# Introduction to Web-based Analyses for Bacterial Typing

Erkison Ewomazino Odih

SEQAFRICA

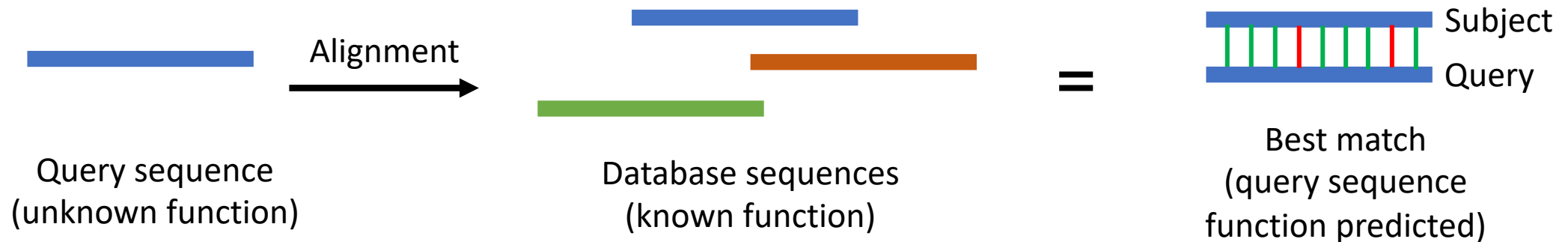# Outline

- General concepts
- Species identification
- Bacterial subtyping (MLST, serotyping)
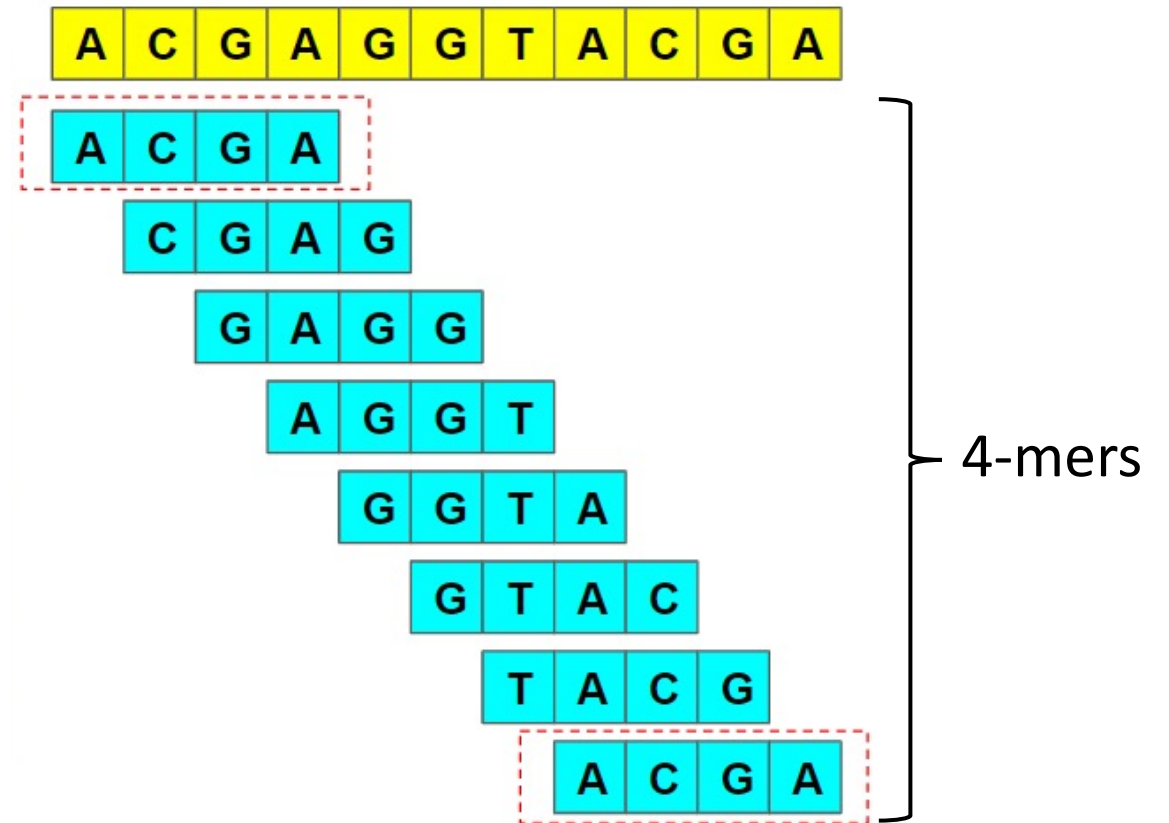- Practical

