

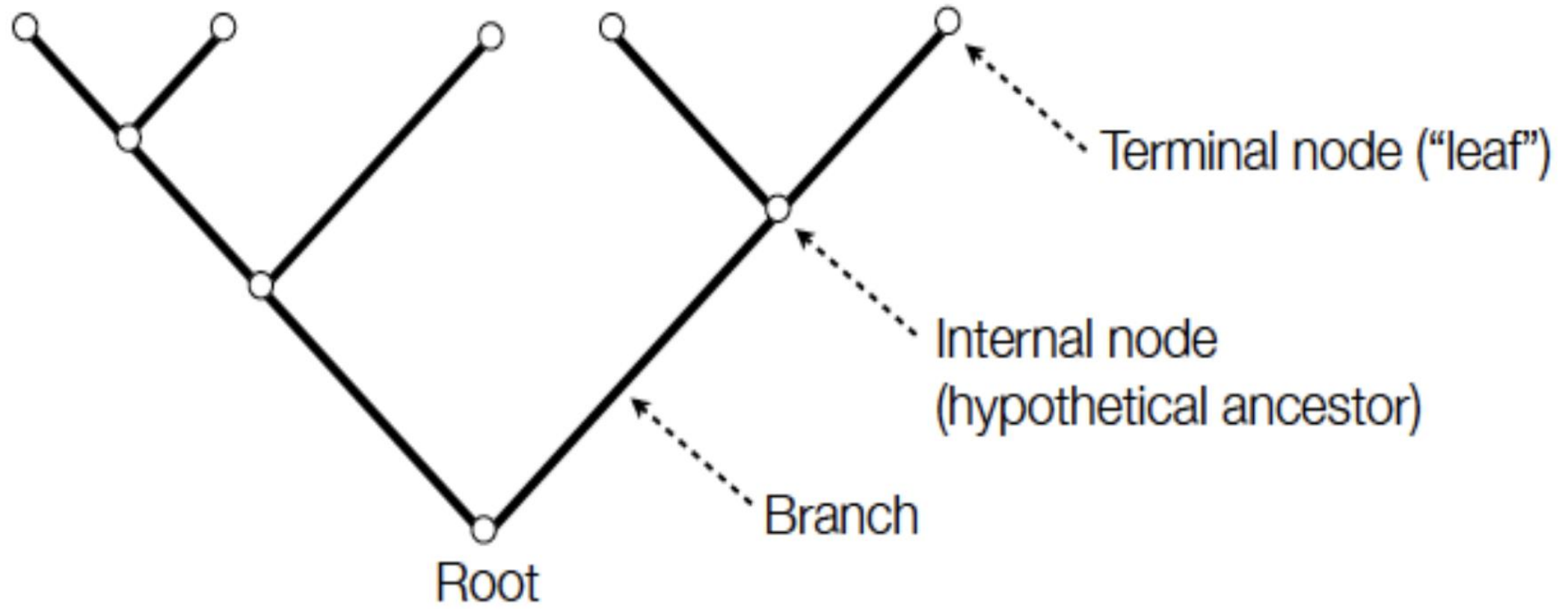
# Using Nextstrain and Nextclade to Analyse Dengue Sequences

