different approaches

  - Removing of reads mapping to human genome (e.g. Metagen-FastQC)

  - Retain reads mapping to viral genome

  - Combine human & viral genome, retain viral reads

- Galaxy allows using other workflows as you please

  - https://workflowhub.eu/

  - https://workflowhub.eu/workflows/99  (Metagen-FastQC)

# WorkflowHub : a FAIR workflow registry

- Workflow management system agnostic

- Registry & repository functionality

- Workflows may remain in their native repositories in their native form

- Open to workflows from all disciplines and any country

- Based on community standards

- Attribution and credit as a central feature

- Added value services: curation, monitoring, …

- WorkflowHub Club open community

- Perpetual Development in the open

  - Registering on behalf of makers

  - Regular releases of new features

  - Agile revisions of features

https://workflowhub.eu

# Integration & development of standards

**Canonical workflow description**

Native or Abstract CWL

**Metadata for registration and discovery**
ComputationalWorkflow and FormalParameter
ComputationalTool
Schema.org profile and types

**API Exchanging Tools and Workflows**
GA4GH Tool Registry Service (TRS) API

Specialisation of RO-Crate to package an executable workflow with all necessary documentation.

**Exchange format for Workflow Hub.**

# WorkflowHub - usegalaxy.eu integration



NOT an execution platform, but can be coupled to execution platform.
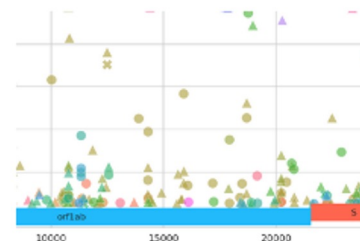
# covid19.galaxyproject.org

- Workflows for different disciplines doing COVID-19 research

- COVID-19 Galaxy webinars

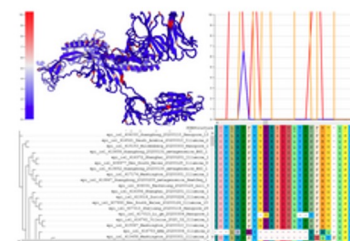🔗 elixir-europe.org/events/2nd-galaxy-elixir-webinar-series

### Genomics
**Assembly and intra-host variation**

- Assembly
- MRCA timing
- Variation analysis
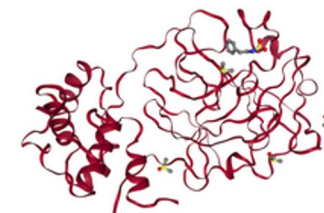- Selection and recombination

### Evolution
**Sites under selection**

- Natural Selection Analysis
- Analysis
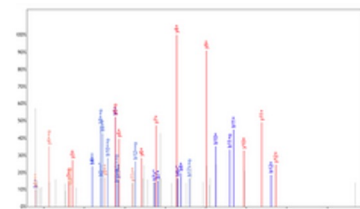- Visualizations
- Observable Notebooks

### Cheminformatics
**Screening of the main protease**

- Compound enumeration
- Generation of 3D conformations
- Docking
- Scoring
- Selection of compounds for synthesis

### Proteomics
**Mass Spectrometry**

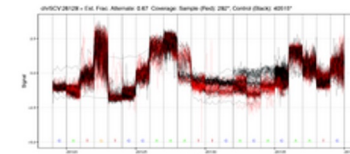- Reanalysis of PXD018117
- Reanalysis of PXD018241

### Artic
**Amplicon based data analysis**

ARTICnetwork
http://artic.network | @NetworkArtic

### direct RNA-seq
**direct RNA-seq data analysis**

- Pre-Processing
- RNA Epigenetics

https://covid19.galaxyproject.org    https://covid19beacon.crg.eu

# Beyond Genomics

# Alignment with other projects and initiatives

# Sensitive Data Infrastructure
## linking genotype & phenotype



European Health Data Space
Belgian Genome Biobank (FWO - EWI)
Genomic Data Infrastructure (Digital Europe Programme)

