

# **NGI-Stockholm -- Science For Life Laboratory**

## **Best-practice analysis for quality checking report**

### **A.Blomberg\_16\_07 -- P4453\_101**

For sample P4453\_101 belonging to the project A.Blomberg\_16\_07 NGI-Stockholm best-practice analysis for quality checking has been performed. For mate pair libraries produced with Nextera, best-practice analysis described at this address has been performed:

[http://res.illumina.com/documents/products/technotes/technote\\_nextera\\_matepair\\_data\\_processing.pdf](http://res.illumina.com/documents/products/technotes/technote_nextera_matepair_data_processing.pdf)

The following tools have been employed (tools are listed in order of execution):

- i trimmomatic : /proj/a2010002/nobackup/sw/mf/bioinfo-tools/misc/trimmomatic/0.30/trimmomatic-0.30.jar
- ii fastqc : /sw/apps/bioinfo/fastqc/0.11.2/milou/fastqc
- iii abyss : /sw/apps/bioinfo/abyss/1.3.5/milou/bin/

The results from each tool is reported in the following sections. Moreover you will find all the results and commands that have been run in the delivery folder on Uppmax

## Trimmomatic

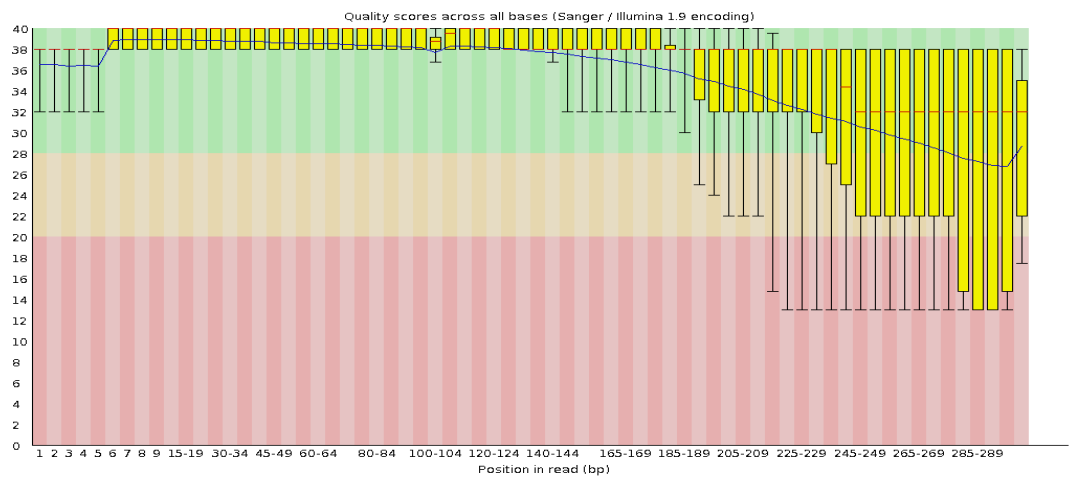
Reads (both paired and mate pairs) can contain parts of the adapter sequence or, in the case of mate pairs, part of the linker sequence. Illumina recommends to remove the adapter before use of the reads in any downstream analysis (this is mandatory for mate pairs).

Adapter sequences removed are:

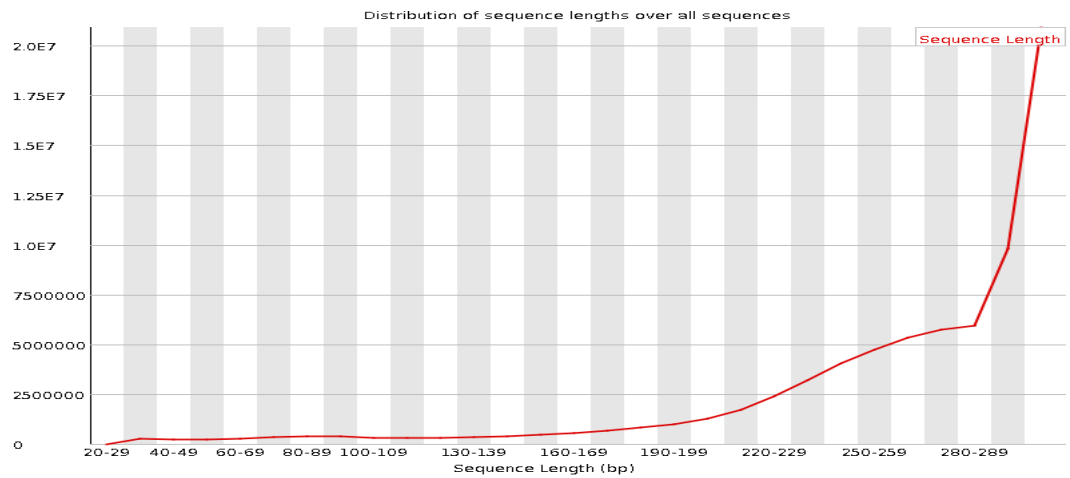
```
i  TACTACTCTTTCCCTACACGACGCTCTTCCGATCT
ii GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT
iii AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT
iv  CAAGCAGAAGACGGCATACGAGATCGGTCTCGGCATTCTGCTGAACCGCTCTTCCGATCT
v   AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT
```

P4453_101	#orig_pairs	#survived_pairs
P4453_101_S1_L002_R1_001_trimmomatic.stdErr	75846689	73276544 (97%)
P4453_101_S1_L001_R1_001_trimmomatic.stdErr	77042638	74309678 (96%)
total	152889327	147586222 (97%)

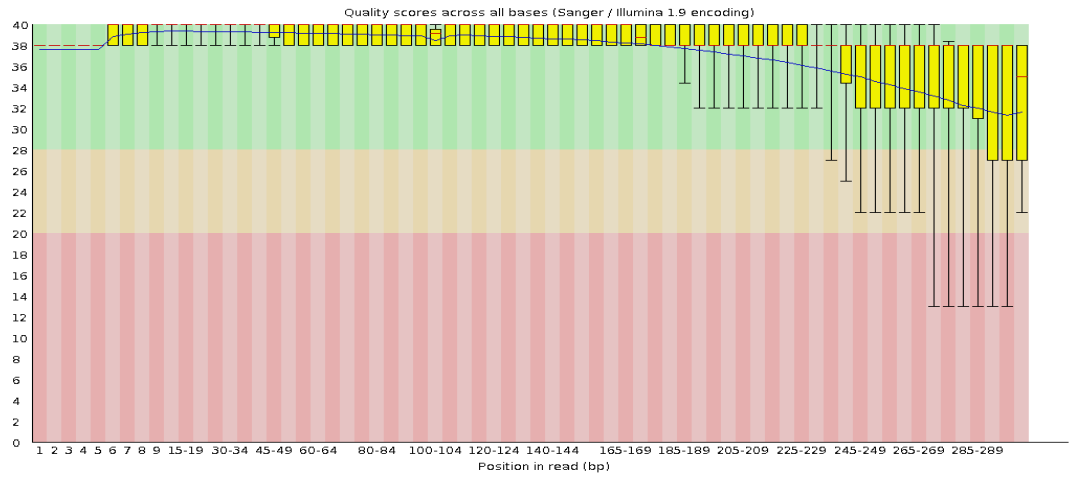
P4453_101	#survived_fw_only	#survived_rv_only	#discarded
P4453_101_S1_L002_R1_001_trimmomatic.stdErr	2278516 (3%)	229750 (0%)	61879 (0%)
P4453_101_S1_L001_R1_001_trimmomatic.stdErr	2435825 (3%)	234266 (0%)	62869 (0%)
total	2435825 (3%)	234266 (0%)	62869 (0%)



