

# Creating a Nextclade dataset

- Bedford Lab Meeting -

**Jennifer Chang, Ph.D.**

Bioinformatic Analyst III

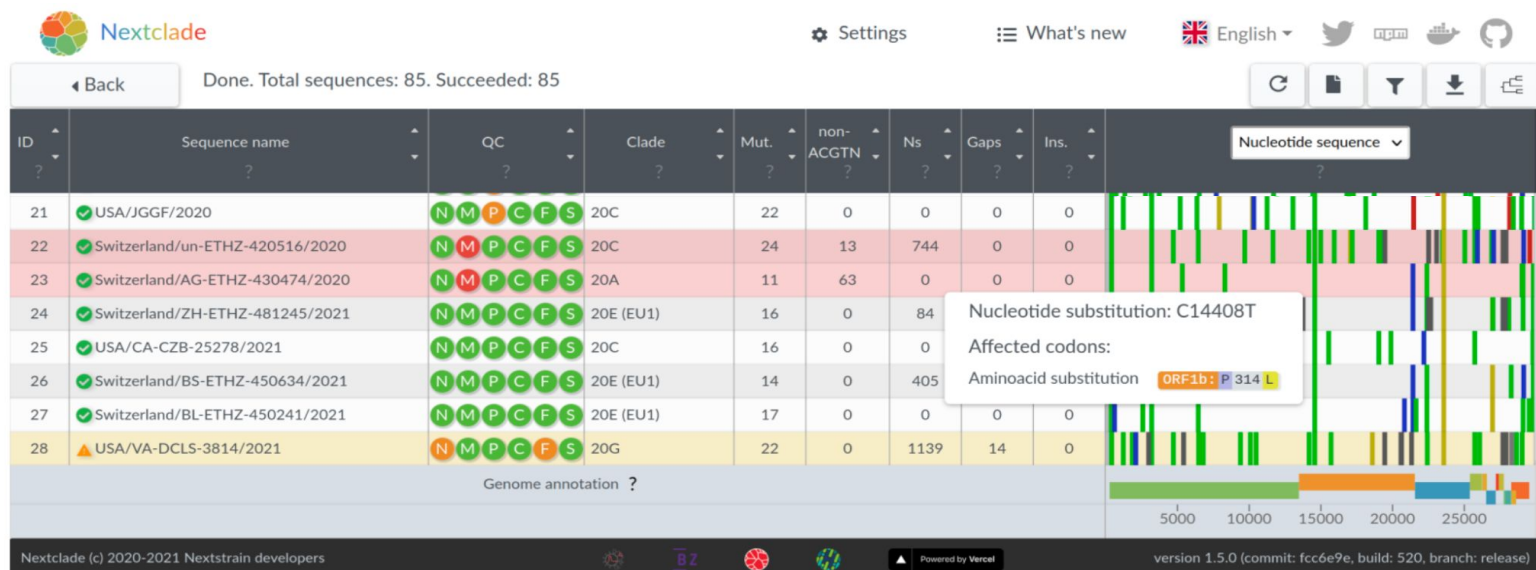
Fred Hutchinson Cancer Center

# Agenda

- What is Nextclade?
  - Statement of need
  - What is a Nextclade dataset?
- How do we create a Nextclade dataset?
  - Nextclade template script
  - Nextalign vs Augur align
- Checking a Nextclade dataset
  - validation and basic checks
- Conclusions

# What is Nextclade?

- (1) Assess the quality of the sequence
- (2) Assign it to a known clade or type
- (3) Compare it to a reference sequence to detect evolutionary changes



**Figure 1:** Overview of the results page with clade assignments, QC metrics, and the nucleotide mutation view. The results can be explored interactively and exported in standard tabular file formats.

# What is a Nextclade dataset?

## Nextclade datasets

tag.json

qc.json

primers.csv

virus\_properties.json

reference.fasta

genemap.gff

sequences.fasta

tree.json

## Nextclade datasets

Nextclade dataset is a set of input data files required for Nextclade to run the analysis:

- reference (root) sequence ( `reference.fasta` )
- reference tree ( `tree.json` )
- quality control configuration ( `qc.json` )
- gene map ( `genemap.gff` )
- PCR primers ( `primers.csv` )
- virus properties ( `virus_properties.json` )

See also: [Input files](#)

Dataset might also include example sequence data ( `sequences.fasta` ) - typically a diverse set of query sequences that represents major clades, used for demonstration and highlights analysis features of Nextclade. Most of the time you want to analyze your own sequence data.

Dataset also includes a file `tag.json` which contains version tag and other properties of the dataset. This file is currently not used by Nextclade and serves only for informational purposes.

An instance of a dataset is a directory containing the dataset files or an equivalent zip archive.

[Nextclade docs: datasets](#)

# What is a Nextclade dataset?

## Nextclade datasets

tag.json

qc.json

primers.csv

virus\_properties.json

reference.fasta

genemap.gff

sequences.fasta

tree.json

The screenshot shows the GitHub repository for `nextstrain/nextclade_data`. The left sidebar displays the file structure, with the `data/datasets` directory expanded. A pink box highlights the `flu_h1n1pdm_ha` directory, which contains files like `genemap.gff`, `primers.csv`, `qc.json`, `reference.fasta`, `sequences.fasta`, `tag.json`, `tree.json`, and `virus_properties.json`. The main content area shows a list of datasets with columns for Name, Last commit message, and Last commit date.

