Pandemic Research Infrastructure (PaRI)

⁴⁴ The one-year PaRI project focused on facilitating Nordic research on pandemics and especially the COVID-19 pandemic"

1 nov 2020–31 okt 2021

Presented by Wolmar Nyberg Åkerström, wolmar.n.akerstrom@uu.se

www.neic.no/pari







A Nordic e-infrastructure project ⊗∩eIC

- 6 partners from Denmark, Estonia, Germany, Norway and Sweden
- **3 observers** from Finland, Norway and Sweden
- Liaised with Nordic & European activities
- Diverse and active reference group







Enable easy and seamless access ⊗∩elC

WP1 Data management

Support for data sharing Access to data resources Index of infrastructure services North Covid-19 Data Portal







Enable easy and seamless access ⊗∩elC

- WP1 Data management Support for data sharing Access to data resources Index of infrastructure services
- **WP2** Analysis platform

Custom analysis portal (Galaxy) Selection of tools & workflows Proof of concept deployment





Enable easy and seamless access ⊗∩elC

- WP1 Data management Support for data sharing Access to data resources Index of infrastructure services
- **WP2 Analysis platform** Custom analysis portal (Galaxy) Selection of tools & workflows Proof of concept deployment
- **WP3 Scaling the analysis** Workflows for data sharing Local Galaxy for sensitive data Nordic pandemic dashboard







- WP1 Data management Support for data sharing Access to data resources Index of infrastructure services
- WP2 Analysis platform
 Custom analysis portal (Galaxy)
 Selection of tools & workflows
 Proof of concept deployment
- **WP3 Scaling the analysis** Workflows for data sharing Local Galaxy for sensitive data Nordic pandemic dashboard

Focusing on pandemic data





- WP1 Data management Support for data sharing Access to data resources Index of infrastructure services
- WP2 Analysis platform
 Custom analysis portal (Galaxy)
 Selection of tools & workflows
 Proof of concept deployment
- WP3 Scaling the analysis
 Workflows for data sharing
 Local Galaxy for sensitive data
 Nordic pandemic dashboard

Focusing on pandemic data



Combining Nordic & EU resources

- → No development of analysis tools
- → No transfer of data ownership
- → No operational services
- ➔ No new hardware



«≫∩PIC

Sensitive personal data can be directly *or* <u>indirectly</u> linked to an individual, such as

- Patient information
- Personal and phenotypic information
- Pathogen genomic information with fragments of human genetic information
- Sample and process information



⊗∩eic

Sensitive personal data can be directly *or* <u>indirectly</u> linked to an individual, such as

- Patient information
- Personal and phenotypic information
- Pathogen genomic information with fragments of human genetic information
- Sample and process information

- Ethics approval and/or contracts required for use and for sharing, including partners, service providers and data sharing platforms.
- Except when it's no longer possible to link the data to an individual.

What agreements & guidance can be prepared for the next outbreak?



FAIR by design for pandemic data ⊗∩eIC



Study & data design

Procedures

data protection, ethics permit, infrastructure, standards, protocols, data dictionaries, data access, ...





FAIR by design for pandemic data ⊗∩eIC



Procedures Biosamples and instruments

data protection, ethics permit, infrastructure, standards, protocols, data dictionaries, data access, ...

Information

system

populations (statistical) and inclusion criteria, physical processing steps, working storage conditions, long-term storage location, sample quality assessment, sample annotations,

reagents, ...





FAIR by design for pandemic data ⊗∩eIC



Procedures Biosamples and instruments

data protection, ethics permit, infrastructure, standards, protocols, data dictionaries, data access, ...

Information

system

populations (statistical) and inclusion criteria, physical processing steps, working storage conditions, long-term storage location, sample quality assessment,

reagents, ...

sample annotations.



Data and computational workflows

digital processing steps, working storage conditions, long-term storage location, data quality assessment, sample/data annotations, reference data, ...





