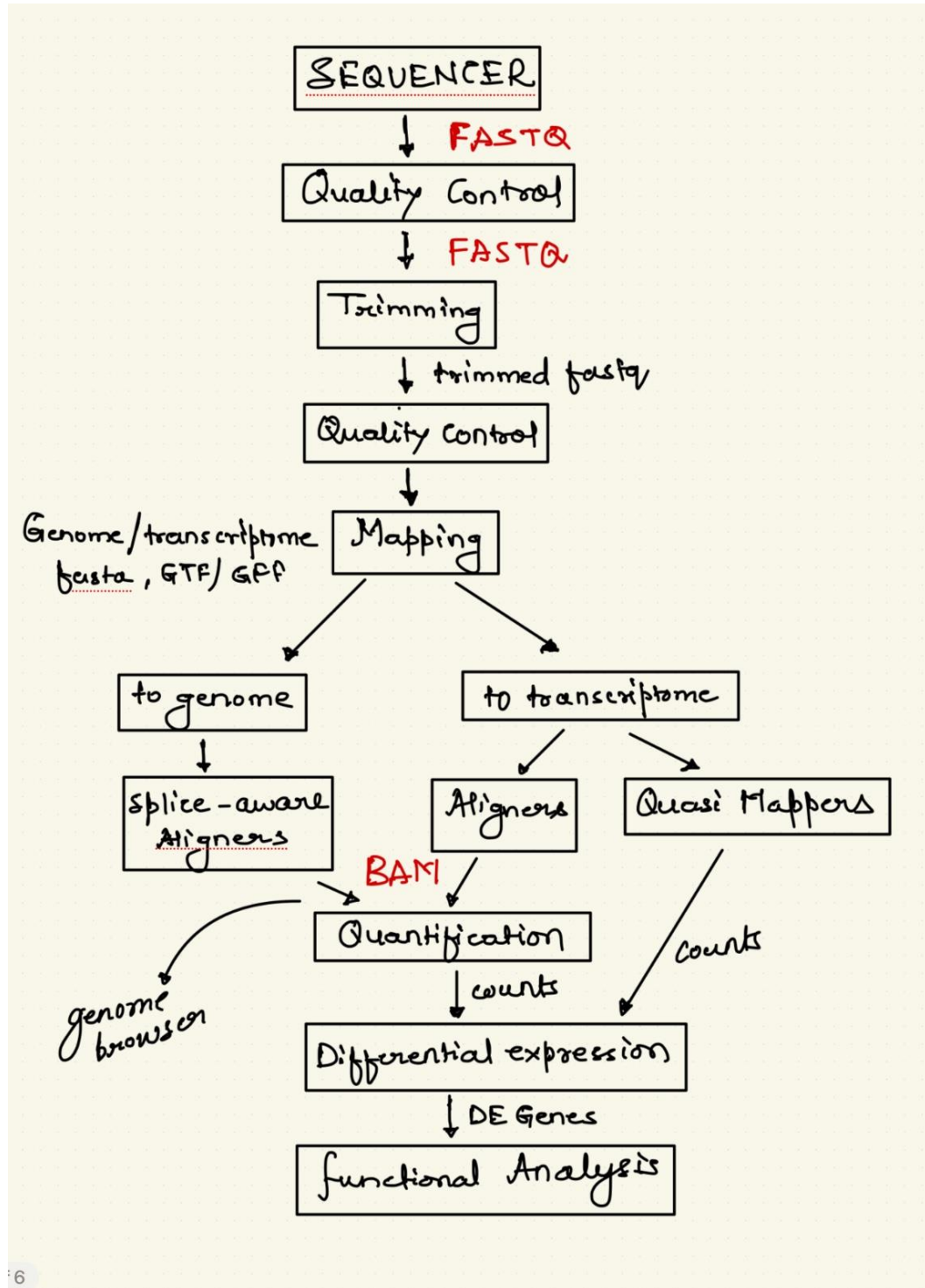


Bio-informatics Tasks

Tasks:

- ❖ Downloading a FASTQ File - https://drive.google.com/file/d/1DGHjbhcRy_zTm6H9C_AUpkzBML-JhtA3/view
- ❖ Converting FASTQ File to SAM and BAM Format.
- ❖ SAM format is known as **Sequence Alignment Map** and BAM files are **Binary representation of SAM Files** which are compressed SAM Files.

General Workflow:



We did work till converting FASTQ to a SAM & BAM file.

