# Introduction to online tools: Kmers, MLST and Serotyping

## SEQAFRICA Module 1

**The Fleming Fund**
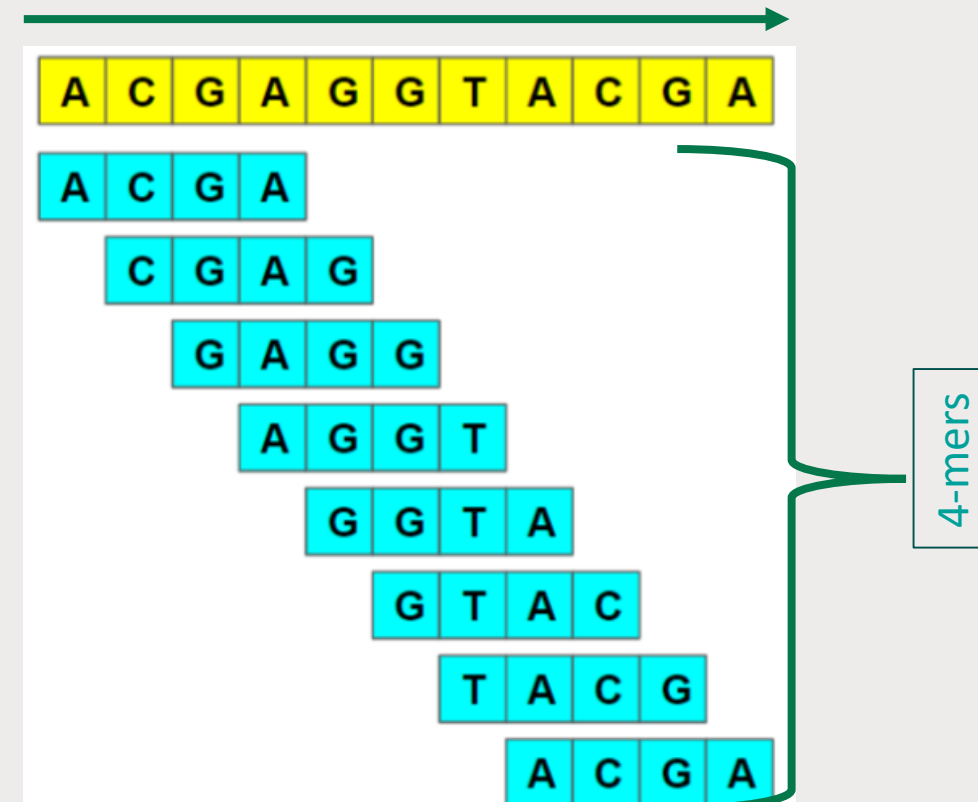*Regional Grants*

19 February 2021

**Stanford Kwenda**
**Lead Data Analyst**

**National Institute for Communicable Diseases**
**South Africa**

# Bioinformatics web-based tools

- Variety of methods and tools are available to analyze bacterial pathogens
- Most bioinformatics tools are implemented in Unix environments
  - Require at least some bioinformatics expertise for usage
- Web-based bioinformatics tools
  - Often free for use
  - Do not require computational power from the user
  - Limited bioinformatics knowledge
  - In some cases demand that that users deposit the analyzed data in public repositories
- Always make an effort to browse through the documentation of web based platforms
  - Helps with choice of parameters and interpretation of results

Uelze *et al.* One Health Outlook (2020)

# Introduction to Kmers

- A k-mer typically refers to all the possible substrings of length k that are contained in a string
- E.g. For sequence with length N
  - No. of K-mers = N-k+1
- K-mer counting
  - Total count, distinct count and unique count
- Computing k-mers is much faster than producing alignments
- Applications using k-mers:
  - Error correction
    - Rare k-mers are more likely due to sequence errors
  - Classification
    - Certain k-mers may uniquely identify genomes
  - Pseudo-alignment
    - New pseudo-aligners can match reads to locations based solely on the presence of common k-mers

# Web-based species identification using Kmers



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

# Web-based species identification using Kmers



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

# Multi-Locus Sequence Typing (MLST)

- A universal, portable, and precise means of typing bacteria
  - Traditionally involved PCR amplification and DNA sequencing of PCR fragments

- Indexes sequences (~500 bp) of representative housekeeping genes
  - Usually from seven loci
  - Each unique allele is assigned an arbitrary integer identifier
  - Sequence type (ST) – unique combinations of the alleles at each locus

- NGS sequence data can be matched to a database of allelic profiles
  - A single nucleotide variation at any of these loci defines a different allele and informs the ST

- MLST often not discriminative enough for outbreak detection or distinguishing highly related strains

Belen et al. Methods Mol Bio (2009)
https://pubmlst.org/multilocus-sequence-typing
Jolley and Maiden. BMC Bioinformatics (2010)

# Web-based MLST typing tools: MLST2.0

https://cge.cbs.dtu.dk/services/MLST/

Must select MLST scheme prior to launching the tool

- Input: raw reads or assembled genomes
- Can only take data for a single isolate at a time
- Some species will have more than one MLST scheme

# Web-based MLST typing tools: PathogenWatch



https://pathogen.watch/

# Web-based MLST typing tools: PathogenWatch



https://pathogen.watch/

# Web-based MLST typing tools: MLST2.0



https://pubmlst.org/

## Submit data to PubMLST

We welcome submissions to the databases hosted on PubMLST. Each organism-specific set of databases is overseen by one or more curators from all over the world who are usually researchers working on that organism and who, therefore, have an understanding of the biology and environment of the organism. Instructions for how to submit to each database can be found linked from the front page for each organism.