Section FAIR by design for pandemic data Section (Section 2) (Sec



Procedures Biosamples and instruments

data protection, ethics permit, infrastructure, standards, protocols, data dictionaries, data access, ...

Information

system

populations (statistical) and inclusion criteria, physical processing steps, working storage conditions, long-term storage location, sample quality assessment, sample annotations,

reagents, ...

Sample workspace Biobank

Data and computational workflows

Outputs

digital processing steps, working storage conditions, long-term storage location, data quality assessment, sample/data annotations, reference data, ...



publications, data, tools, workflows, reports, dashboards, ...





Support for data sharing

- Guidance to projects, labs and other organisations producing or commissioning viral sequencing data
- What constitutes a good and ultimately (re)usable pathogen genome data record with a Nordic perspective?
- National COVID-19 Data Portals as channels to get support





OISU





«»()PIC

 COVID-19 data portal (ENA)
 Open & direct access to assembled genomes and raw reads

SARS-CoV-2 Contextual

Manually curated metadatabase with links to sequences

Federated EGA

Controlled access to sensitive data stored nationally



COVID-19 data portal (ENA) Open & direct access to assembled genomes and raw reads

SARS-CoV-2 Contextual

Manually curated metadatabase with links to sequences

Federated EGA

Controlled access to sensitive data stored nationally

GISAID EpiCoV[™] Contract/license-based access to assembled genomes

ECDC TESSy/EpiPulse Controlled access to surveillance and reporting coordinated by national public health authorities

National registers Often controlled access to health data, statistics, etc.



COVID-19 data portal (ENA) Open & direct access to assembled genomes and raw reads

SARS-CoV-2 Contextual

Manually curated metadatabase with links to sequences

Federated EGA

Controlled access to sensitive data stored nationally

GISAID EpiCoV™ Contract/license-based access to assembled genomes

ECDC TESSy/EpiPulse Controlled access to surveillance and reporting coordinated by national public health authorities

 National registers
 Often controlled access to health data, statistics, etc.

What data should a fully operational pandemic research infrastructure cover?



Analysis platform using Galaxy ⊗∩eIC

- Open platform connecting data, computational workflows, visualisations and other services
- Supports **reproducible**, **and transparent** computational analysis
- Community committed to improving tools & workflows

What can be done to promote FAIR for workflows?







Custom analysis portal (Galaxy) ⊗∩eIC

- Convenient access to Nordic SARS-CoV-2 genomic data and analysis workflows
- Managed collection of tools, versions and documentation
- Bring your own data and build your own workflows
- Automated configuration to allow replication and sharing across environments

Tools	COVID-19 research! Want to learn the best practices for the analysis of SARS-CoV-2 data using Galaxy? Visit the Galaxy SARS-CoV-2 portal at covid19.galaxyproject.org. We mirror all nordic public			History
search tools				search dataset
2. Upload Data				Unnamed hist
Get Data				
Send Data				1 This history is
Collection Operations	Welcome to covid19.usegalaxy.no			your own da an external s
Lift-Over	This subdomain of covid19.usegalaxy.no is devoted to research on coronavirus disease 2019 (COVID-19), which has recently been declared a pandemic by the World Health Organization. This subdomain servers as a companion to our study describing the analysis of early COVID-19 data, the goal of which is to underscore the importance of access to raw data and demonstrate that existing community efforts in curation and deployment of biomedical software can reliably support rapid reproducible research during global crises. This subdomain contains the exact versions of all software used. Our analysis was divided into six parts listed below. Each part has a dedicated page that provides links to input datasets, intermediate and final results, workflows, and Galaxy histories that list all details for each analysis. These workflows can be re-run on this subdomain. 1. Pre-processing of raw read data 2. Assembly of SARS-CoV-2 genome 3. Estimation of timing for most recent common ancestor (MRCA) 4. Analysis of variation within Individual isolates 5. Functional annotation: Analysis of Spike protein substitutions 6. Analysis of recombination and selection			
Text Manipulation				
Convert Formats				
Filter and Sort				
Join, Subtract and Group				
Fetch Alignments/Sequences				
Operate on Genomic Intervals				
Statistics				
Graph/Display Data				
Phenotype Association				
Multiple Alignments				
Assembly	Our Data Policy			
у	Registered Users	Unregistered Users	GDPR Compliance	
	Liser data on will be available as	Processed data will only be	The Galaxy service	



