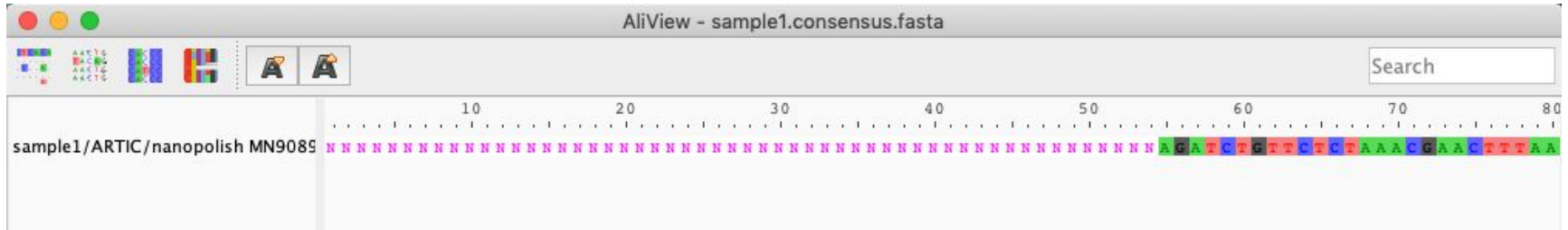# Agenda

- Viewing FASTA files

- Downloading publicly available genomes

- Concatenating consensus genomes

- Sequence Alignment

- Tree inference

# Viewing FASTA files

# How to view FASTA files



- Open in your favorite text editor:
  - TextEdit (Mac)
  - Notepad++ (Windows)
  - gedit (Ubuntu)

# How to view FASTA files



- You can also open FASTA files in specialized software for working with genomic data
- e.g., AliView: https://ormbunkar.se/aliview/ (free software)

- Or view directly on command line using less

# Downloading Publicly Available Data

# Selecting background data for phylogenetics

- For pathogens with limited data, use all available sequences

- Selecting background data is not a trivial task

- Considerations include:

  - Geography

  - Time period

  - Hosts / sample source

  - Subsampling of highly similar sequences

*There is a lack of broadly useful sampling guidance for phylogenetics*

# Downloading publicly available genomes

# Downloading publicly available genomes

# Downloading publicly available genomes

# Downloading publicly available reads

# Concatenating consensus genomes

# Concatenating consensus genomes



- Two methods for concatenating sequences:
  - Command line using cat
  - Copying and pasting in a text editor

- Both methods produce a multi-fasta
- FASTA file format: (.fa, .fas, .fasta); sometimes .mfa is used to indicate multi-fasta

# Concatenating consensus genomes

# Sequence Alignment

# Multiple sequence alignment

# Types of multiple sequence alignment

**Alignment with a reference genome**
- Faster and less computationally intensive
- Ensures consistent coordinates

**Alignment without a reference genome**

# Selecting an appropriate reference genome

- Use the same reference genome as others working on the same outbreak
    - e.g., Wuhan-Hu-1 for SARS-CoV-2
- Use the earliest sequence from the outbreak, if available
- Use a sequence from a prior outbreak in the same location
- Use NCBI RefSeq (https://www.ncbi.nlm.nih.gov/refseq/)

    *The reference genome should always have a date **before** the earliest of your samples*

# MAFFT Server



MAFFT version 7
Multiple alignment program for amino acid or nucleotide sequences

Download version
Mac OS X
Windows
Linux
Source
**Online version**
**Alignment**
mafft --add
Merge
Phylogeny
Rough tree
Merits / limitations
Algorithms
Tips
Benchmarks
Feedback

Hardware was upgraded, Jan 16, 2022. There should be no change in user interface. If you notice any unexpected changes, then please let us know.

To avoid overload, try a light-weight option, for MSA of full-length SARS-CoV-2 genomes (2020/Apr).

For a large number of short sequences, try an experimental service.

