

# Where to start data analysis: quality check, assembly and alignment

Valeria Michelacci

NGS course, June 2015

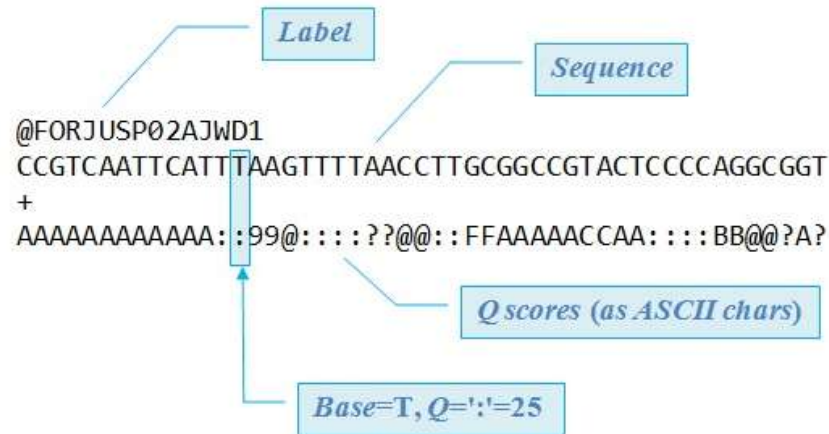


Istituto Superiore di Sanità, Dip. Sanità Pubblica Veterinaria e Sicurezza Alimentare,  
Laboratorio Europeo e Nazionale di Riferimento per *E. coli*



# Quality check

Output of NGS  
sequencers



Input for  
quality check

.fastq file

Sequencing errors would impact every following application

Unreliability of following results (and difficulty to detect the existence of problems!)

# Parameters to control

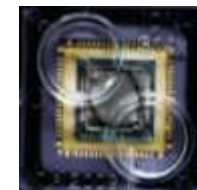
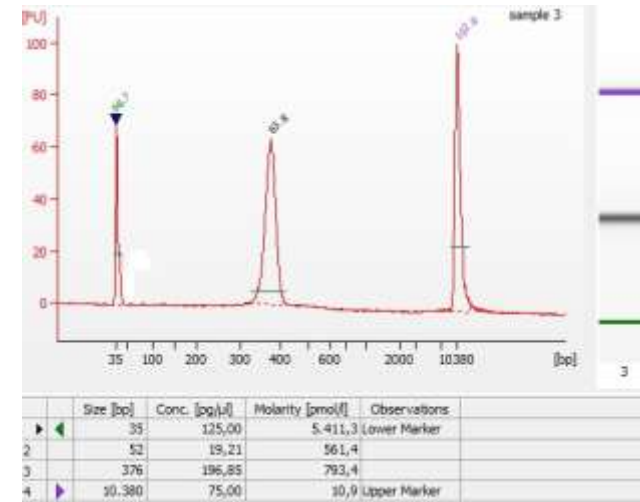
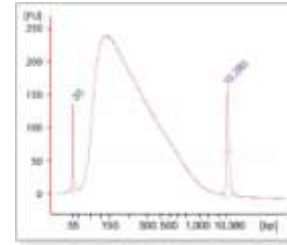
---

- Phred score
- GC content distribution over all sequences
- Distribution of undetermined bases (N)
- Sequence Duplication Levels
- ★ • Length of the reads
- ★ • Coverage

**Adoption of corrective actions is possible to minimize some of these problems**

# Corrective actions

- Check the **quality** and **concentration** of template DNA
- Adjust the **shearing** time and **size selection** condition to optimize the reads length
- Adjust the **concentration of the library** to obtain maximum number of reads (decreasing polyclonality)
- Optimize the **number of samples** to run in each type of run to obtain the desired coverage