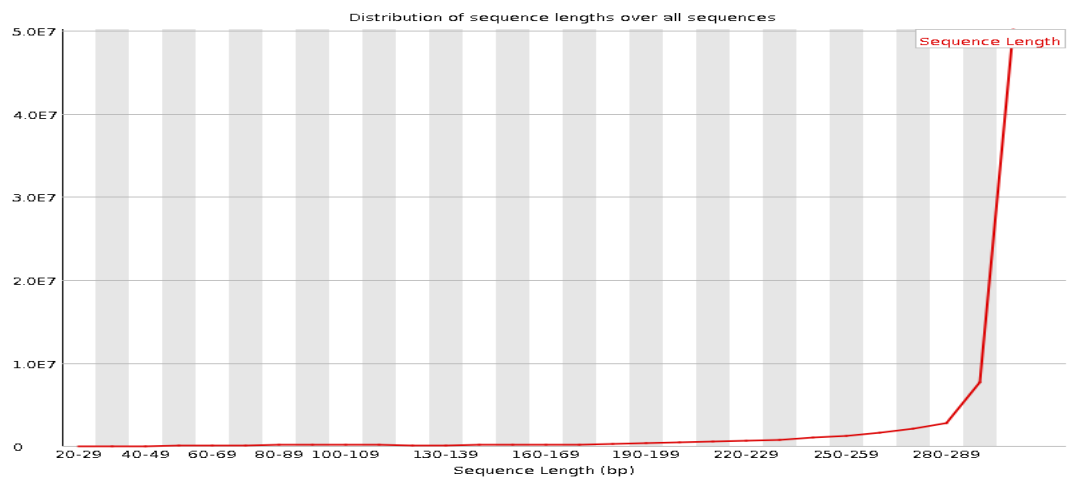
P4453\_101\_S1\_L002\_R2\_001\_fastqc -- Per Base Quality



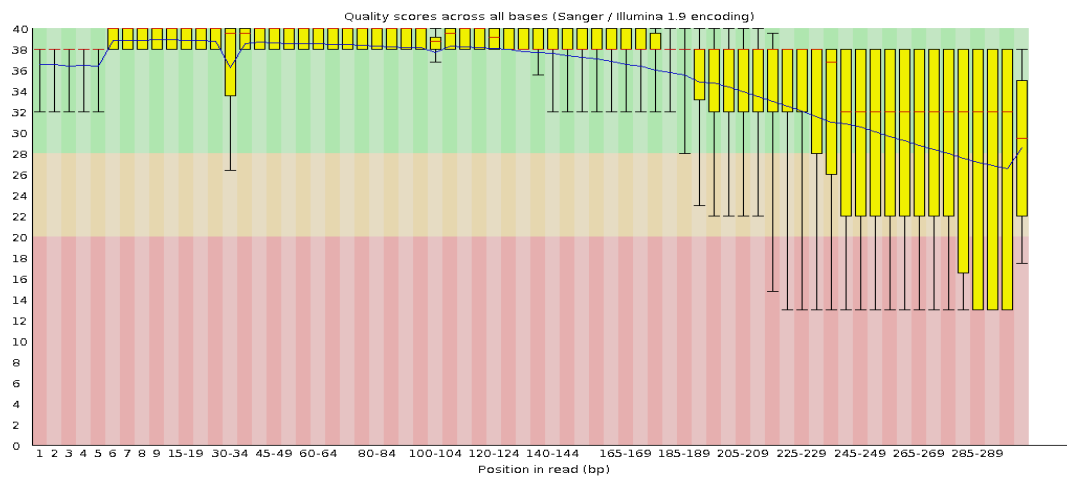
P4453\_101\_S1\_L002\_R2\_001\_fastqc -- Sequence Length Distribution



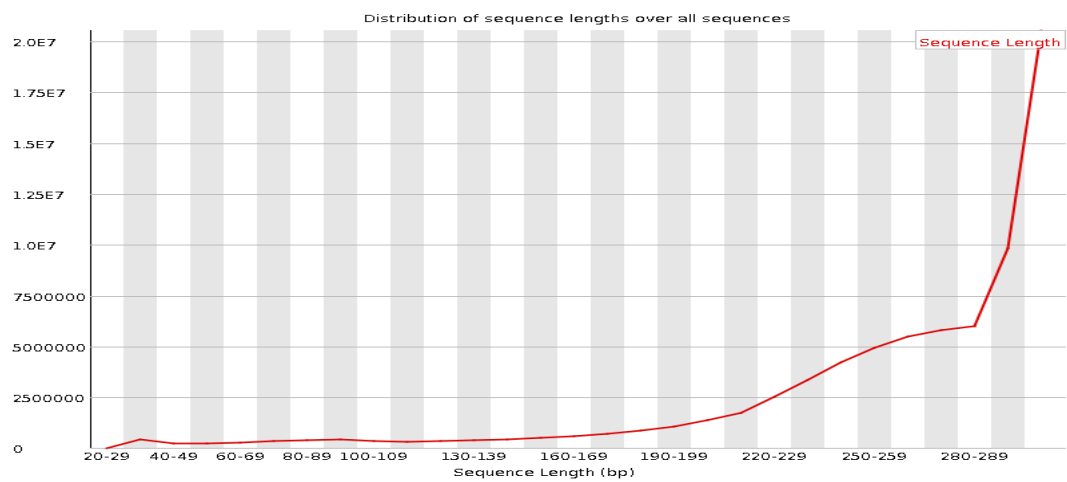
P4453\_101\_S1\_L002\_R1\_001\_fastqc -- Per Base Quality



