# Training course in Bioinformatic tools
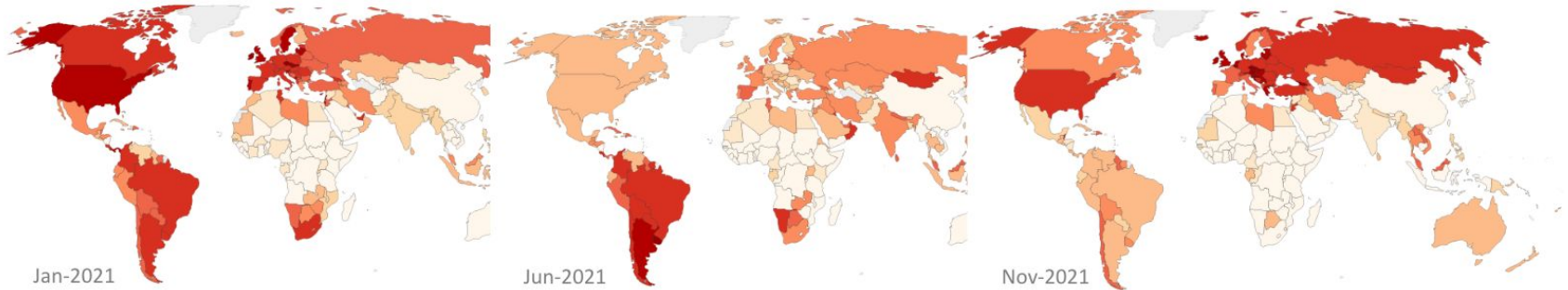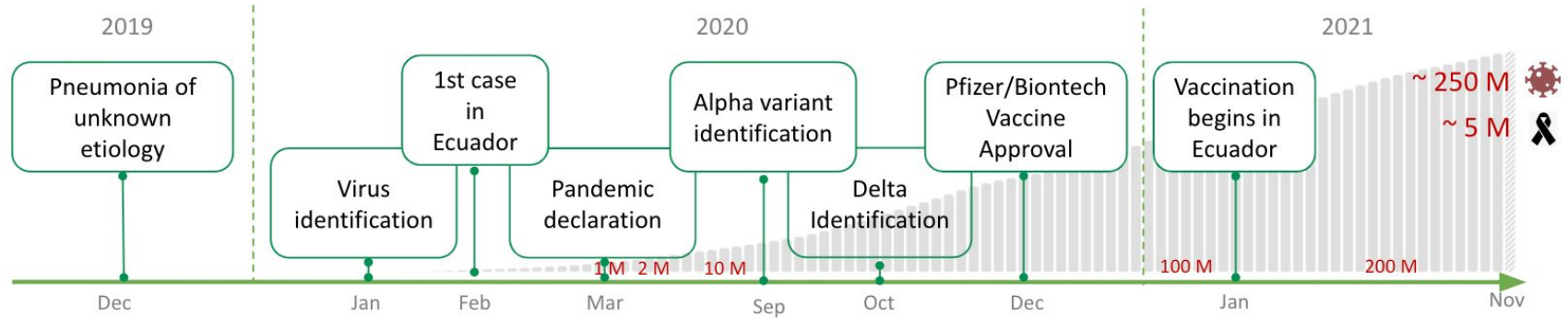
Instituto de Microbiología
Centro de Bioinformática

**Virtual training**
Module 3: SARS-CoV-2 analysis
 May 2023
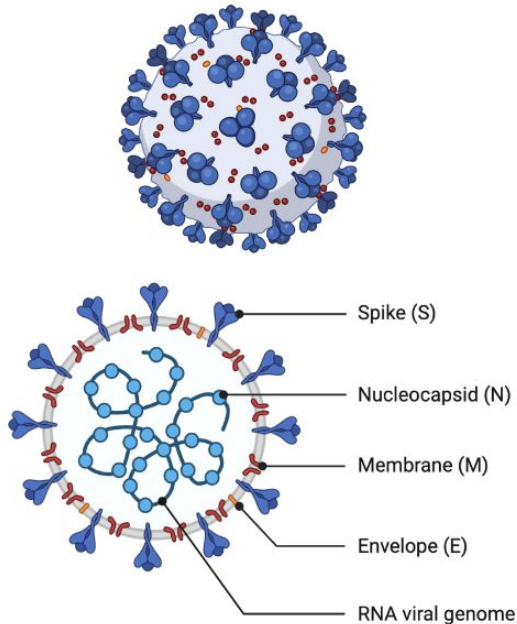
# COVID-19 pandemics



