# Genome Assembly

**STEPS:**

1. Navigate to the directory containing our data, using the following commands:
   **cd /home**
   **cd ngs_hands_on**

2. The conda environment containing the software needed for analysis is activated:
   **conda activate eco_evo_lab**

3. A directory is made to store the outputs we obtain from FastQC quality analysis:
   **mkdir sanchitha_fastqc_raw_output**

```
(base) x@AU-CSLAB24:~$ cdcd /home

Command 'cdcd' not found, but can be installed with:

sudo apt install cdcd

(base) x@AU-CSLAB24:~$ cd /home
(base) x@AU-CSLAB24:/home$ cd x
(base) x@AU-CSLAB24:~$ cd ngs_hands_on
(base) x@AU-CSLAB24:~/ngs_hands_on$ conda activate eco_evo_lab
(eco_evo_lab) x@AU-CSLAB24:~/ngs_hands_on$ ls
TruSeq3-PE.fa                  output_forward_paired_7m.fastq  output_reverse_paired_7m.fastq  quast-5.2.0.tar.gz
output_forward_paired.fastq    output_forward_unpaired.fastq  output_reverse_unpaired.fastq   sanchitha_fastqc_raw_output
output_forward_paired_13m.fastq output_reverse_paired.fastq    quast-5.2.0                     sub_SRR22507560_R1_25m.fastq
```

4. Then, we analyze the quality of the given 25 million reads of DNA using FastQC. Two sets of 25 million reads each - the forward and the reverse reads - are analyzed:
   **fastq sub_SRR22507560_R1_25m.fastq -o sanchitha_fastqc_raw_output**
   **fastq sub_SRR22507560_R2_25m.fastq -o sanchitha_fastqc_raw_output**

```
output_forward_paired_7m.fastq   output_reverse_unpaired.fastq   summarystats
output_forward_paired_7m2.fastq  quast-5.2.0                     trimmomatic-0.39.jar
(eco_evo_lab) x@AU-CSLAB24:~/ngs_hands_on$ mkdir sanchitha_fastqc_rawoutput
(eco_evo_lab) x@AU-CSLAB24:~/ngs_hands_on$ fastqc sub_SRR22507560_R1_25m.fastq -o sanchitha_fastqc_rawoutput
null
Started analysis of sub_SRR22507560_R1_25m.fastq
Approx 5% complete for sub_SRR22507560_R1_25m.fastq
Approx 10% complete for sub_SRR22507560_R1_25m.fastq
Approx 15% complete for sub_SRR22507560_R1_25m.fastq
Approx 20% complete for sub_SRR22507560_R1_25m.fastq
Approx 25% complete for sub_SRR22507560_R1_25m.fastq
Approx 30% complete for sub_SRR22507560_R1_25m.fastq
Approx 35% complete for sub_SRR22507560_R1_25m.fastq
Approx 40% complete for sub_SRR22507560_R1_25m.fastq
Approx 45% complete for sub_SRR22507560_R1_25m.fastq
Approx 50% complete for sub_SRR22507560_R1_25m.fastq
Approx 55% complete for sub_SRR22507560_R1_25m.fastq
Approx 60% complete for sub_SRR22507560_R1_25m.fastq
Approx 65% complete for sub_SRR22507560_R1_25m.fastq
Approx 70% complete for sub_SRR22507560_R1_25m.fastq
Approx 75% complete for sub_SRR22507560_R1_25m.fastq
Approx 80% complete for sub_SRR22507560_R1_25m.fastq
Approx 85% complete for sub_SRR22507560_R1_25m.fastq
Approx 90% complete for sub_SRR22507560_R1_25m.fastq
Approx 95% complete for sub_SRR22507560_R1_25m.fastq
Approx 100% complete for sub_SRR22507560_R1_25m.fastq
Analysis complete for sub_SRR22507560_R1_25m.fastq
```

```
(eco_evo_lab) x@AU-CSLAB24:~/ngs_hands_on$ fastqc sub_SRR22507560_R2_25m.fastq -o sanchitha_fastqc_rawoutput
null
Started analysis of sub_SRR22507560_R2_25m.fastq
Approx 5% complete for sub_SRR22507560_R2_25m.fastq
Approx 10% complete for sub_SRR22507560_R2_25m.fastq
Approx 15% complete for sub_SRR22507560_R2_25m.fastq
Approx 20% complete for sub_SRR22507560_R2_25m.fastq
Approx 25% complete for sub_SRR22507560_R2_25m.fastq
Approx 30% complete for sub_SRR22507560_R2_25m.fastq
Approx 35% complete for sub_SRR22507560_R2_25m.fastq
Approx 40% complete for sub_SRR22507560_R2_25m.fastq
Approx 45% complete for sub_SRR22507560_R2_25m.fastq
Approx 50% complete for sub_SRR22507560_R2_25m.fastq
Approx 55% complete for sub_SRR22507560_R2_25m.fastq
Approx 60% complete for sub_SRR22507560_R2_25m.fastq
Approx 65% complete for sub_SRR22507560_R2_25m.fastq
Approx 70% complete for sub_SRR22507560_R2_25m.fastq
Approx 75% complete for sub_SRR22507560_R2_25m.fastq
Approx 80% complete for sub_SRR22507560_R2_25m.fastq
Approx 85% complete for sub_SRR22507560_R2_25m.fastq
Approx 90% complete for sub_SRR22507560_R2_25m.fastq
Approx 95% complete for sub_SRR22507560_R2_25m.fastq
Approx 100% complete for sub_SRR22507560_R2_25m.fastq
Analysis complete for sub_SRR22507560_R2_25m.fastq
```

