1. Data for analysis
>Download roommate's results (Illumina single-end sequencing run)
Filename:
http://ftp.sra.ebi.ac.uk/vol1/fastq/SRR170/001/SRR1705851/

>Make fasta file with reference sequence
Filename: reference.fasta
https://www.ncbi.nlm.nih.gov/nuccore/KF848938.1?report=fasta

>Download fastq data for the three controls (from sequencing of isogenic reference samples) from SRA FTP:
using
wget <link>
SRR1705858: ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR170/008/SRR1705858/SRR1705858.fastq.gz  saved as C58.fastq.gz
SRR1705859: ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR170/009/SRR1705859/SRR1705859.fastq.gz  saved as C59.fastq.gz
SRR1705860: ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR170/000/SRR1705860/SRR1705860.fastq.gz  saved as C60.fastq.gz


Check the length of referense FASTA
        > wc reference.fasta
        1794 smb  = about 550 AA


**Count the lines for initial  fastq.gz files**

**zcat {filename} | wc -l**

```
(bioinf_practice) maria@DESKTOP-VKN7MM4:~/bioinformatics_practice/IB_Practice_Project2/rawdata$ zcat C58.fastq.gz | wc -l
1026344
(bioinf_practice) maria@DESKTOP-VKN7MM4:~/bioinformatics_practice/IB_Practice_Project2/rawdata$ zcat C59.fastq.gz | wc -l
933308
(bioinf_practice) maria@DESKTOP-VKN7MM4:~/bioinformatics_practice/IB_Practice_Project2/rawdata$ zcat C60.fastq.gz | wc -l
999856
(bioinf_practice) maria@DESKTOP-VKN7MM4:~/bioinformatics_practice/IB_Practice_Project2/rawdata$ zcat roommate.fastq.gz | wc -l
1433060
```

## Check the reads quality for roommate and 3 control fastq files
fastqc -o . {file1} {file2} {file3} {file4}
fastqc -o . C58.fastq.gz C59.fastq.gz C60.fastq.gz roommate.fastq.gz


3 Control have 35-38 Mbp, satisfied 'per base' and 'per sequence' quality and look very similar
failed Checkpoints:
[FAIL]Per base sequence content
[FAIL]Per sequence GC content
[FAIL]Sequence Duplication Levels
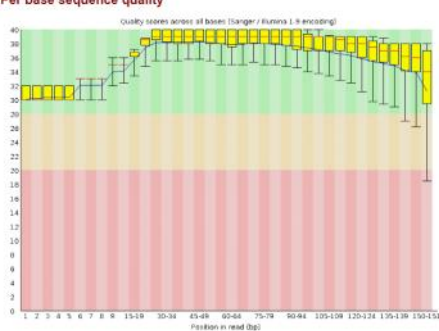[FAIL]Overrepresented sequences
C58 Warning checkpoint:
[WARNING]Sequence Length Distribution

## Control data

### Control C58 File

**Basic Statistics**

| Measure | Value |
|---|---|
| Filename | C58.fastq.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 250580 |
| Total Bases | 38.1 Mbp |
| Sequences flagged as poor quality | 0 |
| Sequence length | 35-151 |
| %GC | 41 |

