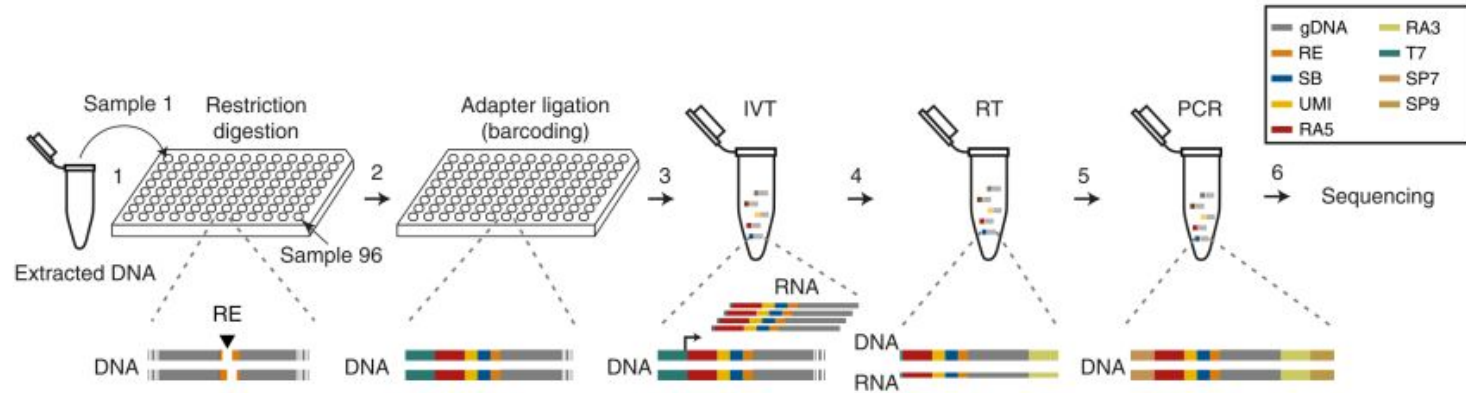


Genomic surveillance of SARS-CoV-2 using COVseq: making it FAIR

Luuk Harbers
2021-12-09

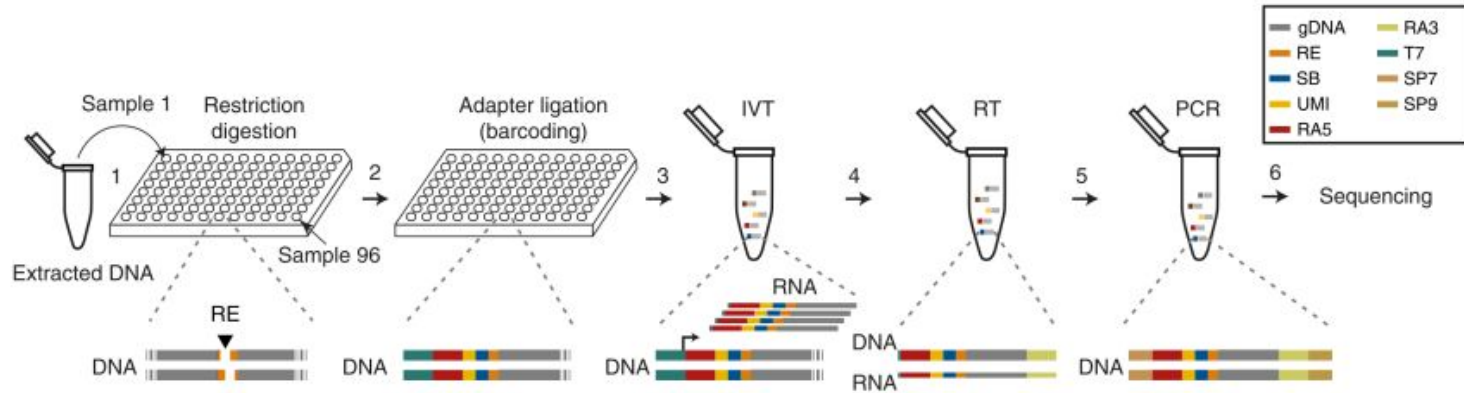
Background

- CUTseq



Background

- CUTseq



- Pandemic