2019       2020       2021

Pneumonia of unknown etiology

1st case in Ecuador

Alpha variant identification

Pfizer/Biontech Vaccine Approval

Vaccination begins in Ecuador

~ 250 M

~ 5 M

Virus identification

Pandemic declaration

Delta Identification

1 M   2 M    10 M        100 M       200 M

Dec    Jan   Feb   Mar    Sep   Oct    Dec     Jan     Nov
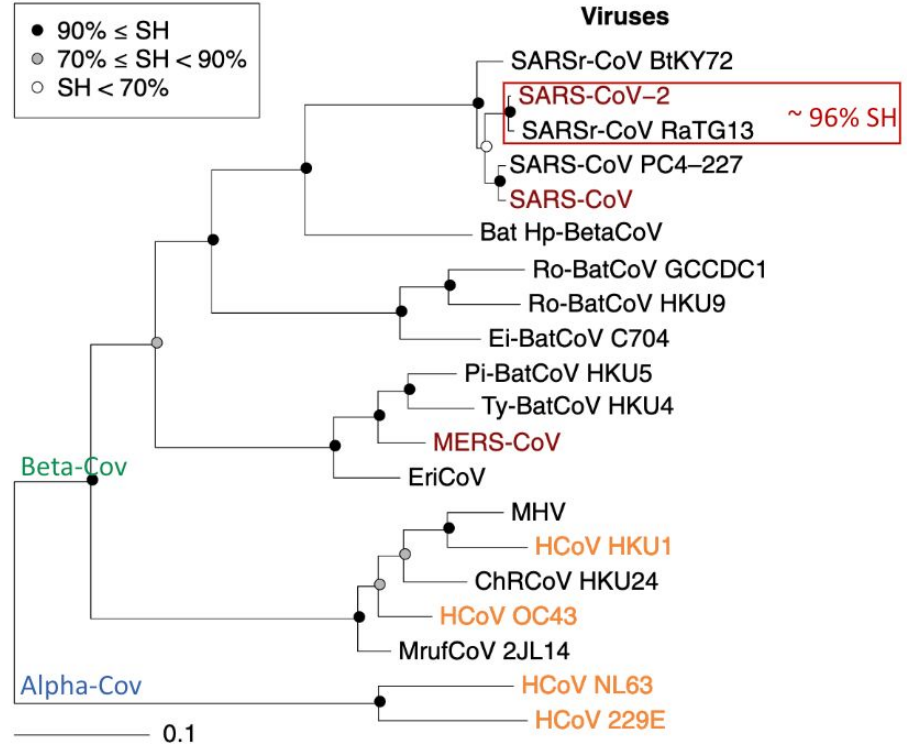
Jan-2021      Jun-2021      Nov-2021

[1] American Society for Microbiology. (2020). *COVID-19 Resources*. [2] WHO. (2021). *COVID-19 Dashboard*. [3] Ritchie *et al.* (2021). *"Coronavirus Pandemic (COVID-19)"* - *OurWorldInData.org.*

2

# SARS-CoV-2 virus



Family: *Coronaviridae*
Genre: *Betacoronavirus*
Subgenre: *Sarbecovirus*

Spike (S)
Nucleocapsid (N)
Membrane (M)
Envelope (E)
RNA viral genome

Source: BioRender.com

**Viruses**

- 90% ≤ SH
- 70% ≤ SH < 90%
- SH < 70%

SARSr-CoV BtKY72
SARS-CoV-2
SARSr-CoV RaTG13    ~ 96% SH
SARS-CoV PC4-227
SARS-CoV
Bat Hp-BetaCoV
Ro-BatCoV GCCDC1
Ro-BatCoV HKU9
Ei-BatCoV C704
Pi-BatCoV HKU5
Ty-BatCoV HKU4
MERS-CoV
EriCoV
Beta-Cov
MHV
HCoV HKU1
ChRCoV HKU24
HCoV OC43
MrufCoV 2JL14
Alpha-Cov
HCoV NL63
HCoV 229E
0.1