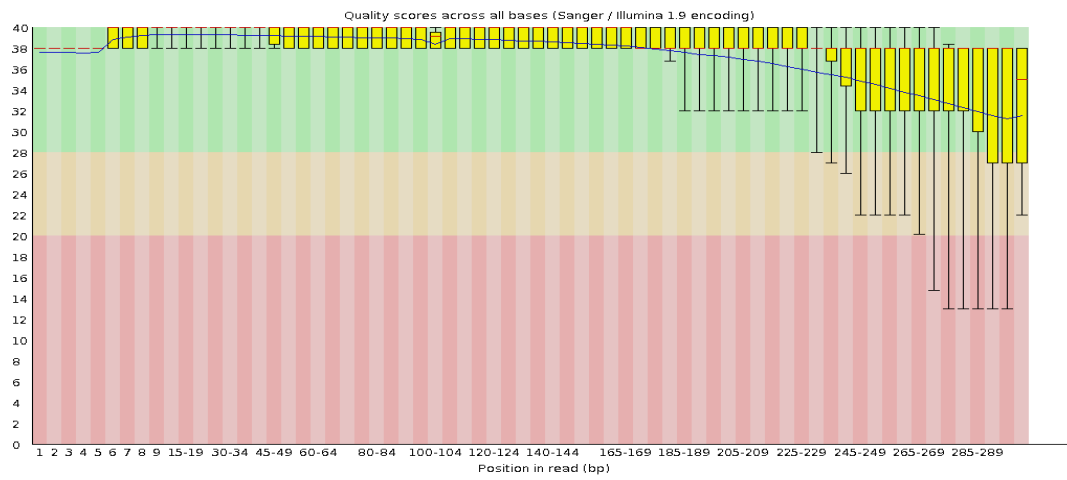
P4453\_101\_S1\_L002\_R1\_001\_fastqc -- Sequence Length Distribution



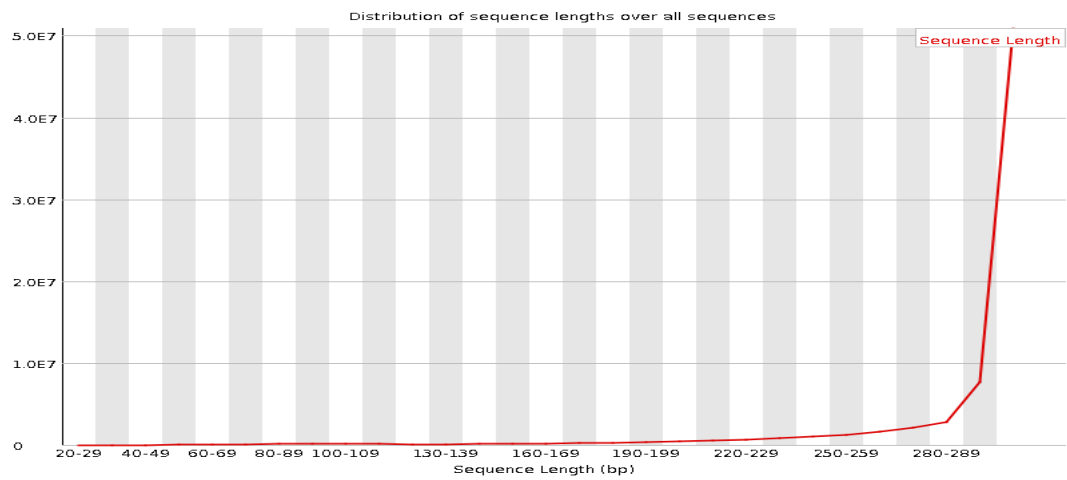
P4453\_101\_S1\_L001\_R2\_001\_fastqc -- Per Base Quality



