

# SARS-CoV-2 at the ENA (COVID-19 Data Platform)

PHA4GE Webinar 21/09/21

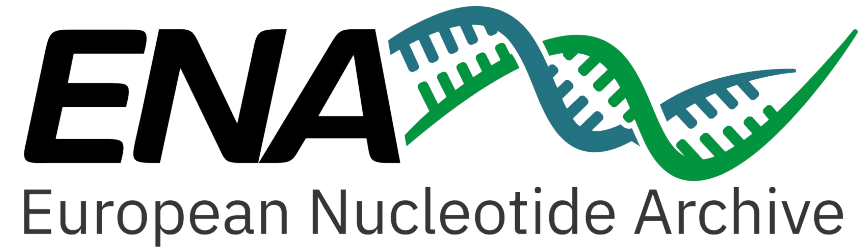
Colman O’Cathail

Bioinformatician, European Nucleotide Archive



# What We Do

- Submitting to the ENA
  - How we relate to the INSDC
  - Our metadata model
  - Submission routes
  - Helpdesk
- The European COVID-19 Data Platform (VEO partners)
  - The COVID-19 Data Portal tour
  - Systematic analyses of public data
  - Data Hubs



# Submitting to ENA

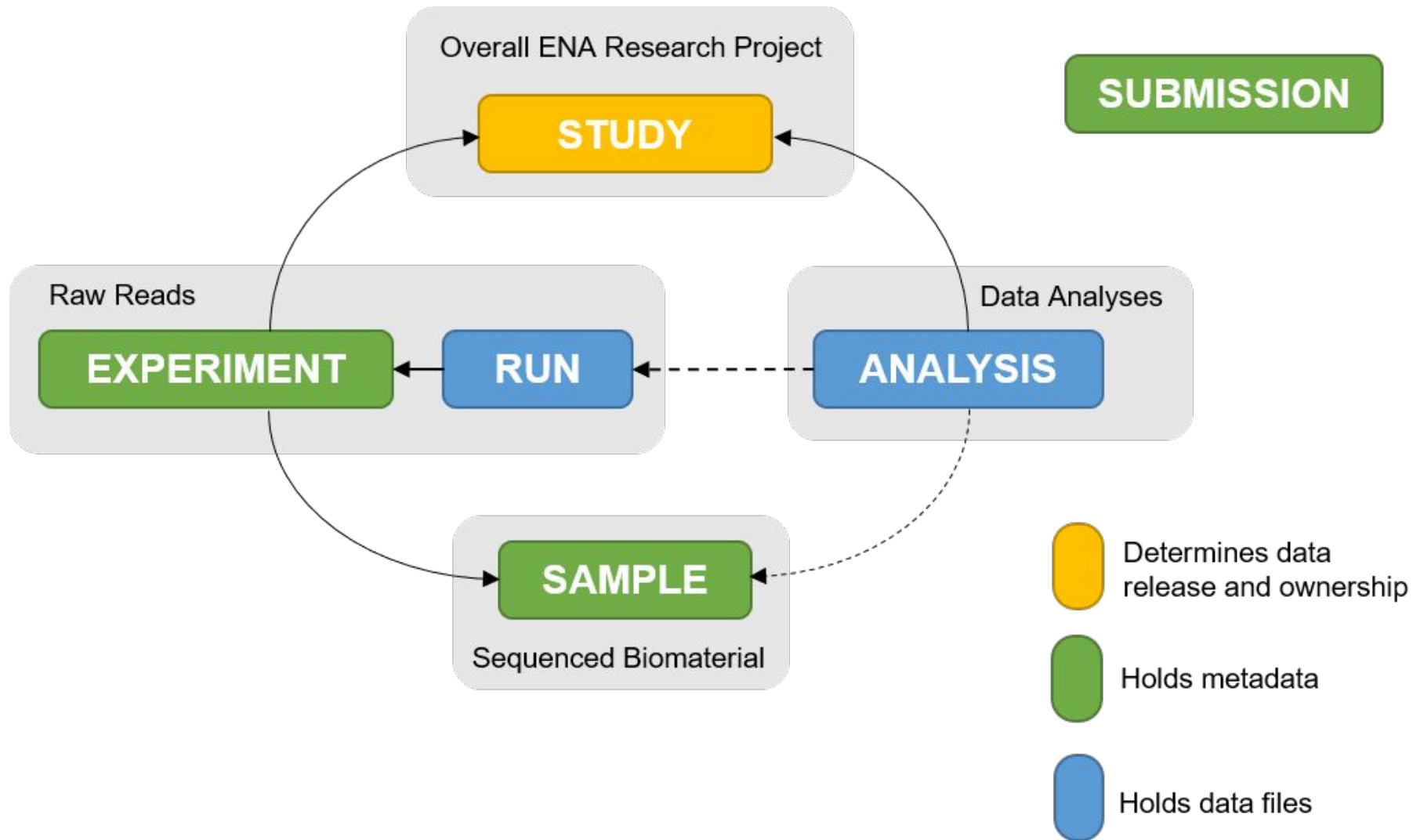
What, why, how

# The INSDC

Data type	DDBJ	EMBL-EBI	NCBI
Next generation reads	<a href="#">Sequence Read Archive</a>	European Nucleotide Archive ( <a href="#">ENA</a> )	<a href="#">Sequence Read Archive</a>
Capillary reads	<a href="#">Trace Archive</a>		<a href="#">Trace Archive</a>
Annotated sequences	<a href="#">DDBJ</a>		<a href="#">GenBank</a>
Samples	<a href="#">BioSample</a>		<a href="#">BioSample</a>
Studies	<a href="#">BioProject</a>		<a href="#">BioProject</a>

[www.insdc.org](http://www.insdc.org)

# The ENA Metadata Model



# Submission methods

	Interactive	Webin-CLI	Programmatic
Study	Y	N	Y
Sample	Y	N	Y
Read data	Y	Y*	Y
Genome Assembly	N	Y	N
Transcriptome Assembly	N	Y	N
Template Sequence	Y	Y	N
Other Analyses	N	N	Y

<https://ena-docs.readthedocs.io/en/latest/index.html>

# Interactive Mode

## Studies (Projects)



Register Study



Studies Report



Submit XMLs (advanced)

## Samples



Register Samples



Samples Report



Register Novel Taxonomy



Submit XMLs (advanced)

## Raw Reads (Experiments and Runs)

*Raw reads can also be submitted using [Webin-CLI](#)*



Submit Reads



Runs Report



Submit XMLs (advanced)