Steps taken:

- As we cannot directly convert the FASTQ Files into SAM and BAM formats, we have to first create a pipeline to do this task where we first get the FASTQ file upon which we do the **“Quality Control”** of reads from sequencer which are in FASTQ files. To perform we used FASTQC as a tool to check quality control reports. Goal of this step: We did quality control to look for bases from the reads that are of poor quality and also to check if any adapter sequences which are present.
- To remove such poor-quality reads and to remove adapter sequences we performed “Trimming”. For this task we used **Trimmomatic tool** to perform this task.
- After getting good quality reads from the above step, we performed **“Mapping”** to a genome file which was downloaded from this link - https://genome-index.s3.amazonaws.com/hisat/grch38_genome.tar.gz
- Using HISAT we mapped this genome files to the Trimmed FASTQ file which contains RNA Sequences, here using HISAT2 and samtools we outputted the BAM and SAM file using Bash pipeline script.
- This study helped us to understand about the genome file system and on this basis, we also understand that using this BAM file we can do quantification such as sequence counts and also this can be used to create Differential Expressions which are further used to do functional analysis depending on the user-based research scenarios.

Shell File Commands: -

```
SECONDS=0
```

```
# change working directory
```

```
cd /Desktop/Fastq/
```

```
# STEP 1: Running fastqc
```

```
fastqc data/demo.fastq -o data/
```

```
# run trimmomatic to trim reads with poor quality
```

```
java -jar ~/Trimmomatic-0.39/trimmomatic-0.39.jar SE -threads 4 data/demo.fastq data/demo_trimmed.fastq  
TRAILING:10 -phred33
```

```
echo "Trimmomatic finished running!"
```

```
fastqc data/demo_trimmed.fastq -o data/
```

```
# STEP 2: Running HISAT2
```

```
mkdir HISAT2
```

```
# get the genome indices
```

```
wget https://genome-index.s3.amazonaws.com/hisat/grch38\_genome.tar.gz
```

```
# run alignment
```

```
hisat2 -q --rna-strandness R -x HISAT2/grch38/genome -U data/demo_trimmed.fastq | samtools sort -o  
HISAT2/demo_trimmed.bam
```

```
hisat2 -q --rna-strandness R -x HISAT2/grch38/genome -U data/demo_trimmed.fastq -S  
HISAT2/demo_trimmed.sam
```

```
echo "HISAT2 finished running!"
```

```
duration=$SECONDS
```

