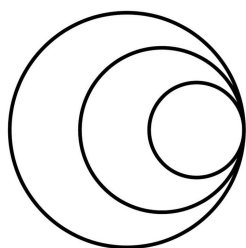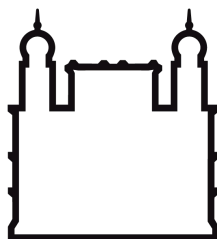# Conceitos sobre montagem e anotação de variantes virais

wellcome
connecting
science

Ministério da Saúde

FIOCRUZ
**Fundação Oswaldo Cruz**
Instituto Gonçalo Moniz
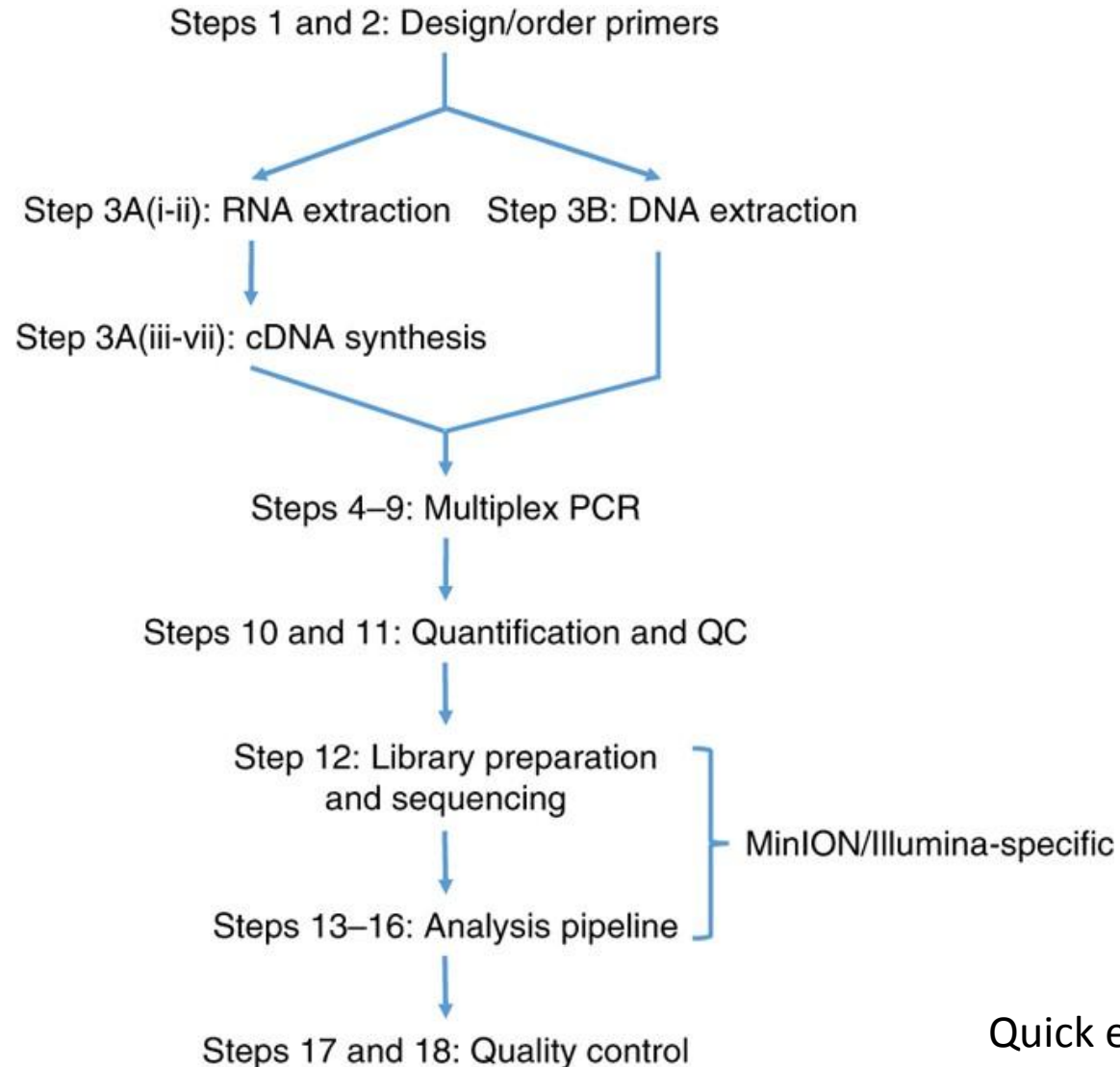
**Túlio Campos, PhD**

**Instituto Aggeu Magalhães**

**Fiocruz Pernambuco**

**tulio.campos@fiocruz.br**

# Informações contidas dos genomas

- De onde surgiu o vírus?

- Como adquiriu a capacidade de infectar humanos?

- Como definir diferentes variantes/cepas?

- Existem variantes/cepas mais infecciosas ou patogênicas?

- Como se espalhou pelo mundo (epidemiologia)?

- Como identificar epítopos para vacinas eficazes?