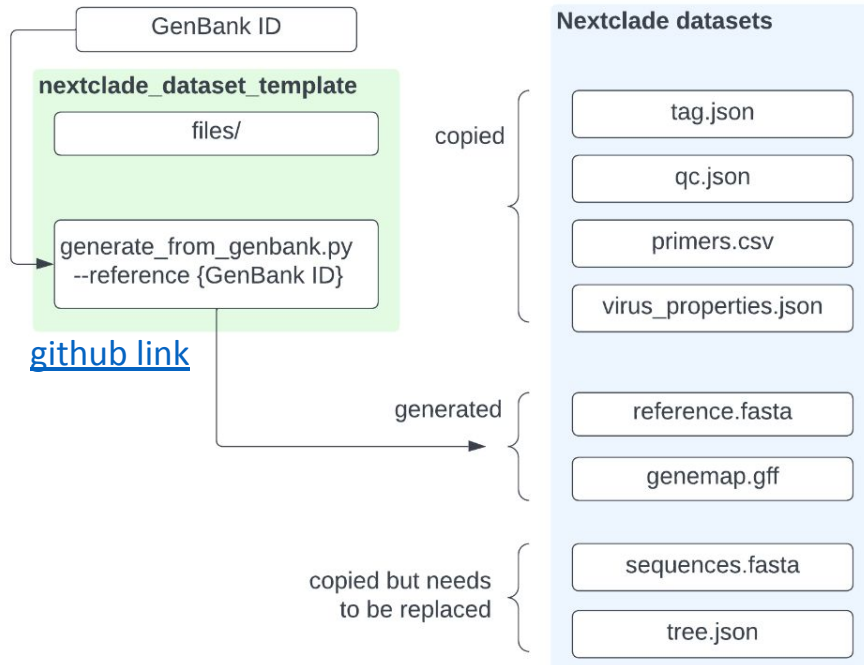
Name	Last commit message	Last commit date
..		
MPXV	Monkeypox dataset update 2023-01-26	4 months ago
flu_h1n1pdm_ha	update flu datasets and fix B/vic annotation	last month
flu_h1n1pdm_na	update flu datasets and fix B/vic annotation	last month
flu_h3n2_ha	update flu datasets and fix B/vic annotation	last month
flu_h3n2_na	update flu datasets and fix B/vic annotation	last month
flu_vic_ha	update flu datasets and fix B/vic annotation	last month
flu_vic_na	update flu datasets and fix B/vic annotation	last month
flu_yam_ha	fix: date in tag wrong month	10 months ago
hMPXV	Monkeypox dataset update 2023-01-26	4 months ago
hMPXV_B1	Monkeypox dataset update 2023-01-26	4 months ago
rsv_a	rsv: update data sets	3 months ago
rsv_b	rsv: update data sets	3 months ago
sars-cov-2-21L	Add placementMaskRanges and new labeled muts (23B)	4 days ago
sars-cov-2-no-recomb	sc2: Update datasets, now containing placement_prior	2 months ago
sars-cov-2	Add placementMaskRanges and new labeled muts (23B)	4 days ago

[https://github.com/nextstrain/nextclade\\_data/tree/master/data/datasets](https://github.com/nextstrain/nextclade_data/tree/master/data/datasets)

# Agenda

- What is Nextclade?
  - Statement of need
  - What is a Nextclade dataset?
- **How do we create a Nextclade dataset?**
  - Nextclade template script
  - Nextalign vs Augur align
- Checking a Nextclade dataset
  - validation and basic checks
- Conclusions

# Picking a reference



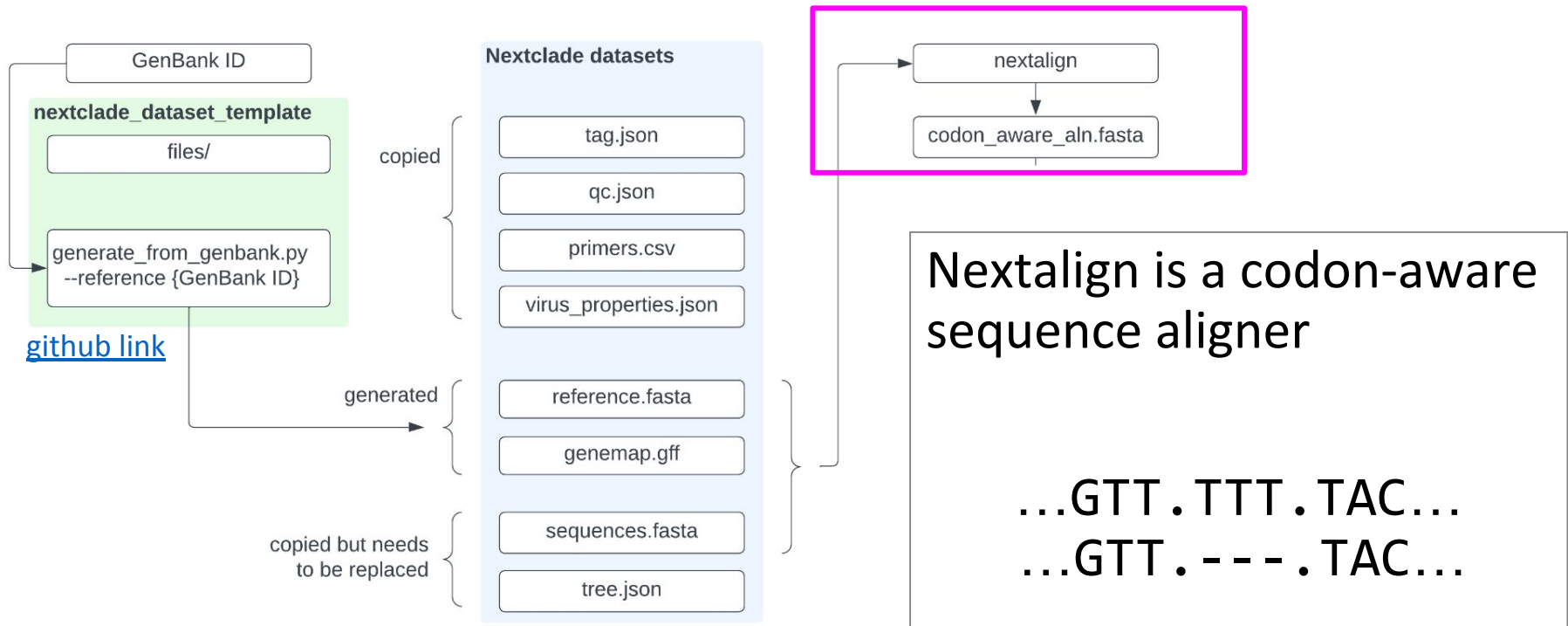
Pick a reference that is close to the base of the tree.

Refseq reference  
(e.g. [NCBI Virus: Dengue](#))

```
git clone https://github.com/nextstrain/nextclade_dataset_template.git
cd nextclade_dataset_template

python generate_from_genbank.py \
  --reference NC_001477 \
  --output-dir denv1_dataset
```

# Nextalign or augur align



```
cat sequences.fasta \  
| nextalign run \  
--jobs=`nproc` \  
--reference reference.fasta \  
--genemap genemap.gff \  
--output-translations translations_{gene}.txt \  
--output-fasta aln.fasta
```

[Nextclade docs: nextalign-cli](#)




# Nextalign or augur align

Nextclade can use a genome annotation to make the alignment more interpretable. Sometimes, the placement of a sequence deletion or insertion is ambiguous as in the following example. The gap could be moved forward or backward by one base with the same number of matches:

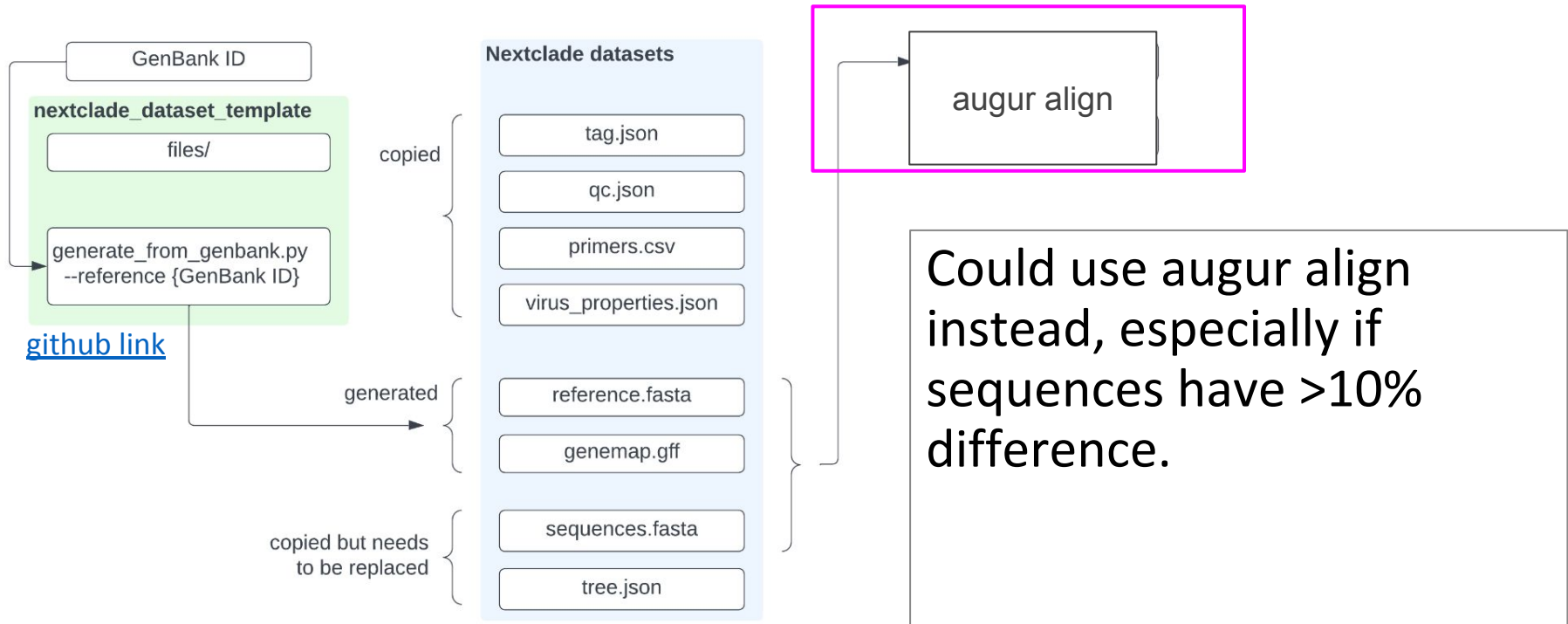
```
Reference   : ... | GTT | TAT | TAC | ...
Alignment 1: ... | GTT | --- | TAC | ...
Alignment 2: ... | GT- | --T | TAC | ...
Alignment 3: ... | GTT | T-- | -AC | ...
```

If a genome annotation is provided, Nextclade will use a lower gap-open-penalty at the beginning of a codon (delimited by the `|` characters in the schema above), thereby locking a gap in-frame if possible. Similarly, nextalign preferentially places gaps outside of genes in case of ambiguities.

```
denv1_dataset > files >  genemap.gff
1  ##gff-version 3
2  ##sequence-region NC_001477.1 1 10735
3  NC_001477.1 feature gene 95 10273 . + . codon_start=1;gene=POLY;gene_name=POLY
4  NC_001477.1 feature gene 710 934 . + . codon_start=1;gene=M;gene_name=M
5  NC_001477.1 feature gene 935 2419 . + . codon_start=1;gene=E;gene_name=E
```

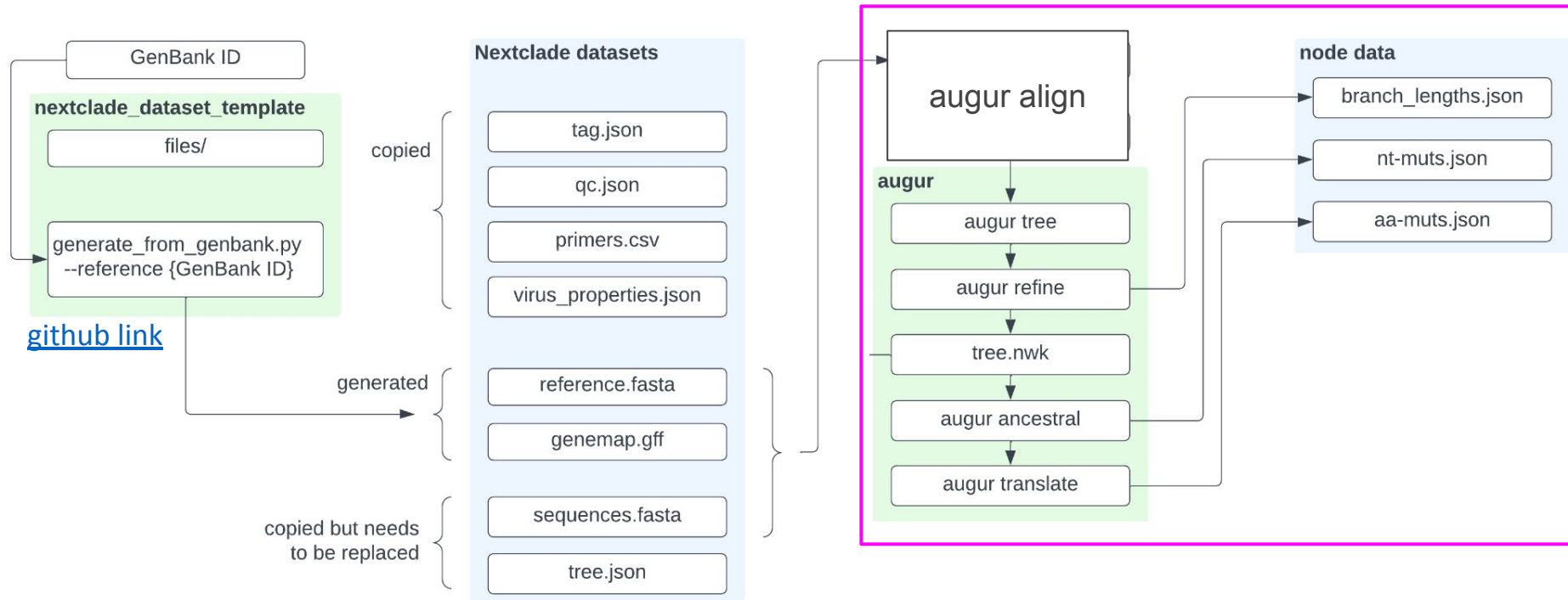
[Nextclade docs: algorithm/01-sequence-alignment.html](https://nextclade.org/docs/algorithm/01-sequence-alignment.html)