Run Files Report



Run Processing Report



Unsubmitted Files Report

## Data Analyses

*Assemblies and annotated sequences must be submitted with [Webin-CLI](#). Other analyses can be submitted as XMLs.*



Generate Annotated Sequence Spreadsheet



Analyses Report



Submit XMLs (advanced)



Analysis File Report



Analysis Processing Report

<https://www.ebi.ac.uk/ena/submit/webin/login>

# Webin-CLI

	Webin-CLI
Study	N
Sample	N
Read data	Y*
Genome Assembly	Y
Transcriptome Assembly	Y
Template Sequence	Y
Other Analyses	N

<https://ena-docs.readthedocs.io/en/latest/index.html>



- Webin-CLI is a java

ENA

```
(base) [ocathail@hh-yoda-08-01 ocathail]$ java -jar webin-cli-3.7.0.jar -help
Usage: java -jar webin-cli-3.7.0.jar [-ascp] [-fields] [-help] [-quick]
                                     [-submit] [-test] [-validate] [-version]
                                     [-centerName=CENTER] -context=TYPE
                                     [-inputDir=DIRECTORY] -manifest=FILE
                                     [-outputDir=DIRECTORY]
                                     [-password=PASSWORD] [-passwordEnv=VAR]
                                     [-passwordFile=FILE] -userName=USER

Description:
Validate and submit files to ENA using the Webin submission service. Use the
-fields option to see supported manifest fields for all contexts or for a
specific -context. Detailed instructions are available from:
https://ena-docs.readthedocs.io/en/latest/cli.html

Options:
  -context=TYPE          Submission type: genome, transcriptome, sequence, reads,
                        taxrefset
  -manifest=FILE         Manifest text file containing file and metadata fields.
  -userName, -username=USER
                        Webin submission account name or e-mail address.
  -password=PASSWORD    Webin submission account password.
  -passwordFile=FILE    File containing the Webin submission account password.
  -passwordEnv=VAR      Environment variable containing the Webin submission
                        account password.
  -inputDir, -inputdir=DIRECTORY
                        Root directory for the files declared in the manifest
                        file. By default the current working directory is used
                        as the input directory.
  -outputDir, -outputdir=DIRECTORY
                        Root directory for any output files written in
                        <context>/<name>/<validate,process,submit> directory
                        structure. By default the manifest file directory is
                        used as the output directory. The <name> is the unique
                        name from the manifest file. The validation reports are
                        written in the <validate> sub-directory.
  -centerName, -centername=CENTER
                        Mandatory center name for broker accounts.
  -validate              Validate files without uploading or submitting them.
  -quick                Validates submitted read files (BAM, CRAM, Fastq) within
                        a fixed time period (5 minutes). All CRAM reference
                        sequence md5 checksums are always validated. When this
                        option is used files may only be partially validated
                        and may fail post-submission processing.
  -submit               Validate, upload and submit files.
  -test                 Use the test submission service.
  -ascp                 Use Aspera (if Aspera Cli is available) instead of FTP
                        when uploading files. The path to the installed "ascp"
                        program must be in the PATH variable.
  -help                 Show this help message and exit.
  -fields               Show manifest fields for all contexts or for the given
                        -context.
  -version              Print version information and exit.
Exit codes: 0=SUCCESS, 1=INTERNAL ERROR, 2=USER ERROR, 3=VALIDATION ERROR
(base) [ocathail@hh-yoda-08-01 ocathail]$
```

# Webin-CLI – Genome Manifest Example

```
STUDY PRJEBXXXX  
SAMPLE SAMEAXXXXX  
ASSEMBLYNAME example_assembly  
COVERAGE 60  
PROGRAM TODO  
PLATFORM TODO  
MINGAPLENGTH TODO  
MOLECULETYPE genomic DNA  
FASTA genome.fasta.gz
```

# Programmatic Submission

	Programmatic
Study	Y
Sample	Y
Read data	Y
Genome Assembly	N
Transcriptome Assembly	N
Template Sequence	N
Other Analyses	Y

<https://ena-docs.readthedocs.io/en/latest/index.html>

# Programmatic Submission

- Most advanced submission method
- Requires users to create and submit XMLs directly
- Good idea to utilise support if going this route
- Can be very powerful when done right

# Submitting to SARS-CoV-2 Data

Recommendations, tools & support

# SARS-CoV-2 Submissions

- Dedicated read the docs page
  - [https://ena-browser-docs.readthedocs.io/en/latest/help\\_and\\_guides/sars-cov-2-submissions.html](https://ena-browser-docs.readthedocs.io/en/latest/help_and_guides/sars-cov-2-submissions.html)
- Dedicated support – [virus-dataflow@ebi.ac.uk](mailto:virus-dataflow@ebi.ac.uk)

# SARS-CoV-2 Submissions

- Study registration as normal
- Sample registration – please select the right checklist
  - ERC000033 – ENA Virus Pathogen Reporting Standard Checklist
- We offer a tool to convert GISAID metadata into ENA metadata

([https://github.com/enasequence/ena-content-dataflow/blob/master/scripts/gisaid\\_to\\_ena.py](https://github.com/enasequence/ena-content-dataflow/blob/master/scripts/gisaid_to_ena.py))

# ERC000033 Sample Checklist

- Mandatory fields:

geographic location (country and/or sea)

host common name

host subject id

host health state

host sex

host scientific name

collector name

collecting institution

isolate

- Recommended fields:

- Organism name = *Severe acute respiratory syndrome coronavirus 2*
  - Taxon ID = 2697049
  - Collection date = *<insert\_here>*
- Sample capture status = “active surveillance in response to outbreak”



# Submitting reads and assemblies

- Reads can be submitted by the users preferred routes
- Assembly submission is recommended via Webin-CLI or the new SARS-CoV-2

Web API

- Don't forget you can reference any archived runs!