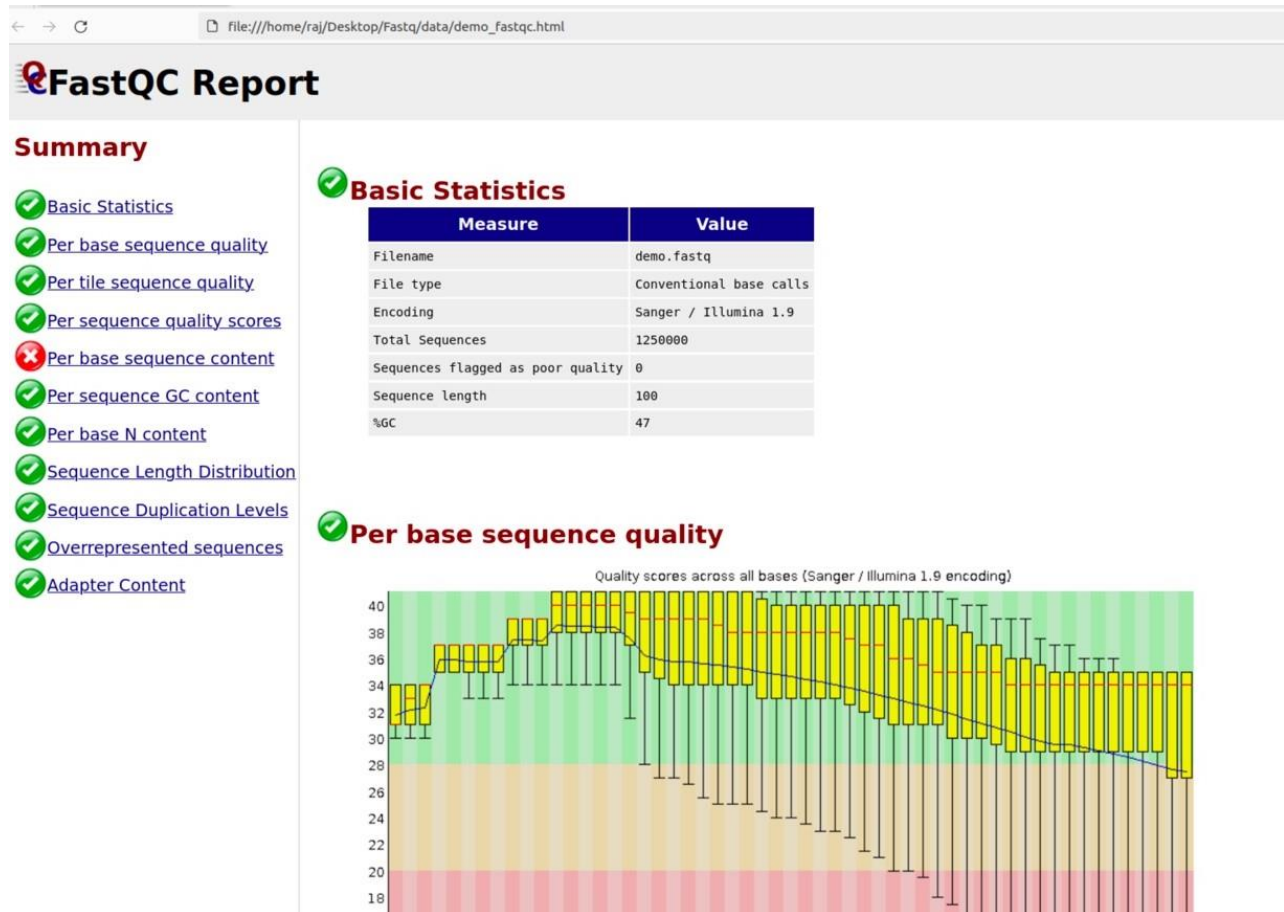
echo "\$((\$duration / 60)) minutes and \$((\$duration % 60)) seconds elapsed."

Working Environment-

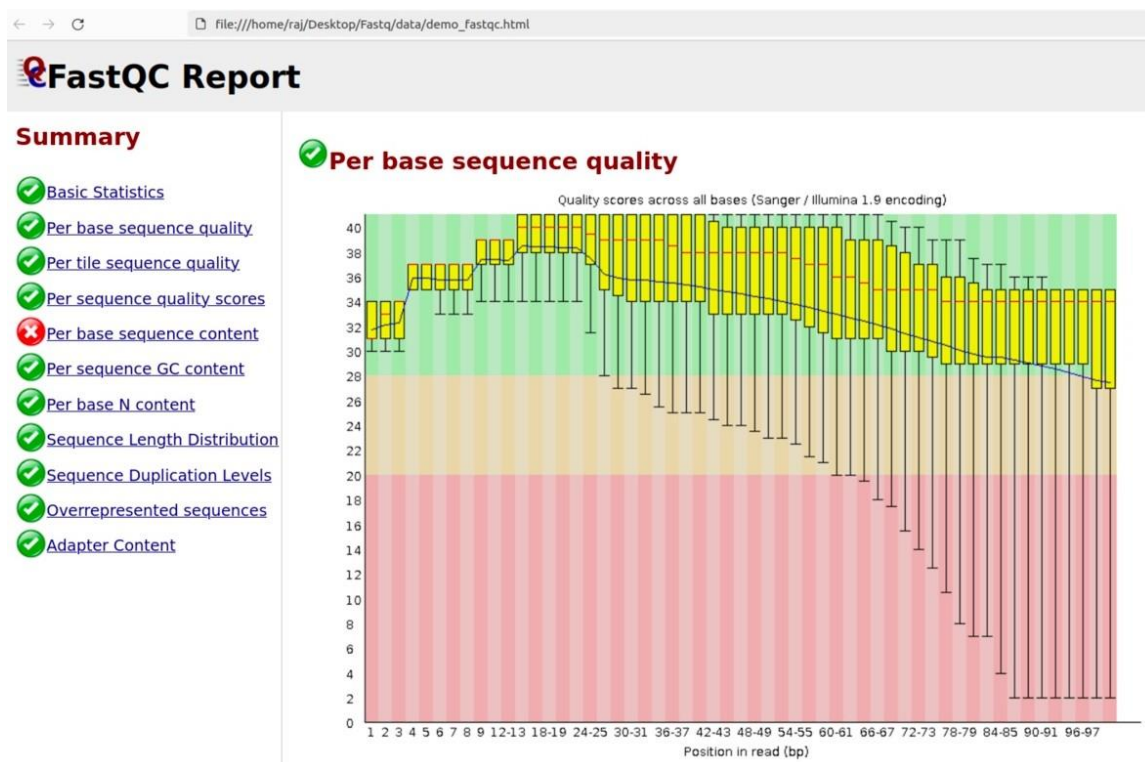
I downloaded Oracle VM VirtualBox on laptop and installed Ubuntu to do this task as my MacBook did not gave me permissions to download third-party apps and was throwing errors with wget commands as I was not able to download the wget on the terminal. Using ubuntu (Linux) terminal the task was successfully performed and executed.