5. The output contains, among other files, a web page containing a thorough analysis of the basis statistics, per base sequence quality, per sequence, per sequence GC content and adaptor quantity.

6. Since it was observed that the adaptor sequences were quite high in the raw data, the reads are trimmed to remove the adapters so that high quality reads could be obtained and further used for assembly. Very low quality reads are also removed using Trimmomatic software, but most of them are conserved.

**java -jar trimmomatic-0.39.jar PE -summary summarystats input_forward.fastq input_reverse.fastq output_forward_paired.fastq output_forward_unpaired.fastq output_reverse_paired.fastq output_reverse_unpaired.fastq ILLUMINACLIP:TruSeq3-PE.fa:2:30:10 SLIDINGWINDOW:4:20 LEADING:3 TRAILING:3 MINLEN:36**

```
(eco_evo_lab) x@AU-CSLAB24:~/ngs_hands_on$ java -jar trimmomatic-0.39.jar PE -summary summarystats sub_SRR22507560_R1_25m.fastq sub_SRR22507560_R2_25m.fastq output_forward_paired2.fastq output_forward_unpaired
2.fastq output_reverse_paired2.fastq output_reverse_unpaired2.fastq ILLUMINACLIP:TruSeq3-PE.fa:2:30:10 SLIDINGWINDOW:4:20 LEADING:3 TRAILING:3  MINLEN:36
TrimmomaticPE: Started with arguments:
 -summary summarystats sub_SRR22507560_R1_25m.fastq sub_SRR22507560_R2_25m.fastq output_forward_paired2.fastq output_forward_unpaired2.fastq output_reverse_paired2.fastq output_reverse_unpaired2.fastq ILLUMINA
CLIP:TruSeq3-PE.fa:2:30:10 SLIDINGWINDOW:4:20 LEADING:3 TRAILING:3 MINLEN:36
Using PrefixPair: 'TACACTCTTTCCCTACACGACGCTCTTCCGATCT' and 'GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT'
ILLUMINACLIP: Using 1 prefix pairs, 0 forward/reverse sequences, 0 forward only sequences, 0 reverse only sequences
Quality encoding detected as phred33
Input Read Pairs: 25000000 Both Surviving: 19817046 (79.27%) Forward Only Surviving: 4226393 (16.91%) Reverse Only Surviving: 509717 (2.04%) Dropped: 446844 (1.79%)
TrimmomaticPE: Completed successfully
```

7. Then, we use FastQC to check if the quality of the trimmed reads is higher than that of the previous reads:

8. Next genome assembly using the trimmed reads is done of the forward and reverse sequences. Genome assembly using 25 million reads would take a long time and might be inefficient. We subsampled 500,000 reads out of the 25 million to perform the assembly.