# SARS-CoV-2 Web API

- [https://ena-browser-docs.readthedocs.io/en/latest/help\\_and\\_guides/Webin-Cli\\_SARS-CoV-2\\_Genome\\_Submission\\_REST\\_API.html](https://ena-browser-docs.readthedocs.io/en/latest/help_and_guides/Webin-Cli_SARS-CoV-2_Genome_Submission_REST_API.html)
- Two service endpoints; Validate & Submit (test and production for both)

```
1 {
2   *"name": "string", # Unique name for the assembly within the Webin submission account
3   *"study": "string", # Study accession number or unique name (alias)
4   *"sample": "string", # Sample accession number or unique name (alias)
5   *"coverage": "number", # The estimated depth of sequencing coverage
6   *"program": "string", # The assembly program
7   *"platform": "string", # The sequencing platform
8   *"sequence": "string", # The assembled genome sequence
9   "description": "string", # Free text description of the genome assembly
10  "minGapLength": "integer", # Minimum length of consecutive Ns to be considered a gap
11  "moleculeType": "genomic DNA",
12  "runRef": "string", # Run accession number containing the raw reads associated with this genome assembly
13  "tpa": boolean, # Set to true for third party assemblies (by default false)
14  "authors": "string" # List of authors associated with this genome assembly (by default authors of the submission account will be used)
15  "address": "string", # Address where this genome was assembled (by default address of the submission account will be used)
16  "submissionTool": "string", # Submission tool that called this endpoint
17  "submissionToolVersion": "string" # Submission tool version that called this endpoint
18 }
```

# SARS-CoV-2 Web API

## Python Example

### Curl Example

```
1 curl -X 'POST' -u Webin-N:password \
2 'https://wwwdev.ebi.ac.uk/ena/submit/webin-cli/api/v1/genome/covid-19' \
3 -H 'accept: application/json' \
4 -H 'Content-Type: application/json' \
5 -d '{
6   "name": "test_1",
7   "study": "PRJEB46468",
8   "sample": "ERS6670887",
9   "coverage": 100,
10  "program": "Illumina",
11  "platform": "Illumina",
12  "sequence": "CTCTCGATCGATCAAATTTGGGTTTAAGGCCCTTGAATT",
13  "description": "test",
14  "minGapLength": 1,
15  "moleculeType": "genomic DNA",
16  "authors": "test",
17  "address": "test"
18 }'
```

```
1 import sys
2 import requests
3 import json
4
5 data = [
6     {
7         "name": "test_1", "study": "PRJEB46811", "sample": "ERS7306048",
8         "coverage": 100, "program": "Illumina", "platform": "Illumina",
9         "sequence": "CTCTCGATCGATCAAATTTGGGTTTAAGGCCCTTGAATT",
10        "description": "test", "minGapLength": 1, "moleculeType": "genomic DNA",
11        "tpa": False, "authors": "test", "address": "test"
12    },
13    {
14        "name": "test_2", "study": "PRJEB46811", "sample": "ERS7306049",
15        "coverage": 100, "program": "Illumina", "platform": "Illumina",
16        "sequence": "CTCTCGATCGATCAAATTTGGGTTTAAGGCCCTTGAATT",
17        "description": "test", "minGapLength": 1, "moleculeType": "genomic DNA",
18        "authors": "test", "address": "test"
19    }
20 ]
21
22 ## Please remove /validate from the URL to submit the genome instead of just validating it
23 server = "https://wwwdev.ebi.ac.uk/ena/submit/webin-cli/api/v1/genome/covid-19/validate"
24
25 for sample in data:
26     sample_json = json.dumps(sample)
27     response = requests.post(
28         server, headers={"accept": "application/json", "Content-Type": "application/json"},
29         data=sample_json, auth=('Webin-XXXXXX', 'password')
30     )
31     status = response.status_code
32     if status != 200:
33         print("Bad REST call : {}".format(status))
34         sys.exit(1)
35     else:
36         receipt = json.loads(response.content)
37         print("{} : {}".format(sample['name'], receipt))
```

# Bulk Webin-CLI tool

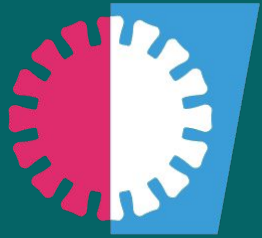
- <https://github.com/enasequence/ena-bulk-webincli>
- Python wrapper that allows users to submit many data with one input spreadsheet
- Containerized using Docker and Singularity
- Create manifests automatically, uploads data and handles submission
- Can be used to validate submissions also before submission

# Bulk Webin-CLI tool

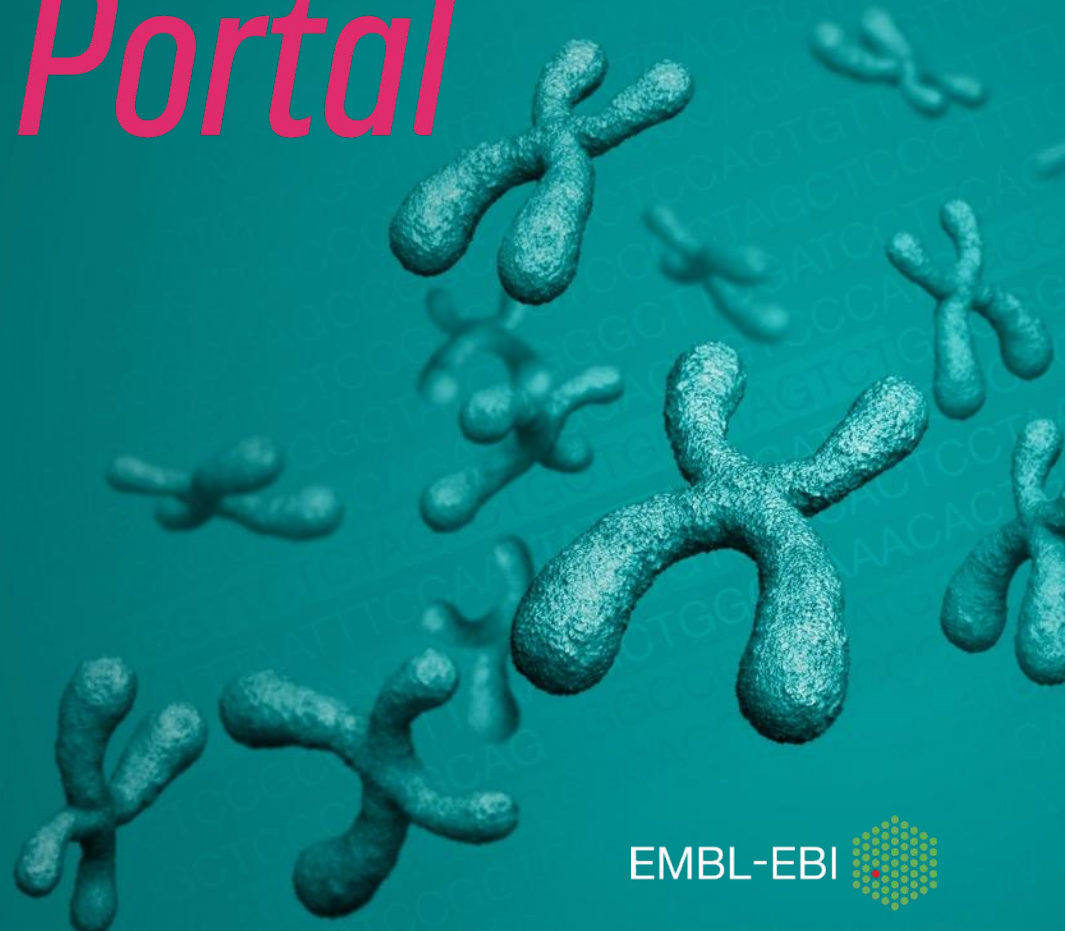
- <https://github.com/enasequence/ena-bulk-webincli>
- Python wrapper that allows users to submit many data with one input spreadsheet
- Containerized using Docker and Singularity
- Create manifests automatically, uploads data and handles submission
- Can be used to validate submissions also before submission