**Per base sequence quality**

C58 failed Checkpoints:
[FAIL]Per base sequence content
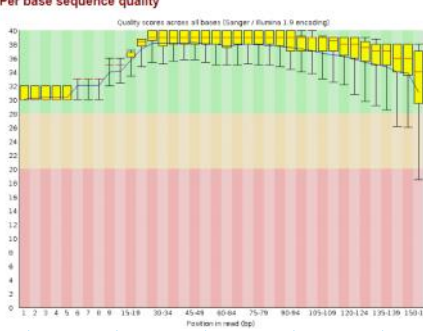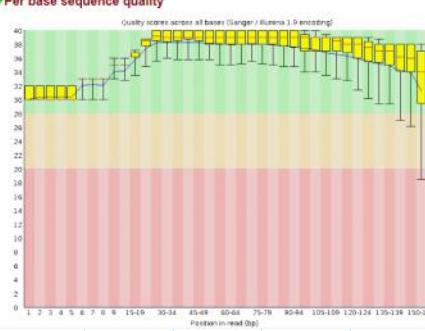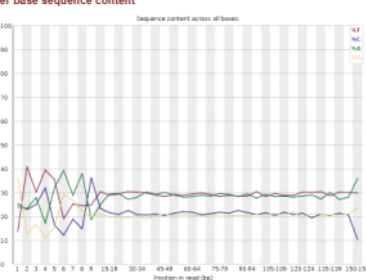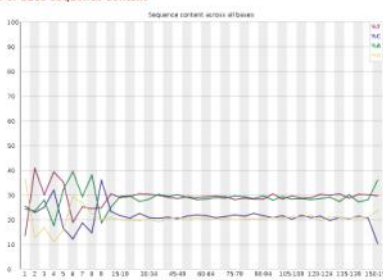[FAIL]Per sequence GC content
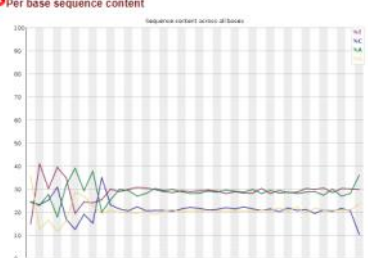[FAIL]Sequence Duplication Levels
[FAIL]Overrepresented sequences

**Per base sequence content**

**Per sequence GC content**

**Sequence Duplication Levels**

### Control C59 File

**Basic Statistics**

| Measure | Value |
|---|---|
| Filename | C59.fastq.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 233527 |
| Total Bases | 34.6 Mbp |
| Sequences flagged as poor quality | 0 |
| Sequence length | 35-151 |
| %GC | 41 |

**Per base sequence quality**

C58 failed Checkpoints:
[FAIL]Per base sequence content
[FAIL]Per sequence GC content
[FAIL]Sequence Duplication Levels
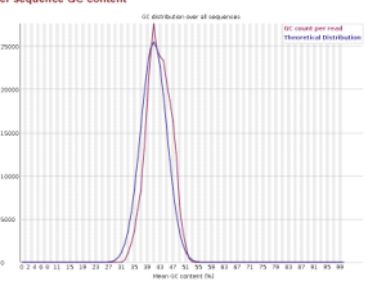[FAIL]Overrepresented sequences

**Per base sequence content**

**Per sequence GC content**

**Sequence Duplication Levels**

### Control C60 File

**Basic Statistics**

| Measure | Value |
|---|---|
| Filename | C60.fastq.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 245984 |
| Total Bases | 37.1 Mbp |
| Sequences flagged as poor quality | 0 |
| Sequence length | 35-151 |
| %GC | 41 |

**Per base sequence quality**

C58 failed Checkpoints:
[FAIL]Per base sequence content
[FAIL]Per sequence GC content
[FAIL]Sequence Duplication Levels
[FAIL]Overrepresented sequences

**Per base sequence content**

**Per sequence GC content**

**Sequence Duplication Levels**

## Overrepresented sequences (left)

Beginning of the list

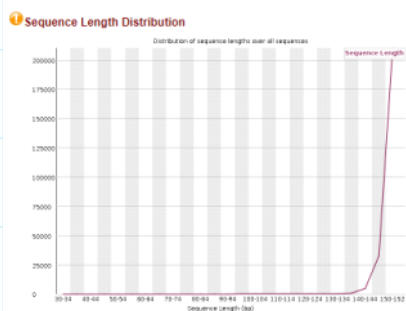## Overrepresented sequences (middle)