# Nextalign or augur align

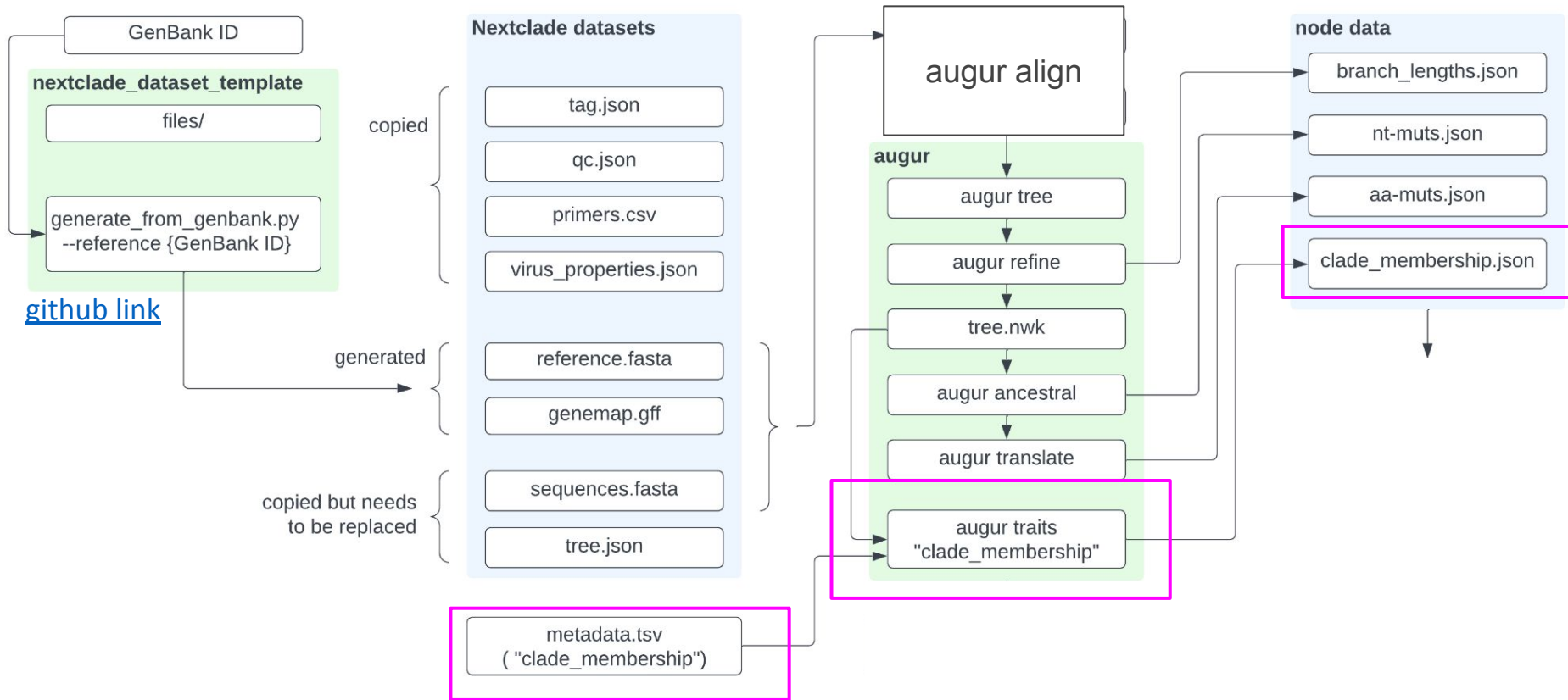


```
augur align \
  --sequences sequences.fasta \
  --reference-sequence reference.fasta \
  --output results/aln.fasta \
  --fill-gaps \
  --nthreads `nproc`
```