Screenshots –

- FASTQC Report after doing initial quality control.

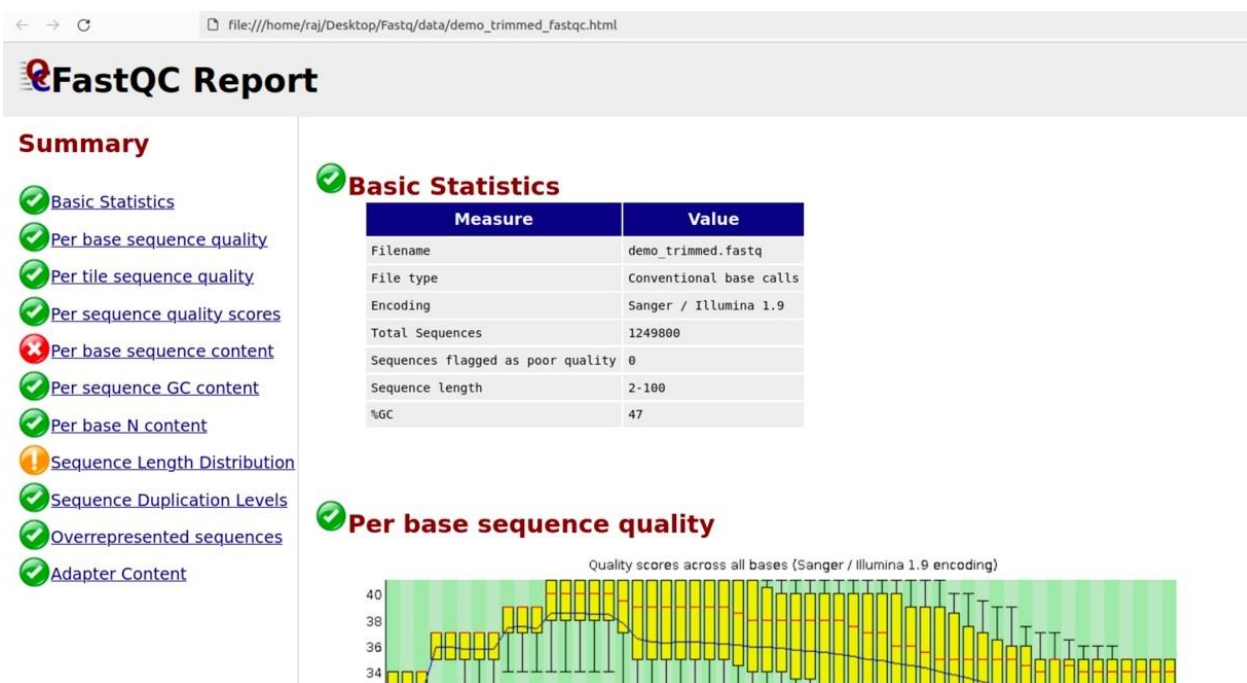


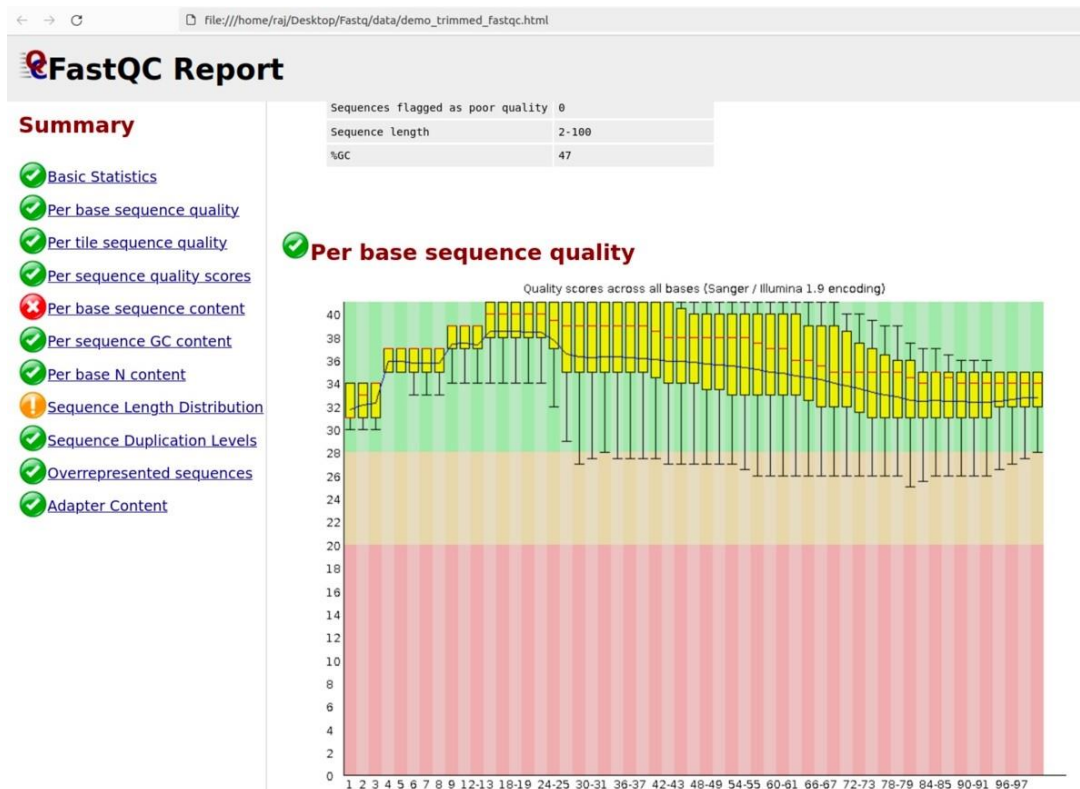
The initial FASTQ file has 1,250,000 Total sequences which we can see in the basic statistics above.



As we can see from the base sequence quality above, the end