Scaling the analysis



- FAIR data sharing requires planning workflows for data mapping and anonymisation
- The analysis portal can be sandboxed to **run on secure infrastructures**
- Proof of concept "real-time" dashboard of geographic distribution of virus variants over time

What common data definitions can we agree on across the Nordics?







Providing Nordic Added Value ⊗∩eIC

- Increased capabilities to support pandemic data sharing and analysis
- Contributed to making more pathogen genome data available for research
- Strengthened contributions to related activities not only in the Nordics
- Work continues in European projects and acoss partners











FAIR Pandemic data sharing

OI90





- Timely and rewarding platform for knowledge exchange between the partners & stakeholders
- Useful to address data sharing guidelines and access from workflows in parallel
- Satisfying to see a complete prototype for managing and executing tools and workflows on Nordic pandemic data sets

→ Increased our capabilities to support pandemic data sharing and analysis

Solution

- → Contributed to making more pathogen genome data available for research
- → Contributed experiences to related activities not only in the Nordics but also worldwide



Professional network

- FAIRification of pandemic related omics data
- Running and administering COVID-19 workflows in Galaxy
- Curating pandemic data for visualisations in Nextstrain / Auspice (possibly other visualisation tools)

PaRI affiliate coordinator:

Wolmar Nyberg Åkerström wolmar.n.akerstrom@uu.se

Outreach

- Expand the professional network to new members
- Share experiences with active projects and new partners

Shared tools & workflows

Populate a GitHub organisation with tools and workflows from the partners



Access to genome data for the virus that causes COVID-19

PaRI provided guidance and hands-on support to deposit genome data and the information collected with the samples according to the schemas and data specifications supported by two different data sharing platforms.

Genome data is acquired by analysing samples collected from individuals. The samples are collected by a diverse set of organisations ranging from public health surveillance initiatives to research projects focusing on different aspects of the virus or the disease. What information is collected with the samples, how the information collected is structured, and in what detail information can be shared varies. To be widely accessible to the public health and research communities the data should be shared through different data sharing platforms with overlapping but different schemas to structure and encode the information. Access to sufficiently detailed data specifications and data mappings is required when developing mediary tools and resources that enable syntactic and semantic interoperability.



Access to tools and workflows to analyse genome data for the virus

PaRI designed data flows and provisioned/hosted workflow systems for running computational workflows on data from data producing partners as well as data available in public data sharing platforms.

There are well-established data formats for storing genome data and the bioinformatics community has developed a wide array of shared practices, tools and computational workflows to automate some aspects of quality control and data analysis. Adopting/developing a computational workflow includes specifications of the expected input and output data as well as considerations on how to document the process to enable peer-review and reproducibility. If all information required for a workflow is available in a data sharing platform, that information could be mapped to the corresponding input data specification resulting in consistently produced outputs for large collections of data.



Access to a dashboard showing geographic distribution of virus variants over time

PaRI designed data flows for regional surveillance and hosted a proof of concept for a dashboard based on Auspice/Nextstrain. Curated public data/data from data producing partners. Mapped to a database hosting contextual data for SARS-CoV-2 sequences. Aggregate data and visualisations presented on regional dashboards can be used to convey information about how a virus/outbreak evolves and spreads and to provide an overview of what data is available for research. The dashboard presents consistent views of underlying data that may have been produced according to varying naming and encoding schemas, such as definitions of geographic regions, age groups, labelling of virus characteristics, etc. A dynamic dashboard requires workflows to be available for mapping information from all the sources of underlying data to the views presented on the dashboard.

Sold (Sector)