Beginning of the list

## Overrepresented sequences (right)
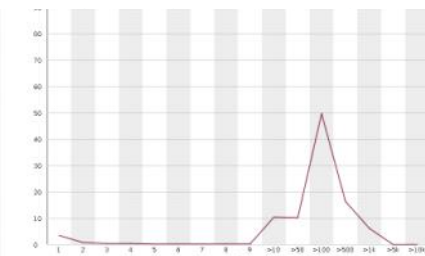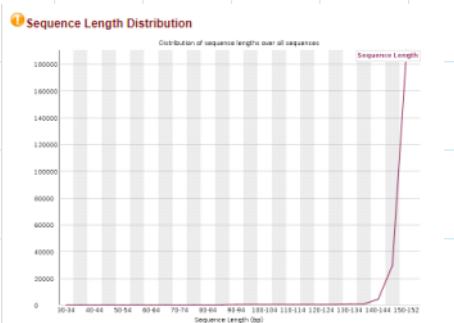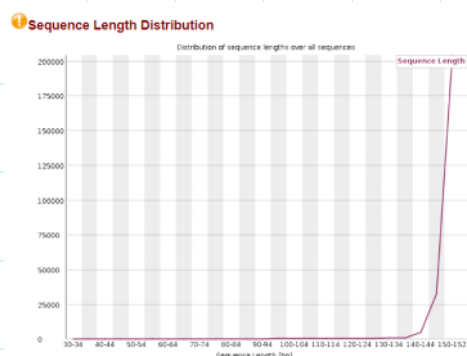
Beginning of the list

C58 Warning checkpoint:
[WARNING]Sequence Length Distribution

C58 Warning checkpoint:
[WARNING]Sequence Length Distribution

C58 Warning checkpoint:
[WARNING]Sequence Length Distribution

### Sequence Length Distribution







# Data from roommate

Data from roommate presents
358265 sequenses and 52.7 Mbp and different from control files quality characteristics. It performs less overrepresented sequences and more normal per sequence GC content

Failed Checkpoints:
[FAIL]Per base sequence content
[FAIL]Sequence Duplication Levels

Warning checkpoints:
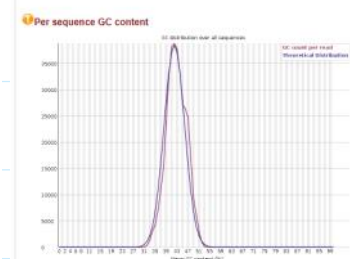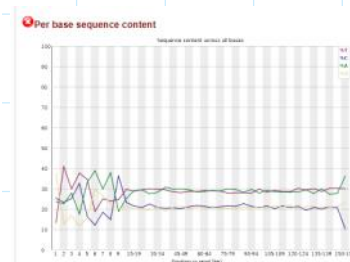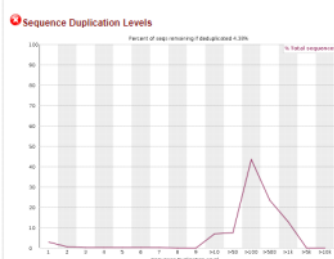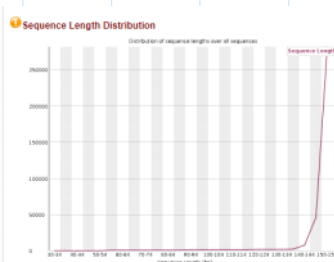[WARNING]Per sequence GC content
[WARNING]Sequence Length Distribution
[WARNING]Overrepresented sequences

### Basic Statistics

| Measure | Value |
| --- | --- |
| Filename | roommate.fastq.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 358265 |
| Total Bases | 52.7 Mbp |
| Sequences Flagged as poor quality | 0 |
| Sequence length | 35-151 |
| %GC | 42 |

### Per base sequence quality



### Sequence Length Distribution



### Per base sequence content



### Sequence Duplication Levels



### Per sequence GC content



### Overrepresented sequences

Beginning of the list

# Points for discussion

1 ноября 2023 г.
19:22
Virus strain:  A/Hong Kong/4801/2014 (H3N2)
//What about vaccination against this strain:
(Need toFind out what strains were in this season's vaccine. Was that one of the flu strains covered by this vaccine?)
// What about epitopes:
Munoz, Deem 2004
https://drive.google.com/file/d/1xe5-4LxIV4bO4mX6jhvrMAqtpOpkWsXm/view

Structure and receptor binding preferences of recombinant human A(H3N2) virus hemagglutinins

https://www.rcsb.org/structure/4WEA