# General concepts

# Sequence Alignment (Applications)

- Discover functional information (annotation)

- Predict molecular structure

- Predict shared ancestry

- Phylogenetic analysis

- Sequence typing and identification (speciation, MLST prediction etc)



Query sequence
(unknown function)

Database sequences
(known function)

Best match
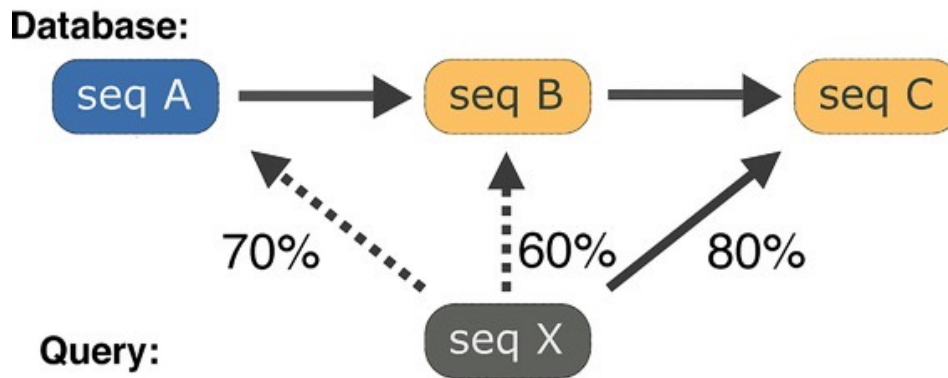(query sequence
function predicted)

# Kmers

- Contiguous substrings of a given sequence, each having length k (positive integer)

