

Assembly, assembly stats, virulotyping, serotyping

Valeria Michelacci

Bioinformatics training,
July 2019



Istituto Superiore di Sanità, Dep. of Food Safety, Nutrition and Veterinary Public Health
European Union and National Reference Laboratory for *E. coli*, Rome, Italy



Assembly (e.g.: SPAdes)

Short sequencing reads

.fastq file

```
@HWI-ST700693:238:B0224ACXX:1:1101:1218:1982
NACACTTGCCTTTGGTGACAGCGGGGCATCCTCAAGC
+
#1=DDDDHHAFF?GEFGIIIIIIIIIIIIIIIIIFI
@HWI-ST700693:238:B0224ACXX:1:1101:1161:1986
NGATTTTGACCTCTCCAGTTTCCTCTTAACACTTTG
+
#1=BDFFFHGHGJJJJIJHJJJJJJJJJJJJJJ
@HWI-ST700693:238:B0224ACXX:1:1101:1193:1989
NTATCCAGCCTGCGGTGCTACTTGGTGAAGAGGAT
+
#1=DDFFFHGHGJJJFGHJJJJJIEGECDFHCC?
@HWI-ST700693:238:B0224ACXX:1:1101:1440:1981
NTCAAGAATCCAAGTGGGGCCAGCATAATGTACGCT
+
#1=DDFFFHGHDFDAEGIIIFGIIICGGHGBFGEFDHI
@HWI-ST700693:238:B0224ACXX:1:1101:1367:1983
NATTAGAACAGATCGCTACTTCCGCCGAAGATACAT
+
#4BDFFFHHHHHJGIIJJJJJJJJJJJJJJJJ
@HWI-ST700693:238:B0224ACXX:1:1101:1395:1988
NTGGAACGTTTTTAAACGCGGAGACAGCGTGGAGT
+
#1=DDFFFHCFHJJJJJJJJJJJJJJJJJJJJJJ
@HWI-ST700693:238:B0224ACXX:1:1101:1285:1994
NCTTTGCTGTATTGACCGTTTGTAGATTTGAATCTT
+
#4=DDFFFHBBBBHHHIGIJFHIJFGGGIGIHIJJII
@HWI-ST700693:238:B0224ACXX:1:1101:1632:1989
NTCTATGAATGTTCAAGCGGTAGCTGAGGAGAGTCC
+
```



Partially assembled genome (contigs)

.fasta file

```
>NODE 1 length 449 cov 4.835189
ATCTTTTCGCGCCTTCCAGCTCCAGCCATTCCGGAACCGTTCGCGAGAAAACGGGCGTAAATC
GGGTAAAGACATAGCGCGGTTTGTACGGCGCATGACCTTCAAACATATCGCAGATTACACC
TTTCATCCAGCGCGCGCGGGGCTTCGCGAGGAAGCTGTGGTAAAGCAGATTGTTTTCTGC
TTCCAGTGCCAGAAAATGGCGCTTCTGCTCCGGGCTAAGCACTGGGCTGGTACAATTTG
CTGGCAACGTTGTTGCAGTGCATTTTCATGAGAAGTGGGCATCTTCTTTCTTTTATGC
CGAAGGTGATGCGCCATTGTAAGAAGTTTCGTGATGTTCACTTTGATCCTGATGCGTTTG
CCACCCTGACGCATTCAATTTGAAAGTGAATTTTGAACCAGATCGCATTACAGTGATG
CAAACCTGTAAGTAGATTTCTTAATTGTGATGTGATCGAAGTGTGTTGCGG
>NODE 2 length 309 cov 4.686084
ACTGGTCAGTGCGGTATCCTTGACAAATGGCCGATTGGACGTCTGGCGGATAAGTTTGG
TCGACTGCTGGTGTGCGTGTTCAGGTCTTGTGCGTATTCTCGGCAGTATCGCGATGCT
TAGCCAGCGCGGATGCCCCAGCGTTATTCATCCTCGGTGCCGCTGGCTTACGCTATA
TCCGGTGGCGATGGCATGGGCTTGCAGAAAAGTTGAACATCATCAACTGGTGGCGATGAA
CCAGGCCCTTACTGTTGAGCTACTGTGGGAAGTCTGCTTGGCCCGTCATTTACCGCTAT
GCTAATGCAGAAATTTCTCCGATAATTTATTGTT
>NODE 3 length 101 cov 3.346535
AGCGCATGAGCGCGCAGCGCGCGGTTACGTGGTGCATCAGCATGATGTTGGCCGGAGAG
TACAGAGACTCCCCTTCATCCATGATGCCCTTTTACCAGCAGTTCTTCAATCATCACC
AGACC
>NODE 4 length 311 cov 3.610933
CATCAACGCTAAAAGCCAAGATGACGCAGACCGCAAGCTTCCGGTCCGCTGGGTGTTCCG
GCGGGAACGGAAATGAGAAAAGCTCAATCACATATTGCCCATTAAGCGCCAAATCCCCTT
TCCATGAGTCCGCGGCTTCGCGATAGACTTCGCTTTCGACGCGTAAAACCAAGAAATCGC
AGTAGAAAAGCTTGTCCAGGCATAATCCGTGCATATCGCAATATGGTGAACCTGTT
TTAAACCCAGCATAAAGTCTCCTTTATTTGTTAACAGCACGTTACTCGCCCGAAGCCG
TCTGGCAAGTTATCCCGCATTTTGGAGTCTGTA
>NODE 5 length 186 cov 4.973118
CGAAGATATAAGAAAAGCGAACCAGAAAAGAATGCCGGAGAATTCATCAATTCATCACCTG
CATTGAGCAGATTTGCAAGTCTCAATAACCGGTAAATCCAGCCCCAACGTTGGTGTGAT
AGAGGAATTTACGCCCGGATTTTCCGCCGATAAACGCAACTGATGGTAGTAAATCCATCG
ACGAGGTGTTGGCCTTTTGTTCGCGCTGA
```