- Porque é que as vacinas são mais/menos eficazes para algumas variantes/cepas?
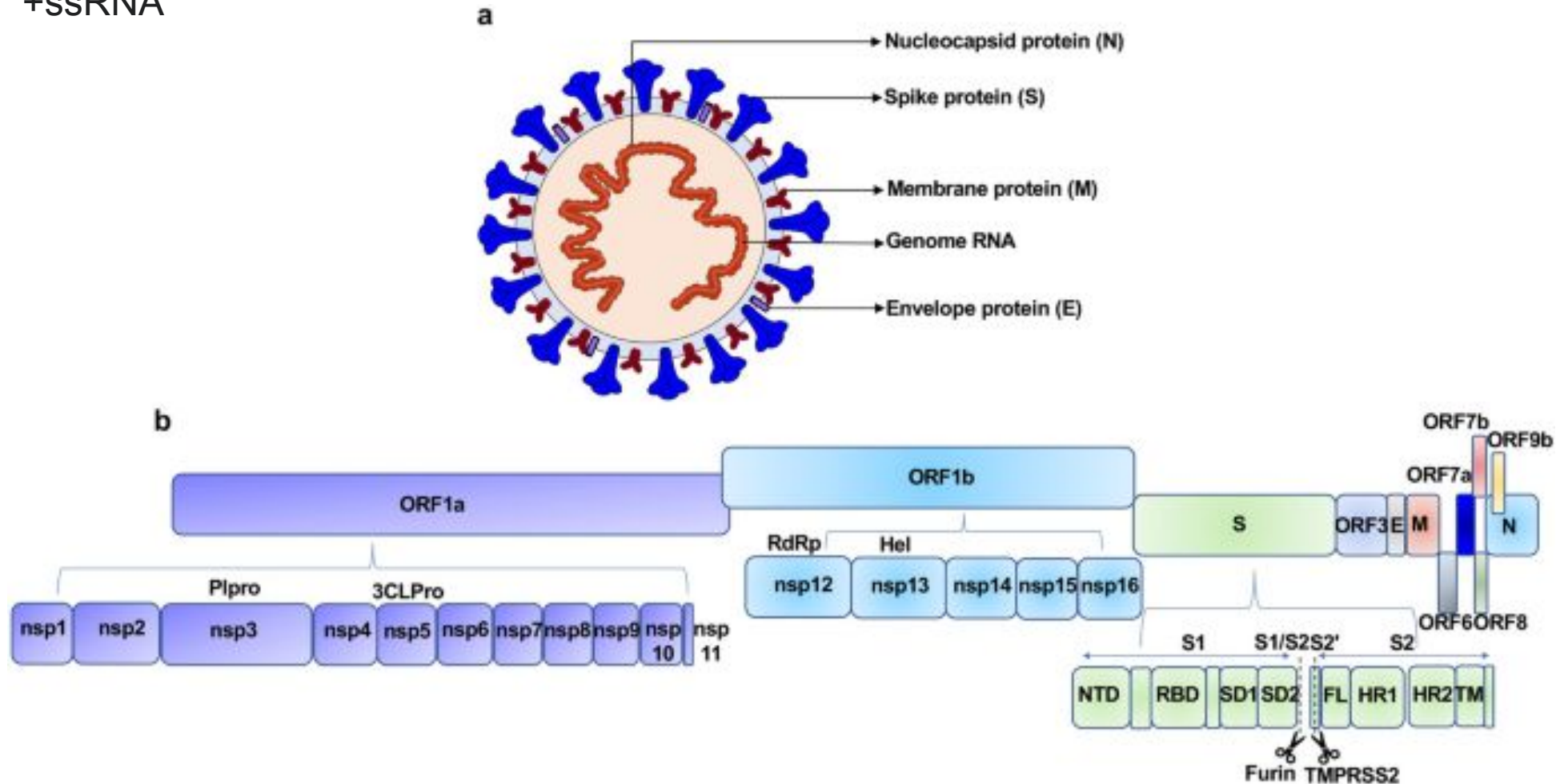
# Das amostras aos genomas