# Standard augur commands

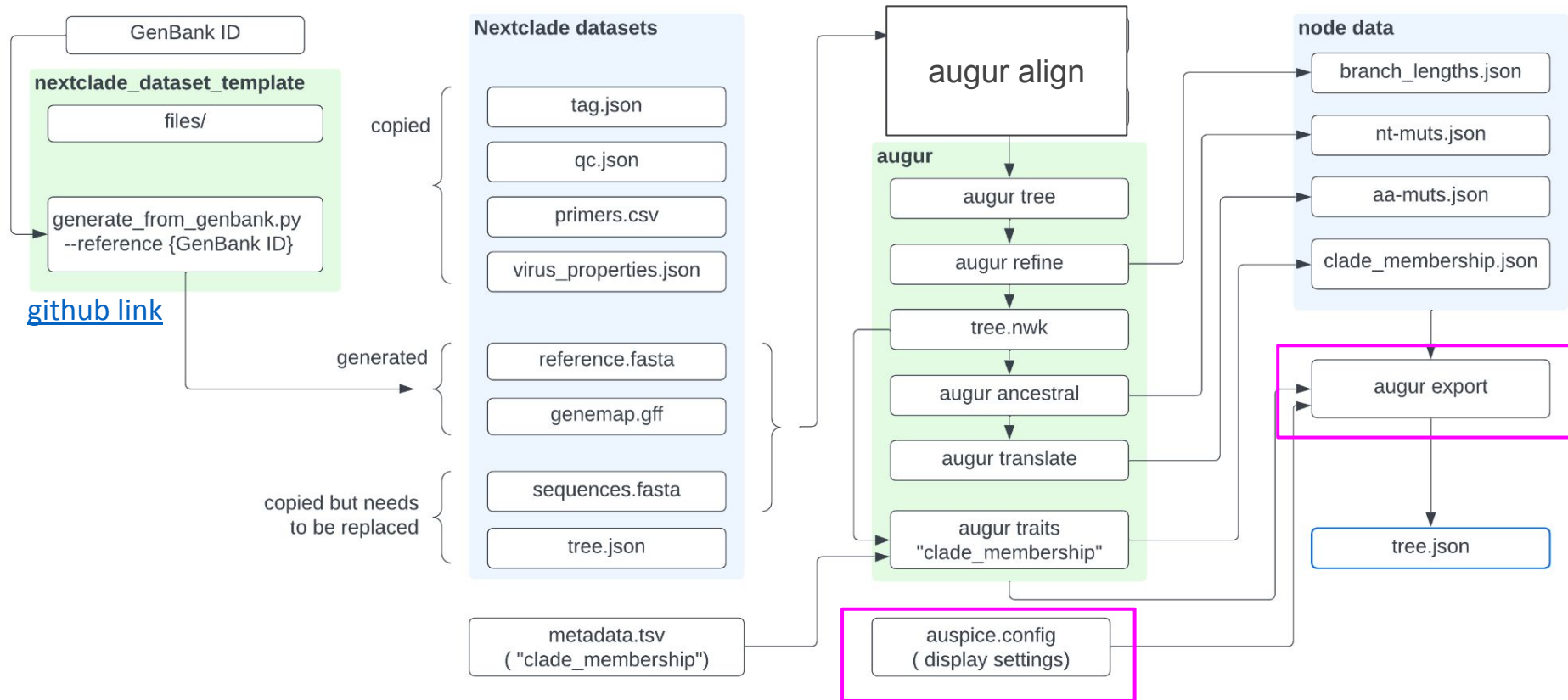


# Clade membership

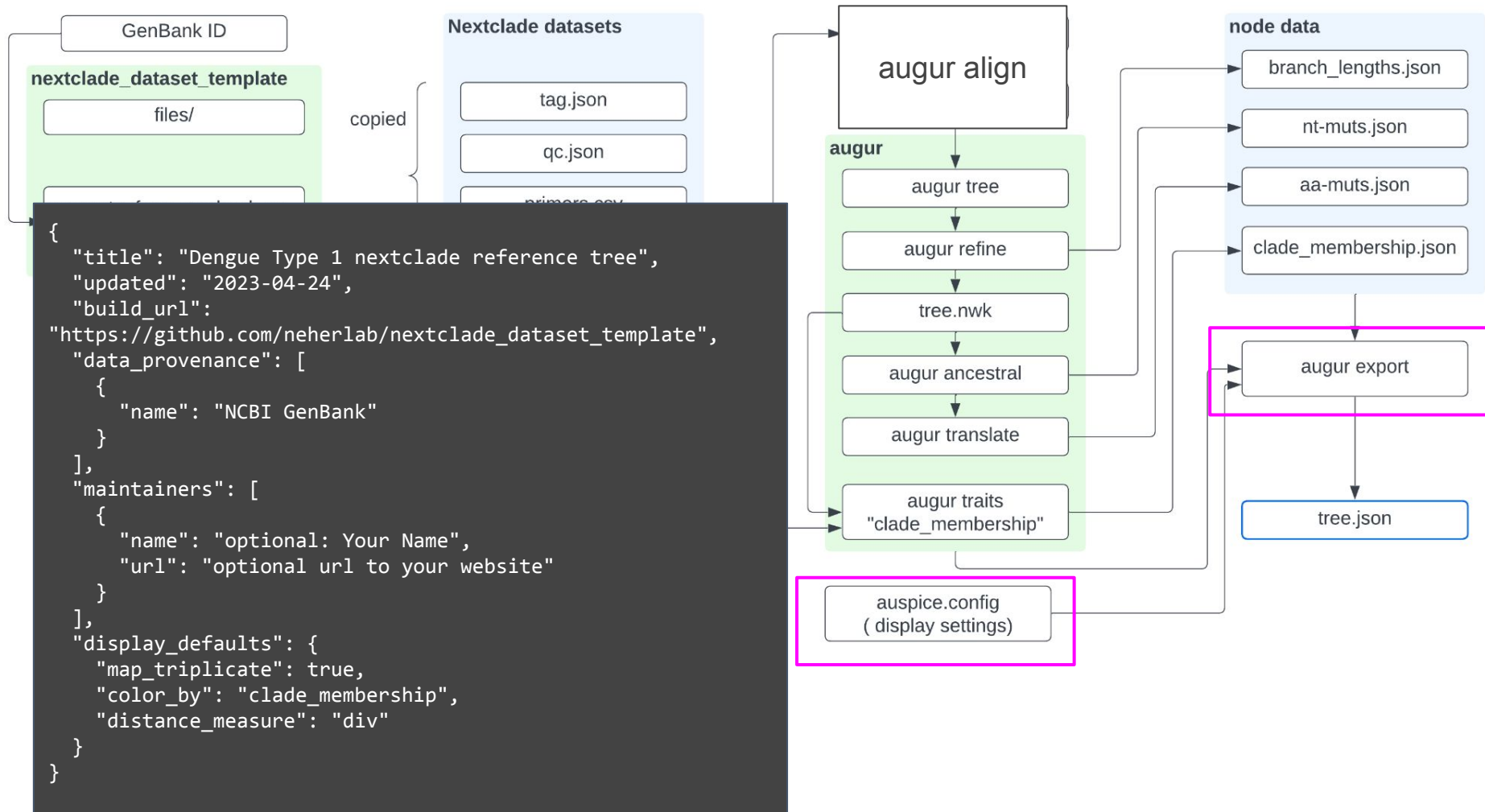


metadata.tsv should have "strain" and "clade\_membership" columns

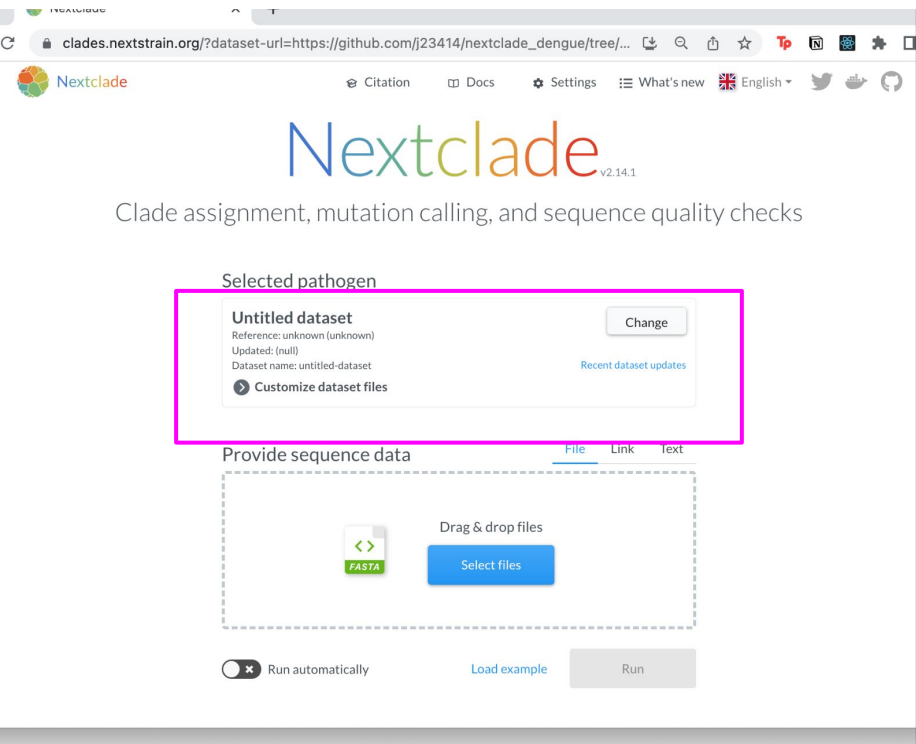
# Augur export



# Augur export



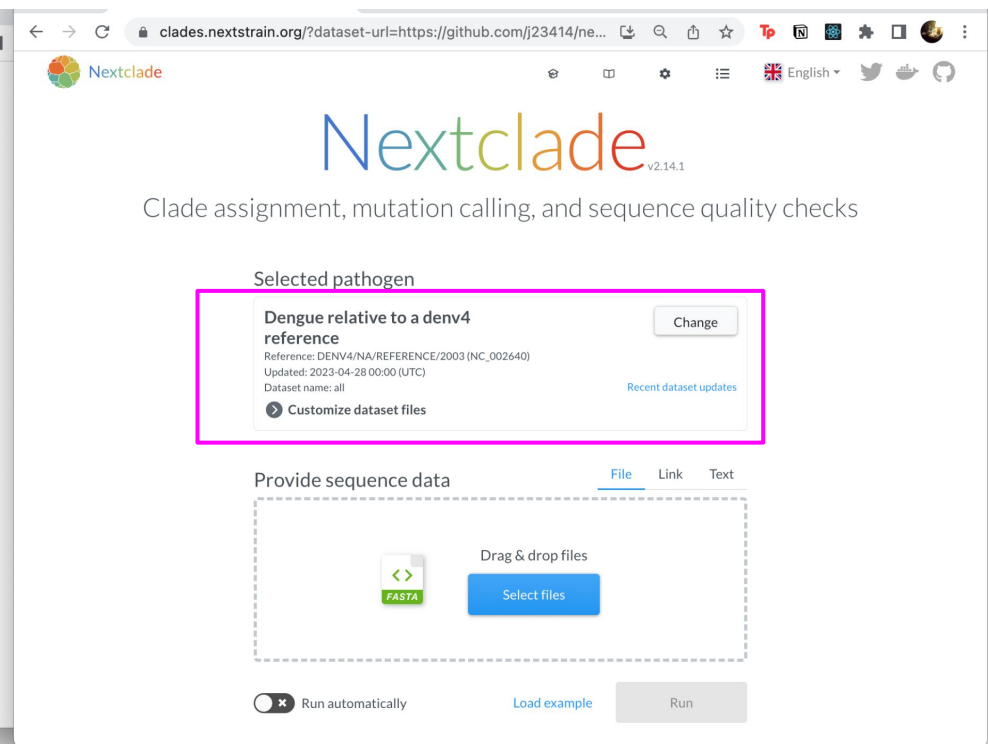
# Augur export



The screenshot shows the Nextclade web interface at the URL `clades.nextstrain.org/?dataset-url=https://github.com/j23414/nextclade_dengue/tree/...`. The page features the Nextclade logo (v2.14.1) and the text "Clade assignment, mutation calling, and sequence quality checks". Below this, the "Selected pathogen" section is highlighted with a pink box and contains the following information:

- Untitled dataset** (with a "Change" button)
- Reference: unknown (unknown)
- Updated: (null)
- Dataset name: untitled-dataset
- [Recent dataset updates](#)
- [Customize dataset files](#)

The "Provide sequence data" section is below, with tabs for "File", "Link", and "Text". It includes a dashed box for "Drag & drop files" with a "FASTA" icon and a "Select files" button. At the bottom, there is a "Run automatically" toggle (disabled), a "Load example" link, and a "Run" button.