FastqSize \approx GenomeSize x Coverage x 2

At least 0.5 GB per genome

FastaSize for *E. coli* contigs

~5.5 MB



Filter SPAdes repeats/1

Input: contigs file & file with contigs stats

Coverage cut-off ratio:

This is the average coverage ratio cutoff. For example: if the average coverage is 100 and a coverage cut-off ratio of 0.5 is used, then any contigs with coverage lower than 50 will be eliminated.

Repeat cut-off ratio:

This is the coverage ratio cutoff to determine repeats in contigs. For example: if the average coverage is 100 and a repeat cut-off ratio of 1.75 is used, then any contigs with coverage more than or equal to 175 will be marked as repeats.

Length cut-off

Length for average coverage calculation (default = 5000)



Filter SPAdes repeats/2

What does it do?

Using the output of SPAdes (a fasta and a stats file, either from contigs or scaffolds), it filters the fasta files, discarding all sequences that are under a given length or under a calculated coverage. Repeated contigs are detected based on coverage.

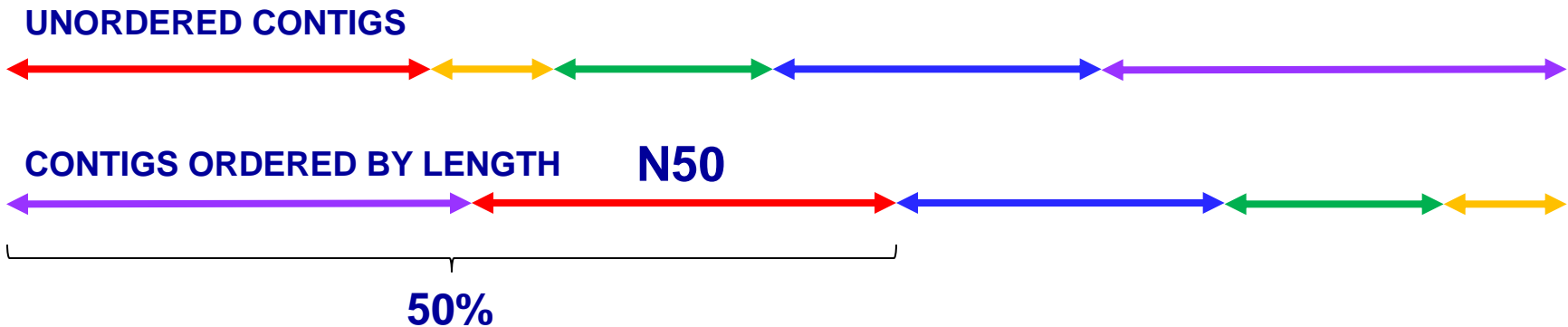
Output

- **Filtered sequences (with repeats)**
 - Will contain the filtered contigs/scaffolds including the repeats. These are the sequences that passed the length and minimum coverage cutoffs.
 - For workflows, this output is named **output_with_repeats**
- **Filtered sequences (no repeats)**
 - Will contain the filtered contigs/scaffolds excluding the repeats. These are the sequences that passed the length, minimum coverage and repeat cutoffs.
 - For workflows, this output is named **output_without_repeats**
- **Repeat sequences**
 - Will contain the repeated contigs/scaffolds only. These are the sequences that were excluded for having high coverage (determined by the repeat cutoff).
 - For workflows, this output is named **repeat_sequences_only**
- **Discarded sequences**
 - If selected, will contain the discarded sequences. These are the sequences that fell below the length and minimum coverage cutoffs, and got discarded.
 - For workflows, this output is named **discarded_sequences**
- **Results summary** : If selected, will contain a summary of all the results.