—  
Ammar Aziz – VIDRL

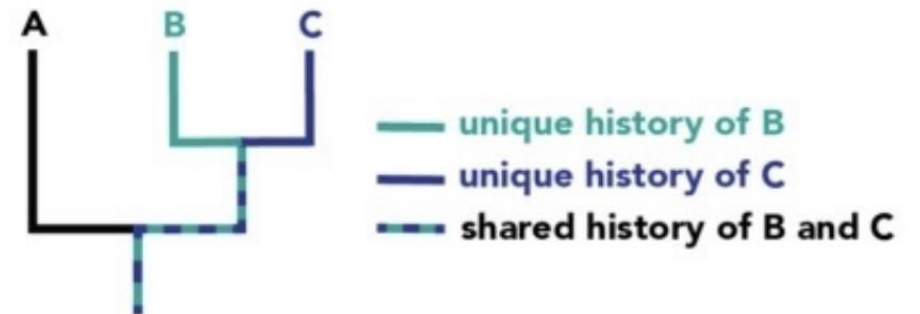
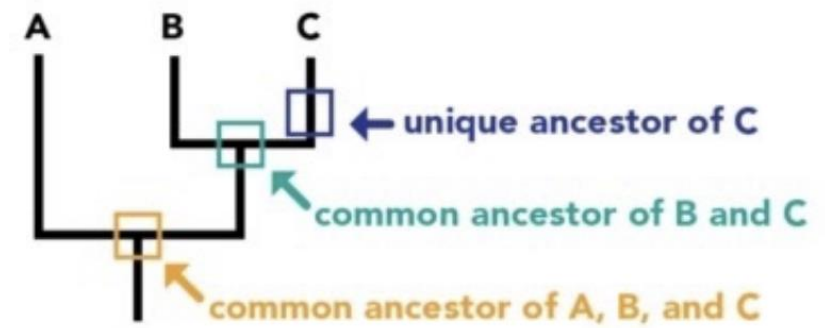
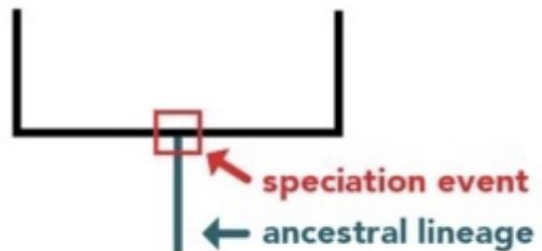
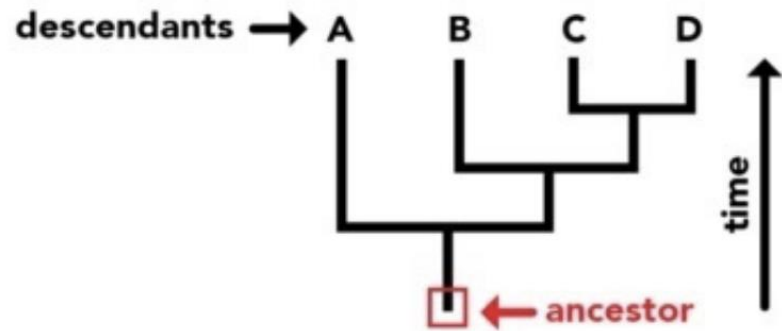
# Objectives

- Refresher on Phylogenies
- Dengue global surveillance using Nextstrain
- Clade/Lineage analysis using Nextclade