# ORCID Data Claiming

- <https://www.ebi.ac.uk/ebisearch/orcidclaimdocumentation.ebi>
- Allows researchers to claim credit for your deposited public data records



# **COVID-19** *Data Portal*

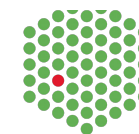


# The who behind the portal



Eötvös Loránd  
University

EMBL-EBI







# Portal Home



About ▾ News Partners Related resources FAQ Bulk downloads Submit data

[Viral Sequences](#) [Host Sequences](#) [Expression](#) [Proteins](#) [Biochemistry](#) [Imaging](#) [Literature](#)

## Accelerating research through data sharing

[Read and sign our letter in support of open COVID-19 data >](#)

### Viral sequences →

Raw and assembled sequence and analysis of SARS-CoV-2 and other coronaviruses.

[4,483,812 records >](#)

### Expression →

Gene and protein expression data of human genes implicated in the virus infection of the host cells. Identifying cell types and genes with highest expression in SARS-CoV-2 infections.

[112 records >](#)

### Biochemistry →

COVID-19 pathways, interactions, complexes, targets and compounds

### Host sequences →

Raw and assembled sequence and analysis of human and other hosts.

[15,450 records >](#)

### Proteins →

Curated functional and classification data on the SARS-CoV-2 protein entries and associated protein receptors.

[2,207 records >](#)

### Imaging →

Biological images from microscopy and other platforms

### Latest news →



7 Sep 2021

[Latest VEO report published](#)

28 Jul 2021

[July 2021 VEO report available to read](#)

30 Jun 2021

[New bulk downloader tool](#)

17 Jun 2021

[VEO...](#)

SUPPORT & FEEDBACK

# Data exploration

Showing 15 of 1,196,055 in Viral sequences > Sequences



## Data types

- All (4,591,150)
- Sequences (1,196,055)
- Reference sequences (2)
- Raw reads (1,408,392)
- Sequenced samples (1,372,557)
- Systematic Analyses (600,964)
- Studies (466)
- Genes (22)
- Browser (1)
- Variants (12,691)

## Release Date ⓘ

## Collection date ⓘ

## Last modification date ⓘ

## Organisms

- Severe acute respiratory syndrome coronavirus 2 (1,193,858)
- Scotophilus bat coronavirus 512 (442)

< ..... > Edit table view

<input type="checkbox"/>	Accession	Lineage	Cross-references <span>ⓘ</span>	Collection date	Country	Center name
<input type="checkbox"/>	OU517059		ENA Study (1) <a href="#">↗</a>	Jul 28, 2021	Germany	Robert Koch Institute
<input type="checkbox"/>	MW751646	B.1.298	Coding (Standard) (12) <a href="#">↗</a> <span>See all ▾</span>	Nov 9, 2020	USA	
<input type="checkbox"/>	OU132334	B.1.619	ENA Study (1) <a href="#">↗</a>		Germany	Robert Koch Institute
<input type="checkbox"/>	OU411290	B.1.617.2 <span>Delta</span>	Viral sequences > Systematic Analyses (1) <span>See all ▾</span>	Jun 29, 2021	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	OU429775	B.1.617.2 <span>Delta</span>	Viral sequences > Systematic Analyses (1) <span>See all ▾</span>	Jul 5, 2021	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	MZ252427	B.1.95	Coding (Standard) (11) <a href="#">↗</a>	May 5, 2021	USA	
<input type="checkbox"/>	MT612232	B.1.338	Proteins > Protein sequences (1) <span>See all ▾</span>	May 9, 2020	Australia	The Peter Doherty Institute for Infection and
<input type="checkbox"/>	BS001066	B.1.1.214	Coding (Standard) (12) <a href="#">↗</a>	Jan 30, 2021	Japan	
<input type="checkbox"/>	FR998245	B.1.1.83	BioSamples (1) <a href="#">↗</a> <span>See all ▾</span>	Apr 6, 2020	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	HG993784	B.1.1.7 <span>Alpha</span>	ENA Study (1) <a href="#">↗</a> <span>See all ▾</span>	Feb 13, 2021	France	UMR 1283/8199
<input type="checkbox"/>	MZ371189	B.1.617.2 <span>Delta</span>	Coding (Standard) (11) <a href="#">↗</a>	May 26, 2021	USA	
<input type="checkbox"/>	OU306801	B.1.617.2 <span>Delta</span>	ENA Study (1) <a href="#">↗</a> <span>See all ▾</span>	May 28, 2021	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	OU023582	B.1.1.7 <span>Alpha</span>	Viral sequences > Systematic Analyses (1) <span>See all ▾</span>	Feb 21, 2021	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	FR998350	B.1.1.1	BioSamples (1) <a href="#">↗</a> <span>See all ▾</span>	May 11, 2020	United Kingdom	COVID-19 Genomics UK Consortium

SUPPORT & FEEDBACK

# Data exploration

## Organisms

- Severe acute respiratory syndrome coronavirus 2 (1,193,819)
- Scotophilus bat coronavirus 512 (442)
- Murine coronavirus (407)

[More... >](#)

## Center name

- COVID-19 Genomics UK Consortium (525,489)
- Robert Koch Institute (126,189)
- CDC-OAMD (47,470)

[More... >](#)

## Country

- United Kingdom (527,776)
- USA (458,557)
- Germany (126,361)

[More... >](#)

## Coverage %



1195801 in this range

[Apply](#)