   Subsampling:
   **seqtk sample -s100 output_forward_paired.fastq 500000 >
   output_forward_paired_halfm.fastq
   seqtk sample -s100 output_reverse_paired.fastq 500000 >
   output_reverse_paired_halfm.fastq**



9. Megahit is used to perform the genome assembly using the 2 sets of 500,000 reads:

   **megahit -1 output_forward_paired_halfm.fastq -2
   output_reverse_paired_halfm.fastq -o output_directory_halfm**

10. We then do the genome assembly using 25 million reads. This took quite a bit of time.

**megahit -1 output_forward_paired.fastq -2 output_reverse_paired.fastq -o output_dir_25m**



11. The two genome assemblies are analyzed using QUAST software. We go to the directory where the QUAST software is installed and use the following commands to do quality analysis of the half million reads' assembly and the 25 million reads' assembly.
**./quast.py /home/x/ngs_hands_on/output_directory_halfm/final.contigs.fa -o quast_output_halfm_sanchitha**

**./quast.py /home/x/ngs_hands_on/output_dir_25m/final.contigs.fa -o quast_output_25m_sanchitha**





**RESULTS:**

1. FastQC Results of Raw Reads:
   a. Forward Raw Reads:

# Basic Statistics

| Measure | Value |
|---|---|
| Filename | sub_SRR22507560_R1_25m.fastq |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 25000000 |
| Total Bases | 4 Gbp |
| Sequences flagged as poor quality | 0 |
| Sequence length | 161 |
| %GC | 38 |

## Per base sequence quality



Quality scores across all bases (Sanger / Illumina 1.9 encoding)

## Per sequence quality scores



Quality score distribution over all sequences

## Per sequence GC content



GC distribution over all sequences

**Adapter Content**

b. Reverse Raw Reads:

**Basic Statistics**

| Measure | Value |
|---|---|
| Filename | sub_SRR22507560_R2_25m.fastq |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 25000000 |
| Total Bases | 4 Gbp |
| Sequences flagged as poor quality | 0 |
| Sequence length | 161 |
| %GC | 38 |

## Per base sequence quality



Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Position in read (bp)

## Per sequence quality scores



Quality score distribution over all sequences

Average Quality per read

Mean Sequence Quality (Phred Score)

**⚠Per sequence GC content**



**❌Adapter Content**



2. FastQC Results of Trimmed Reads:
   a. Forward Paired Reads:

# Basic Statistics

| Measure | Value |
|---|---|
| Filename | output_forward_paired.fastq |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 19817046 |
| Total Bases | 3 Gbp |
| Sequences flagged as poor quality | 0 |
| Sequence length | 36-161 |
| %GC | 37 |

## Per base sequence quality



Quality scores across all bases (Sanger / Illumina 1.9 encoding)

## ✅ Per sequence quality scores



Quality score distribution over all sequences

## ⚠ Per sequence GC content



GC distribution over all sequences

**Adapter Content**



b. Reverse Paired Reads:

**Basic Statistics**

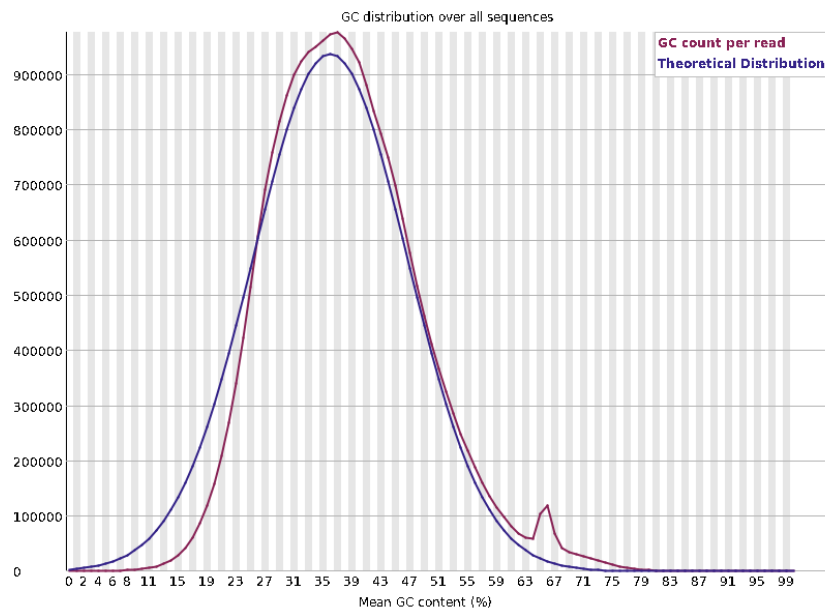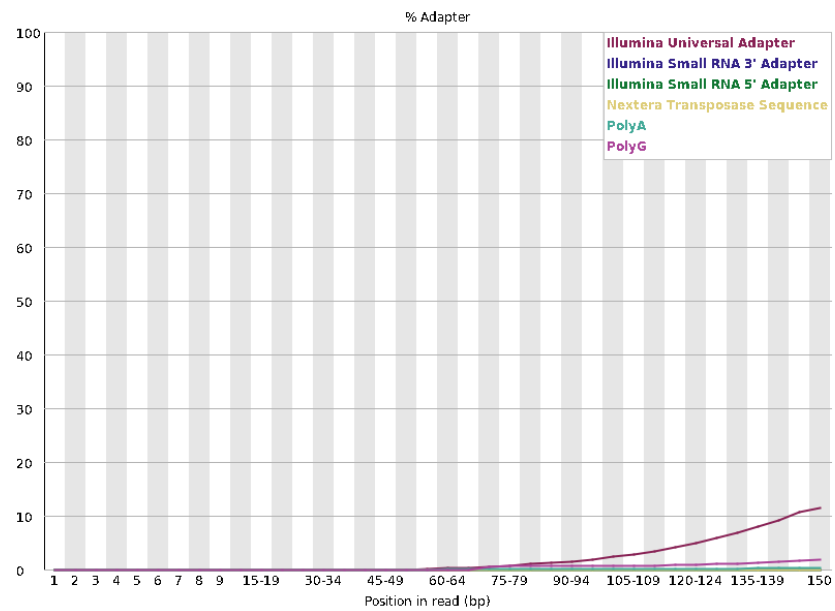| Measure | Value |
|---|---|
| Filename | output_reverse_paired.fastq |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 19817046 |
| Total Bases | 3 Gbp |
| Sequences flagged as poor quality | 0 |
| Sequence length | 36-161 |
| %GC | 37 |

## ✅ Per base sequence quality



Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Position in read (bp)

## ✅ Per sequence quality scores



Quality score distribution over all sequences

Average Quality per read

Mean Sequence Quality (Phred Score)

## ✅ Per sequence GC content
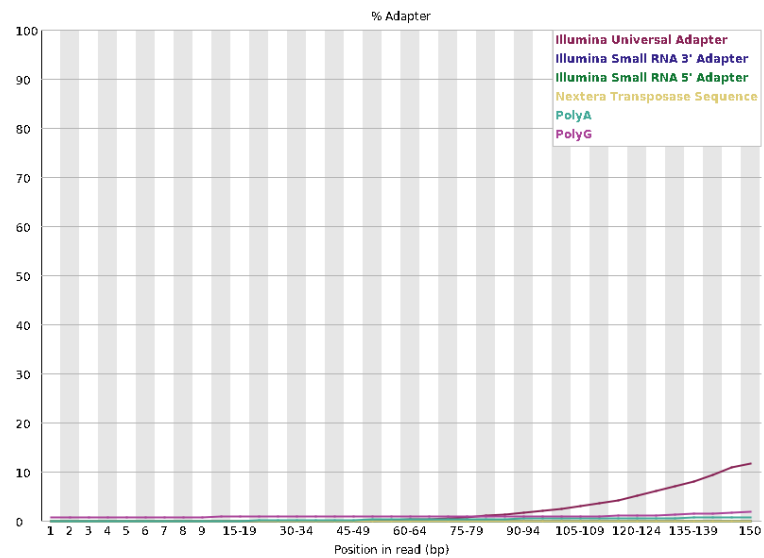


## ✅ Adapter Content



- Thus, comparing the above FastQC results from the raw and trimmed data, we can see that the number of reads decreases from 25 million to 19.8 million, indicating that around 5 million low quality reads were removed. We also observe that, on average, the per base sequence quality scores for the trimmed reads are higher than that of the raw data. The Adaptor Content in the raw data is around 10%, while that of the trimmed reads are about 0%, showing that the software has successfully removed the Adaptor sequences.