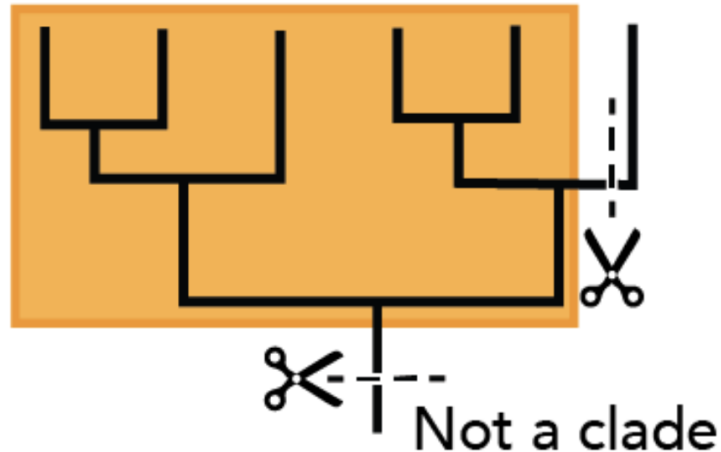
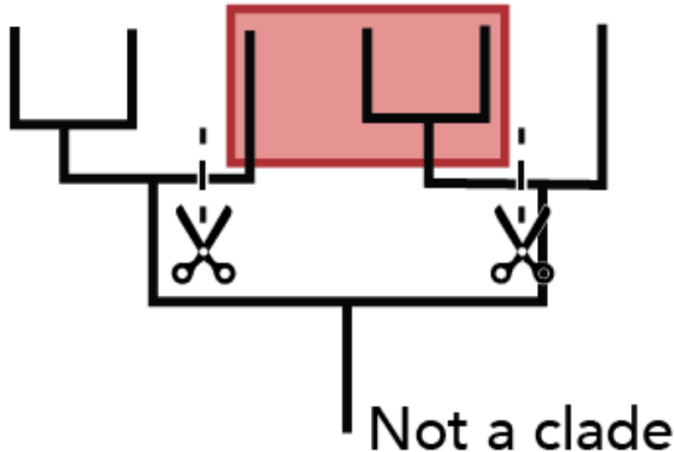
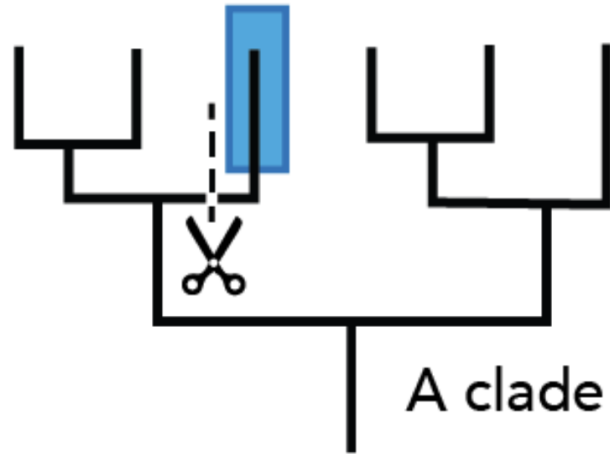
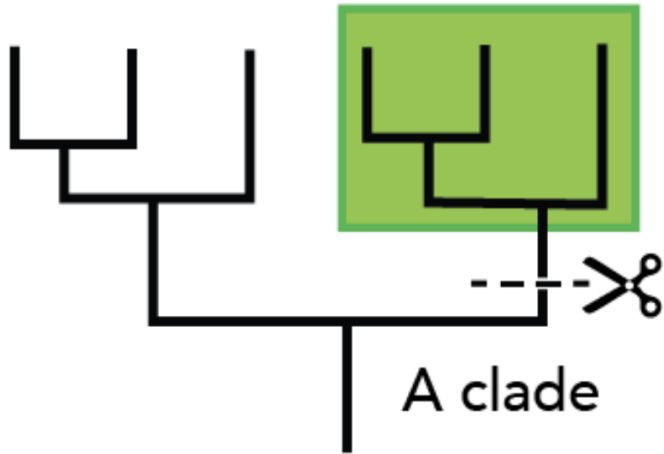
# Trees - Terminology



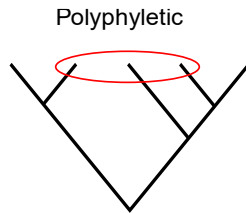
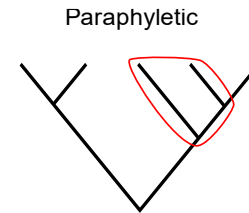
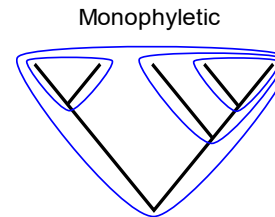
# Reading Trees



# What is a Clade?



- Monophyletic group - All and only the descendants of a common ancestor (including that ancestor).
- Paraphyletic group - Does not include all of the descendants of the common ancestor.
- Polyphyletic group - Does not include the common ancestor.





# What is Nextstrain?



Nextstrain is a project to harness the scientific and public health potential of pathogen genome data. Our goal is to aid epidemiological understanding of pathogen spread and evolution and improve outbreak response.