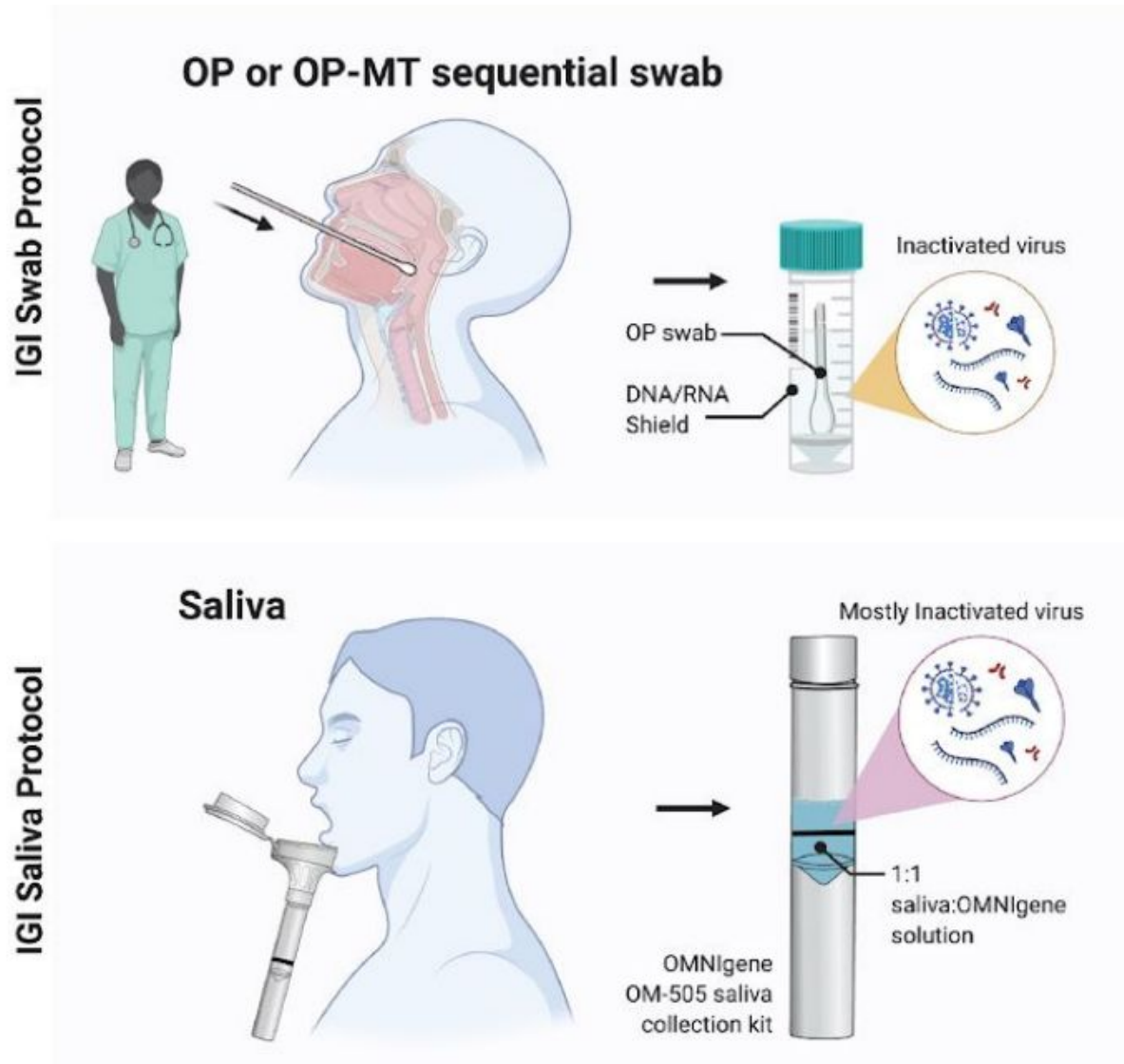Quick et al. (2017)

# SARS-CoV-2 e suas características
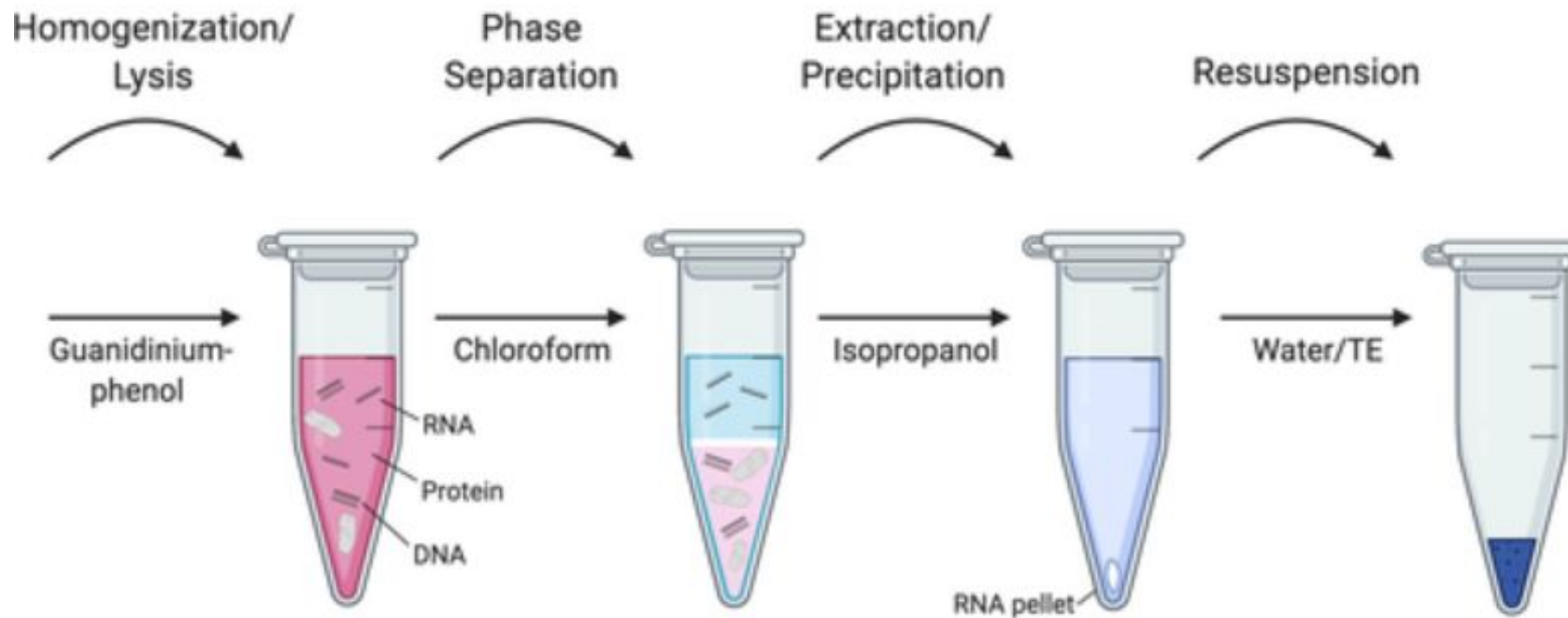
**S**evere **A**cute **R**espiratory **S**yndrome **Co**rona**V**irus 2