- Similar sequences expected to share more kmers in common



4-mers

# Databases

- Publicly accessible repository of [annotated] sequences
- Any prediction tool is only as good as its queried database
- Curated vs non-curated databases
- Find best matches for a query sequence in the database of sequences
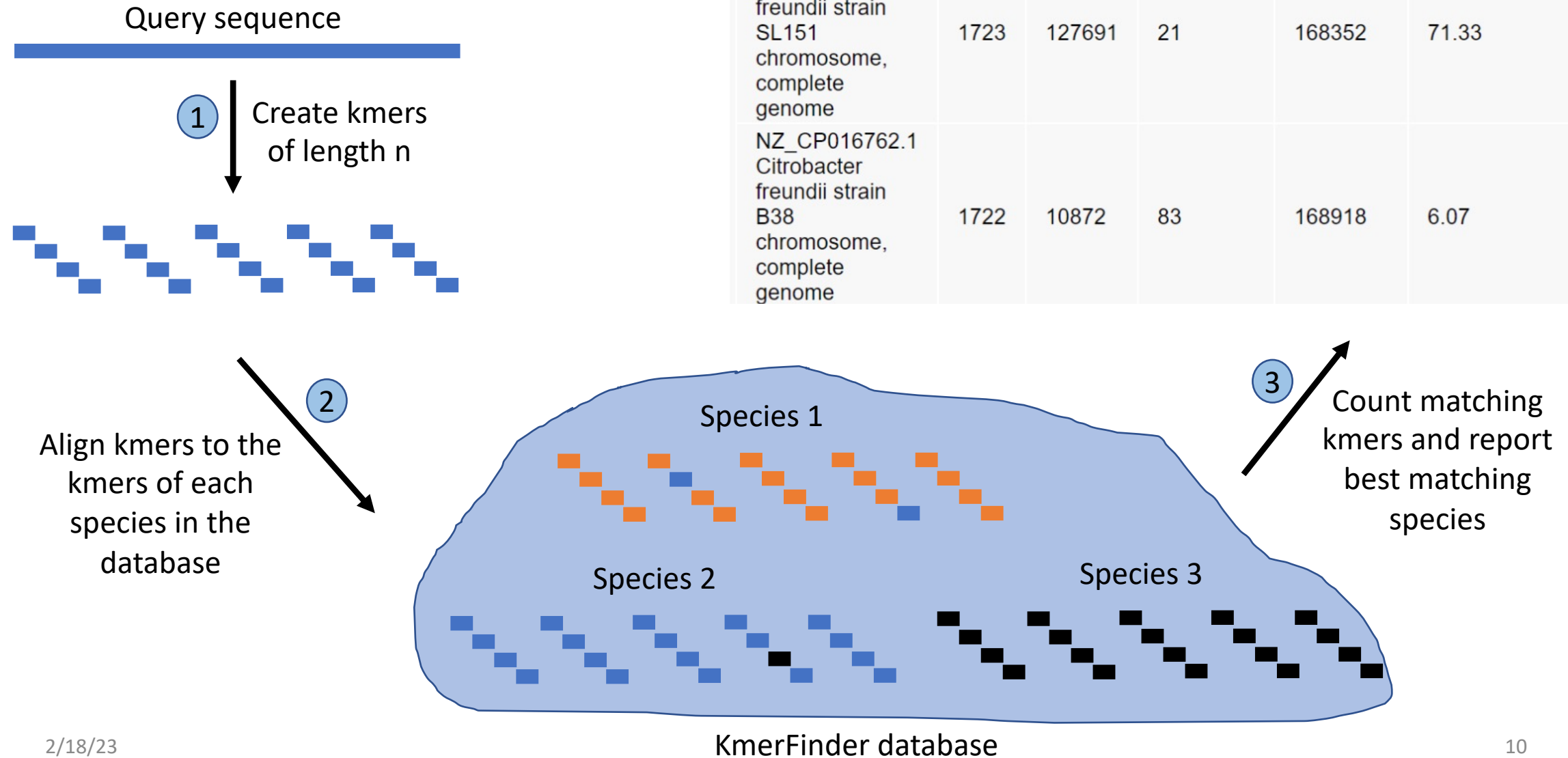
# Species identification

SEQAFRICA

# KmerFinder Tool

https://cge.food.dtu.dk/services/KmerFinder/

Query sequence

① Create kmers of length n

② Align kmers to the kmers of each species in the database

③ Count matching kmers and report best matching species

| Template | Num | Score | Expected | Template length | query_coverage |
|---|---|---|---|---|---|
| NZ_CP016952.1 Citrobacter freundii strain SL151 chromosome, complete genome | 1723 | 127691 | 21 | 168352 | 71.33 |
| NZ_CP016762.1 Citrobacter freundii strain B38 chromosome, complete genome | 1722 | 10872 | 83 | 168918 | 6.07 |