- Adaptation of CUTseq → COVseq



ARTICLE



<https://doi.org/10.1038/s41467-021-24078-9>

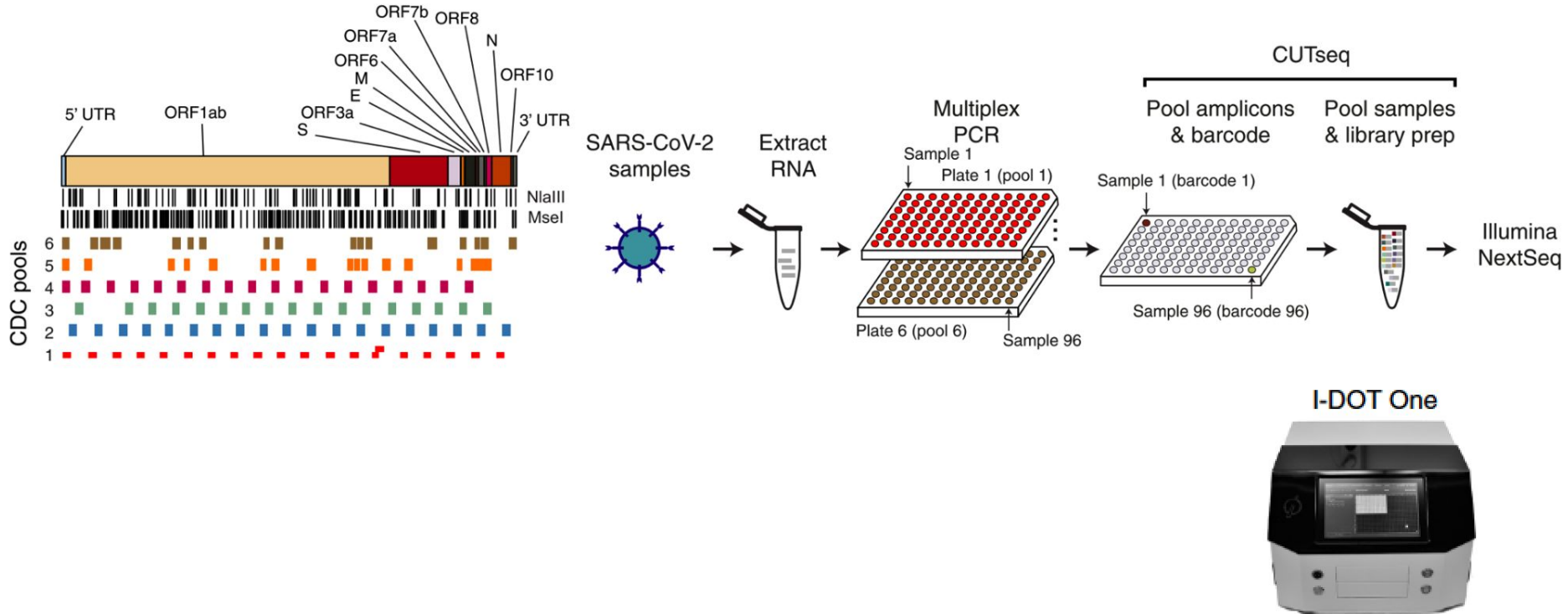
OPEN

COVseq is a cost-effective workflow for mass-scale SARS-CoV-2 genomic surveillance