Source: Gorbalenya *et al.* (2020). *Nature Microbiology.*

3

# SARS-CoV-2 virus



**Mutation rate of SARS-CoV-2**
$3 \times 10^{-6}$/nt/cycle[2]
$2 - 3$ nt/month

Source: Duffy *et al.* (2018). *PLOS Biology;* [2] Borges *et al.* (2021).
*Evolutionary Biology*



Source: Nextstrain. (2021). *Genomic epidemiology of novel coronavirus.*

# SARS-CoV-2 virus

| WHO | PANGO Lineage | RBD | Furin site | Other | Feature |
|------|------|------|------|------|------|
| - | B.1 | - | - | D614G | ↑ transmission |
| Alpha | B.1.1.7 | (E484K), (S494P), N501Y | P681H | D614G, A570D, T716I, S982A, D1118H, (K1991N) | ↑ transmission, ↑ viral load, moderate ab evasion |
| Beta | B.1.351 | E484K, K417N, N501Y | - | D614G, A701V | ↑ transmission, moderate ab evasion |
| Gamma | P.1 | E484K, K417T, N501Y | - | D614G, H655Y, T1027I | ↑ transmission, moderate ab evasion, ↓ vaccine efficiency |
| Delta | B.1.617.2 | (K417N), L452R, T487K | P681R | D614G, D950N | ↑ transmission, moderate ab evasion, ↓ vaccine efficiency |

Source: Jackson *et al.* (2021). *Nature Reviews.*



Spike protein
RBD
S1
Furin site
S2



Source: Hoffman *et al.* (2020). *Molecular Cell.*

5

# SARS-CoV-2 lineages and variants

# SARS-CoV-2 lineages and variants



Source: Hodcroft, 2022


United Kingdom

| 10 Apr 2023 - 24 Apr 2023 | | |
|---|---|---|
| Variant | Num seq | Freq |
| 23A (Omicron) | 466 | 0.64 |
| 22F (Omicron) | 170 | 0.24 |
| 22D (Omicron) | 56 | 0.08 |
| 23B (Omicron) | 17 | 0.02 |

# Genomic surveillance



**What is GENOMIC SURVEILLANCE?**

It involves constantly monitoring pathogens...

**AND**

Analyzing their similarities and differences

**HELPING US TO:**

Monitor diseases

Control pathogens

Tailor interventions and recommendations for the public

Develop countermeasures, like vaccines

Stamp out disease

**The Global Genomic Surveillance Strategy**
for Pathogens with Pandemic and Epidemic Potential

**World Health Organization**

8

# Workflow



**Collection of samples**

- Private and public
- All provinces in Ecuador
- With and without symptoms

- NP swabs
- DNA/RNA Shield
- Informed consent

**Transport**

- At 4 °C

Instituto de Microbiología
Universidad San Francisco de Quito

**Diagnostic**

- RT-qPCR
- Testing N and S genes
- Interpretation based on Ct-value

**Storage**

- Positive samples
- At -80 °C

# Workflow



**a. RNA extraction**
- Storage in DNA/RNA Shield
- Extraction of viral RNA
- RNA to cDNA

**b. PCR and sequencing library**

29.8kb amplified
- Multiplex-PCR: 98 pairs of primers (V1 Y V3)
- Electrophoresis
- Qubit quantification
- Library: barcodes and adapters

**c. Sequencing**
- MinION flow-cell priming
- MinKNOW: base-calling, demultiplexing.
- Epi2Me - Medaka: consensus sequence

# RNA to cDNA

# Amplicon-based sequencing



**Genomic or Template DNA**

**Gene-specific Primer**

Gene X

Primer

**PCR Amplification**

...TGAACCATTGTTCAATATCG...
T
T
T
T
T
T

**Alignment**

**Sequencing**

12

# Amplicon-based sequencing

# SARS-CoV-2 whole-genome primers

New England Biolabs

# SARS-CoV-2 whole-genome primers

ARTIC network

# Agarose gel electrophoresis

# Workflow



a. RNA extraction
- Storage in DNA/RNA Shield
- Extraction of viral RNA
- RNA to cDNA

b. PCR and sequencing library
29.8kb amplified
- Multiplex-PCR: 98 pairs of primers (V1 Y V3)
- Electrophoresis
- Qubit quantification
- Library: barcodes and adapters

c. Sequencing
- MinION flow-cell priming
- MinKNOW: base-calling, demultiplexing.
- Epi2Me - Medaka: consensus sequence