Species 1

Species 2

Species 3

KmerFinder database

2/18/23

10

# Pathogenwatch (Speciator)
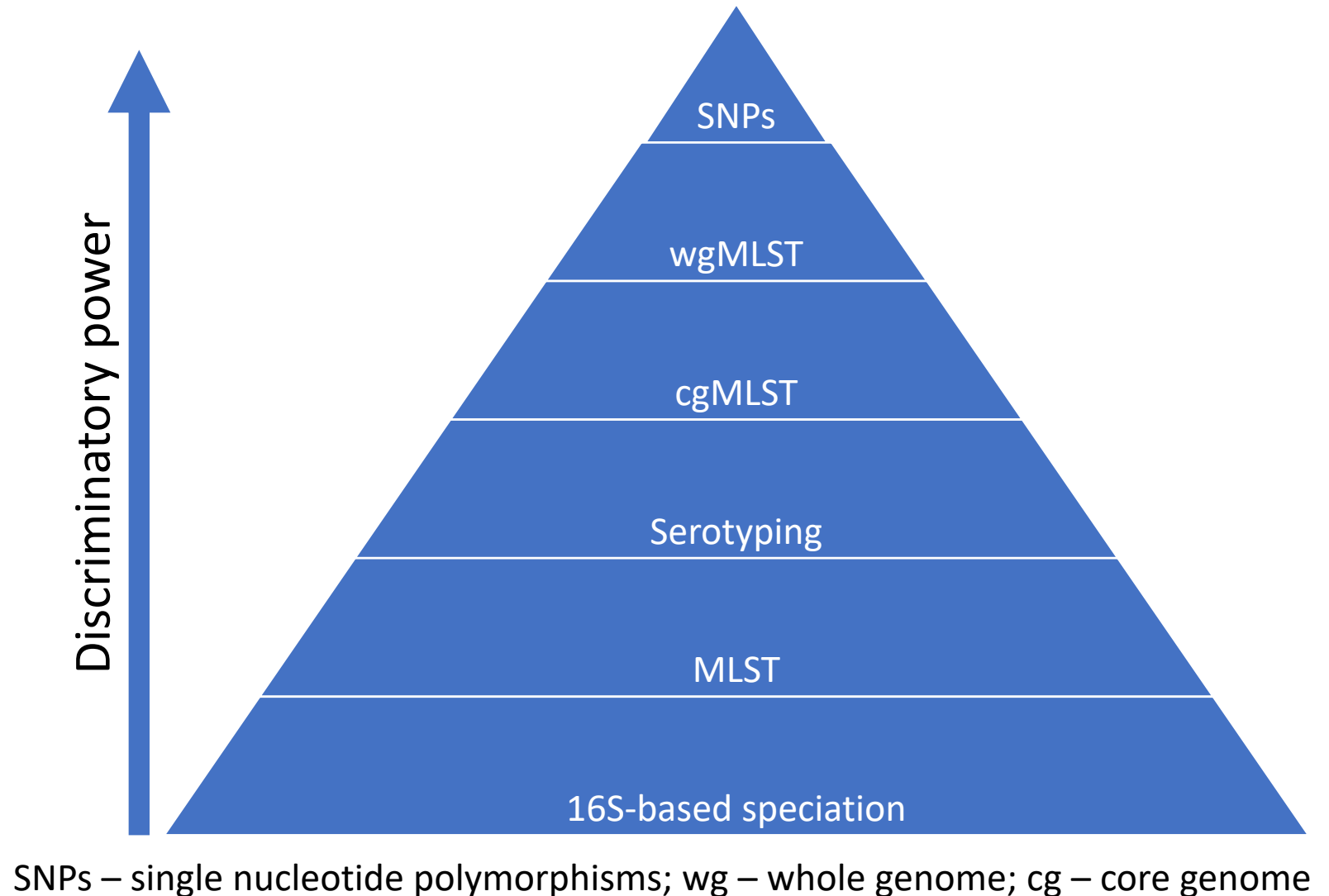
# Bacterial subtyping

# Why subtype?

Sub-classification of separate bacterial strains <u>within</u> the same species

- Delineate virulent subtypes (e.g., *E. coli* O157:H7)