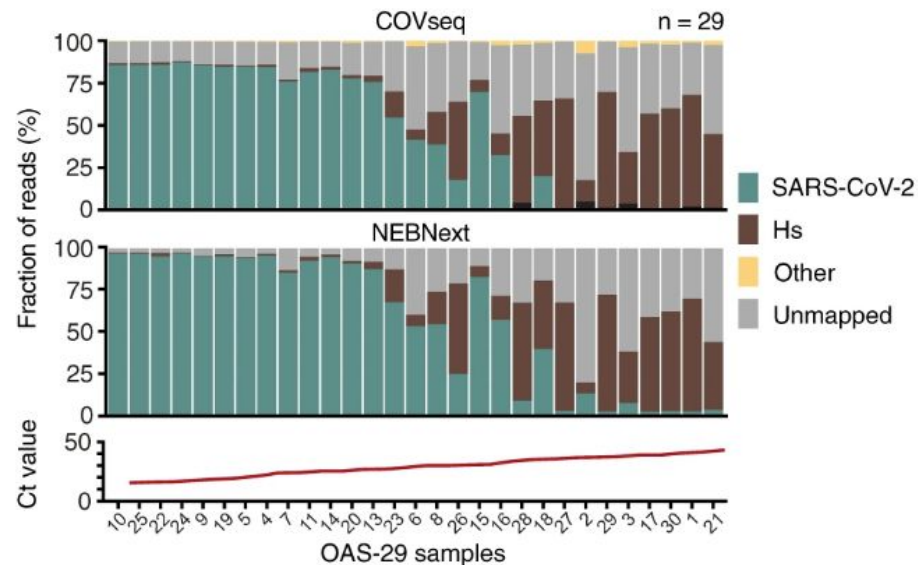
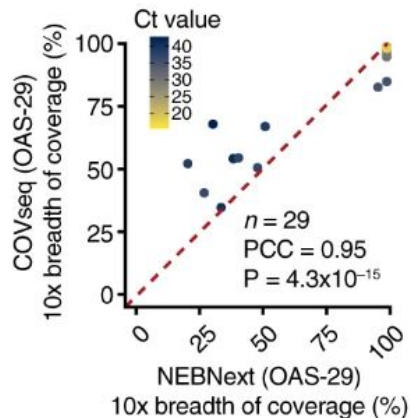
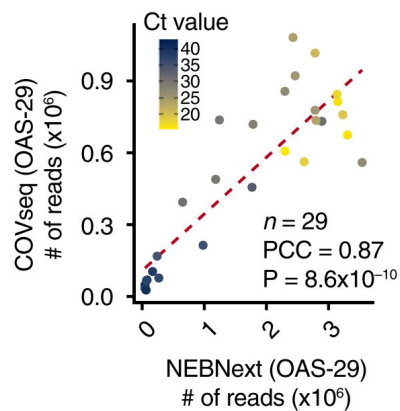
Michele Simonetti ^{1,2,7}, Ning Zhang ^{1,2,3,7}, Luuk Harbers ^{1,2,7}, Maria Grazia Milia⁴, Silvia Brossa⁵, Thi Thu Huong Nguyen ^{1,2}, Francesco Cerutti ⁴, Enrico Berrino^{5,6}, Anna Sapino^{5,6}, Magda Bienko ^{1,2}, Antonino Sottile⁵, Valeria Ghisetti ⁴✉ & Nicola Crosetto ^{1,2}✉

