# Bacterial species identification using sequencing data

**Introduction to Bioinformatics - Module 1**

Public Health Alliance for Genomic Epidemiology

Ifeoluwa  J. Akintayo

# Bacteria Identification

The process of determining the taxon a bacteria isolate belongs.

*It is important as a bacteriologist to identify pathogens in order to give the appropriate treatment options.*

# Bacteria Identification Methods

- Phenotypic Method (Macroscopy & Microscopy)
  - ✓ Appearance
  - ✓ Gram stains etc

- Biochemical test Method
  - ✓ Acid production
  - ✓ Gas production
  - ✓ Enzyme production
  - ✓ Sugar fermentation etc

- **Genotypic Method –** Genetic material of the organism
  - ✓ Pattern-Fingerprint based technology
  - ✓ **Sequence based technology**

# Sequence-based technique

The use of DNA sequences of the unknown species, compared to a comprehensive sequence database (e.g *National Center for Biotechnology Information* - NCBI) of known species -*Many organisms have similarities to a known species.*

with the use of

**Sequence alignment tool**– *arrangement of DNA nucleotides of an unknown species to a set of known species to identify regions of similarities*.

# Methods for species identification using sequencing data

- Basic Local Alignment Search Tool – BLAST

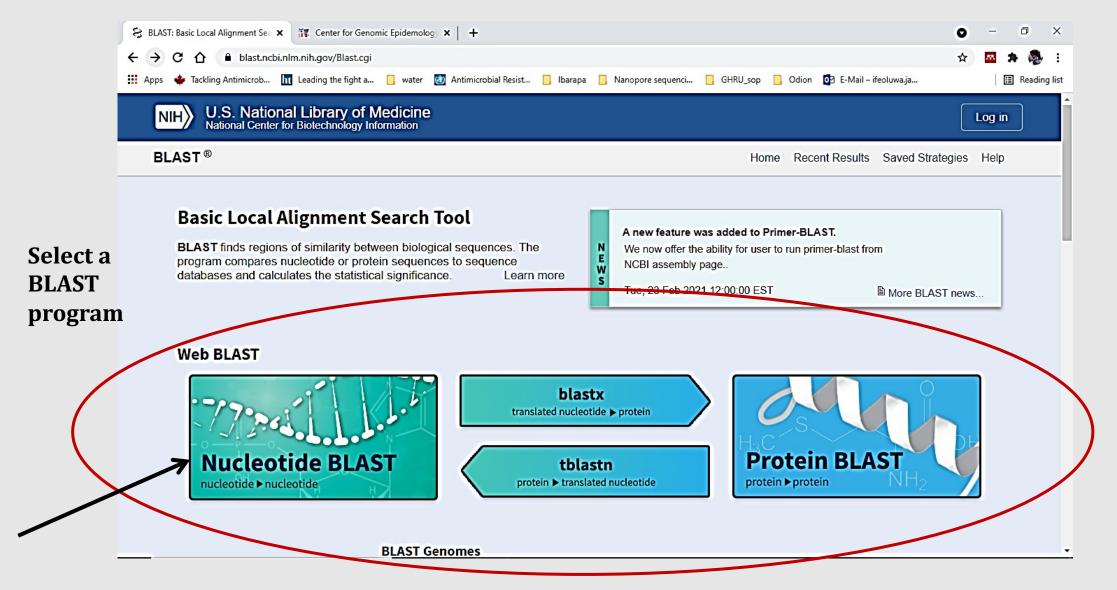- SpeciesFinder

- K-merFinder

- Pathogenwatch

# Basic Local Alignment Search Tool - BLAST

- Matches unknown sequence to known published sequences

- Compares genes and protein sequences against public database

- Search database for maximum alignment
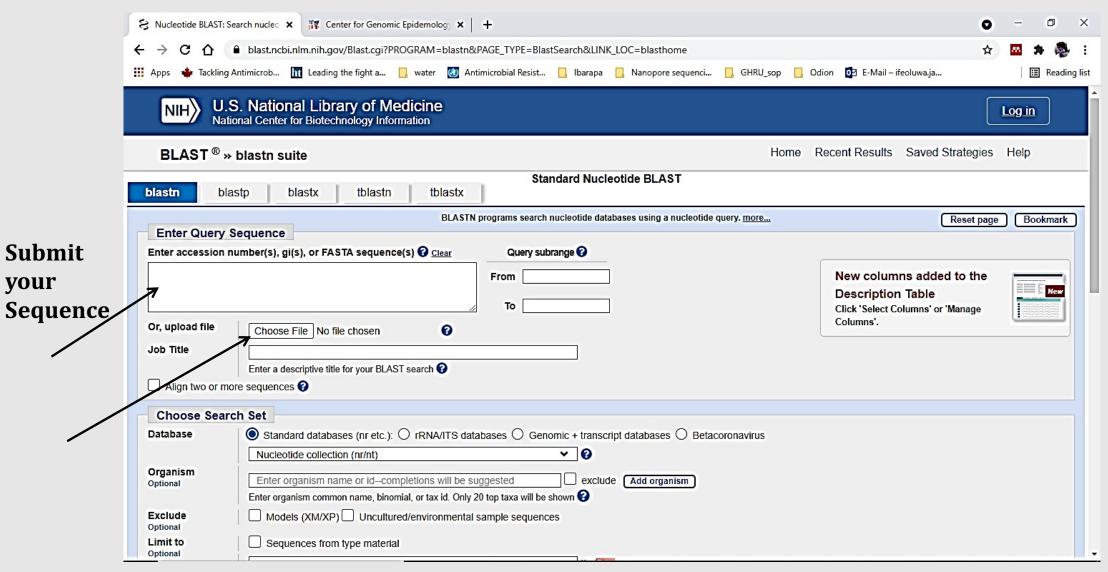
- Returns most similar sequences from the database