+ssRNA



Zhang et al. (2021)

# Passo 1: coleta de amostras



Hamilton et al. (2021)

LACEN - Recife/Pernambuco
http://portal.saude.pe.gov.br

# Passo 2: Extração de RNA total



Adapted from: https://www.addgene.org/protocols/kit-free-rna-extraction/

Ambrosi et al., 2021

# Passo 3: rRT-PCR (quantificação e cDNA)



Adaptado de: https://microbeonline.com/rt-pcr-principles-applications/

# Passo 4: PCR multiplex ou enriquecimento?



Quick et al. (2017)

# Passo 5: preparação de biblioteca Illumina

Nextera XT vs. COVID-seq

# Passo 6: Sequenciamento Illumina

**1** **Library Preparation**

Fragment DNA
Repair ends
Add A overhang
Ligate adapters
Purify

**2** **Cluster Generation**

Hybridize to flow cell
Extend hybridized template
Perform bridge amplification
Prepare flow cell for sequencing

**3** **Sequencing**

Perform sequencing
Generate base calls

**4** **Data Analysis**

Images
Intensities
Reads
Alignments

Flowcell

# Passo 6: Sequenciamento Illumina

DNA

Single molecule array

3' 5'

5'

**Library Preparation**

**Cluster Growth**

**Sequencing**

1 2 3 4 5 6 7 8 9

T G T A C G A T...

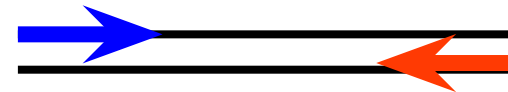**Image Acquisition**

**Base Calling**

# Passo 6: Illumina sequencing

**Leituras single-end**

Read1

ATGTTCCATAAGC…

**Leituras paired-end**

Read1

ATGTTCCATAAGC…

Read2

CCGTAATGGCATG…

| | iSeq 100 | MiniSeq | MiSeq Series ⊕ | NextSeq 550 Series ⊕ | NextSeq 1000 & 2000 |
|---|---|---|---|---|---|
| Run Time | 9.5–19 hrs | 4–24 hours | 4–55 hours | 12–30 hours | 11–48 hours |
| Maximum Output | 1.2 Gb | 7.5 Gb | 15 Gb | 120 Gb | 360 Gb* |
| Maximum Reads Per Run | 4 million | 25 million | 25 million † | 400 million | 1.2 billion* |
| Maximum Read Length | 2 × 150 bp | 2 × 150 bp | 2 × 300 bp | 2 × 150 bp | 2 × 150 bp |

# Análise de dados: FASTA/FASTQ

FASTA – genes e genomas

```
>NG_008679.1:5001-38170 Homo sapiens paired box 6 (PAX6)
ACCCTCTTTTCTTATCATTGACATTTAAACTCTGGGGCAGGTCCTCGCGTAGAACGCGGCTGTCAGATCT
GCCACTTCCCCTGCCGAGCGGCGGTGAGAAGTGTGGGAACCGGCGCTGCCAGGCTCACCTGCCTCCCCGC
CCTCCGCTCCCAGGTAACCGCCCGGGCTCCGGCCCCGGCCCGGCTCGGGGCCCGCGGGGCCTCTCCGCTG
CCAGCGACTGCTGTCCCCAAATCAAAGCCCGCCCCAAGTGGCCCCGGGGCTTGATTTTTGCTTTTAAAAG
GAGGCATACAAAGATGGAAGCGAGTTACTGAGGGAGGGATAGGAAGGGGGGTGGAGGAGGGACTTGTCTT
TGCCGAGTGTGCTCTTCTGCAAAAGTAGCAAAATGTTCCACTCCTAAGAGTGGACTTCCAGTCCGGCCCT
GAGCTGGGAGTAGGGGGCGGGAGTCTGCTGCTGCTGTCTGCTAAAGCCACTCGCGACCGCGAAAAATGCA
GGAGGTGGGGACGCACTTTGCATCCAGACCTCCTCTGCATCGCAGTTCACGACATCCACGCTTGGGAAAG
TCCGTACCCGCGCCTGGAGCGCTTAAAGACACCCTGCCGCGGGTCGGGCGAGGTGCAGCAGAAGTTTCCC
GCGGTTGCAAAGTGCAGATGGCTGGACCGCAACAAAGTCTAGAGATGGGGTTCGTTTCTCAGAAAGACGC
```