Michele Simonetti, Ning Zhang and Luuk Harbers are equally contributing authors

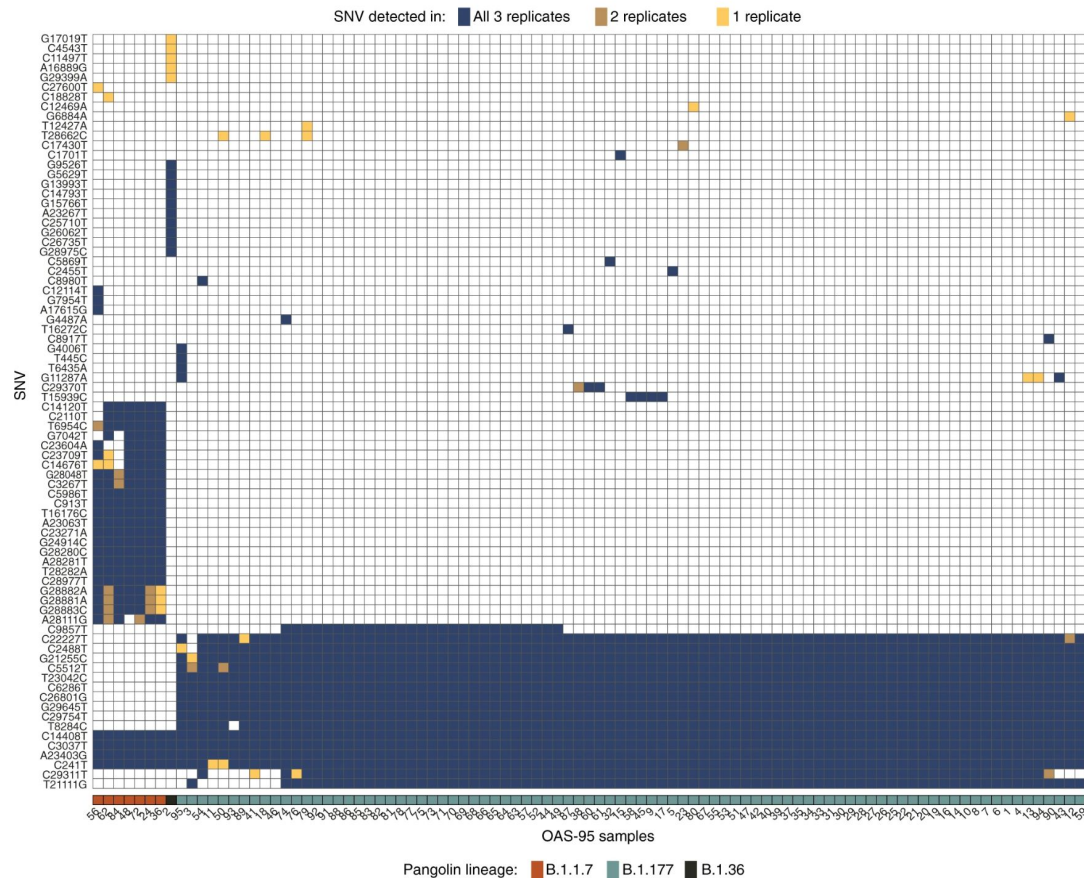
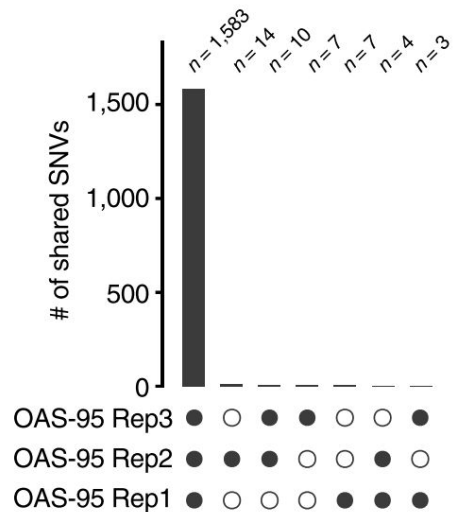
COVseq workflow