- Epidemiological surveillance

- Outbreak investigation

- Identify emerging pathogenic strains (e.g., hybrid pathotypes)
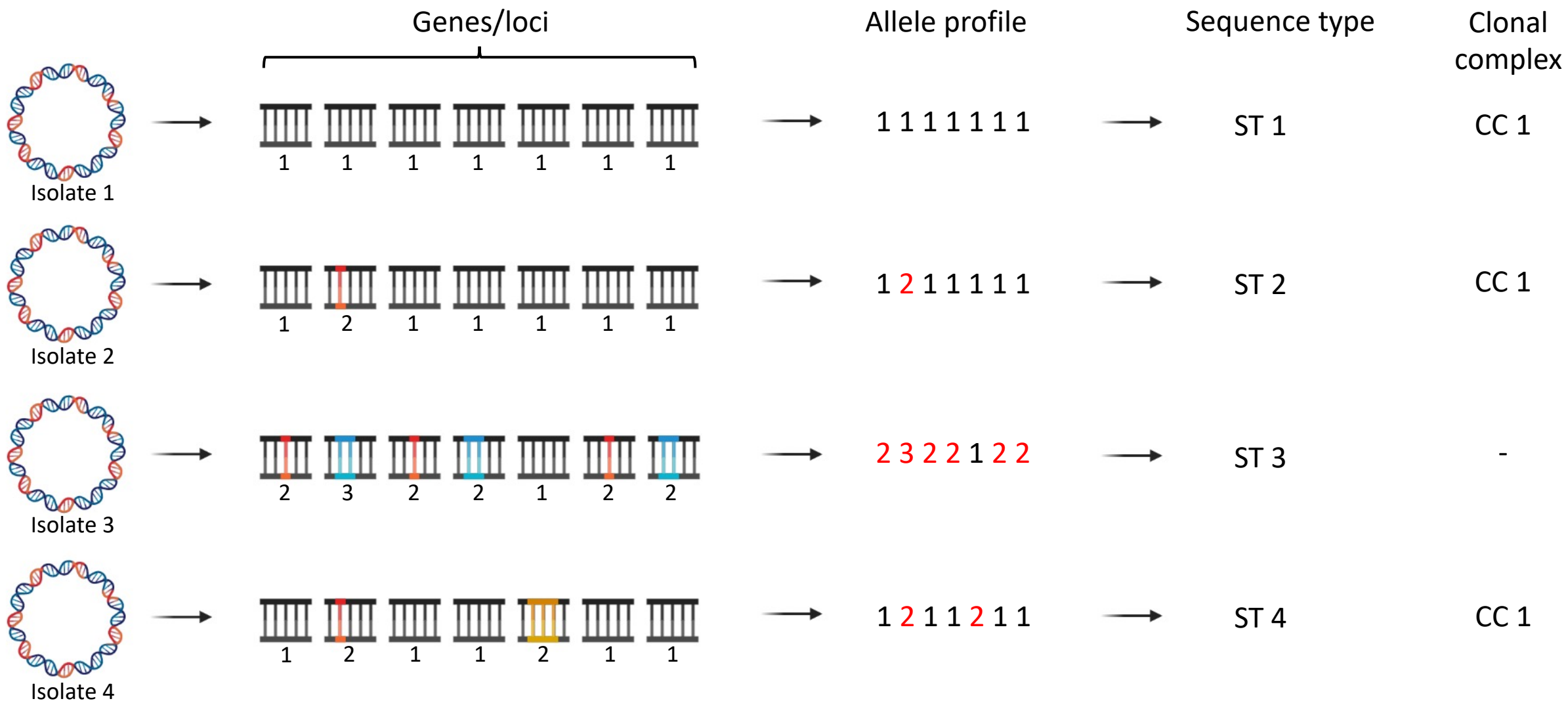
# Sequence-based subtyping

Merits:

- Not all genetic changes translate to morphological change

- Morphological convergence

- More "characters" compared

- Increased WGS access

Discriminatory power

SNPs

wgMLST

cgMLST

Serotyping

MLST

16S-based speciation

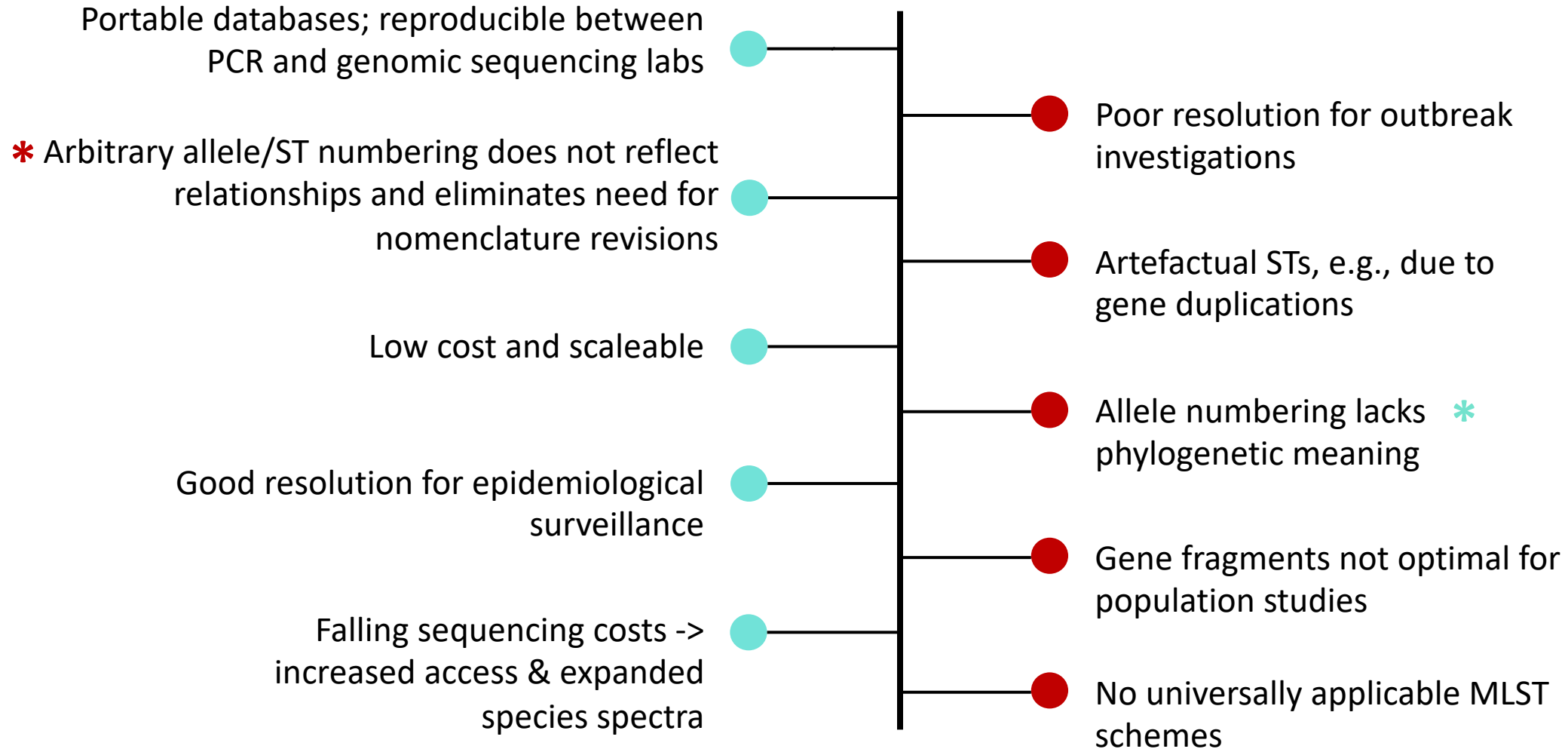SNPs – single nucleotide polymorphisms; wg – whole genome; cg – core genome

# Multi-locus sequence typing (MLST)

- Whole genome- or PCR-based sequencing of selected (n) housekeeping genes

- Genetic changes in these genes are constrained but discriminatory

- Multiple loci (usually 7) studied to address lack of congruence in bacteria

- Typically compares sequences of gene fragments (300 – 550 bp)

- Allele numbers for each loci assigned in order of discovery

- All allele numbers form an <u>allelic profile</u>

- Each ST corresponds to a unique allelic profile; STs are also assigned in order of discovery; database updated with new alleles and STs

- Clonal complexes typically defined as clusters with 1 or 2 varying loci

Genes/loci     Allele profile     Sequence type     Clonal complex

Isolate 1 → 1 1 1 1 1 1 1 → 1 1 1 1 1 1 1 → ST 1 → CC 1

Isolate 2 → 1 2 1 1 1 1 1 → 1 2 1 1 1 1 1 → ST 2 → CC 1

Isolate 3 → 2 3 2 2 1 2 2 → 2 3 2 2 1 2 2 → ST 3 → -

Isolate 4 → 1 2 1 1 2 1 1 → 1 2 1 1 2 1 1 → ST 4 → CC 1

# MLST: Pros and cons

Portable databases; reproducible between PCR and genomic sequencing labs

Poor resolution for outbreak investigations

* Arbitrary allele/ST numbering does not reflect relationships and eliminates need for nomenclature revisions

Artefactual STs, e.g., due to gene duplications

Low cost and scaleable

Allele numbering lacks phylogenetic meaning *

Good resolution for epidemiological surveillance

Gene fragments not optimal for population studies

Falling sequencing costs -> increased access & expanded species spectra

No universally applicable MLST schemes

# PubMLST

Public databases for molecular typing and microbial genome diversity

A collection of open-access, curated databases that integrate population sequence data with provenance and phenotype information for over 100 different microbial species and genera.

