# Enhancing NCBI Pathogen Detection cluster surveillance with ncbi-cluster-tracker

**Samuel Baird**

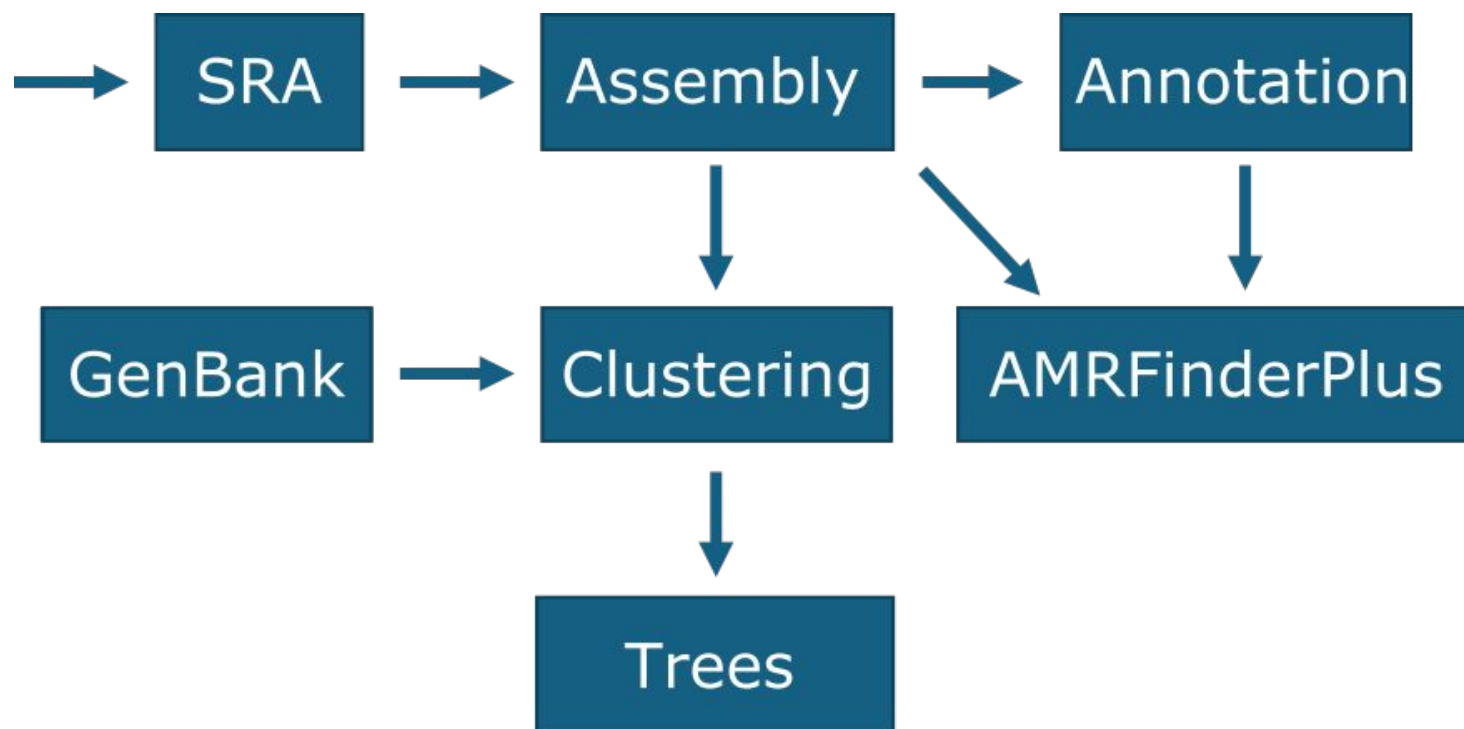**April 22, 2025**

**AMD Mountain Regional Bioinformatics Conference**

COLORADO
Department of Public
Health & Environment

# NCBI Pathogen Detection

Automated system that clusters related bacterial and fungal pathogen genome sequences submitted to NCBI to help identify possible links between cases for genomic epidemiology