**314v2**  
**30–100Mb**



**316v2**  
**300Mb–1Gb**



**318v2**  
**600Mb–2Gb**

# Mapping

Alignment of the sequencing reads on a reference sequence

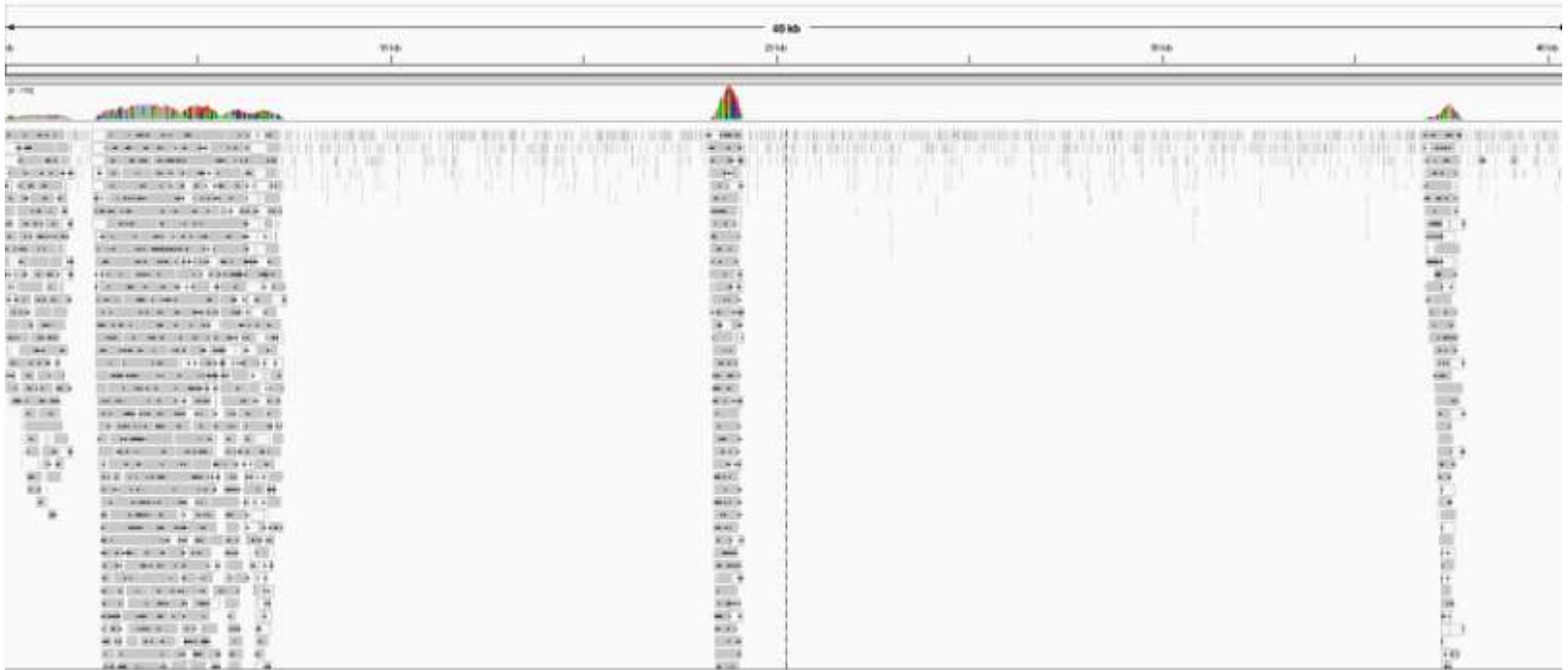


Ref seq

Possibility to directly inspect the **presence/absence of a target sequence** and the presence of **SNPs at interesting positions**

# Mapping: an example

Mapping of the reads deriving from the WGS of a VTEC strain on the reference sequence of a bacteriophage (e.g. BP933W harbouring *vtx2* in EDL933)



If no reads are mapping only on a few conserved spots (e.g. integrase and shiga toxins coding genes), no similar phage is present in the test strain

With this approach we can check the presence of features, but we can't determine the complete sequence of the corresponding bacteriophage