The screenshot shows the Nextclade web interface at the URL `clades.nextstrain.org/?dataset-url=https://github.com/j23414/ne...`. The page features the Nextclade logo (v2.14.1) and the text "Clade assignment, mutation calling, and sequence quality checks". Below this, the "Selected pathogen" section is highlighted with a pink box and contains the following information:

- Dengue relative to a denv4 reference** (with a "Change" button)
- Reference: DENV4/NA/REFERENCE/2003 (NC\_002640)
- Updated: 2023-04-28 00:00 (UTC)
- Dataset name: all
- [Recent dataset updates](#)
- [Customize dataset files](#)

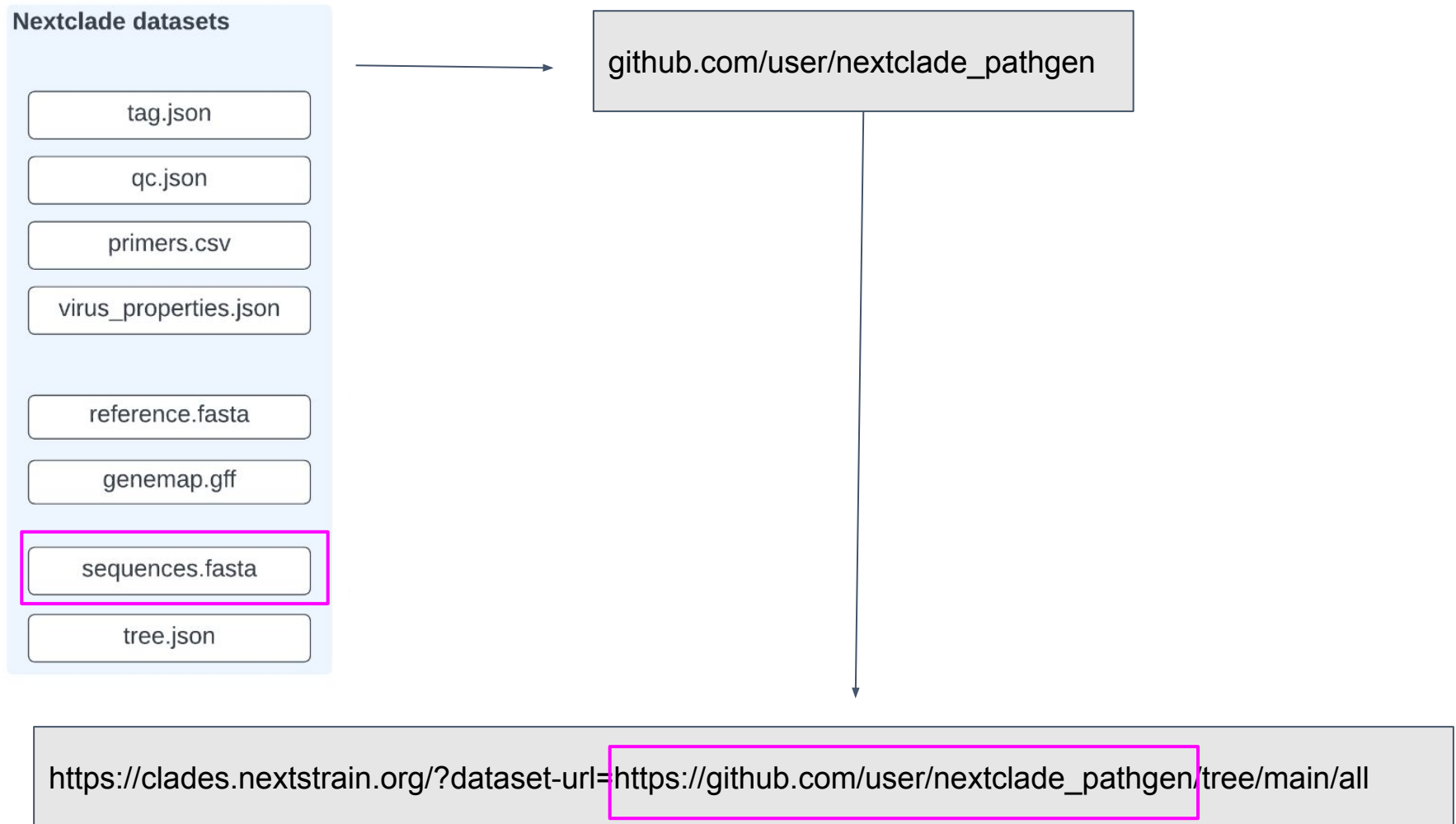
The "Provide sequence data" section is below, with tabs for "File", "Link", and "Text". It includes a dashed box for "Drag & drop files" with a "FASTA" icon and a "Select files" button. At the bottom, there is a "Run automatically" toggle (disabled), a "Load example" link, and a "Run" button.