COLORADO
Department of Public
Health & Environment

# NCBI Pathogen Detection

Isolates and associated clusters can be searched on the NCBI Isolate Browser website and phylogenetic trees investigated using the SNP Tree Viewer
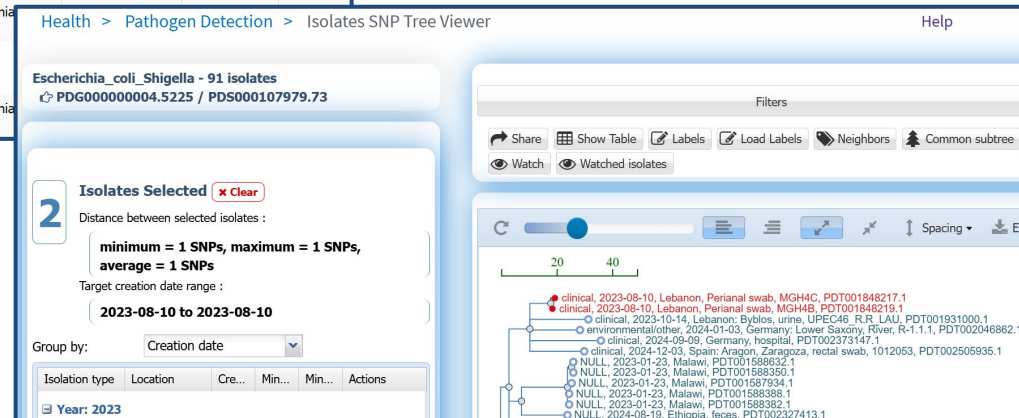


**Isolates Browser**

**SNP Tree Viewer**

https://www.ncbi.nlm.nih.gov/pathogens/

# NCBI Pathogen Detection

Pathogen Detection data can be also be accessed through the FTP file server and BigQuery data warehouse



## Index of /pathogen/Results/Klebsiella

| Name | Last modified | Size |
|------|---------------|------|
| Parent Directory | | - |
| AMR/ | 2025-04-16 23:52 | - |
| Clusters/ | 2025-04-16 23:52 | - |
| Exceptions/ | 2025-04-16 23:56 | - |
| Metadata/ | 2025-04-16 22:01 | - |
| SNP_trees/ | 2025-04-16 23:33 | - |
| PDG000000030.960.final.descriptor.xml | 202 | |
| PDG000000030.960.kmer.descriptor.xml | 202 | |
| PDG000000030.960.snps.descriptor.xml | 202 | |

**FTP Server**

**BigQuery data warehouse**

https://ftp.ncbi.nlm.nih.gov/pathogen/Results/
https://www.ncbi.nlm.nih.gov/pathogens/docs/getting_started_bigquery/

COLORADO
Department of Public
Health & Environment

# Original system for tracking HAI clusters

Every week: submit new sequencing data to NCBI, run BigQuery SQL query, copy results to Google sheet, review new clusters with epi at weekly meeting

# Goals for a new system and report format

1. Indicate new isolates and new clusters from the previous week
2. Annotate with additional metadata not published on NCBI (exact collection dates, internal lab IDs...)
3. Additional data visualizations (isolate counts over time, pairwise SNP distance matrix)
4. Report can be securely shared and explored
5. Agnostic to CDPHE data systems, could be deployed and adopted by other agencies

COLORADO
Department of Public
Health & Environment

# ncbi-cluster-tracker

## Inputs

**Previous cluster data for comparison** (optional)

Required    Optional metadata

```
biosample,id,collection_date,...
SAMN1010,482841,2024-08-02,...
SAMN1063,482896,2025-02-01,...
SAMN1029,482002,2024-12-18,...
```
**Sample sheet CSV**

**Browser CSV/TSV** (alternative to BigQuery)

## NCBI Pathogen Detection

**Isolates Browser**

BigQuery

**pdbrowser Data Warehouse**

FTP

**File server**

## ncbi-cluster-tracker

**Clusters**

| Cluster | Count | Change |
|---------|-------|--------|
| PDS123 | 3 | new |
| **PDS289** | **18** | **+2** |
| PDS661 | 6 | +0 |
| PDS180 | 5 | +0 |

**Isolates**

| BioSample | Source | |
|-----------|--------|----|
| SAMN1010 | internal | new |
| SAMN3122 | external | new |
| SAMN1063 | internal | |
| SAMN281 | external | |

## Outputs

**HTML**

**Cluster report**

**Clusters and isolates CSV**

snps/

**VCFs, Newicks**

# Demo

**Clusters table:** View clusters associated with internal isolates and changes to isolate counts since the previous report was created.