FASTQ – leituras Illumina



| Phred Quality Score | Probability of incorrect base call | Base call accuracy |
|---|---|---|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.9% |
| 40 | 1 in 10,000 | 99.99% |
| 50 | 1 in 100,000 | 99.999% |
| 60 | 1 in 1,000,000 | 99.9999% |

# Análise de dados: qualidade

# Análise de dados: filtragem/trimagem

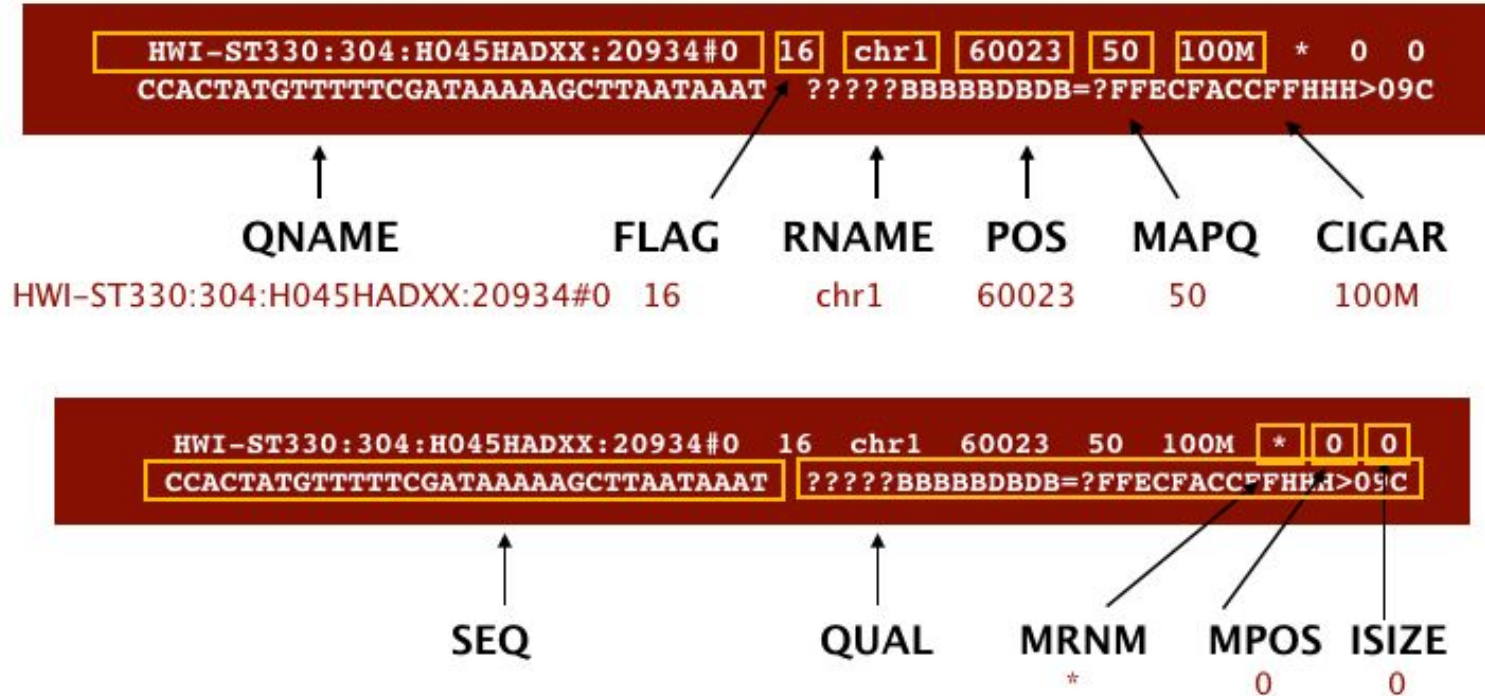

Example: fastp – Chen et al. (2018)

# Análise de dados: alinhamento



DNA/cDNA

Fragment DNA (PCR amplify)