## How to submit

Some smaller databases currently just have an E-mail link to the curator where data can be sent directly. The larger or newer databases use an automated submission system that allows you to upload data via the website. Submissions are then routed automatically to a curator. We are gradually rolling out the automated system for all databases. To submit, please:

1. Access the site for the organism that you are submitting for.
2. Click the 'Submit' link on the organism page and follow the instructions.

https://pubmlst.org/submit-data

## Submitting data using the submission system

The automated submission system allows users to submit data (new alleles, profiles, or isolates) to the database curators for assignment and upload to the database. The submission system is enabled on a per-database basis so will not always be available.

If the system is enabled, new submissions can be made by clicking the 'Manage submissions' link on the database front page.



https://bigsdb.readthedocs.io/en/latest/submissions.html

# **Bacterial serotyping**

- Serotypes (serovars) – groups within a single species of microorganisms, (e.g. bacteria) which share distinctive surface structures e.g. surface antigens
  - Serotyping – classifying species at a sub-species level based on e.g. antigen properties into serogroups


- Strain – single isolates from pure cultures of a given species with distinct phenotypic/genotypic traits


- Virulence and pathogenicity tend to correlate well with subtype assignments

https://www.microscopemaster.com/serotype.html
Uelze *et al.* One Health Outlook (2020)
https://www.cdc.gov/salmonella/reportspubs/salmonella-atlas/serotyping-importance.html

# Serotyping in *Salmonella*

- *Salmonella* can be separated into many serotypes based on:
  - O antigen – outermost portion of the bacteria's surface covering
  - H antigen – slender threadlike structure that is part of the flagella

- Serotypes determined based on the distinct combination of O and H antigens
  - >2500 serotypes described for *Salmonella*
  - <100 serotypes account for most human infections

- Each known antigen is assigned a number and letter code
  - Combined into a seroformula
  - For example the White-Kauffmann-Le Minor scheme for *Salmonella*

https://www.cdc.gov/salmonella/reportspubs/salmonella-atlas/serotyping-importance.html
Uelze *et al.* One Health Outlook (2020)

Web-based tools for Salmonella serotyping

http://www.denglab.info/SeqSero2

seqsero@gmail.com
to me ▾

SeqSero2 Output:

Input_files    Predicted_antigenic_profile    Predicted_serotype(s)
ecoli.fasta    -:-:-
salmonella.fasta    3,10:y:1,5

Individual Prediction Details:

Output directory:    ecoli.fasta_directory
Input files:   ecoli.fasta
O antigen prediction:   -
H1 antigen prediction(fliC):   -
H2 antigen prediction(fljB):   -
Predicted subspecies:   -
Predicted antigenic profile:   -:-:-
Predicted serotype:    - -:-:-
Note:   The input genome cannot be identified as Salmonella. Check the input for taxonomic ID, contamination, or sequencing quality.

Output directory:    salmonella.fasta_directory
Input files:   salmonella.fasta
O antigen prediction:   3,10
H1 antigen prediction(fliC):   y
H2 antigen prediction(fljB):   1,5
Predicted subspecies:   I
Predicted antigenic profile:   3,10:y:1,5
Predicted serotype:    Orion
Note:

http://www.denglab.info/SeqSero2

# Web-based tools for *Salmonella* serotyping



- Input: assembled genome
- Can accept multiple assemblies

https://lfz.corefacility.ca/sistr-app/

# Web-based tools for *Salmonella* serotyping

# Web-based tools for *Salmonella* serotyping



https://cge.cbs.dtu.dk/services/SalmonellaTypeFinder-1.4/

# Serotyping in *E. coli*

- WGS of *E. coli* is replacing established subtyping methods such as PFGE
  - Traditional methods such as use of antibodies to test for surface antigens e.g. O, H, and K antigens

- Approx. 190 known *E. coli* serovars
  - Currently ~186 different O-groups and 53 H-types
  - Serotyping is highly complex

Fratamico et al. Front Microbiol (2016)
Uelze *et al.* One Health Outlook (2020)

# Web-based tools for *E. coli* serotyping



https://cge.cbs.dtu.dk/services/SerotypeFinder/

# Web-based tools for *E. coli* serotyping

https://cge.cbs.dtu.dk/services/SerotypeFinder/