## Pathogen surveillance

Our website, [nextstrain.org](https://nextstrain.org), provides real-time snapshots of evolving pathogen populations such as [SARS-CoV-2](#), [influenza](#), and [Ebola](#). We use interactive visualizations to enable exploration of curated datasets and analyses which are continually updated when new genomes are available. This offers a powerful pathogen surveillance tool to virologists, epidemiologists, public health officials, and community scientists. In many cases old snapshots of these analyses are able to be easily accessed, see [viewing previous analyses](#) for more.

## Open-source software

The software we write to power [all parts](#) of Nextstrain—bioinformatics, visualizations, analysis pipelines, data management, and more—is entirely [open-source](#) and [available to the public](#). We aim to empower the wider genomic epidemiology and public health communities to tweak our analyses, create new ones, and communicate what we do.

# Nextstrain – Reading trees

- <https://nextstrain.org/narratives/trees-background/>

# Nexstrain Demonstration

- Dengue E tree: <https://nextstrain.org/dengue/all/E>
- Filtering:
  - By country
  - By serotype
- Change colors:
  - Date
  - Country
  - Region
- Interactively explore tree
- Tanglegram (WGS – E tree)



# Nextclade: analysis of viral genetic sequences

Nextclade is an open-source project for viral genome alignment, mutation calling, clade assignment, quality checks and phylogenetic placement.

Nextclade consists of a set of related tools:

- Nextclade Web - a web application available online at [clades.nextstrain.org](https://clades.nextstrain.org)
- Nextclade CLI - a command-line tool

Both tools are powered by the same algorithms, they consume the same inputs and produce the same outputs, but they differ in the user interface, the features included, and the degree of customization. It is recommended to start with Nextclade Web and later proceed to CLI tools if you have more advanced use-cases (for example, repeated batch processing, bioinformatics pipelines).

# Download sequencing data

- Go to XXXXXX
- Download your data and a country specific dataset
- Launch [Nextclade.org](https://nextclade.org)


# Nextclade<sup>v3.10.0</sup>

Clade assignment, mutation calling, and sequence quality checks

Data input here

Provide sequence data

[File](#) [Link](#) [Text](#) [Example](#) ▾




Drag & drop files or folders

Select files

Selected reference dataset [i](#)

☐ Suggest automatically [Reset](#) [Suggest](#)



**DENV-2**  
**community**

Reference: Thailand/CDC-16681/1964 (NC\_001474.2)  
Updated at: 2024-10-17 16:48:48 (UTC)  
Dataset name: community/v-gen-lab/dengue/denv2

[Open tree](#) [Load example](#)























[Change reference dataset](#) [Run](#)

Select reference dataset here

# Nextclade supports many viruses

Change reference dataset



	<b>SARS-CoV-2</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>SARS-CoV-2 (Mature protein)</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>SARS-CoV-2 (BA.2)</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>SARS-CoV-2 (XBB)</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>SARS-CoV-2 (BA.2.86)</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza A H3N2 pdm HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza A H3N2 pdm HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza A H3N2 pdm HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza A H3N2 pdm HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza A H3N2 pdm HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza A H3N2 pdm HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza A H3N2 pdm HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza B Victoria HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza B Victoria HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza B Victoria HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Influenza B Yamagata HA</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>RSV-A</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>RSV-B</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Mpox virus (All clades)</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Mpox virus (Clade I)</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Mpox virus (Clade Ibb)</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)
	<b>Mpox virus (Lineage B.1)</b> Reference: NC_045886.2 (2019-12-01) Accession: NC_045886.2 (2019-12-01) Dataset: SARS-CoV-2 (2019-12-01)

# Nextclade<sup>v3.10.0</sup>

Clade assignment, mutation calling, and sequence quality checks

If unsure which ref dataset, click 'Suggest'

Add more sequence data

Add more sequence data

File

Link

Text

Example ▾

<>

FASTA

Drag & drop files or folders

Select files

Sequence data you've added

☆ Examples for 'nextstrain/dengue/all'

Remove all

Selected reference dataset ⓘ

☐ Suggest automatically

Reset

Suggest

D

DENV-2

community

Reference: Thailand/CDC-16681/1964 (NC\_001474.2)  
Updated at: 2024-10-17 16:48:48 (UTC)  
Dataset name: community/v-gen-lab/dengue/denv2