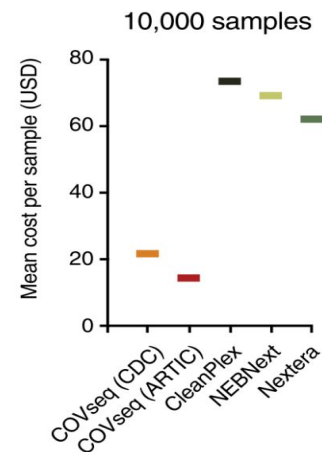
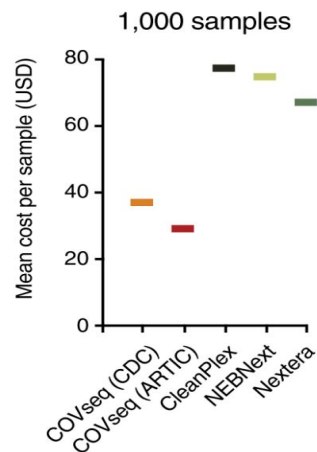
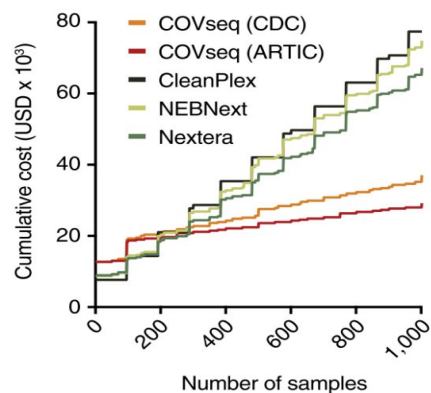
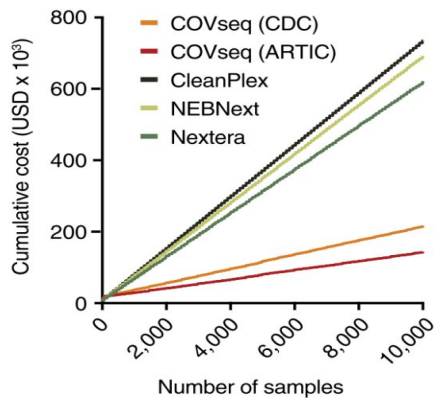
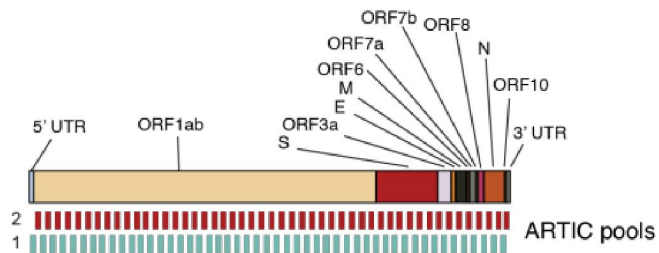
Results



Results



Results



FAIR principles in COVseq

- **F**indable
- **A**ccessible
- **I**nteroperable
- **R**eusable

FAIR principles in COVseq

Data availability

The BAM files used to generate all the plots in the main Figures and Supplementary Figures have been deposited to the European Nucleotide Archive (ENA) and are available at the following link: <https://www.ebi.ac.uk/ena/browser/view/PRJEB42601>. All reference sequences used in this study are listed in Supplementary Table [2](#). All the GISAID data used in this study are described in Supplementary Data [7](#) and are available at <https://www.gisaid.org>.

FAIR principles in COVseq

- 691 bam files from 274 samples
 - Sequencing related metadata
 - Sample related metadata