# Agenda

- What is Nextclade?
  - Statement of need
  - What is a Nextclade dataset?
- **How do we create a Nextclade dataset?**
  - Nextclade template script
  - Nextalign vs Augur align
- Checking a Nextclade dataset
  - validation and basic checks
- Conclusions



# Validate Nextclade dataset



# Validate Nextclade Dengue - env1

## Nextclade datasets

tag.json

qc.json

primers.csv

virus\_properties.json

reference.fasta

genemap.gff

sequences.fasta

tree.json

github.com/j23414/nextclade\_dengue

j23414 / nextclade\_dengue Public

<> Code Issues 2 Pull requests Zenhub Actions Projects Wiki Security Insights

main 3 branches 0 tags Go to file Add file <> Code

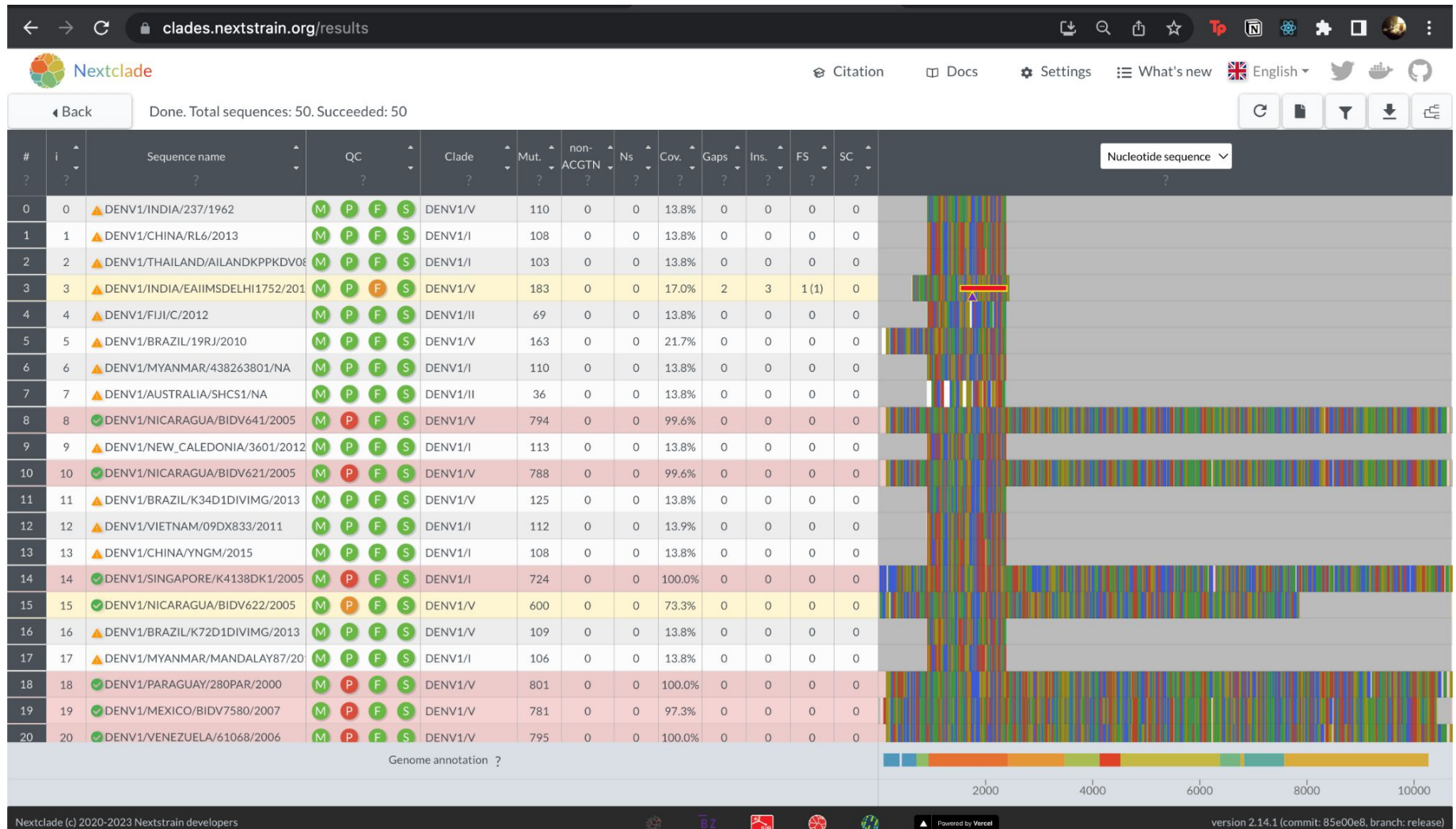
**Your main branch isn't protected**  
Protect this branch from force pushing or deletion, or require status checks before merging. [Learn more](#) [Protect this branch](#)