Open tree

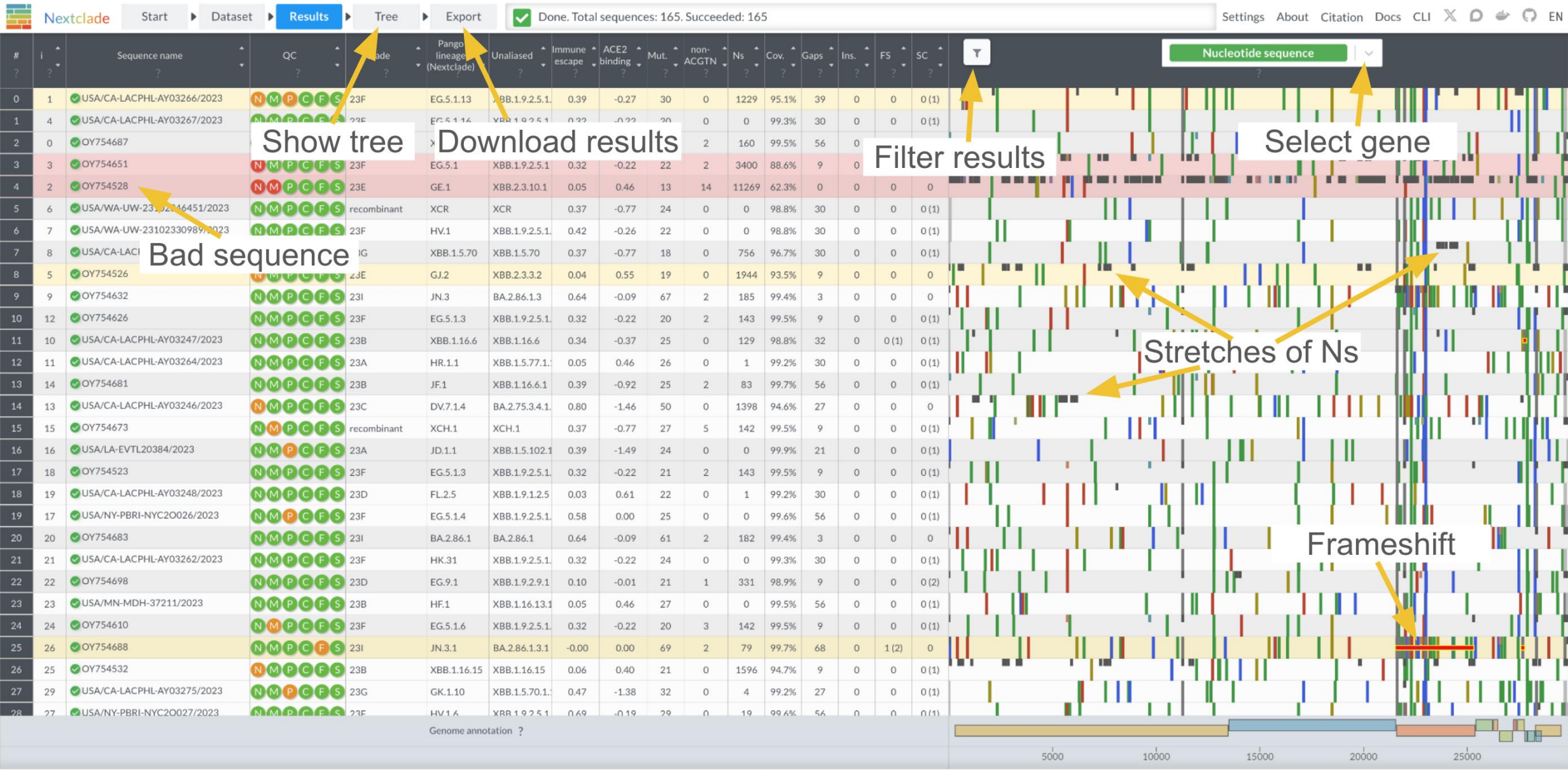
Load example

Change reference dataset

Run

Click Run to start analysis





Sequence name	QC	Clade	Mut.	non-ACGTN	Ns	Cov.	Gaps	Ins.	FS	SC
?	?	?	?	?	?	?	?	?	?	?