| 28,636,800 | 932,444 | 673,526 |
|:---:|:---:|:---:|
| ALLELES | ISOLATES | GENOMES |

pubmlst.org

# *Escherichia coli* MLST schemes (pubmlst.org)

### Achtman scheme

| ST | adk | fumC | gyrB | icd | mdh | purA | recA |
|----|-----|------|------|-----|-----|------|------|
| 1 | 4 | 2 | 2 | 4 | 4 | 4 | 4 |
| 2 | 5 | 3 | 2 | 6 | 5 | 5 | 4 |
| 3 | 6 | 4 | 3 | 7 | 7 | 7 | 6 |
| 4 | 6 | 5 | 4 | 8 | 8 | 8 | 2 |
| 5 | 7 | 6 | 5 | 9 | 9 | 8 | 2 |
| 6 | 8 | 7 | 1 | 1 | 10 | 8 | 6 |
| 7 | 9 | 8 | 5 | 1 | 11 | 8 | 7 |
| 8 | 10 | 9 | 5 | 10 | 12 | 9 | 2 |
| 9 | 6 | 4 | 3 | 7 | 7 | 7 | 8 |
| 10 | 10 | 11 | 4 | 8 | 8 | 8 | 2 |

### Pasteur scheme

| ST | dinB | icdA | pabB | polB | putP | trpA | trpB | uidA |
|----|------|------|------|------|------|------|------|------|
| 1 | 1 | 1 | 2 | 1 | 1 | 2 | 3 | 1 |
| 2 | 8 | 2 | 7 | 3 | 7 | 1 | 4 | 2 |
| 3 | 3 | 8 | 5 | 11 | 8 | 3 | 5 | 3 |
| 4 | 2 | 4 | 6 | 4 | 1 | 6 | 1 | 1 |
| 5 | 5 | 3 | 3 | 10 | 5 | 8 | 2 | 5 |
| 6 | 1 | 7 | 1 | 9 | 2 | 20 | 1 | 6 |
| 7 | 6 | 6 | 4 | 2 | 6 | 7 | 2 | 4 |
| 8 | 23 | 9 | 8 | 12 | 9 | 11 | 7 | 13 |
| 9 | 9 | 20 | 15 | 7 | 4 | 9 | 6 | 9 |
| 10 | 4 | 18 | 10 | 5 | 2 | 4 | 1 | 6 |

https://pubmlst.org/bigsdb?db=pubmlst_escherichia_seqdef
https://enterobase.readthedocs.io/en/latest/mlst/mlst-legacy-info-ecoli.html

# MLST 2.0 (CGE)

Platform: Web-based

Input: Sequence reads (.fastq) <u>OR</u> assemblies (.fasta)

URL: https://cge.food.dtu.dk/services/MLST/

Comments: Available for many clinically important species; regularly updated DB



2/18/23 SEQAFRICA 36

# Enterobase

Platform:     Web-based

Input:     Sequence reads (fastq)

URL:     https://enterobase. warwick.ac.uk/

Comments:     Specific for *E. coli*, *Salmonella*, etc. Also generates cgMLST, serotyping results, etc.

# MLST

Platform:  CLI (Linux/OS X)

Input:  Assembled and/or annotated genomes (fasta/GenBank/EMBL)

URL:  https://github.com/tseemann/mlst

Comments:  Preferred for large datasets; manually update database once installed

# *E. coli* serotyping

- Serological typing based on differences in the lipopolysaccharide <u>O antigen</u>, capsular <u>K antigen</u>, and flagellar <u>H antigen</u>.

- Serotype information very useful epidemiologically because it is directly linked to antigenic response

- Poor phylogenetic correlation due to high propensity for recombination in these genes

- WGS serotyping based on sequence similarity

- O:H typing (standard serotyping) from WGS data:

  - O group:
    - O-antigen processing genes: *wzx, wzy, wzm*, and *wzt*
  - H group:
    - Flagellin-encoding genes – *fliC, flkA, fllA, flmA* and *flnA*

