**EMBL-EBI**

**ENA** European Nucleotide Archive

# Automating a SARS-CoV-2 ENA submission tool with Nextflow

Zahra Waheed and Ahmad Zyoud

European Nucleotide Archive
European Molecular Biology Laboratory - European Bioinformatics Institute

## Who are the European Nucleotide Archive (ENA)?
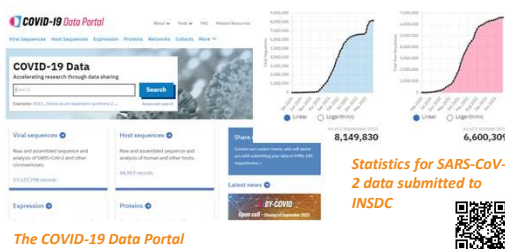
The ENA:

- Is a global, open-access nucleotide sequence repository.

- Covers non-sensitive raw sequence data, sequence assembly information and functional annotation for all organisms.

- Is the European arm of the International Nucleotide Sequence Database Collaboration (INSDC), which includes NCBI and DDBJ

- Mirrors data between NCBI and DDBJ

## The ENA & SARS-CoV-2

During the COVID-19 pandemic an unprecedented volume of SARS-CoV-2 data was submitted to the ENA. SARS-CoV-2 submissions now make up ~24% of all ENA raw sequence reads, which have been shared from 105 countries.

## The COVID-19 Data Portal

- All public SARS-CoV-2 sequences and raw reads submitted to the ENA feed into the COVID-19 Data Portal

- This was launched by EMBL-EBI in April 2020 in response to the pandemic

- The Portal brings together a variety of SARS-CoV-2 related biological datatypes and visualisation tools for global access and analysis

*Statistics for SARS-CoV-2 data submitted to INSDC*

8,149,830      6,600,309

*The COVID-19 Data Portal*

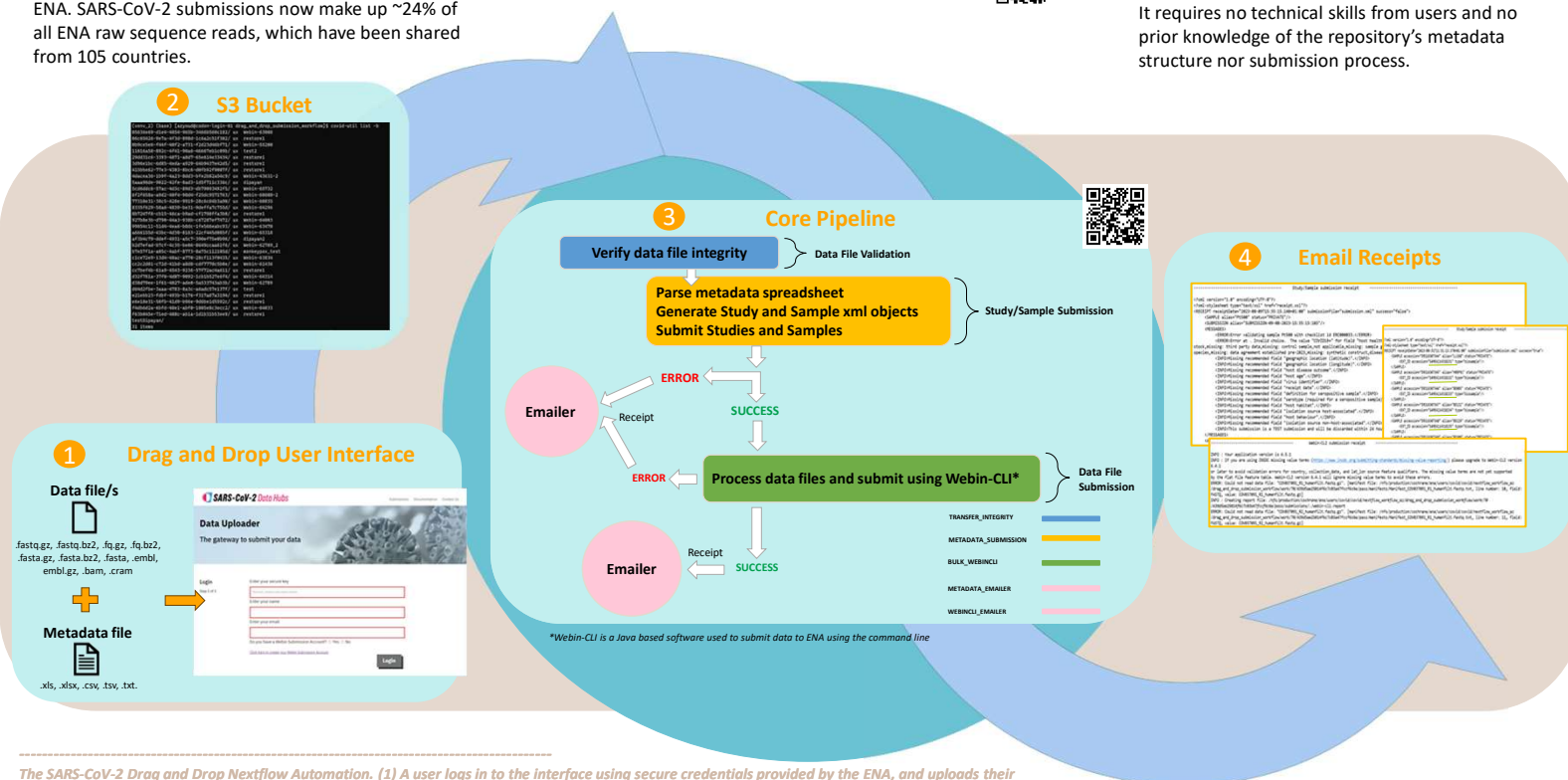## The SARS-CoV-2 Drag and Drop Submission Tool