Clusters | Isolates | Cluster details

Command:

```
ncbi-cluster-tracker tests/data/sample_sheet.csv --compare-dir tests/data/20240826/ --browser-file tests/data/pdbrowser_
```

⟷ Comparing to tests/data/20240826/

☰ 11 rows   ▯ 10 columns   ⊞ 110 cells                                    ▾ Run SQL Query   Export ∨

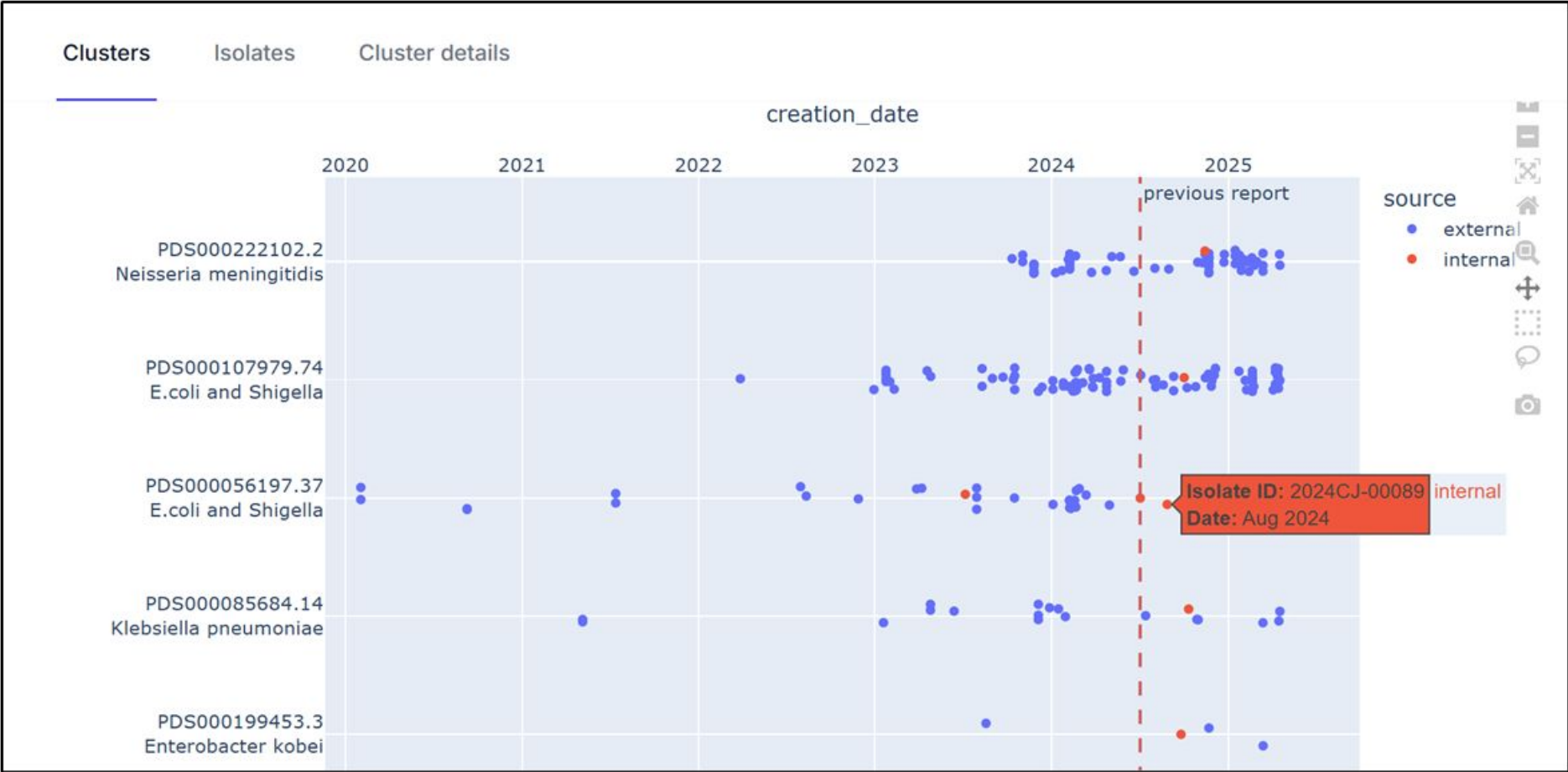| | cluster | taxgroup_name | internal_count | external_count | change | latest_added |
|---|---|---|---|---|---|---|
| 1 | PDS000205465.1 | Acinetobacter baumannii | 1 | 1 | new cluster | 2024-11-05 |
| 2 | PDS000107979.74 | E.coli and Shigella | 1 | 91 | new cluster | 2025-04-17 |
| 3 | PDS000056197.37 | E.coli and Shigella | 3 | 38 | +1 / +11 | 2025-04-17 |
| 4 | PDS000218192.1 | E.coli and Shigella | 1 | 2 | new cluster | 2025-02-06 |
| 5 | PDS000056207.36 | E.coli and Shigella | 1 | 23 | new cluster | 2025-01-30 |
| 6 | PDS000139993.9 | Enterobacter hormaechei | 1 | 13 | new cluster | 2024-08-27 |