- “Mut.”: number of mutations with respect to the reference sequence
- “non-ACGTN”: number of ambiguous nucleotides that are not *N*
- “Ns”: number of missing nucleotides indicated by *N*
- “Gaps”: number of nucleotides that are deleted with respect to the reference sequence
- “Ins.”: number of nucleotides that are inserted with respect to the reference sequence
- “FS”: Number of uncommon frame shifts (total number, including common frame shifts are in parentheses)
- “SC”: Number of uncommon premature stop codons (total number, including common premature stops are in parentheses)



QC Field

Hover mouse over highlighted sample to show popup.

QC

?

Warning

Failed

N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S

N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S

Overall QC score: 56  
Overall QC status: mediocre  
Detailed QC assessment:

N

Missing Data: good  
No issues

P

Private Mutations: good  
No issues

F

Frame shifts: mediocre  
Unexpected 1 frame shift(s) detected:  
NS5:404-900. QC score: 75

S

Stop codons: good  
No issues

N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S
N	P	F	S

Overall QC score: 773  
Overall QC status: bad  
Detailed QC assessment:

N

Missing Data: bad  
Too much missing data found. Total Ns: 2880  
(1100 allowed). QC score: 278

P

Private Mutations: good  
No issues

F

Frame shifts: good

S

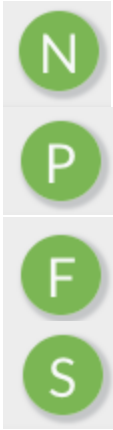
Stop codons: good  
No issues



N	P	F	S	2III_C.1.1	508	0	2880	68.0%	0	0	0	0
---	---	---	---	------------	-----	---	------	-------	---	---	---	---

Why was it marked as fail? Too much missing data (N or -)

# QC Metrics



- Missing Data – threshold 1000
  - Private Mutations – cutoff 216, typical 72
  - Frame Shift
  - Stop Codon
- 
- Nextclade implements a variety of quality control metrics to quickly spot problems in your sequencing/assembly pipeline.
  - Bad sequences are colored red, mediocre ones yellow and good ones white. You can view detailed results of the QC metrics by hovering your mouse over a sequences QC entry:

## Change Genetic Feature to Nucleotide to show full genome

The diagram illustrates a two-step process for changing the genetic feature. The top panel shows the initial state with 'Genetic feature' set to 'CDS' and 'Relative to' set to 'Reference'. An orange arrow points down to the bottom panel, where 'Genetic feature' has been changed to 'Nucleotide sequence' (highlighted with a blue border) and 'Relative to' remains 'Reference'.