# Assembly

## Short sequencing reads

**.fastq file**

```
@HWI-ST700693:238:B0224ACXX:1:1101:1218:1982
NACACTTGCTTTGGTGACAGCGGGGCATCCTCAAGC
+
#1=DDDDHHAFF?GEFGIIIIIIIIIIIIIIIIIFI
@HWI-ST700693:238:B0224ACXX:1:1101:1161:1986
NGATTTTGACCTCTCCAGTTTCCTCTTAACACTTTC
+
#1=BDFFFHGHGJJJJIJHIJJJJJJJJJJJJJJ
@HWI-ST700693:238:B0224ACXX:1:1101:1193:1989
NTATCCAGCCTGCGGTGCTACTTGGTGGAAGAGGAT
+
#1=DDFFFHGHGJJJFGHJJJJJJIEGECDFHCC?
@HWI-ST700693:238:B0224ACXX:1:1101:1440:1981
NTCAAGAATCCAAGTGGGGCCAGCATAATGTACGCT
+
#1=DDFFFHGHDFDAEGIIFGIICGGHGBFGEFDHI
@HWI-ST700693:238:B0224ACXX:1:1101:1367:1983
NATTAGAACAGATCGCTACTTCGCCGAAGATACAT
+
#4BDFFFHGHGJJJJIJJJJJJJJJJJJJJJJJJ
@HWI-ST700693:238:B0224ACXX:1:1101:1395:1988
NTGGAACGTTTTTAAACGCGGAGACAGCGTGGAGT
+
#1=DDFFFHCFHJJJJJJJJJJJJJJJJJJJJJJ
@HWI-ST700693:238:B0224ACXX:1:1101:1285:1994
NCTTTGCTGTATGACCGTTTGTAGATTTGAATCCT
+
#4=DDFFFHBBBBHHHIGIJFHIJFGGGIGIHIJJII
@HWI-ST700693:238:B0224ACXX:1:1101:1632:1989
NTCTATGAATGTTCAAGCGGTAGCTGAGGAGAGTCC
+
```



## Partially assembled genome (contigs)

**.fasta file**

```
>NODE 1 length 449 cov 4.835189
ATCTTTTCGCGCCTTCCAGCTCCAGCCATTCCGGAACCGTTCCGCAGAAAACGGGGCGTAAATC
GGGTAAAGACATAGCGCGGTTTGTACGGCGCATGACCTTCAAACATATCGCAGATTACACC
TTCATCCAGCGCGCGGGGCTTCCGCAGGAAGCTGTGGTAAAGGCAGATTGTTTTCTGC
TTCAGTCCAGAAAATGGCGCTTCTGCTCCGGCTAAGCACTGGGCTGGTGACAATTTG
CTGGCAACGTTGTTGCAGTGCATTTTATGAGAAGTGGGCATCTTCTTTCTTTTATGC
CGAAGGTGATGCGCCATTGTAAAGAGTTTTCGTGATGTTCACTTTGATCCTGATGCGTTTG
CCACCCTGACGCATTATTTGAAAGTGAATTTTGAACCAGATCGCATTACAGTGATG
CAAACCTGTAAAGTAGATTTCTTAATTGTGATGTGATCGAAGTGTGTTGCGG
>NODE 2 length 309 cov 4.686084
ACTGGTCAGTGCGGTATCCTTGACAAATGGCCGATTGGACGTCTGGCGGATAAGTTTGG
TCGACTGCTGGTGTGCGTGTTCAGGTCTTTGTGTCATTCTCGGCAGTATCGCGATGCT
TAGCCAGCGCGGATGCCCCAGCGTTATTCATCCTCGGTGCCGCTTTCGCTATA
TCCGGTGGCGATGGCATGGGCTTGCAGAAAAGTTGAACATCATCAACTGGTGGCGATGAA
CCAGGCCCTTACTGTTGAGCTACTGTGGGAAGTCTGCTTGGCCCGTCATTTACCGCTAT
GCTAATGCAGAAATTTCTCCGATAATTTATTGTT
>NODE 3 length 101 cov 3.346535
AGCGCATGAGCGCGCAGCGCGCGTTCAGTGGTGCATCAGCATGATGTTGGCCGGAGAG
TACAGAGACTCCCCTTCATCCATGATGCCCTTTTACCAGCAGTTCTTCAATCATCACC
AGACC
>NODE 4 length 311 cov 3.610933
CATCAACGCTAAAAGCCAAGTACGCGAGACCGCAAGCTTCCGGTCCGCTGGGTGTTCCG
GCGGAAACGGAAATGAGAAAAGCTCAATCACATATTGCCATTAAAGCGCAAAATCCCCTT
TCCATGAGTCCGCGGCTTCCGATAGACTTCGCTTTCGACGCGTAAACCAAGAAATATCGC
AGTAGAAAAGCTTGTCCAGCATATCCGTGCATATCGCAATATCGCAATATGGTGAACCTGTT
TTAAACCCAGCATAAAGTCTCCTTTATTTGTTAAACAGCACGTTACTCGCCCGAAGCCG
TCTGGCAAGTTATCCCGCATTTTGGAGTCTGTA
>NODE 5 length 186 cov 4.973118
CGAAGATATAAGAAAAGCGAACCAGAAAAGAAATGCCGGAGAACTTCAATCAATTCACCTG
CATTGAGCAGATTTGCAAGTCTCAATAACCGGTAAATCCAGCCCCAACGTTGGTGTGAT
AGAGGAATTTACGCCCGGATTTTCCGCCGATAACGCAACTGATGGTAGTAAATCCATCG
ACGAGGTGTTGGCCTTTTGTTCGGCGTGA
```