<input type="checkbox"/>	FR995092	B.1.1.29	BioSamples (1) <a href="#">↗</a>	See all <a href="#">v</a>	Apr 4, 2020	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	FR990449	B.1.160	BioSamples (2) <a href="#">↗</a>	See all <a href="#">v</a>	Jan 29, 2021	Switzerland	Swiss Pathogen Surveillance Platform (S
<input type="checkbox"/>	OD944682	B.1.1.1.7	<span style="border: 1px solid blue; padding: 2px;">Alpha</span> Viral sequences > Systematic Analyses (1)	See all <a href="#">v</a>	Jan 10, 2021	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	OU377872	B.1.617.2	<span style="border: 1px solid blue; padding: 2px;">Delta</span> BioSamples (1) <a href="#">↗</a>	See all <a href="#">v</a>	Jun 26, 2021	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	OU113468	B.1.1.7	<span style="border: 1px solid blue; padding: 2px;">Alpha</span> BioSamples (2) <a href="#">↗</a>	See all <a href="#">v</a>	Mar 24, 2021	Germany	Robert Koch Institute
<input type="checkbox"/>	MZ187102	B.1.1.7	<span style="border: 1px solid blue; padding: 2px;">Alpha</span> Coding (Standard) (12) <a href="#">↗</a>		Apr 15, 2021	USA	
<input type="checkbox"/>	FR994316	B.1.1.220	BioSamples (1) <a href="#">↗</a>	See all <a href="#">v</a>	Apr 5, 2020	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	FR997260	B.1.1.279	BioSamples (1) <a href="#">↗</a>	See all <a href="#">v</a>	Apr 18, 2020	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	FR995150	B.1	BioSamples (1) <a href="#">↗</a>	See all <a href="#">v</a>	Apr 2, 2020	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	MZ232211	B.1.429			Mar 22, 2021	USA	
<input type="checkbox"/>	OU556889				Aug 10, 2021	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	OD991475	B.1.1.7	<span style="border: 1px solid blue; padding: 2px;">Alpha</span>		Feb 15, 2021	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	OU103567	B.1.1.7	<span style="border: 1px solid blue; padding: 2px;">Alpha</span>		Mar 16, 2021	United Kingdom	COVID-19 Genomics UK Consortium
<input type="checkbox"/>	MZ342566	B.1.1.7	<span style="border: 1px solid blue; padding: 2px;">Alpha</span>		May 24, 2021	USA	
<input type="checkbox"/>	MZ152854	B.1.1.7	<span style="border: 1px solid blue; padding: 2px;">Alpha</span>		May 1, 2021	USA	Minnesota Department of Health
<input type="checkbox"/>	MW495237	B.1			Jan 5, 2021	USA	Quest Diagnostics
<input type="checkbox"/>	MW749952	B.1.517			Jan 6, 2021	USA	Broad Institute of MIT and Harvard
<input type="checkbox"/>	MW838286	B.1			Jan 18, 2021	USA	
<input type="checkbox"/>	OU263732	B.1.1.7	<span style="border: 1px solid blue; padding: 2px;">Alpha</span>		Apr 21, 2021	United Kingdom	COVID-19 Genomics UK Consortium

# Data exploration

## Instrument platform

- ILLUMINA (1,187,210)
- OXFORD\_NANOPORE (147,740)
- PACBIO\_SMRT (140,255)

## Library selection

- PCR (1,097,326)
- RT-PCR (323,582)
- unspecified (32,282)

[More... >](#)

## Library strategy

- AMPLICON (1,421,244)
- WGA (24,197)
- WGS (17,582)

## Instrument model

- Illumina NovaSeq 6000 (831,870)
- Sequel II (140,251)
- Illumina MiSeq (138,722)

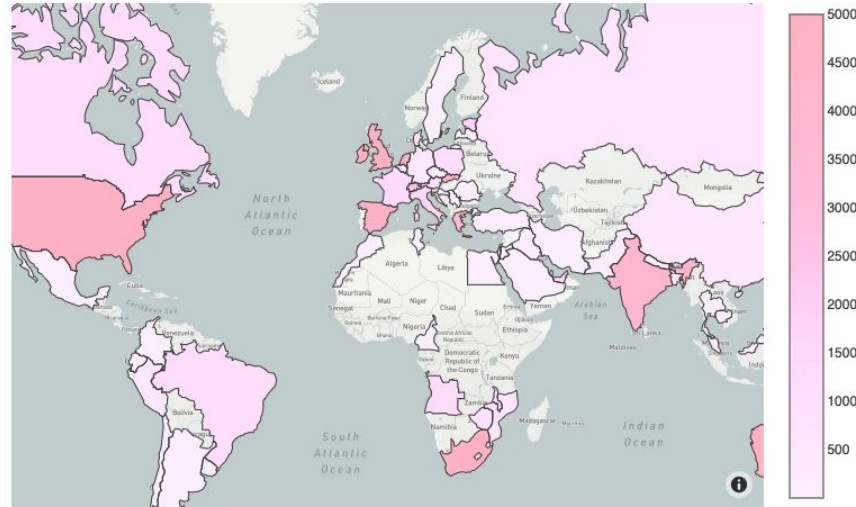
[More... >](#)

<input type="checkbox"/>	DRX276349	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280768	Jap
<input type="checkbox"/>	DRX276352	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280771	Jap
<input type="checkbox"/>	ERX5962095	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	Illumina NovaSeq 6000 paired end sequencing; COG-UK/ALDP-18D5E87/SANG:210714_A00495_0594_AHFTYKDRXY/1t47	Unit
<input type="checkbox"/>	DRX276354	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280773	Jap
<input type="checkbox"/>	ERX5962097	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	Illumina NovaSeq 6000 paired end sequencing; COG-UK/QEUEH-18D71B9/SANG:210714_A00495_0594_AHFTYKDRXY/1t50	Unit
<input type="checkbox"/>	DRX276356	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280775	Jap
<input type="checkbox"/>	ERX5962098	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	Illumina NovaSeq 6000 paired end sequencing; COG-UK/QEUEH-18D708F/SANG:210714_A00495_0594_AHFTYKDRXY/1t52	Unit
<input type="checkbox"/>	DRX276358	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280777	Jap
<input type="checkbox"/>	DRX276359	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280778	Jap
<input type="checkbox"/>	ERX5962105	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	Illumina NovaSeq 6000 paired end sequencing; COG-UK/PLYM-18D784B/SANG:210714_A00495_0594_AHFTYKDRXY/1t68	Unit
<input type="checkbox"/>	DRX276362	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280781	Jap
<input type="checkbox"/>	DRX276363	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280782	Jap
<input type="checkbox"/>	ERX5962107	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	Illumina NovaSeq 6000 paired end sequencing; COG-UK/PLYM-18D76AB/SANG:210714_A00495_0594_AHFTYKDRXY/1t70	Unit
<input type="checkbox"/>	DRX276365	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280784	Jap
<input type="checkbox"/>	ERX5962108	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	Illumina NovaSeq 6000 paired end sequencing; COG-UK/ALDP-18D5DF3/SANG:210714_A00495_0594_AHFTYKDRXY/1t71	Unit
<input type="checkbox"/>	DRX276366	BioSamples (1) <a href="#">↗</a>	<a href="#">See all</a> ▼	NextSeq 550 paired end sequencing of SAMD00280785	Jap