Experimental service for aligning raw reads (2019/Aug)

Multiple sequence alignment and NJ / UPGMA phylogeny

**Input**:
Paste protein or DNA sequences in fasta format. Example

or upload a **plain text** file: Choose File  No file chosen
☐ Use DASH to add homologous structures (protein only)  *New! 2018/Dec/23*
   ◉ Ouput original plus DASH sequences  ○ Output original sequences only
☐ Give structural alignment(s) externally prepared
☐ Allow unusual symbols (Selenocysteine "U", Inosine "i", non-alphabetical characters, *etc.*)  Help
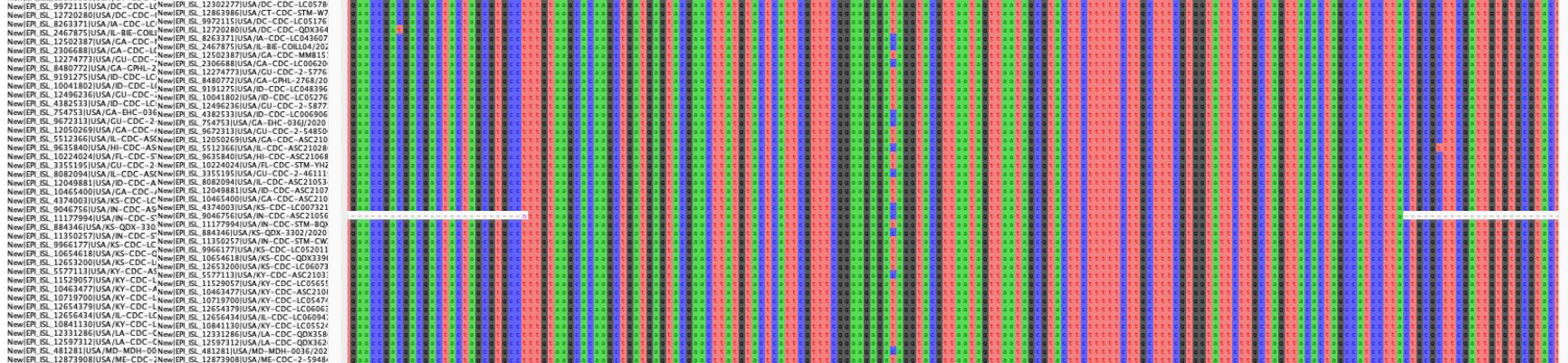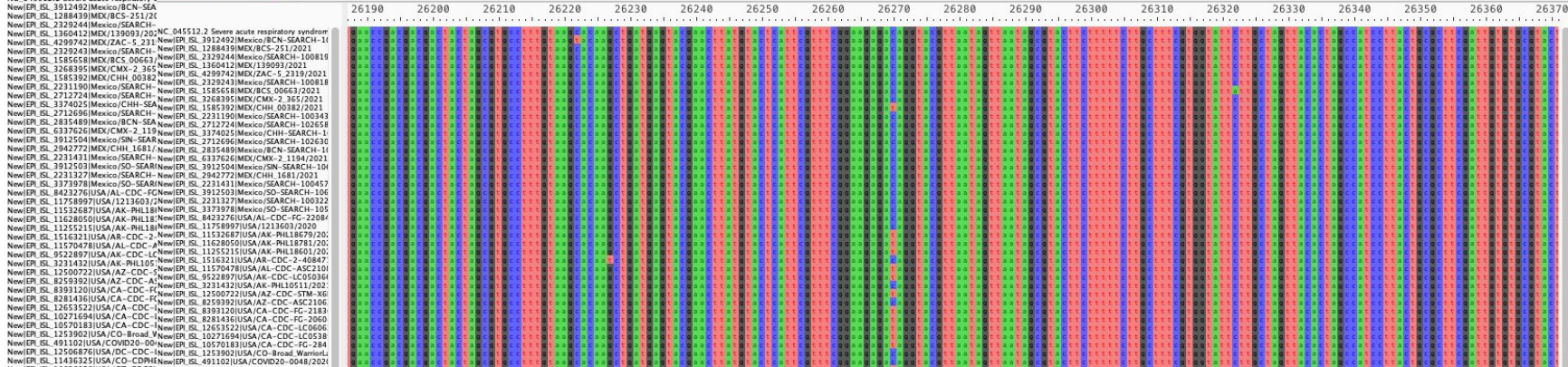
# MAFFT Commandline

mafft --auto **consensus.fasta** > **consensus_aln.fasta**

# Phylogenetic Tree Inference

# IQ-Tree Server

http://iqtree.cibiv.univie.ac.at/

# IQ-Tree Commandline

iqtree2 -s **consensus_aln.fasta** -T AUTO -m TEST -B 1000

# What questions do you have?