Sequence DNA

Unaligned sequence

Aligned sequences

Reference genome

# Análise de dados: SAM/BAM/GFF
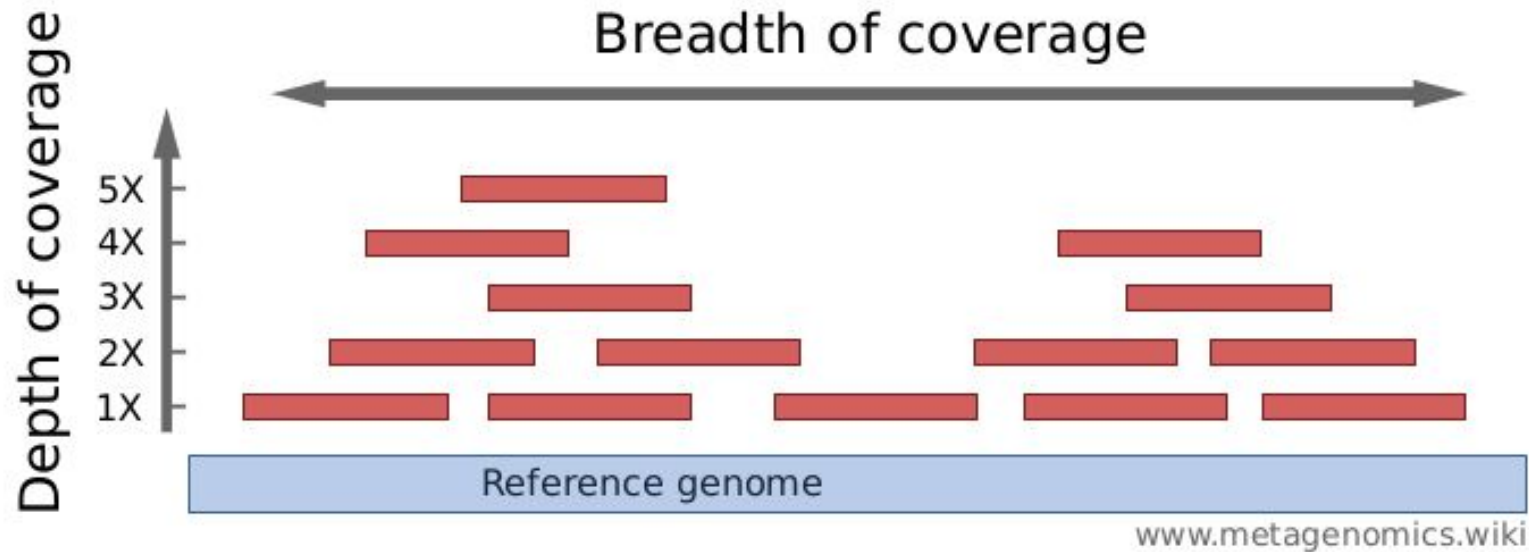
SAM/BAM



GFF

```
0  ##gff-version 3.2.1
1  ##sequence-region ctg123 1 1497228
2  ctg123 . gene            1000  9000  . + .  ID=gene00001;Name=EDEN
3  ctg123 . TF_binding_site 1000  1012  . + .  ID=tfbs00001;Parent=gene00001
4  ctg123 . mRNA            1050  9000  . + .  ID=mRNA00001;Parent=gene00001;Name=EDEN.1
5  ctg123 . mRNA            1050  9000  . + .  ID=mRNA00002;Parent=gene00001;Name=EDEN.2
6  ctg123 . mRNA            1300  9000  . + .  ID=mRNA00003;Parent=gene00001;Name=EDEN.3
7  ctg123 . exon            1300  1500  . + .  ID=exon00001;Parent=mRNA00003
```

# Análise de dados: profundidade vs. cobertura



Exemplo:
Genoma: 10 Mbp
Sequenciamento: 5 milhões de leituras x 100bp = 50Mbp (sequenciados totais)

Portanto, a **profundidade** média esperada em cada posição é 5x. No entanto, isso deve ser calculado em cada posição.

Já a **cobertura** refere-se à porcentagem do genoma suportada por leituras de sequenciamento

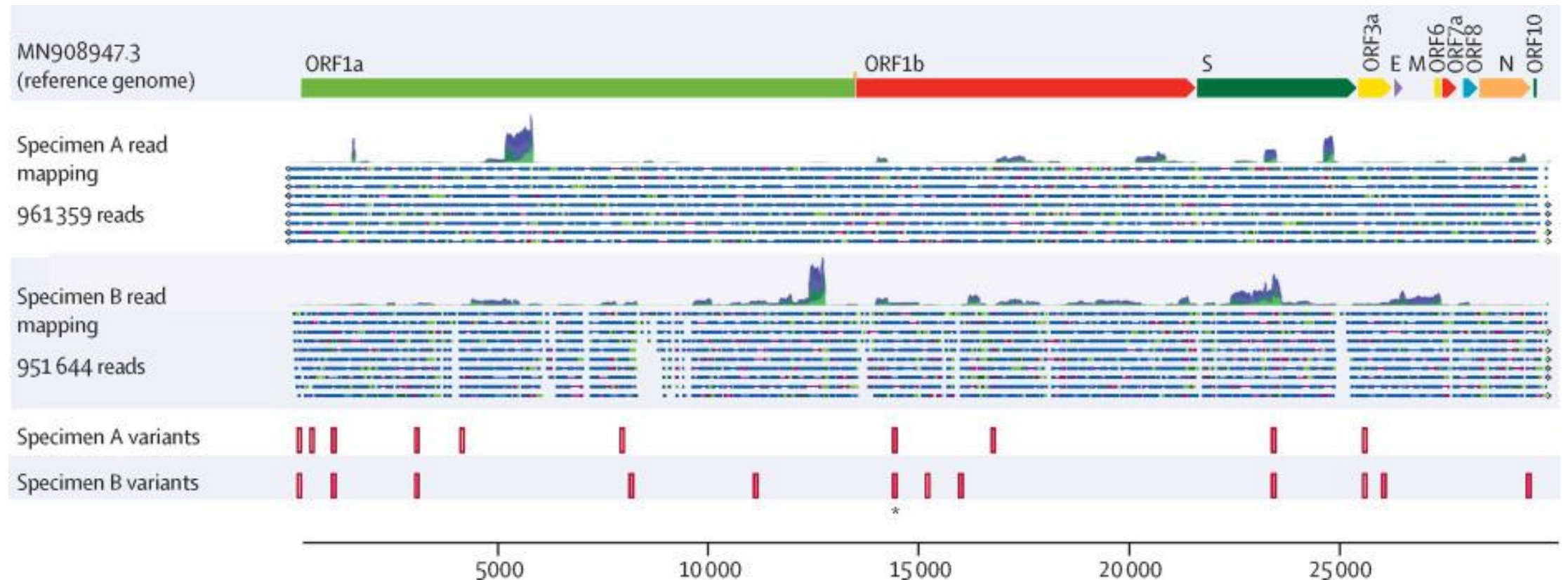# Análise de dados: profundidade vs. cobertura

| Sequencing applications | Recommended Coverage |
|---|---|
| Whole genome sequencing (WGS) | 15X to 60X |
| Whole exome sequencing (WES) | 100X |
| RNA sequencing (RNA-seq) | 5 to 100 M reads per sample depending on target study |
| ChIP-Seq | 100X |