Hide Column Selection

<input type="checkbox"/> base_count	<input type="checkbox"/> broker_name	<input type="checkbox"/> center_name
<input type="checkbox"/> cram_index_aspera	<input type="checkbox"/> cram_index_ftp	<input type="checkbox"/> cram_index_galaxy
<input checked="" type="checkbox"/> experiment_accession	<input type="checkbox"/> experiment_alias	<input type="checkbox"/> experiment_title
<input type="checkbox"/> fastq_aspera	<input type="checkbox"/> fastq_bytes	<input type="checkbox"/> fastq_ftp
<input type="checkbox"/> fastq_galaxy	<input type="checkbox"/> fastq_md5	<input type="checkbox"/> first_created
<input type="checkbox"/> first_public	<input type="checkbox"/> instrument_model	<input type="checkbox"/> instrument_platform
<input type="checkbox"/> last_updated	<input type="checkbox"/> library_layout	<input type="checkbox"/> library_name
<input type="checkbox"/> library_selection	<input type="checkbox"/> library_source	<input type="checkbox"/> library_strategy
<input type="checkbox"/> nominal_length	<input type="checkbox"/> nominal_sdev	<input type="checkbox"/> read_count
<input checked="" type="checkbox"/> run_accession	<input type="checkbox"/> run_alias	<input checked="" type="checkbox"/> sample_accession
<input type="checkbox"/> sample_alias	<input type="checkbox"/> sample_title	<input checked="" type="checkbox"/> scientific_name
<input type="checkbox"/> secondary_sample_accession	<input type="checkbox"/> secondary_study_accession	<input type="checkbox"/> sra_aspera
<input type="checkbox"/> sra_bytes	<input checked="" type="checkbox"/> sra_ftp	<input type="checkbox"/> sra_galaxy
<input type="checkbox"/> sra_md5	<input checked="" type="checkbox"/> study_accession	<input type="checkbox"/> study_alias
<input type="checkbox"/> study_title	<input type="checkbox"/> submission_accession	<input type="checkbox"/> submitted_aspera
<input type="checkbox"/> submitted_bytes	<input type="checkbox"/> submitted_format	<input checked="" type="checkbox"/> submitted_ftp
<input type="checkbox"/> submitted_galaxy	<input type="checkbox"/> submitted_md5	<input checked="" type="checkbox"/> tax_id

Select default

Download report: [JSON](#) [TSV](#) [Download Files as ZIP](#) [Download selected files](#)

Study Accession	Sample Accession	Experiment Accession	Run Accession	Tax Id	Scientific Name	Submitted files: FTP	SRA fi
PRJEB42601	SAMEA8620010	ERX3486380	ERR5779464	2697049	Severe acute respiratory syndrome coronavirus 2	<input type="checkbox"/> MS45.trim.sorted.bam	h..
PRJEB42601	SAMEA8620011	ERX3486381	ERR5779465	2697049	Severe acute respiratory syndrome coronavirus 2	<input type="checkbox"/> MS47.trim.sorted.bam	h..

AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ
geographic location (latitude)	geographic location (longitude)	geographic location (region and locality)	host disease outcome	host age	virus identifier	definition for seropositive sample	serotype (required for seropositive sample)	host habitat	isolation source host-associated
R	R	R	R	R	R	R	R	R	R
52.2053° N	0.1218° E	Hinxton-Cambridgeshire	recovered	10	181 60 JJUL			wild	kidney cell line vsgp e6
DD	DD			years					

FAIR principles in COVseq

Code availability

All the custom code used for processing COVseq sequencing data and the custom MATLAB code used in the Cost Analysis (see [Supplementary Notes](#)) is available at <https://github.com/ljwharbers/COVseq> and the repository is linked to Zenodo at the following link: <https://doi.org/10.5281/zenodo.4776499>.

FAIR principles in COVseq

- All code available on github with a release on zenodo
- Readme file including extra information with how to run the preprocessing and analyses

For COVseq libraries make sure all the paths to the required scripts in the config are correct. Build the python/cython library using `$ python setup.py build_ext --inplace`, dependencies for demultiplexing are: pandas, argparse and pysam.

Preparation for preprocessing COVseq libraries

For the demultiplexing of fastq files a custom python script is used. Main input required for this script is the fastq file, a list of barcodes used (no column names, just one barcode per row), the length of the barcode (default 8) and the number of mismatches allowed. To filter out reads that map too far from cutsites you also need a bed file with all the cut site locations in the genome and the read length. For the ivar pipeline the only extra file that is required is the file with the primer locations used for the amplicon sequencing.

Demultiplexing COVseq libraries

Since COVseq libraries are multiplexed libraries (multiple samples in one library). These libraries need to be demultiplexed. To demultiplex COVseq fastq files you can run the script in the `Demultiplex` folder.

Build the python/cython library using `$ python setup.py build_ext --inplace` and make sure the following dependencies are met: `pandas`, `argparse` and `pysam`.