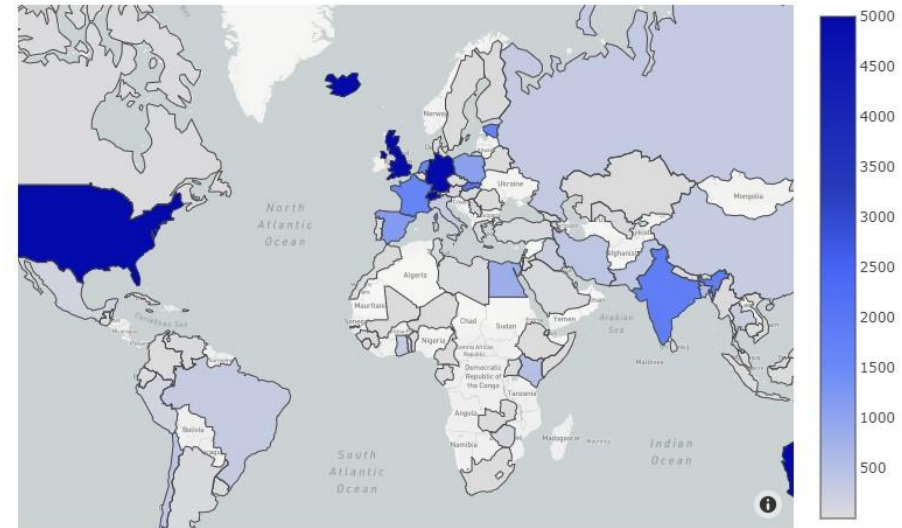
# Improved Stats Page



Raw sequences submitted by Country



Nucleotide sequences submitted by Country



Country	Sequences submitted	Raw sequences submitted
United Kingdom	527,776	886,601
USA	458,557	385,978
Germany	126,362	320
Switzerland	41,050	2,473
Australia	13,311	15,629
Iceland	5,365	N/A

[Download summary data \(.csv\)](#)

Summary data columns: Sequences submitted, Raw Sequences submitted

SUPPORT & FEEDBACK

# Systematic Analysis



- Workflows automatically pick up SARS-CoV-2 raw sequencing reads
- Variant calling, filtered variant calls & consensus sequences

This data is produced as part of the [VEO project](#). For more information, see the [ENA project](#).

Showing 4 of 64 in [Viral sequences](#) > [Systematic Analyses](#)



## Data types

- All (4,591,150)
- Sequences (1,196,055)
- Reference sequences (2)
- Raw reads (1,408,392)
- Sequenced samples (1,372,557)
- Systematic Analyses (600,964)
- Studies (466)
- Genes (22)
- Browser (1)
- Variants (12,691)

[Download](#) [Statistics](#) [Phylogeny](#)

[Edit table view](#)

<input type="checkbox"/>	Run accession	Sample accession	Collection date	Country	Full VCF <a href="#">i</a>	Filtered VCF <a href="#">i</a>	Consensus Sequences <a href="#">i</a>
<input type="checkbox"/>	ERR6134035	ERS7006985	Jan 19, 2021	Germany	Complete <a href="#">+</a>	Complete <a href="#">+</a>	Pending
<input type="checkbox"/>	ERR6134036	ERS7006986	Jan 18, 2021	Germany	Complete <a href="#">+</a>	Complete <a href="#">+</a>	Pending
<input type="checkbox"/>	ERR6134037	ERS7015866	Jan 22, 2021	Germany	Complete <a href="#">+</a>	Complete <a href="#">+</a>	Pending
<input type="checkbox"/>	SRR13649823	SRS8211031	Jun 1, 2020	Germany	Complete <a href="#">+</a>	Pending	Pending

Showing 15 results

[Previous](#) Page 5 of 5 [Next](#)

# Pipelines

- **COVID-19 Sequence Analysis Workflow (ELTE)**

- Processing Illumina raw reads

- Output (per sample):

- Unfiltered VCF
- Filtered VCF (cutoff AF=0.25)
- Consensus sequence

- **Nanopore Analysis Workflow (EMC)**

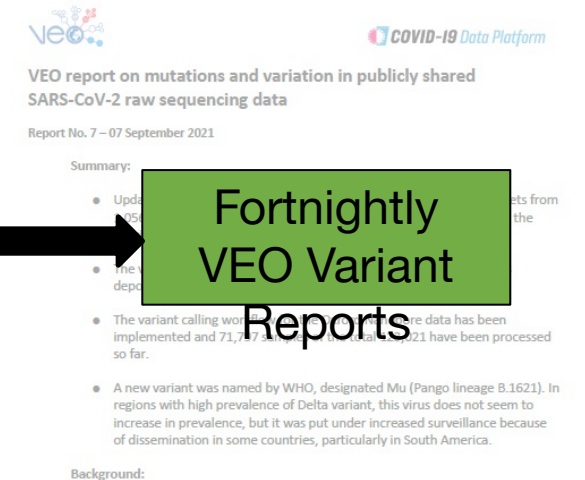
- Processing Nanopore raw reads

- Output (per sample):

- Unfiltered VCF

Fortnightly  
Snapshots

Fortnightly  
VEO Variant  
Reports



# Fortnightly Snapshots and Archival

- Latest: 14<sup>th</sup> Snapshot
  - Analysed: 601,136 samples
- Archival:
  - Analysis umbrella project: [PRJEB45555](#)
    - Unfiltered VCF analysis - [PRJEB43947](#)
    - Filtered variant calls – [PRJEB45554](#)
    - Consensus sequences – [PRJEB45619](#)

The screenshot shows the ENA European Nucleotide Archive interface. At the top, there is a search bar with the text "Enter text search terms" and a search icon. Below it, there are examples of search terms: "h1sone, BN000955" and "Enter accession" with a "View" button. The main content area displays the project title "Project: PRJEB45555" and a description: "All public SARS-CoV-2 INSDC raw read data is systematically analysed to produce a set of uniform variant calls (VCF format) and consensus sequences, using the COVID Sequence Analysis Workflow and the Nanopore Analysis Workflow. Briefly, reads are mapped to the SARS-CoV-2 reference genome, variants are called using LoFreq and annotated using SnpEff. Variant sets are filtered and used to construct consensus sequences." Below this, it states "This project is an umbrella project holding together 3 child projects, each consisting of a different output from our workflow." A list of collaborating institutions is provided: EMBL-EBI, Erasmus Medical Center (EMC), Netherlands, Eötvös Loránd University (ELTE), Hungary, and Technical University of Denmark, Denmark. A "Show Less" button is visible. On the right side, there is a sidebar with options for "View" (XML, XML (STUDY)), "Download" (XML, XML (STUDY)), "Navigation" (Show), and "Component Projects" (Hide).