How does BLAST work?

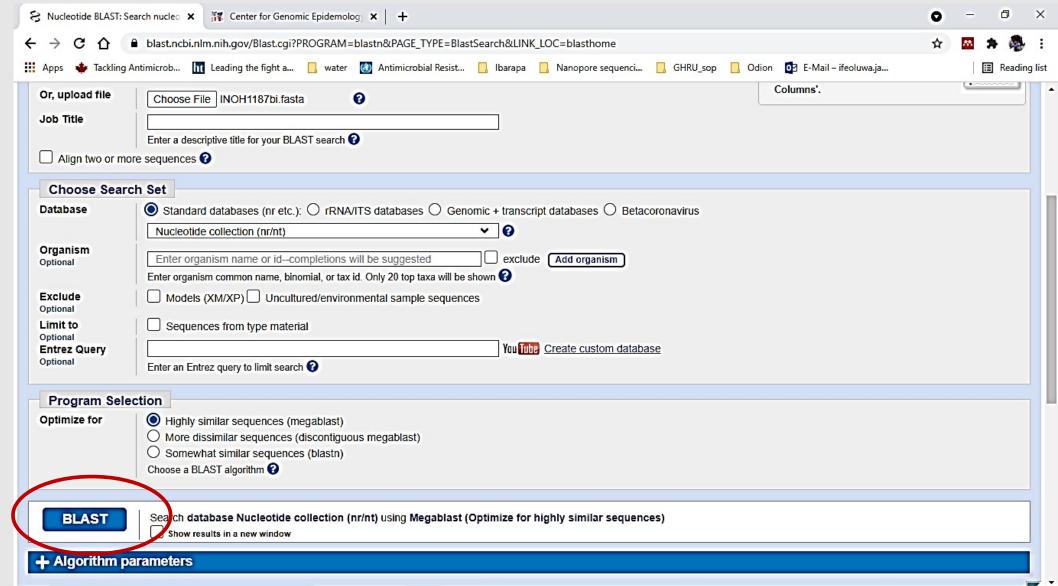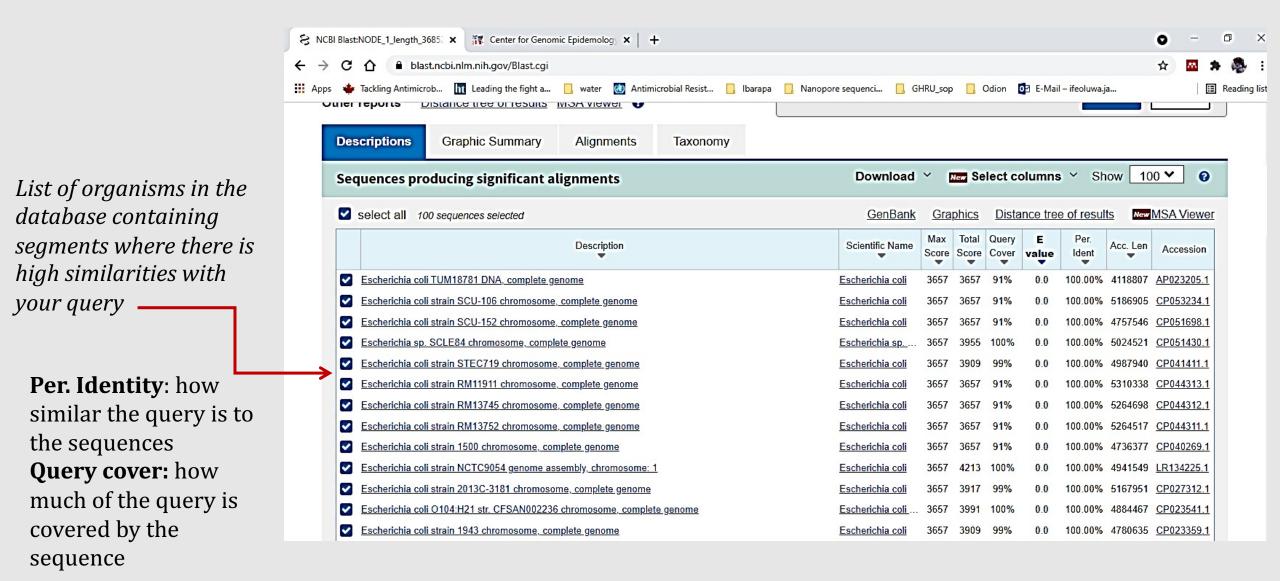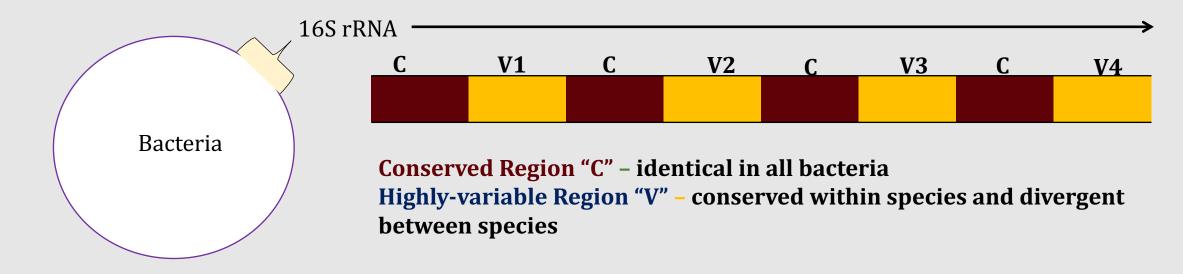https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE_TYPE=BlastSearch