Source: Illumina and genohub

References:

- Sims D, Sudbery I, Ilott NE, Heger A, Ponting CP. Sequencing depth and coverage: key considerations in genomic analyses. Nature Reviews Genetics. 2014 Feb;15(2):121-32.

# Análise de dados: alinhamento



Tillett et al. (2021)

# Análise de dados: chamada de variantes (VCF)

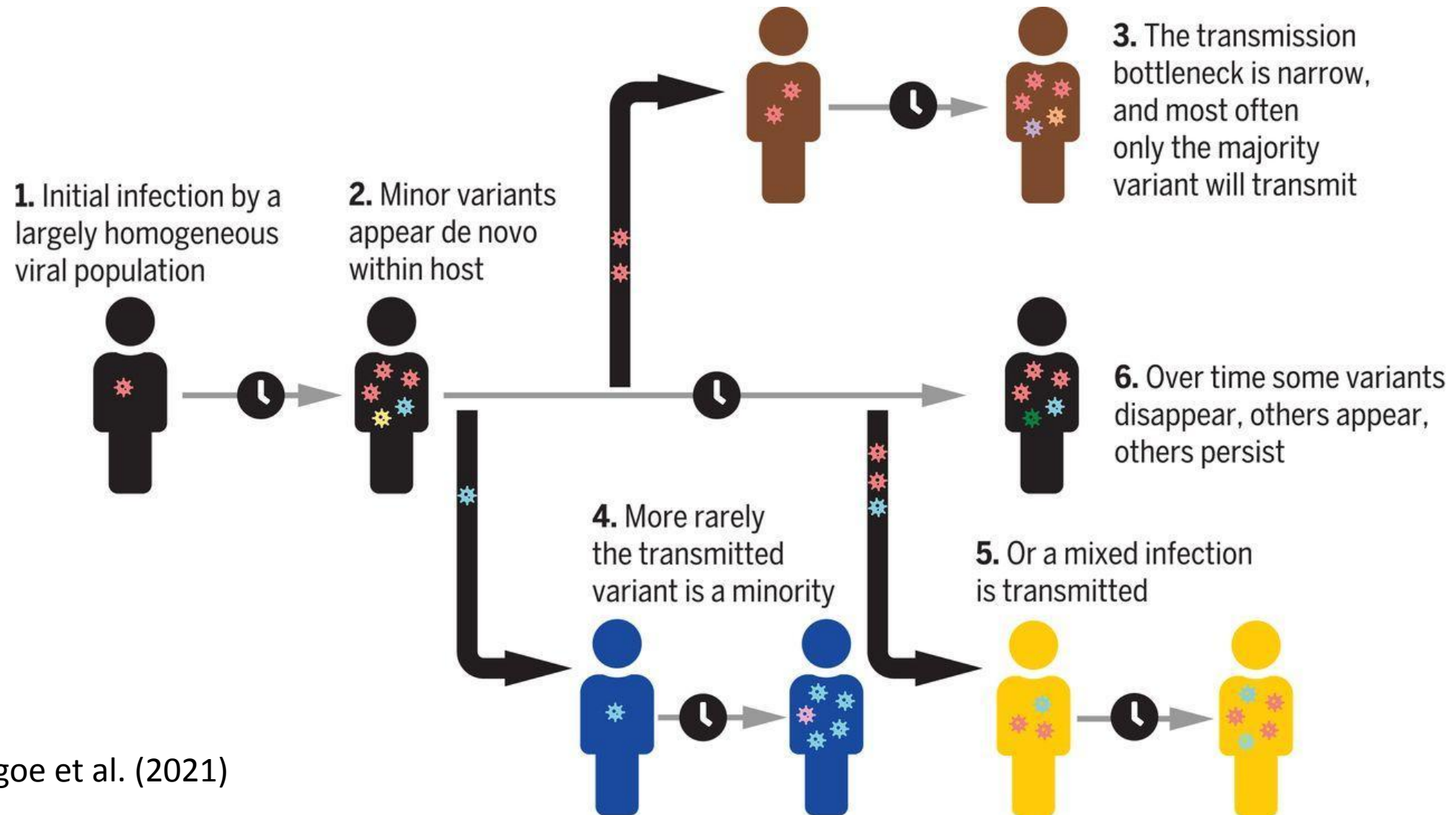| | | | | | | | |
|---|---|---|---|---|---|---|---|
| CP000819.1 | 1521 | . | C | T | 207 | . | DP=9;VDB=0.993024;SGB=-0.662043;MQSB=0.974597;MQ0F=0;AC=1;AN=1;DP4=0,0,4,5;MQ=60 |
| CP000819.1 | 1612 | . | A | G | 225 | . | DP=13;VDB=0.52194;SGB=-0.676189;MQSB=0.950952;MQ0F=0;AC=1;AN=1;DP4=0,0,6,5;MQ=60 |
| CP000819.1 | 9092 | . | A | G | 225 | . | DP=14;VDB=0.717543;SGB=-0.670168;MQSB=0.916482;MQ0F=0;AC=1;AN=1;DP4=0,0,7,3;MQ=60 |
| CP000819.1 | 9972 | . | T | G | 214 | . | DP=10;VDB=0.022095;SGB=-0.670168;MQSB=1;MQ0F=0;AC=1;AN=1;DP4=0,0,2,8;MQ=60   GT:PL |
| CP000819.1 | 10563 | . | G | A | 225 | . | DP=11;VDB=0.958658;SGB=-0.670168;MQSB=0.952347;MQ0F=0;AC=1;AN=1;DP4=0,0,5,5;MQ=60 |
| CP000819.1 | 22257 | . | C | T | 127 | . | DP=5;VDB=0.0765947;SGB=-0.590765;MQSB=1;MQ0F=0;AC=1;AN=1;DP4=0,0,2,3;MQ=60   GT:PL |
| CP000819.1 | 38971 | . | A | G | 225 | . | DP=14;VDB=0.872139;SGB=-0.680642;MQSB=1;MQ0F=0;AC=1;AN=1;DP4=0,0,4,8;MQ=60   GT:PL |
| CP000819.1 | 42306 | . | A | G | 225 | . | DP=15;VDB=0.969686;SGB=-0.686358;MQSB=1;MQ0F=0;AC=1;AN=1;DP4=0,0,5,9;MQ=60   GT:PL |
| CP000819.1 | 45277 | . | A | G | 225 | . | DP=15;VDB=0.470998;SGB=-0.680642;MQSB=0.95494;MQ0F=0;AC=1;AN=1;DP4=0,0,7,5;MQ=60 |
| CP000819.1 | 56613 | . | C | G | 183 | . | DP=12;VDB=0.879703;SGB=-0.676189;MQSB=1;MQ0F=0;AC=1;AN=1;DP4=0,0,8,3;MQ=60   GT:PL |
| CP000819.1 | 62118 | . | A | G | 225 | . | DP=19;VDB=0.414981;SGB=-0.691153;MQSB=0.906029;MQ0F=0;AC=1;AN=1;DP4=0,0,8,10;MQ=59 |
| CP000819.1 | 64042 | . | G | A | 225 | . | DP=18;VDB=0.451328;SGB=-0.689466;MQSB=1;MQ0F=0;AC=1;AN=1;DP4=0,0,7,9;MQ=60   GT:PL |