# Cluster timelines: View when isolates were added to each cluster.

**Isolates table**: View metadata and associated clusters for specific isolates. Table can be sorted, filtered, and exported for further analysis.

**Cluster details:** Look up details about specific clusters using the dropdown menu or search bar.

Clusters    Isolates    **Cluster details**

[NEW] PDS000056197.37 - E.coli and Shigella    ▼

*E.coli and Shigella* cluster PDS000056197.37

| Internal isolates | External isolates | Total isolates |
|---|---|---|
| **3** +1 | **38** +11 | **41** +12 |

New internal isolates added:

- lab_id_25 / 2024CJ-00089

New external isolates added:

- SAMN39500079 / EC33821
- SAMN39500114 / EC23022

**Tree links and labels**: Load the Pathogen Detection tree with defined isolates selected and download label file to display custom metadata on the tree.

# SNP matrix: View pairwise SNP distances between internal and external isolates.



| | isolate_id | collection_date | geo_loc_name | SAMN18319160 | SAMN41036368 | SAMN41036328 | SAMN41036353 | SAMN32777951 | SAMN32745605 | SAMN14997961 | SAMN14997952 | lab_id_24 | lab_id_25 | lab_id_23 | SAMN29503575 | SAMEA112938329 | SAMEA112938368 | SAMN28857294 | SAMEA6368822 | SAMN30393072 | SAMN20033248 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SAMN18319160 | rs104 | 2018-05-16 | Banglad...: Dhaka | 0 | 18 | 17 | 16 | 24 | 20 | 21 | 20 | 32 | 32 | 22 | 22 | 20 | 19 | 24 | 13 | 17 | 9 |
| SAMN41036368 | 19AR0778 | 2019 | New Zealand | 18 | 0 | 3 | 2 | 22 | 18 | 19 | 18 | 32 | 32 | 22 | 24 | 22 | 21 | 26 | 23 | 27 | 19 |
| SAMN41036328 | 148450 | 2019 | New Zealand | 17 | 3 | 0 | 1 | 21 | 17 | 18 | 17 | 31 | 31 | 21 | 23 | 21 | 20 | 25 | 22 | 26 | 18 |
| SAMN41036353 | 19AR0675 | 2019 | New Zealand | 16 | 2 | 1 | 0 | 20 | 16 | 17 | 16 | 30 | 30 | 20 | 22 | 20 | 19 | 24 | 21 | 25 | 17 |
| SAMN32777951 | 20200317_MGL_35 | 2020-02-17 | India: ... Sahyog | 24 | 22 | 21 | 20 | 0 | 4 | 25 | 24 | 38 | 38 | 28 | 30 | 28 | 27 | 32 | 29 | 33 | 25 |
| SAMN32745605 | 2020031...L-35_B2 | 2020-02-17 | India: ... Sahyog | 20 | 18 | 17 | 16 | 4 | 0 | 21 | 20 | 34 | 34 | 24 | 26 | 24 | 23 | 28 | 25 | 29 | 21 |
| SAMN14997961 | BA8153 | 2019 | India: Vellore | 21 | 19 | 18 | 17 | 25 | 21 | 0 | 1 | 35 | 35 | 25 | 27 | 25 | 24 | 29 | 26 | 30 | 22 |
| SAMN14997952 | BA33222 | 2018 | India: Vellore | 20 | 18 | 17 | 16 | 24 | 20 | 1 | 0 | 34 | 34 | 24 | 26 | 24 | 23 | 28 | 25 | 29 | 21 |
| ⭐ lab_id_24 | 2024CJ-00073 | 2024-07-03 | USA | 32 | 32 | 31 | 30 | 38 | 34 | 35 | 34 | 0 | 0 | 10 | 38 | 36 | 35 | 40 | 37 | 41 | 33 |
| ⭐ lab_id_25 | 2024CJ-00089 | 2024-08-27 | USA | 32 | 32 | 31 | 30 | 38 | 34 | 35 | 34 | 0 | 0 | 10 | 38 | 36 | 35 | 40 | 37 | 41 | 33 |
| ⭐ lab_id_23 | 2023CJ-00175 | 2023-07-07 | USA | 22 | 22 | 21 | 20 | 28 | 24 | 25 | 24 | 10 | 10 | 0 | 28 | 26 | 25 | 30 | 27 | 31 | 23 |
| SAMN29503575 | 21D20CPO003B | 2020 | Canada | 22 | 24 | 23 | 22 | 30 | 26 | 27 | 26 | 38 | 38 | 28 | 0 | 20 | 21 | 26 | 27 | 31 | 23 |