P4453\_101\_S1\_L001\_R2\_001\_fastqc -- Sequence Length Distribution



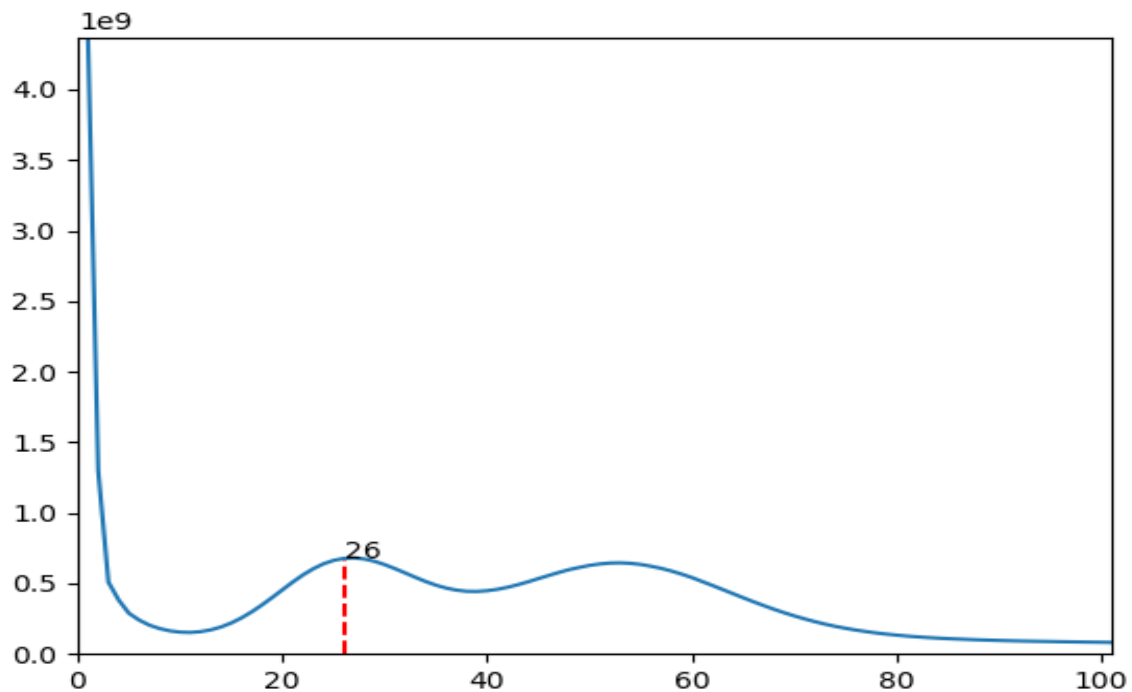
P4453\_101\_S1\_L001\_R1\_001\_fastqc -- Per Base Quality



P4453\_101\_S1\_L001\_R1\_001\_fastqc -- Sequence Length Distribution

## Abyss

A possible way to assess the complexity of a library even in absence of a reference sequence is to look at the kmer profile of the reads. The idea is to count all the kmers (i.e., sequence of length  $k$ ) that occur in the reads. In this way it is possible to know how many kmers occur 1,2,...,  $N$  times and represent this as a plot. This plot tell us for each  $x$ , how many  $k$ -mers (y-axis) are present in the dataset in exactly  $x$ -copies. In an ideal world (no errors in sequencing, no bias, no repeating regions) this plot should be as close as possible to a gaussian distribution. In reality we will always see a peak for  $x=1$  (i.e., the errors) and another peak close to the expected coverage. If the genome is highly heterozygous a second peak at half of the coverage can be expected.



kmer profile with  $k=35$ .