FastqSize  $\approx$  GenomeSize x Coverage x 2

**At least 300 MB per genome**

FastaSize for *E. coli* contigs

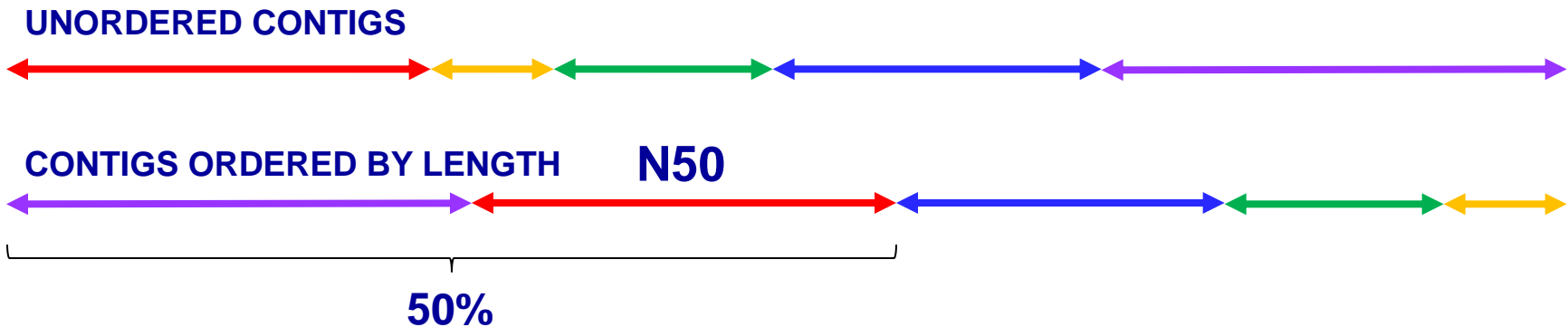
**~10 MB**



# 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



# Presence of genetic features

---

**Need for long contigs to investigate the presence of interesting genetic features by blast analysis**

- Virulence genes
- Serotype associated genes
- Genes part of the MLST scheme
- Presence of pathogenicity islands and investigation of insertion loci
- Any gene of interest