**ENA**  
European Nucleotide Archive

Enter text search terms Search

Examples: h1sone, BN000955

Enter accession View

Examples: Tixon:9608, BN000955, PRJEB402

Home Submit Search Rulespace About Support

**Project: PRJEB45555**

All public SARS-CoV-2 INSDC raw read data is systematically analysed to produce a set of uniform variant calls (VCF format) and consensus sequences, using the [COVID Sequence Analysis Workflow](#) and the [Nanopore Analysis Workflow](#). Briefly, reads are mapped to the SARS-CoV-2 reference genome, variants are called using [LoFreq](#) and annotated using [SnpEff](#). Variant sets are filtered and used to construct consensus sequences.

This project is an umbrella project holding together 3 child projects, each consisting of a different output from our workflow.

This work is carried out as part of the [Versatile emerging infectious disease observatory \(VEO\) project](#), a collaboration between:

- [EMBL-EBI](#)
- [Erasmus Medical Center \(EMC\), Netherlands](#)
- [Eötvös Loránd University \(ELTE\), Hungary](#)
- [Technical University of Denmark, Denmark](#)

Show Less

**Study Title:** Parent project for all VEO SARS-CoV-2 taskforce analysis work

**Center Name:** EMBL-EBI

**ENA-FIRST-PUBLIC:** 2021-06-07

**ENA-LAST-UPDATE:** 2021-09-03

**View:** XML  
XML (STUDY)

**Download:** XML  
XML (STUDY)

**Navigation:** Show

**Component Projects:** Hide



# Bulk Downloader

- Utility based on Java8 to download directly from ENA FTP
- Downloads data in the following formats:
  - XML, FASTA, EMBL, FASTQ, SUBMITTED
- Download accessions
- Download data using FTP or Aspera
- Create scripts to easily download data and keep the data up to date

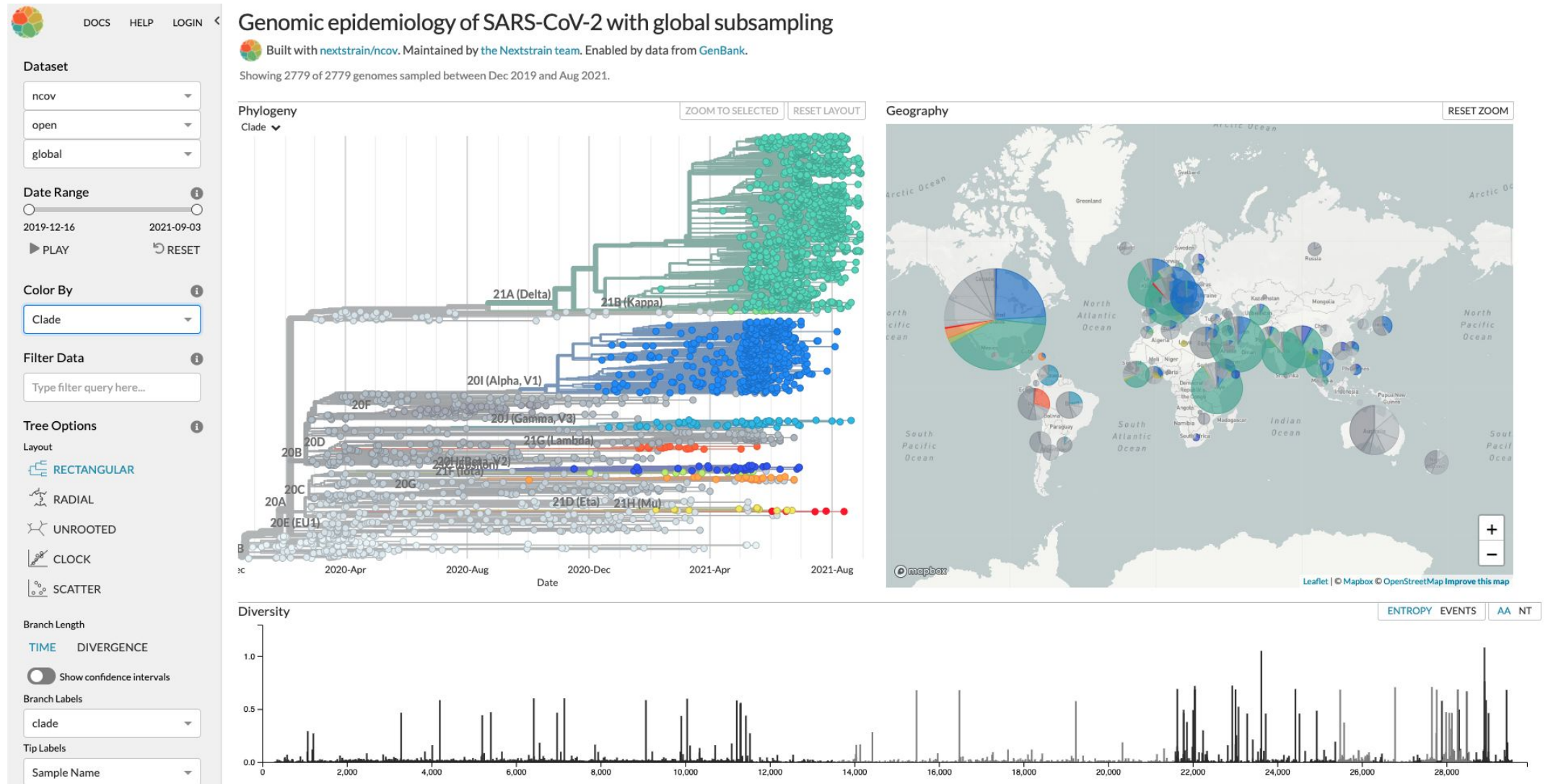
```
#. /////////////// eeeeeee eeeeeeee eeee eee eeee eeeeeeee eee eeeeeee
###. #. /////////////// eeee eeee eee eee eeee eee eeee eeee eeee eeee eeee
## #####. /////////////// eeee eeee eee eee eeee eee eeee eeee eeee eeee eeee
#####. /////////////// eeee eee eee eee eeeeeee eeee eeee eeee eeee eee
(#####. /////////////// eeee eeee eee eee eeeeeee eeee eeeeeeee eeee eeee eee
#####. /////////////// eeeeeee eeeeeeee eeeee eeee eeeeeeee eeee eeeeeee
## #####. /////////////// ##### # ##### # ##### ## ##### # #
### #####. /////////////// # # # # # # # # # # # # # # # # # # #
###. #. /////////////// # # # # # # # ##### # # # # # # # # # #
#. /////////////// ##### # # # # # # # # # # # # # # # # # # # # # #

(Copyright © EMBL 2021)

Welcome to the Covid-19 Data Portal's data downloader utility!

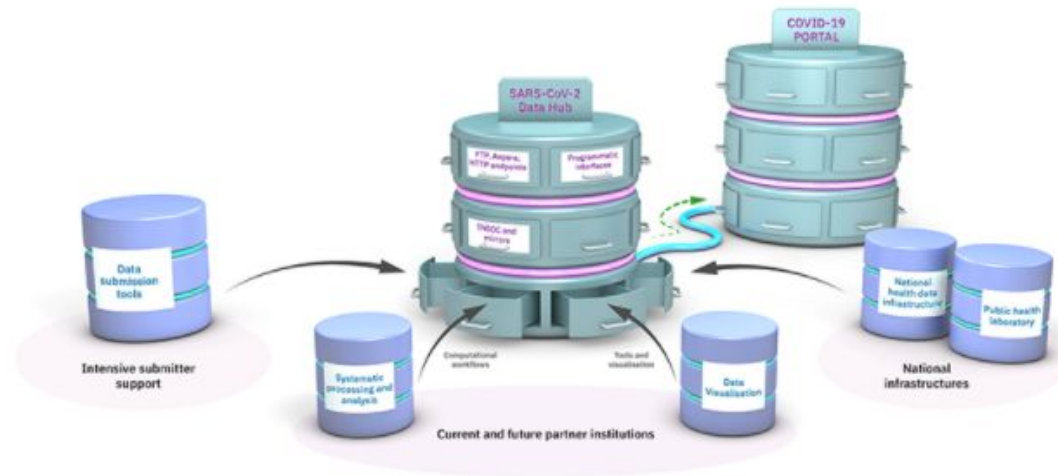
Select from the options below:
-----
For Viral Sequences enter 1
For Host Sequences enter 2
For Help enter 3
For Privacy Notice enter 4
To exit enter 0 (zero)
-----
```