| column | info |
|---|---|
| CHROM | contig location where the variation occurs |
| POS | position within the contig where the variation occurs |
| ID | a . until we add annotation information |
| REF | reference genotype (forward strand) |
| ALT | sample genotype (forward strand) |
| QUAL | Phred-scaled probability that the observed variant exists at this site (higher is better) |
| FILTER | a . if no quality filters have been applied, PASS if a filter is passed, or the name of the filters this variant failed |

Referência

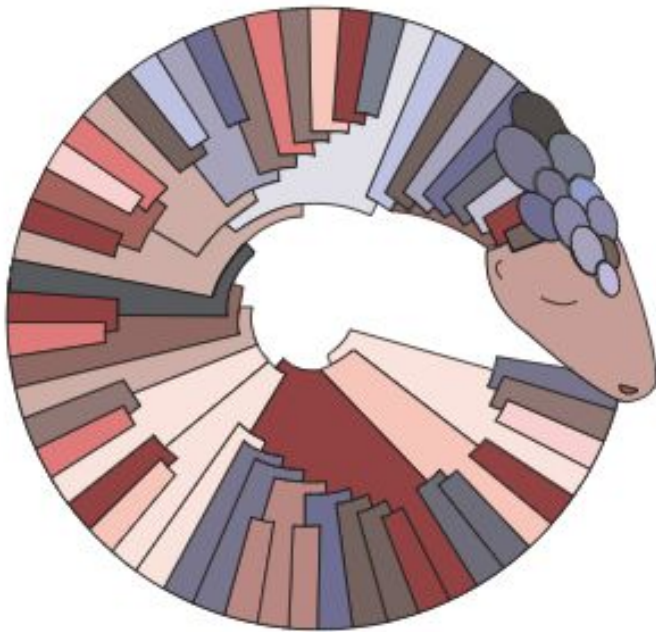

sequencing errors          SNP

# Análise de dados: consenso e variantes



Lythgoe et al. (2021)

# Análise de dados: definição de linhagem (Pangolin/Nextclade)



Command-line tool

GNU General Public License v3.0

Web application

Developed by the Centre for Genomic Pathogen Surveillance.

Rambaut et al. (2020)

# Publicação de genomas: GISAID

# Conda

# Contêineres



https://cloudhelix.io/blog/post/introduction-containerisation-enterprises

# Workflow: ViralFlow