j23414 Mask regions other than Env 1116eb5 6 hours ago 5 commits

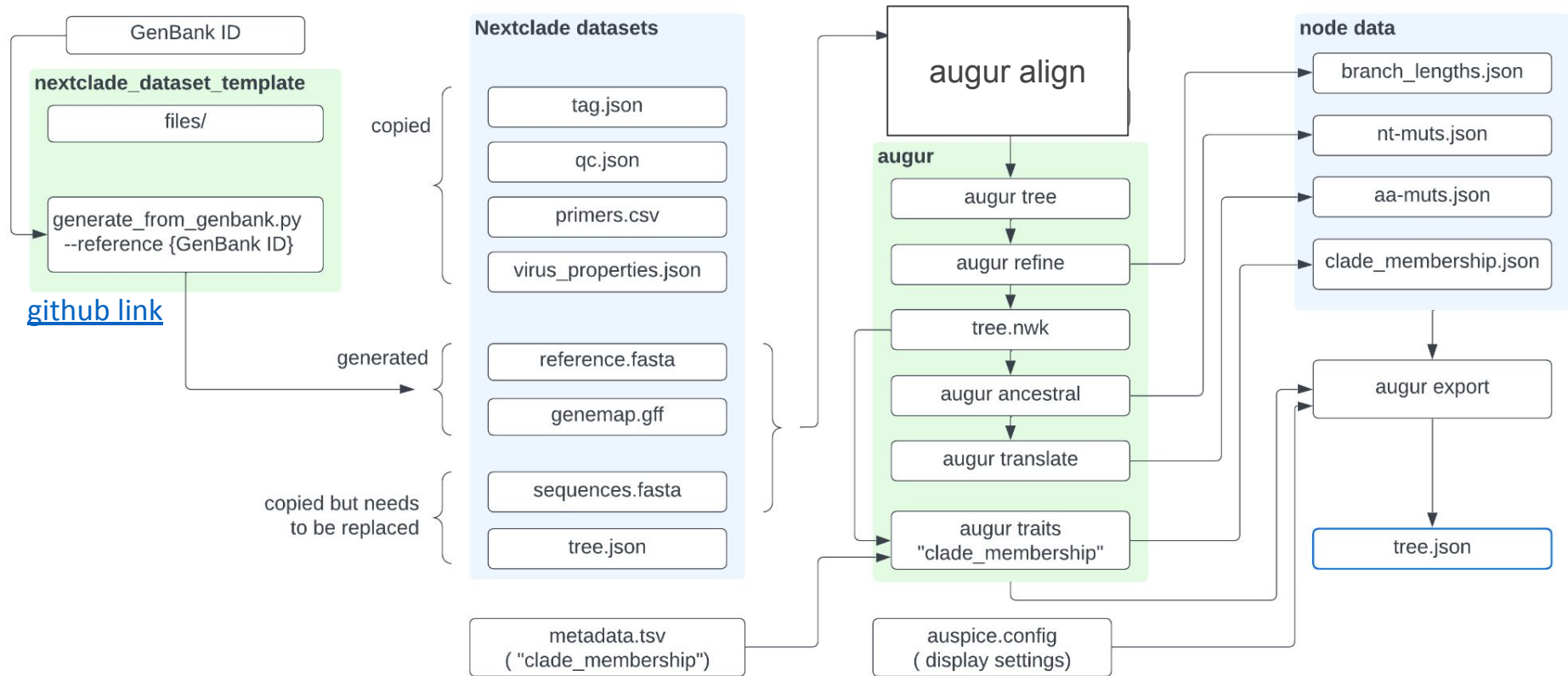
all	Mask regions other than Env	6 hours ago
env1	Mask regions other than Env	6 hours ago
env2	Mask regions other than Env	6 hours ago
env3	Mask regions other than Env	6 hours ago
env4	Mask regions other than Env	6 hours ago
01_run_Nextclade_template.sh	Generate Nextclade dataset template from the following commands:	9 hours ago
02_create_reference_tree.sh	Changes for Dengue dataset	8 hours ago
LICENSE	Initial commit	3 weeks ago
README.md	Generate Nextclade dataset template from the following commands:	9 hours ago

[https://clades.nextstrain.org/?dataset-url=https://github.com/j23414/nextclade\\_dengue/tree/main/env1](https://clades.nextstrain.org/?dataset-url=https://github.com/j23414/nextclade_dengue/tree/main/env1)

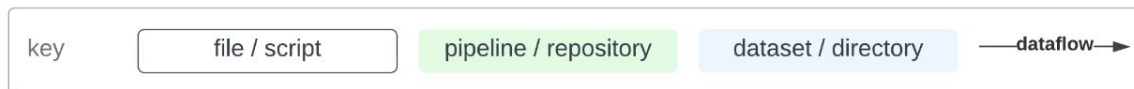
# Validate Nextclade Dengue - denv1



# In Summary



validation: [https://clades.nextstrain.org/?dataset-url=https://github.com/user/nextclade\\_pathgen/tree/main](https://clades.nextstrain.org/?dataset-url=https://github.com/user/nextclade_pathgen/tree/main)



[https://github.com/j23414/nextclade\\_dengue](https://github.com/j23414/nextclade_dengue)

[https://github.com/neherlab/nextclade\\_data\\_workflows](https://github.com/neherlab/nextclade_data_workflows)

# References

- Aksamentov, I., Roemer, C., Hodcroft, E.B. and Neher, R.A., 2021. [Nextclade: clade assignment, mutation calling and quality control for viral genomes](#). Journal of open source software, 6(67), p.3773.
- Hadfield, J., Megill, C., Bell, S.M., Huddleston, J., Potter, B., Callender, C., Sagulenko, P., Bedford, T. and Neher, R.A., 2018. [Nextstrain: real-time tracking of pathogen evolution](#). Bioinformatics, 34(23), pp.4121-4123.
- Huddleston, J., Hadfield, J., Sibley, T.R., Lee, J., Fay, K., Ilcisin, M., Harkins, E., Bedford, T., Neher, R.A. and Hodcroft, E.B., 2021. [Augur: a bioinformatics toolkit for phylogenetic analyses of human pathogens](#). Journal of open source software, 6(57).
- Sayers, E.W., Cavanaugh, M., Clark, K., Pruitt, K.D., Schoch, C.L., Sherry, S.T. and Karsch-Mizrachi, I., 2022. [GenBank](#). Nucleic acids research, 50(D1), p.D161.
- [Formal specifications of GFF3 format from Sequence Ontology](#)
- Reese, M.G., Moore, B., Batchelor, C., Salas, F., Cunningham, F., Marth, G.T., Stein, L., Flicek, P., Yandell, M. and Eilbeck, K., 2010. [A standard variation file format for human genome sequences](#). Genome biology, 11, pp.1-9.
- Arendsee, Z.W., Baker, A.L.V. and Anderson, T.K., 2022. [smot: a python package and CLI tool for contextual phylogenetic subsampling](#). Journal of Open Source Software, 7(80), p.4193.