# Library preparation - Nanopore



End repair and dA-tailing

Ligation of barcodes

BC

Pool samples and clean up

Ligation of sequencing adaptors

Clean up

Legend: RNA, DNA, Barcode (BC), Adaptor, Motor protein

BCxx: Barcode xx

# Flowcell loading

# MinKNOW software

# EPI2ME software

# ARTIC Pipeline

# Workflow



Lineage assignment

Assignment of clades and lineages

Nextclade
GISAID

Mutations and substitutions

Phylogeny

Alignment

MAFFT version 7

- Wuhan-Hu-1 reference genome
- Aliview: check alignment

Phylogenetic tree construction

IQ-TREE

- Maximum Likelihood, GTR, 1000x
- Visualization in iTOL

Statistics

R Studio

# Nextclade - Nextstrain

# Nextclade

# Nextclade

# Nextclade

Frequencies (colored by Clade )



- 20I (Alpha, V1) or B.1.1.7
- 20H (Beta, V2) or B.1.351
- 20J (Gamma, V3) or P.1
- 21A (Delta) or B.1.617.2
- o Omicron : 21K (Omicron) or BA.1, 21L (Omicron) or BA.2, 22A (Omicron) or BA.4, 22B (Omicron) or BA.5, 22C (Omicron) or BA.2.12.1, 22D (Omicron) or BA.2.75, 22E (Omicron) or BQ.1, 22F (Omicron) or XBB, 23A (Omicron) or XBB.1.5, 23B (Omicron) or XBB.1.16

27

# Pangolin COVID-19 Lineage assigner

## pangolin

**Phylogenetic Assignment of Named Global Outbreak Lineages**

Pangolin was developed to implement the dynamic nomenclature of SARS-CoV-2 lineages, known as the Pango nomenclature. It allows a user to assign a SARS-CoV-2 genome sequence the most likely lineage (Pango lineage) to SARS-CoV-2 query sequences.

It is available as a command line tool and a web application. The web application was developed by the Centre for Genomic Pathogen Surveillance. The command line tool is open source software available under the GNU General Public License v3.0.

# Pangolin COVID-19 Lineage assigner

| | File name | Sequence name | Lineage | Assignment probability |
|---|---|---|---|---|
| | Retry Failed Sequences | Reset entries | Upload another file | Help |

| | File name | Sequence name | Lineage | Assignment probability |
|---|---|---|---|---|
| **— FAILED (Click warning icon for more info) 2 sequences** | | | | |
| ⊙ | cluster.fasta | EDB003 | | |
| ⊙ | cluster.fasta | EDB004 | | |
| | | | | |
| **— ANALYSED (Click tick icon for more info) 8 sequences** ⬇ | | | | |
| ✓ | cluster.fasta | EDB001 | B.1.1.65 🏴🌐 | 1.0 |
| ✓ | cluster.fasta | EDB002 | B.1.1.65 🏴🌐 | 1.0 |
| ✓ | cluster.fasta | EDB005 | B 🏴🌐 | 1.0 |
| ✓ | cluster.fasta | EDB006 | B.1.1.65 🏴🌐 | 1.0 |
| ✓ | cluster.fasta | EDB007 | B.1.1.65 🏴🌐 | 1.0 |
| ✓ | cluster.fasta | EDB008 | B.1.1.65 🏴🌐 | 1.0 |
| ✓ | cluster.fasta | EDB009 | B.1.1.65 🏴🌐 | 1.0 |
| ✓ | cluster.fasta | EDB010 | B 🏴🌐 | 1.0 |

Pangolin (version v2.0.7, lineages version 2020-08-29) is built by Áine, JT, Verity, Emily and Andrew. Web Application by
Centre for Genomic Pathogen Surveillance

# GISAID database