# *Salmonella* serotyping

- Phenotypic serotyping
  - Labor-intensive and expensive
  - Requires procurement and storage of multiple antisera
- Serovars designated by names or antigenic formula in the format: **O:H1:H2**
- Over 2500 serotypes in White-Kauffmann-Le Minor scheme
- Specific combinations of O and H antigenic types represent serotypes/serovars
- Genome-based serotyping targets same antigens as phenotypic assays:
  - Somatic (O) group antigen – *rfb* gene cluster (*wzx, wzy,* others)
  - Flagellar antigens – *fliC* and *fljB*

Zhang et al (2015) https://doi.org/10.1128/JCM.00323-15

# SeqSero (CGE)

Platform: Web-based

Input: Sequence reads (.fastq)
OR assemblies (.fasta)

URL: https://cge.food.dtu.dk/services/SeqSero/

Comments: *Salmonella* serotyping

## Center for Genomic Epidemiology

Home          Services          Publications          Contact

### SeqSero 1.2

Service | Instructions | Output | Article abstract | Citations

SeqSero predicts the Salmonella serotype of either the pre-assembled or raw read sequence data provided to the service.
**Note:** This service is hosted by CGE but all credit and scientific questions should be given to the original authors from Deng Lab (SeqSero).

**More info on Salmonella serotypes**
From Deng Labs website

Software version: Available on GitHub here
Download the Salmonella determinants databases from: Deng Labs website (zip file)

**Select Data type**
Assembled Genome/Contigs

# SISTR

Platform:   Web and CLI (Linux / OS X)

Input:   Assemblies (fasta)

URL:   https://github.com/phac-nml/sistr_cmd
https://sistr-app.herokuapp.com/

Comments:   *Salmonella* typing: serovar and serogroup prediction, cgMLST, etc.



| HOME | RESULTS | QUEUE | HISTORY |

## SISTR: Salmonella In Silico Typing Resource

We present the Salmonella In Silico Typing Resource (SISTR) version 1.1.1, a bioinformatics platform for rapidly performing simultaneous in silico analyses for several leading subtyping methods on draft Salmonella genome assemblies. In addition to performing serovar prediction by genoserotyping, this resource integrates sequence-based typing analyses for: Multi-Locus Sequence Typing (MLST), ribosomal MLST (rMLST), and core genome MLST (cgMLST).

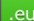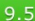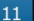Click or Drop assembly file(s) in FASTA format here for typing.

Submit

**Citation:** The *Salmonella In Silico* Typing Resource (SISTR): an open web-accessible tool for rapidly typing and subtyping draft *Salmonella* genome assemblies. Catherine Yoshida, Peter Kruczkiewicz, Chad R. Laing, Erika J. Lingohr, Victor P.J. Gannon, John H.E. Nash, Eduardo N. Taboada. *PLoS ONE 11(1): e0147101. doi: 10.1371/journal.pone.0147101*

**Notes:**
1) Do not submit sensitive private data to this public demo website. Rather deploy SISTR web application privately.
2) Submitted data is stored temporary due to hosting virtual machine shut down after 30 min of inactivity.
3) This site could be deployed locally or on your own infrastructure with source code available at https://github.com/phac-nml/sistr-web-app

© Copyright 2021. Canada's National Microbiology Laboratory

# ECTyper

Platform: Web and CLI

Input: Assemblies (fasta) or raw reads (fastq)

URL: https://github.com/phac-nml/ecoli_serotyping

https://usegalaxy.eu/root?tool_id=ectyper

Comments: *E. coli* serotyping



## ECTyper (an easy typer)

`ECTyper` is a standalone versatile serotyping module for *Escherichia coli*. It supports both *fasta* (assembled) and *fastq* (raw reads) file formats. The tool provides convenient species identification coupled to quality control module giving a complete, transparent and reference laboratories suitable report on E.coli serotyping.

## Dependencies:

- python >= 3.5
- bcftools >= 1.8
- blast == 2.7.1
- seqtk >= 1.2
- samtools >= 1.8
- bowtie2 >= 2.3.4.1
- mash >= 2.0

# Questions

# Practical

- CGE tools: https://cge.food.dtu.dk/services/

  - Species identification

  - MLST determination

  - Serotyping (if applicable)

- Multi-analyses tools:

  - Pathogenwatch [and Kleborate] - https://pathogen.watch/

  - Enterobase - https://enterobase.warwick.ac.uk/