Following this you can run the demultiplexing. An example command would be: `$ demultiplex_withcython.py -f {fastq1} -f2 {fastq2} --paired -o {output} -l {logfile-output} -b {list-of-barcodes} -m {mismatches}` With the {list-of-barcodes} being a text file with 1 COVseq barcode per line (no headers). For more information regarding the different commands you can run `demultiplex_withcython.py --help`.

Running FastQ-Screen

To get information regarding to the mapping percentages of sample specific fastq files (output of demultiplexing step) we used FastQ-Screen (version 0.14.1) https://www.bioinformatics.babraham.ac.uk/projects/fastq_screen/.

Simply add another entry in the database with the SARS-CoV-2 reference genome and run with default settings.

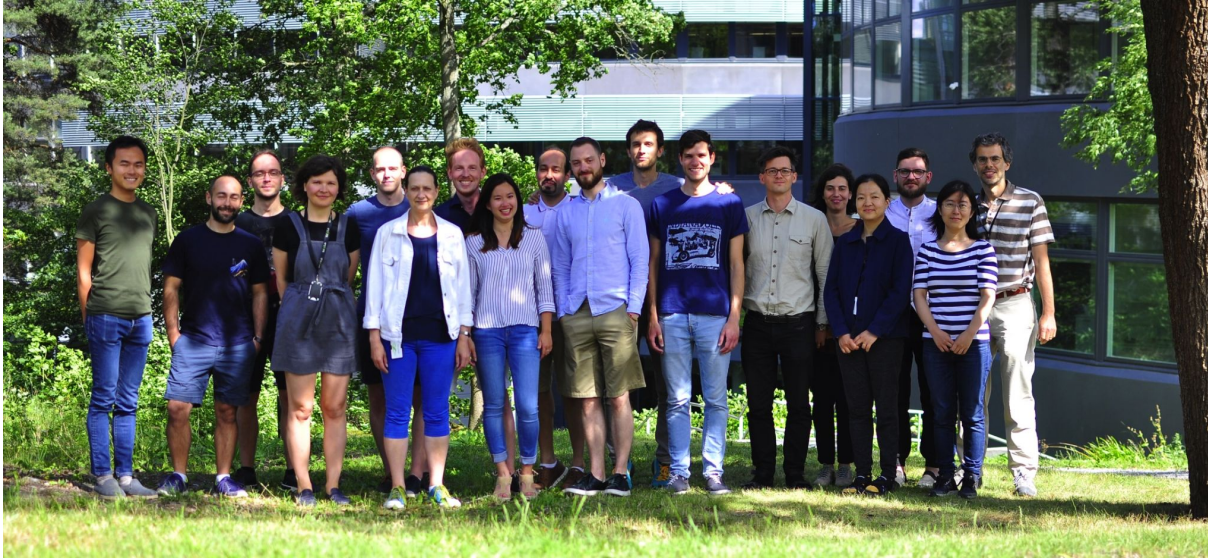
Due to sensitive patient information we do not share any patient specific fastq files.

Running nf-co.re/viralrecon

For further processing of fastq files we used the nextflow based pipeline from nf-core called viralrecon (version 1.1.0) <https://nf-co.re/viralrecon/1.1.0>. For any extra information or troubleshooting please check out their website and/or join the slack channel for the specific pipeline.

Commands we used to run this pipeline are as follows: `$ nextflow run nf-core/viralrecon --input {samplesheet.csv} --genome 'NC_045512.2' --fasta {sarscov2-fastafile} --save_reference --protocol amplicon --amplicon_bed {ampliconbedfile} --skip_assembly --skip_markduplicates --skip_mosdepth --callers ivar --outdir {outdir} -profile docker --max-cpus 40 -r 1.1.0`

Acknowledgements



BiCro Lab:

Nicola Crosetto
Michele Simonetti
Ning Zhang

Collaborators:

Maria Grazia Milia
Silvia Brossa
Francesco Cerutti
Enrico Berrino