# Division of responsibilities

- Hardware Infrastructure
  - Secure storage
  - Secure compute
- Services
  - Secure access to data
  - Trusted Research Environment
- Interoperability
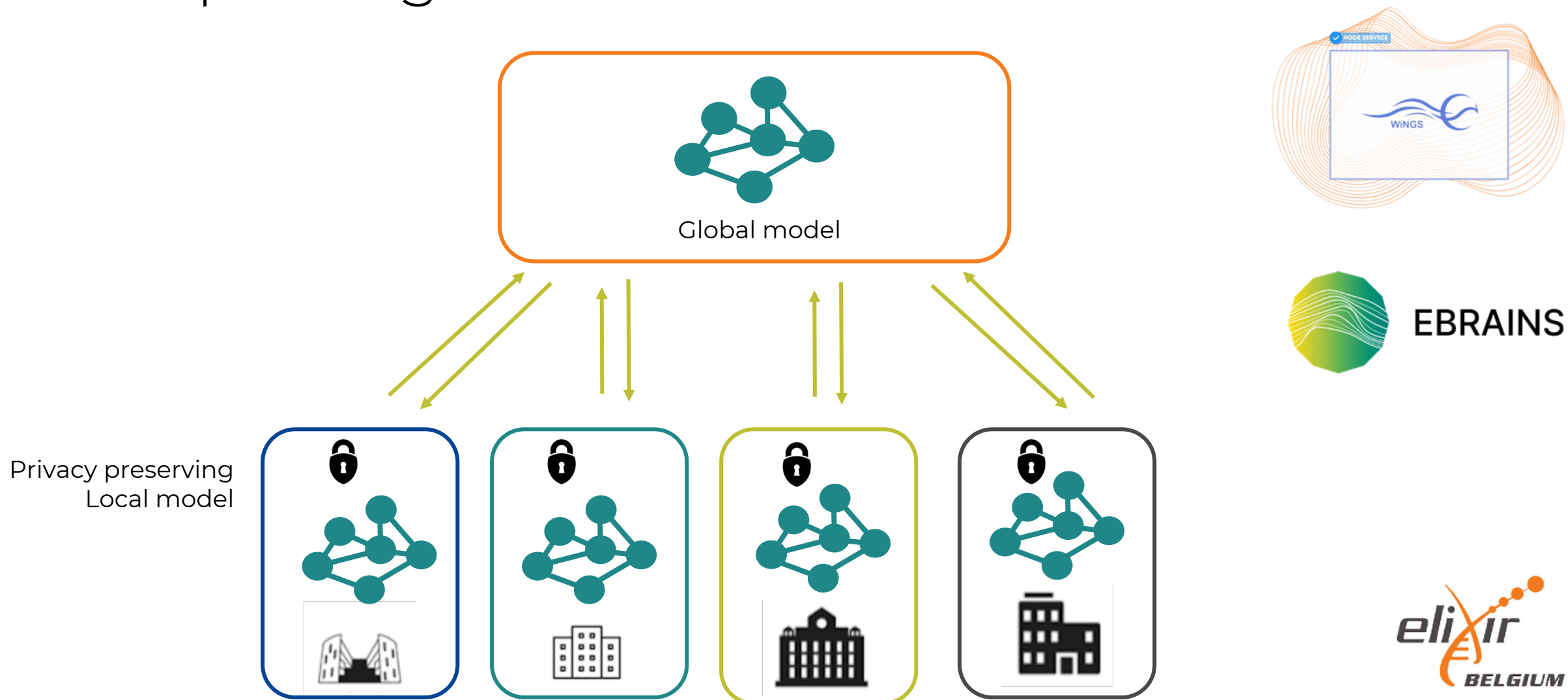  - Standards e.g. GA4GH
  - Alignment with EU developments
- Data management & ELSI

Accelerating research
building tomorrow's services

# Federated learning framework
## underpinning diverse use cases



Global model

Privacy preserving
Local model

# Biodiversity



Collaborative genome annotation

Functional annotation through comparative genomics

Processing & analysis
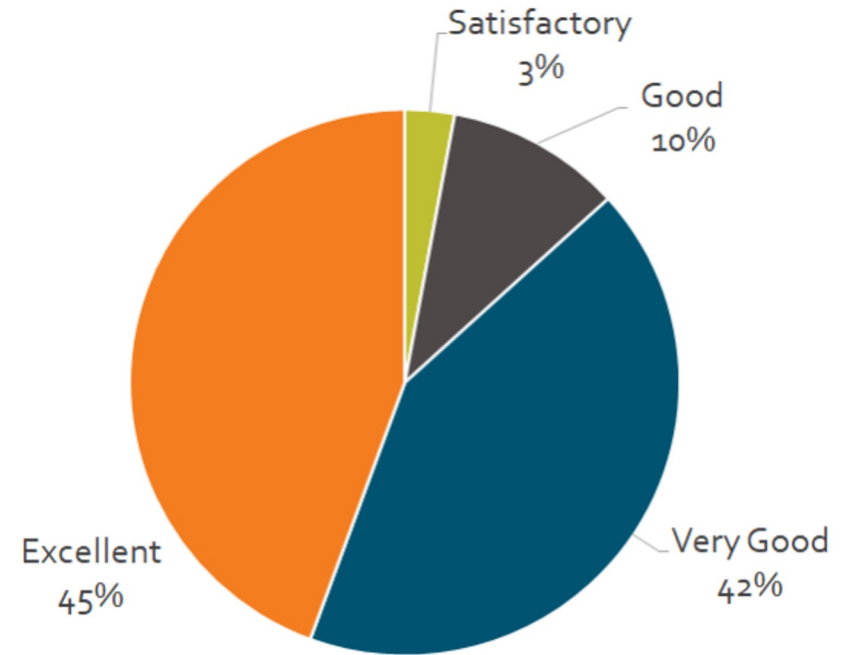
# Bringing services to the users

# Training



e-Learning

F2F training

Hackathon

Knowledge exchange

Satisfactory 3%

Good 10%

Excellent 45%

Very Good 42%

*"Nicely interactive with hands-on exercises"*

*"Taking the time to demo parts that may be relevant to certain researchers is a great plus as it allows you to see the course material 'in the wild'"*

# Outreach & communication

Continue building the ELIXIR community in Flanders & Belgium

Engage with life-sciences community

**Data users**

Academia
Industry
Public health

**Data generators**

Core facilities
Sequencing Centers
…

**Policy makers**

**Funders**

Thank you