

Bacterial WGS training : Exercise 6

Title Chromosome, plasmid, resistance and virulence annotation

Training dataset:

- How many genes there are in my sample?
- Are there virulence and/or antibiotic resistance genes?
- Where are the genes located?
- Which plasmids are present in the sample?
- How do I visualize the results?

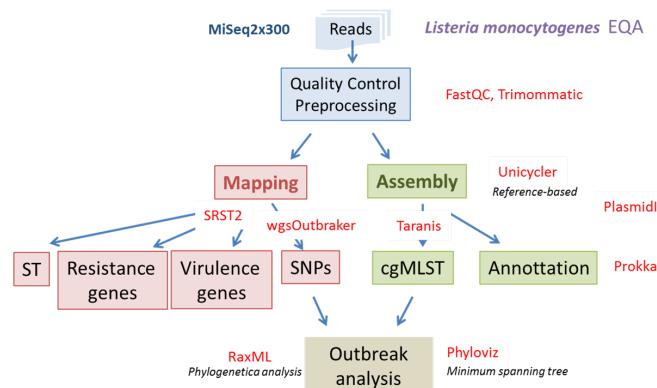
Questions:

- Annotate virulence and ABR genes
- Determine gene variants
- Determine plasmidome
- Locate annotated genes
- Results interpretation

Objectives:

Time estimation: 1 h

- Key points:**
- Comparing annotation using mapping vs assembly
 - Plasmid, virulence and resistance determination



- Bacterial WGS training : Exercise 6
 - Introduction
 - Training dataset description
 - Exercise
 - Mapping based annotation
 - Results should look like that
 - Assembly based annotation
 - Results should look like these

Introduction

In this exercise we are going to determine the genomic content of a multidrug-resistant (MDR) *K. pneumoniae* isolate. First we will use [srst2](#) to assess the resistome and later, we will use [plasmidID](#) to infer biological and positional information to sequences and see where the genes, detected with mapping strategy, are located.

Training dataset description

The sample we are going to analyse is an *in silico* dataset obtained with [wgsim](#) using a sample of [Klebsiella pneumoniae subsp. pneumoniae HS11286](#) available at ncbi.

Exercise

Mapping based annotation

To execute [srst2](#), which maps the reads against a antibiotic resistance genes database (ARGannot), lets execute this command:

```

cd
cd wgs/bacterial_wgs_training_dataset/ANALYSIS
nextflow run ../../bacterial_wgs_training/main.nf \
--reads '../REFERENCES/plasmidid_test/KPN*_R{1,2}.fastq.gz' \
--fasta ../REFERENCES/listeria_NC_021827.1_NoPhages.fna \
-profile conda \

```

```
--gtf ../REFERENCES/listeria_NC_021827.1_NoPhages.gff \
--srst2_resistance ../REFERENCES/ARGannot.r1.fasta \
--srst2_virulence ../REFERENCES/EcOH.fasta \
--step mapAnnotation \
--outdir 07-mapAnnotation \
-resume
```

Results should look like that

| Sample | DB | gene | allele | coverage | depth | diffs | uncertainty | divergence | length | maxMAF | clusterid |
|------------|-------------|-----------------|--------------|----------|--------|--------------|-------------|------------|--------|--------|-----------|
| KPN_TEST_R | ARGannot.r1 | RmtB_AGly | RmtB_1580 | 100.0 | 12.09 | 1snp | | 0.132 | 756 | 0.125 | 309 |
| KPN_TEST_R | ARGannot.r1 | TEM-1D_Bla | TEM-117_968 | 100.0 | 33.386 | 2snp | | 0.262 | 764 | 0.382 | 205 |
| KPN_TEST_R | ARGannot.r1 | KPC-1_Bla | KPC-14_809 | 100.0 | 5.412 | 1indel | | 0.0 | 876 | 0.333 | 184 |
| KPN_TEST_R | ARGannot.r1 | AmpH_Bla | AmpH_634 | 100.0 | 11.373 | 14snp | | 1.206 | 1161 | 0.143 | 86 |
| KPN_TEST_R | ARGannot.r1 | CTX-M-9_Bla | CTX-M-14_102 | 100.0 | 26.676 | 1snp | | 0.114 | 876 | 0.412 | 190 |
| KPN_TEST_R | ARGannot.r1 | StrA_AGly | StrA_1501 | 100.0 | 12.502 | 2snp | | 0.249 | 804 | 0.167 | 263 |
| KPN_TEST_R | ARGannot.r1 | StrB_AGly | StrB_1614 | 100.0 | 9.545 | 1snp | | 0.119 | 837 | 0.167 | 227 |
| KPN_TEST_R | ARGannot.r1 | AadA_AGly | AadA2_1605 | 100.0 | 9.306 | 2snp | | 0.256 | 780 | 0.167 | 229 |
| KPN_TEST_R | ARGannot.r1 | SHV-OKP-LEN_Bla | SHV-11_1287 | 100.0 | 9.401 | | | 0.0 | 861 | 0.143 | 164 |
| KPN_TEST_R | ARGannot.r1 | TetRG_Tet | TetRG_605 | 96.209 | 6.48 | 10snp24holes | edge0.0 | 1.642 | 633 | 0.5 | 373 |
| KPN_TEST_R | ARGannot.r1 | DfrA_Tmt | DfrA12_1089 | 99.799 | 8.389 | 1indel | | 0.0 | 498 | 0.143 | 418 |
| KPN_TEST_R | ARGannot.r1 | TetG_Tet | TetG_632 | 100.0 | 9.963 | | | 0.0 | 1176 | 0.25 | 80 |
| KPN_TEST_R | ARGannot.r1 | SullI_Sul | SullI_1219 | 100.0 | 11.094 | 1snp | | 0.123 | 816 | 0.2 | 256 |

This table is a full report of all the ARG found with all mapping stats.