Genetic feature *i* Relative to *i*

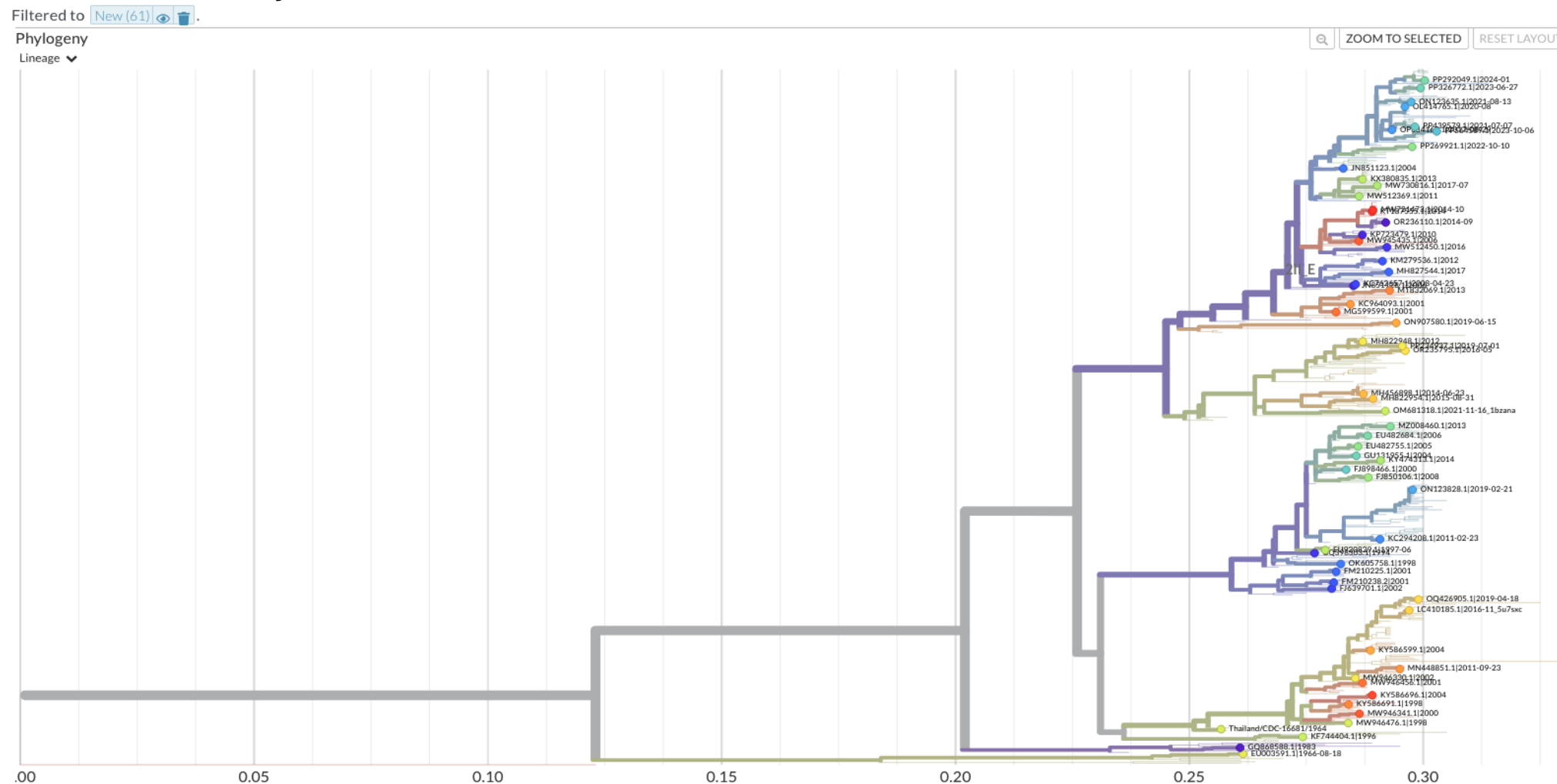
Gene CDS E Reference

Genetic feature *i* Relative to *i*

Nucleotide sequence Reference

Note: In Dengue the E gene is routinely sequenced as it's genetically distinct between serotypes and provides sufficient resolution for clade information and possible lineage information.

The tree is nearly identical to Nextstrain tree, so interact with it.



# Nextclade Demonstration

- Loading consensus sequences
- Selecting datasets
- Interpreting results
- Exporting results

# Nextclade Practice

Use the amplicon data you generate in the lab to answer these questions:

1. What is the clade/lineage of each sample?
2. Check each samples QC metrics, which would you use for:
  - Clade/Lineage calling?
  - Phylogenetic analysis?
  - Speciation?
3. Check for odd genomic features (premature stop, gaps, etc).
4. Which samples would you pass/fail?

# Questions? + Resources

- Nextclade is also a CLI tool and is available in Galaxy.
  - <https://docs.nextstrain.org/projects/nextclade/en/stable/user/nextclade-cli/index.html>
- The web interface is excellent, provides all the features of the CLI and even more.
- It's possible to create your own reference datasets for Nextclade but it requires knowledge of the organism and advanced bioinformatic skills.