**Submit your Sequence**

# BLAST OUTPUT

*List of organisms in the database containing segments where there is high similarities with your query*

**Per. Identity**: how similar the query is to the sequences

**Query cover:** how much of the query is covered by the sequence

# SpeciesFinder

Predicts species by using 16S rRNA genes- *a gene present in all bacteria species*



16S rRNA

**C**  **V1**  **C**  **V2**  **C**  **V3**  **C**  **V4**

Bacteria

**Conserved Region "C"** – **identical in all bacteria**

**Highly-variable Region "V"** – **conserved within species and divergent between species**

- The gene is predicted from the sequence, predicted sequence is aligned against the 16S database

- The species associated with the best hit is the final prediction.

# K-merFinder

- Predicts species by examining the number of overlapping 16-mers
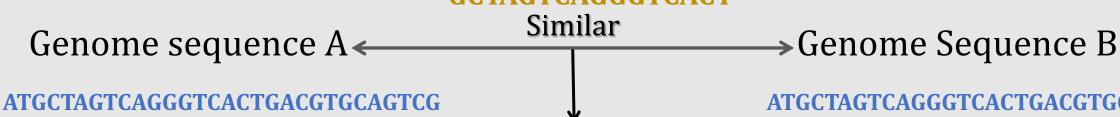
  $K$ is number of bases =16

- KmerFinder tool cuts the unknown species genome into 16mers, this is compared with 16mers of the known species in the database and gives the best matching species.

  Sequence    **ATGCTAGTCAGGGTCACTGACGTGCAGTCG**

  16-mers     **ATGCTAGTCAGGGTCA**

              **TGCTAGTCAGGGTCAC**

              **GCTAGTCAGGGTCACT**

Similar

Genome sequence A ← → Genome Sequence B

ATGCTAGTCAGGGTCACTGACGTGCAGTCG

**Unknown species**

ATGCTAGTCAGGGTCACTGACGTGCAGTCG

**Known species**

**Share the same K-mers**

# Center for Genomic Epidemiology

## Your job is being processed

Wait here to watch the progress of your job, or fill in the form below to get an email message upon completion.

To get notified by email: [                    ] [Notify me via email]

This page will update itself automatically.

# Public Health Alliance for Genomic Epidemiology

## Center for Genomic Epidemiology

| Home | Services | Instructions | Output |

## KmerFinder-3.2 Server - Results

**KmerFinder 3.2 results:**

| Template | Num | Score | Expected | Template_length | Query_Coverage | Template_Coverage | Depth | tot_query_Coverage |
|---|---|---|---|---|---|---|---|---|
| NZ_CP070071.1 Escherichia coli strain FDAARGOS_1287 chromosome, complete genome | 11051 | 156560 | 3 | 169133 | 92.56 | 93.85 | 0.93 | 92.56 |
| NZ_AP022362.1 Escherichia coli strain E302 chromosome, complete genome | 5996 | 2494 | 55 | 173594 | 1.47 | 1.39 | 0.01 | 60.75 |
| NZ_CP015088.1 Escherichia coli O25b:H4 extrachomosomal sequence | 14431 | 101 | 0 | 997 | 0.06 | 8.53 | 0.10 | 0.30 |

# Pathogenwatch

❖SpeciesFinder

❖KmerFinder

❖Pathogenwatch