All submissions to the COVID-19 Data Portal must follow the ENA's metadata structure below:

with multiple submission routes possible for each object.

To help maximise the rate and volume of SARS-CoV-2 data shared to the Portal, the Drag and Drop Uploader Tool was developed by EBI (ENA and Archive Infrastructure and Technology (AIT)) teams as a simpler alternative to existing ENA submission services.
It requires no technical skills from users and no prior knowledge of the repository's metadata structure nor submission process.

### ② S3 Bucket

### ① Drag and Drop User Interface

**Data file/s**

.fastq.gz, .fastq.bz2, .fq.gz, .fq.bz2, .fasta.gz, .fasta.bz2, .fasta, .embl, embl.gz, .bam, .cram

**Metadata file**

.xls, .xlsx, .csv, .tsv, .txt.

### ③ Core Pipeline

**Verify data file integrity** — Data File Validation

**Parse metadata spreadsheet
Generate Study and Sample xml objects
Submit Studies and Samples** — Study/Sample Submission

Emailer — Receipt — **ERROR** / **SUCCESS**

**Process data files and submit using Webin-CLI*** — Data File Submission

Emailer — Receipt — **ERROR** / **SUCCESS**

TRANSFER_INTEGRITY
METADATA_SUBMISSION
BULK_WEBINCLI
METADATA_EMAILER
WEBINCLI_EMAILER

*Webin-CLI is a Java based software used to submit data to ENA using the command line*

### ④ Email Receipts

*The SARS-CoV-2 Drag and Drop Nextflow Automation. (1) A user logs in to the interface using secure credentials provided by the ENA, and uploads their data files and a metadata spreadsheet (with the file extensions specified) to a specific Amazon S3 Bucket (2). (3) All data is then transferred to an EBI High Performance Compute environment and submitted through the core pipeline. (4) Notifications of a successful or unsuccessful ENA submission will be received by email.*

## Pipeline Directory Structure

### Backend Pipeline

- S3 Bucket transfer
- File integrity validation
- Metadata parsing
- Study and sample metadata submission
- Data file submission
- Email notification

### Standalone Pipeline

- Metadata parsing
- Study and sample metadata submission
- Data file submission
- Email notification

## Limitations

- Workflow does not support 're-runs' - cannot be run from a specific stage in case of errors / partial submissions
- The tool submit only a single data type (Raw reads or genome assemblies) per run
- Currently specific to SARS-CoV-2 data

### Feature improvements

- Support re-run of pipeline from specific step, e.g.:
  - Transfer > submit metadata > submit data > emailer
- Submission of both raw reads and genome assemblies in one Nextflow run
- Containerisation with Docker / Singularity
- Refactor hard-coded script variables to support other small pathogen submissions to ENA