Assembly stats

N50

the **length** of the smallest contig among the set of the largest contigs that together cover at least 50% of the assembly



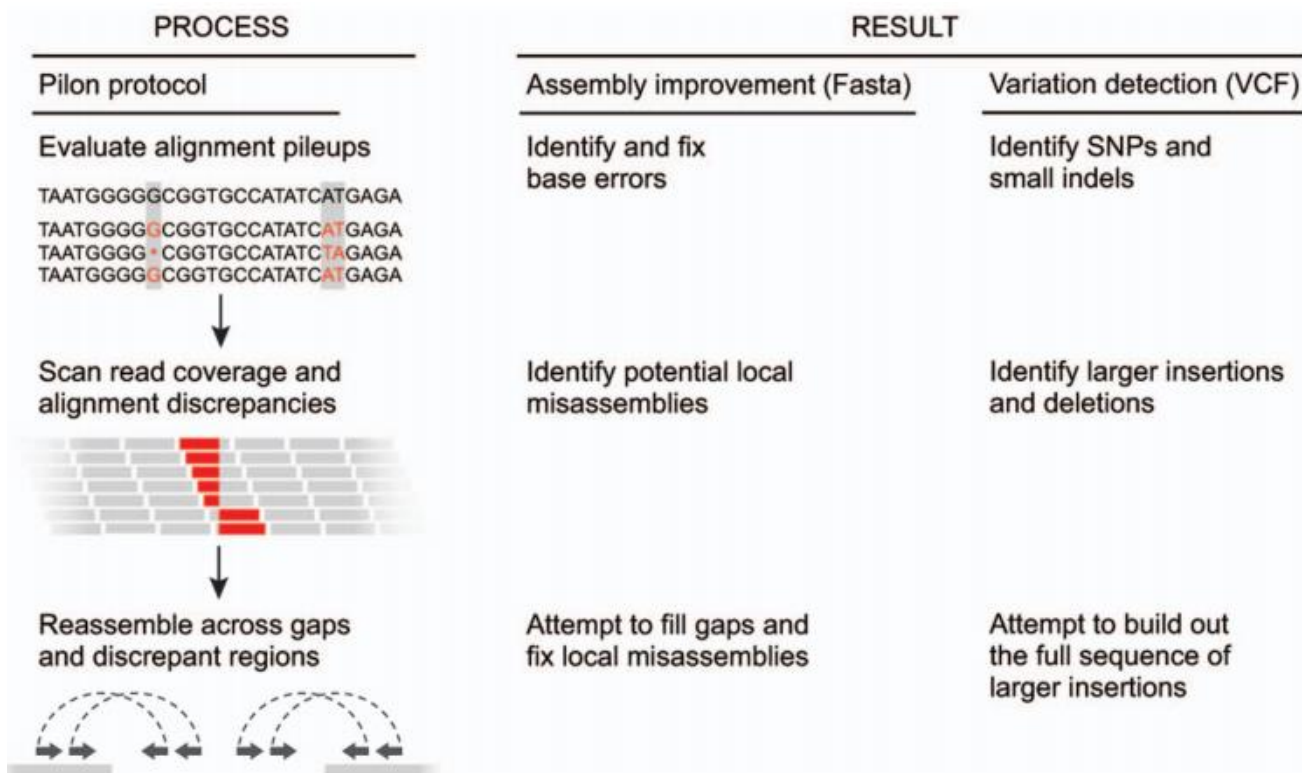
Other intuitive parameters to check:

Maximum contig length

Coverage of the contigs

Consensus length

Pilon – contigs refinement



Realignment of the reads on a «reference sequence»:

we use Bowtie2 as alignment tool and the contigs as ref seq

Pilon uses the result of the alignment to improve the assembly:

it outputs better assembled contigs

Assembly stats: check bacterial contigs

Contigs Evaluator v1.0 on file dataset_126093.dat

Estimated genome size: 5000000 bp

Assembled nucleotides: 5754440 bp

Estimated coverage: 1.15 x

N. contigs: 818

Average contig length: 7035

Median contig length: 457

Maximum contig length: 165129

N. contigs \geq 200 bp: 572 (69.9 %)

N. contigs \geq 2,000 bp: 204 (24.9 %)

N50: 58429

NG50: 72678



Assembly stats: Quast

Plots: [Cumulative length](#) [Nx](#) [NGx](#) [GC content](#)