**ANALYSIS:**

Results for:

Half a Million Reads Assembly          25 Million Reads:

| Statistics without reference | final.contigs | | Statistics without reference | final.contigs |
|---|---|---|---|---|
| # contigs | 15 556 | | # contigs | 10 683 |
| # contigs (>= 0 bp) | 41 875 | | # contigs (>= 0 bp) | 21 459 |
| # contigs (>= 1000 bp) | 1422 | | # contigs (>= 1000 bp) | 8111 |
| # contigs (>= 5000 bp) | 3 | | # contigs (>= 5000 bp) | 3704 |
| # contigs (>= 10000 bp) | 1 | | # contigs (>= 10000 bp) | 2585 |
| # contigs (>= 25000 bp) | 0 | | # contigs (>= 25000 bp) | 1295 |
| # contigs (>= 50000 bp) | 0 | | # contigs (>= 50000 bp) | 503 |
| Largest contig | 13 535 | | Largest contig | 271 056 |
| Total length | 11 060 080 | | Total length | 112 545 635 |
| Total length (>= 0 bp) | 21 124 484 | | Total length (>= 0 bp) | 115 749 172 |
| Total length (>= 1000 bp) | 1 883 079 | | Total length (>= 1000 bp) | 110 731 196 |
| Total length (>= 5000 bp) | 26 995 | | Total length (>= 5000 bp) | 100 418 295 |
| Total length (>= 10000 bp) | 13 535 | | Total length (>= 10000 bp) | 92 455 409 |
| Total length (>= 25000 bp) | 0 | | Total length (>= 25000 bp) | 71 381 048 |
| Total length (>= 50000 bp) | 0 | | Total length (>= 50000 bp) | 43 392 561 |
| N50 | 684 | | N50 | 37 302 |
| N90 | 528 | | N90 | 4604 |
| auN | 822.8 | | auN | 53 124 |
| L50 | 5968 | | L50 | 803 |
| L90 | 13 404 | | L90 | 3887 |
| GC (%) | 39.01 | | GC (%) | 37.83 |
| **Mismatches** | | | **Mismatches** | |
| # N's per 100 kbp | 0 | | # N's per 100 kbp | 0 |
| # N's | 0 | | # N's | 0 |

The two assemblies compared are the one which is done using 25 millions reads, and one which is done by randomly sampling 500,000 out of the 25 million reads. Since the number of reads used is lesser, it will result in a smaller amount of the genome being assembled, as we have an overall smaller amount of data obtained from the sample. Thus, using a lesser number of reads leads to a less accurate assembly. We can see this in the 3 values tabulated below:

| | Half a Million Reads | 25 Million Reads |
|---|---|---|
| Number of Contigs (≥500bp) | 15,556 | 10,683 |
| N50 | 13,535 | 37,302 |
| Length of Assembled Genome | 11,060,080 | 112,545,635 |

1. Number of Contigs: We can see that the assembly which is done with more reads produces fewer contigs. Contigs are contiguous sequences of DNA that are assembled from the reads. With more reads, it is easier for the software to extend contigs and merge overlapping sequences, resulting in a reduced number of contigs and more coverage of the genome. Therefore, with less number of reads(half a million), there are more contigs, indicating that the resulting genome would be more fragmented since enough overlaps were not found.

2. N50 Value: The N50 number is a value that describes the contiguity of an assembly. It represents the contig length such that using equal or longer contigs will produce half of the bases of the assembly. When assembly is done with 25 million reads, the N50 value is more, implying that there are longer and more contiguous sequences in the assembly. Having a high N50 is an indicator of a good quality assembly. However, the N50 of the half million reads assembly is quite low, indicating a lower quality.

3. Length of Genome Assembled: When assembly is done with 25 million reads, it leads to longer contigs and, thus, a more complete assembly. Thus, the resulting assembled genome is longer. When assembly is done with half a million reads, the length of the genome assembled is much less, due to the fact that there is lesser coverage of the genome, resulting in the possibility for more fragmentation of the genome and more gaps being present in the sequence.

Thus, we can see that the genome assembly from the 25 million reads is of better quality than that of the subsampled half a million reads assembly. However, increasing the reads beyond that does not guarantee a higher quality of assembly and thus finding the optimal number of reads for assembly such that the quality is high but the time taken for assembly is not too much is an important step in the assembly process.