Assembly based annotation

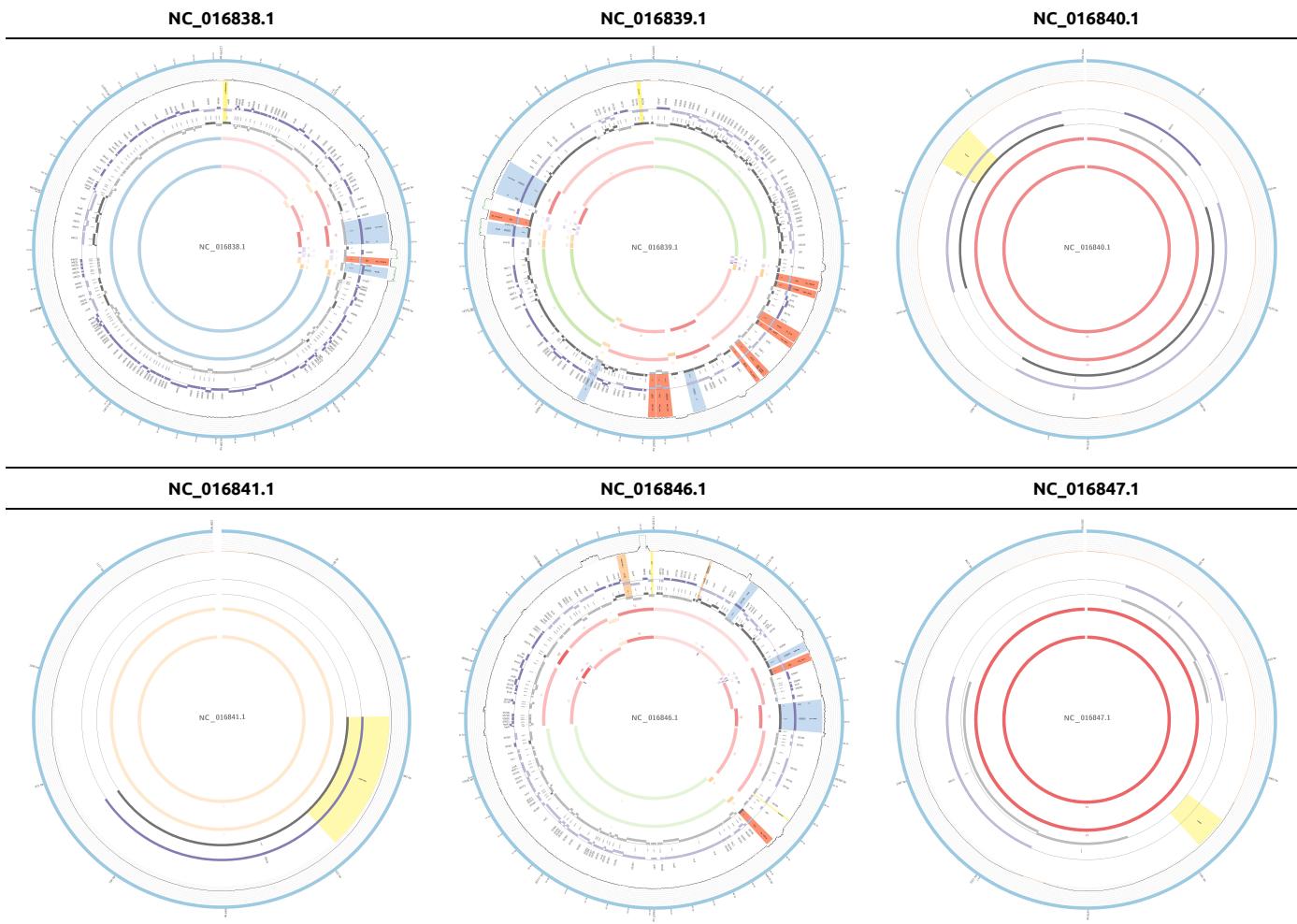
Now, using the contigs assembled using those same reads, we can determine the exact location of those ARG. ARG can be located on the chromosome but mostly on plasmids. In that case, we are going to focus on plasmid derived ARG using the annotation feature of plasmidID. To run the analysis lets use this command:

```
cd
cd wgs/bacterial_wgs_training_dataset/REFERENCES/
cp -r /mnt/ngs_course_shared/bacterial_wgs_training_dataset/REFERENCES/plasmidid_test .
```

Now we can run the nextflow

```
cd ../ANALYSIS/
nextflow run ../../bacterial_wgs_training/main.nf \
--reads '../REFERENCES/plasmidid_test/KPN*_R{1,2}.fastq.gz' \
--fasta ../REFERENCES/listeria_NC_021827.1_NoPhages.fna \
-profile conda \
--gtf ../REFERENCES/listeria_NC_021827.1_NoPhages.gff \
--plasmidid_database ../REFERENCES/plasmidid_test/plasmids_TEST_database.fasta \
--plasmidid_config ../REFERENCES/plasmidid_test/plasmidid_config.txt \
--step plasmidID \
--outdir 08-plasmidID \
-resume
```

Results should look like these



Those are the 6 plasmids that this isolate had, have a look at those pictures and find out if the genes are the same allele.

Are all the genes located with srst2 bound to plasmids?