# Technical note

Report-generating code written in Python (no HTML/CSS/JavaScript) using Arakawa to create interactive tables and plots within the report

```python
import altair as alt
import arakawa as ar
from vega_datasets import data

df = data.iris()

plot_base = alt.Chart(df).mark_point().interactive()

ar.Group(
    "Iris analysis",
    ar.Select(
        ar.DataTable(df, label='Data'),
        ar.Group(
            ar.Plot(plot_base.encode(
                x='sepalLength',
                y='sepalWidth',
                color='species')
            ),
            ar.Plot(plot_base.encode(
                x='petalLength',
                y='petalWidth',
                color='species')),
            columns=2,
            label='Plots'
        )
    )
)
```



https://ninoseki.github.io/arakawa/latest/blocks/layout-blocks/#group

# Limitations of tracking clusters with Pathogen Detection

- Limitations of ncbi-cluster-tracker:
    - Currently does not incorporate AMR results
    - Cannot add notes directly to report (this was easy to do in the Google Sheets system!)
- General limitations with using Pathogen Detection:
    - Clustering distance threshold and timeline are fixed, resulting in sensitivity and specificity issues with detecting potential outbreaks
    - Can be difficult to find useful or actionable information about closely related external isolates
    - Particularly relevant to HAI: inability to detect plasmid clusters to help track the spread of AMR genes

COLORADO
Department of Public
Health & Environment

# Limitations of tracking clusters with Pathogen Detection

- Limitations of ncbi-cluster-tracker:
  - Currently does not incorporate AMR results
  - Cannot add notes directly to report (this was easy to do in the Google Sheets system!)
- General limitations with using Pathogen Detection:
  - Clustering distance threshold and timeline are fixed, resulting in sensitivity and specificity issues with detecting potential outbreaks
  - Can be difficult to find useful or actionable information about closely related external isolates
  - Particularly relevant to HAI: inability to detect plasmid clusters to help track the spread of AMR genes

**Overall:** Clusters identified through ncbi-cluster-tracker and Pathogen Detection can be a good *starting point* (hypothesis generation), and can help confirm existing links identified through epidemiology, but require follow-up analysis and investigation

COLORADO
Department of Public
Health & Environment

# Future directions

- Automated querying of Pathogen Detection without needing BigQuery access
- Include AMR genes
- Ability to add notes to clusters / isolates
- Include annotated tree directly within report
- Incorporate references to any literature that exists about external isolates and clusters
- Provide more options to configure analysis
- Easier installation and usage (Pip/Docker options and cross-platform support)

# THANKS!

Documentation and source code can be found here:

github.com/CDPHE-bioinformatics/ncbi-cluster-tracker

Feature requests, bug reports, and pull requests are welcomed!

_____

**More questions?**

sam.baird@state.co.us
cdphe_bioinformatics_lab@state.co.us