# BLAST-based tools for the characterization of assembled genomes

## EURL CPS
## Marina Cavaiuolo

**Joint Training Course of the inter EURLs Working Group on NGS**

# Characterisation of bacterial genomes

**Search for genetic features**
*e.g.* Antimicrobial resistance genes
Virulence gene

**Typing**
*e.g.* MLST
Serotyping, virulotyping

**Principle of the analyses:**

Alignment of sequences:
*e.g.* whole genomes, genes
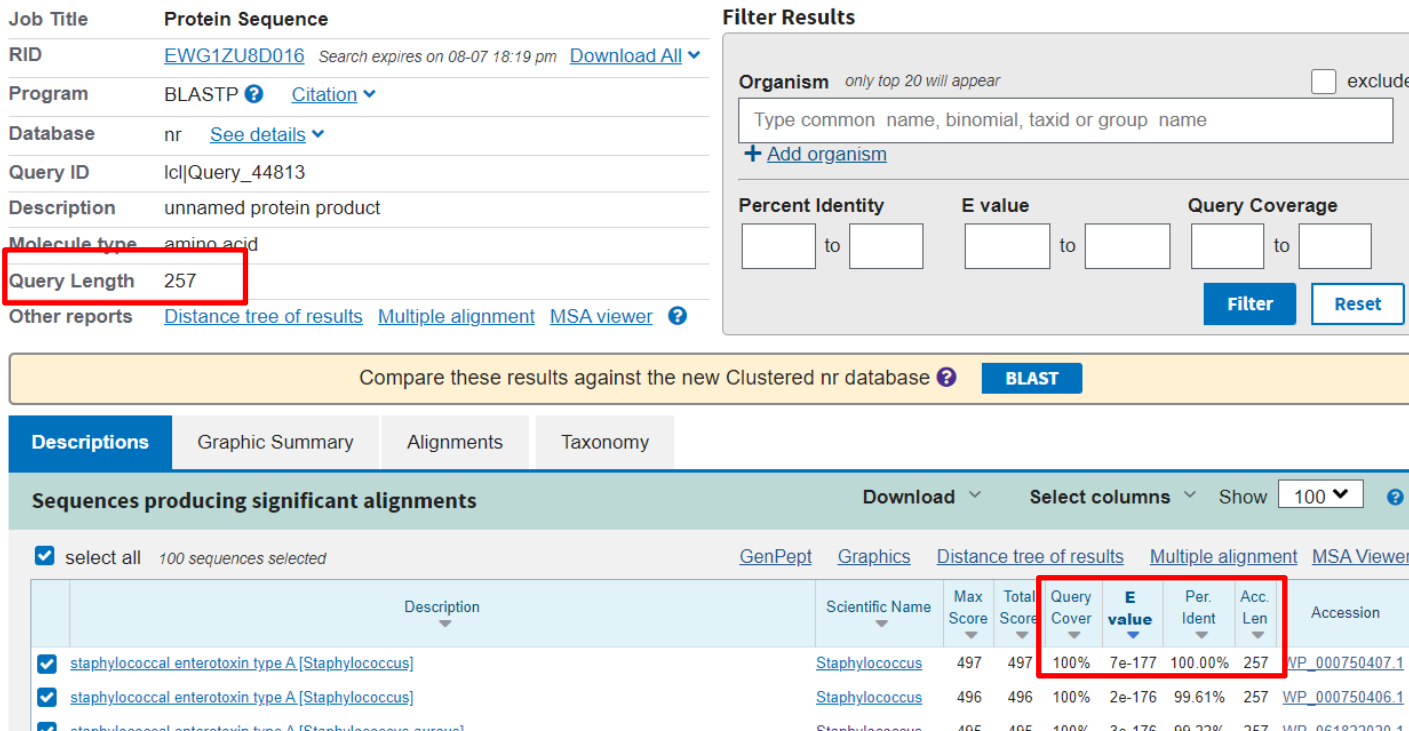
Mapping of NGS reads

**specific algorithms**

**database of reference gene sequences**

**BLAST** algorithm

# Basic Local Alignment Search Tool (BLAST)

BLAST is a sequence comparison <u>algorithm</u> to find homologous genes and proteins

**BLAST** takes a **query** (single or multiple sequences) in fasta format and **compares** it to a **subject** (single or multiple sequences) by constructing local alignments



**BLAST programs**

| Program | Database (Subject) | Query |
|---------|--------------------|-------|
| BLASTN | Nucleotide | Nucleotide |
| BLASTP | Protein | Protein |
| BLASTX | Protein | Nt. ➔ Protein |
| TBLASTN | Nt. ➔ Protein | Protein |
| TBLASTX | Nt. ➔ Protein | Nt. ➔ Protein |

# How to interpret the results of a BLAST search?

Blast output: Alignment. The result of matching up nucleotide/amino acids of two or more biological sequences to achieve maximal levels of identity and conservation (for proteins), for the purpose of assessing the degree of similarity and the possibility of homology.

For NCBI's web-page, the default format for output is HTML



## 1) How good is the match?

**e-value:** the number of expected hits of similar quality (score) that could be found just by chance

$E = 10^{-4}$ is considered the cutoff point

$E = 0$ means that the sequences are identical

## 2) How long is the alignement ?

**Coverage**: the % of the query length that aligns with the subject.

## 3) How similar are the aligned segments?

**Identity**: The extent to which two sequences have the same residues at the same positions in an alignment.

# How to use BLAST?

## 1) BLAST+ standalone suite



https://www.ncbi.nlm.nih.gov/books/NBK52637/

## 2) BLAST online



https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastn&PAGE_TYPE=BlastSearch&LINK_LOC=blasthome

## 3) NCBI BLAST on Galaxy



## 4) Many web-tools and software implement BLAST as search engine:
1. Upload of the assembled genomes
2. Blast analysis against a database

# Search for antimicrobial resistance genes using BLAST web-tools

## Center for Genomic Epidemiology - CGE



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

**Select the species**

**Select the type of input**

## Institut Hospitalier Universitaire Méditerranée Infection
## IHU – Méditerranée Infection



**Paste the query sequence**

https://ifr48.timone.univ-mrs.fr/blast/arg-annot_v6.html

No registration needed

# Search for antimicrobial resistance genes using BLAST web-tools

The Comprehensive Antibiotic Resistance Database: a rigorously curated collection of characterized, peer-reviewed resistance determinants and associated antibiotics, organized by the Antibiotic Resistance Ontology (ARO) and AMR gene detection models.



**Download the whole database**
**If you want to run Blast locally**

https://card.mcmaster.ca/

# Search for virulence factors using BLAST web-tools

The virulence factor database (VFDB) is an integrated and comprehensive online resource for curating information about virulence factors of bacterial pathogen

http://www.mgc.ac.cn/cgi-bin/VFs/genus.cgi?Genus=Staphylococcus

**Download the whole database if you want to run the blast locally**

**Paste query sequence**

No registration needed

# Search for virulence factors using BLAST web-tools

## CGE tools



Select species

Define coverage and % identity

Type of input

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

No registration needed

# Search for antimicrobial resistance genes and virulence factors

**ABRICATE on Galaxy**

# Thank you