# NextStrain Visualisation

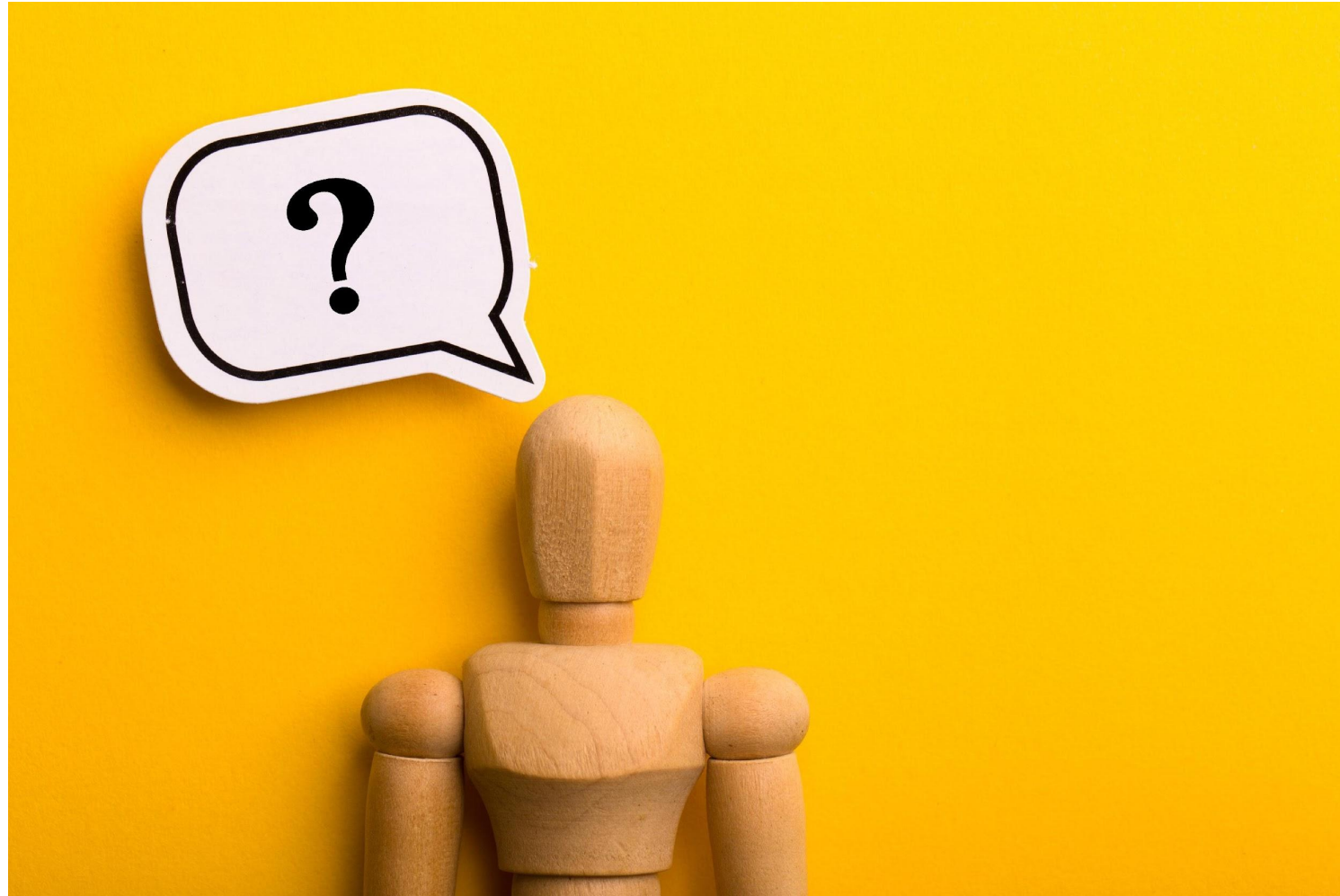


# Private Data Hubs

- Reference-based mapping workflows (CSAW, NAW)
- Evergreen phylogenetic tree (without visualisation)
- Nextstrain reports (including phylogeny and other plots)
- Lineage classifications



# Questions?



[ocathail@ebi.ac.uk](mailto:ocathail@ebi.ac.uk)

[virus-dataflow@ebi.ac.uk](mailto:virus-dataflow@ebi.ac.uk)