sequence quality scores or phred are not that good, hence we did trimming.

- After trimming



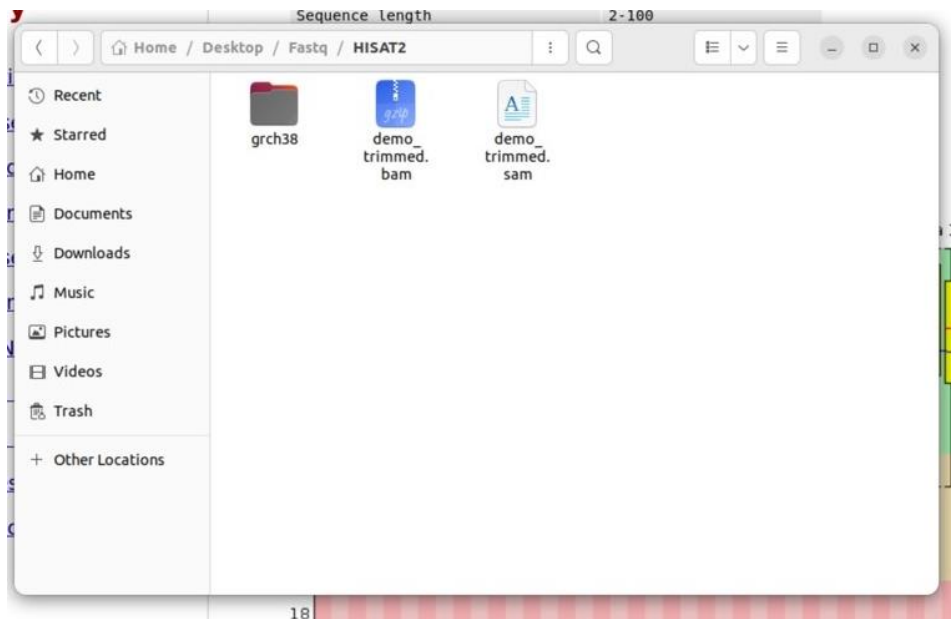


After trimming we again did quality control to check the results of the trimming, as you can see above all sequences are of good quality and it got improved from the initial fastq file.

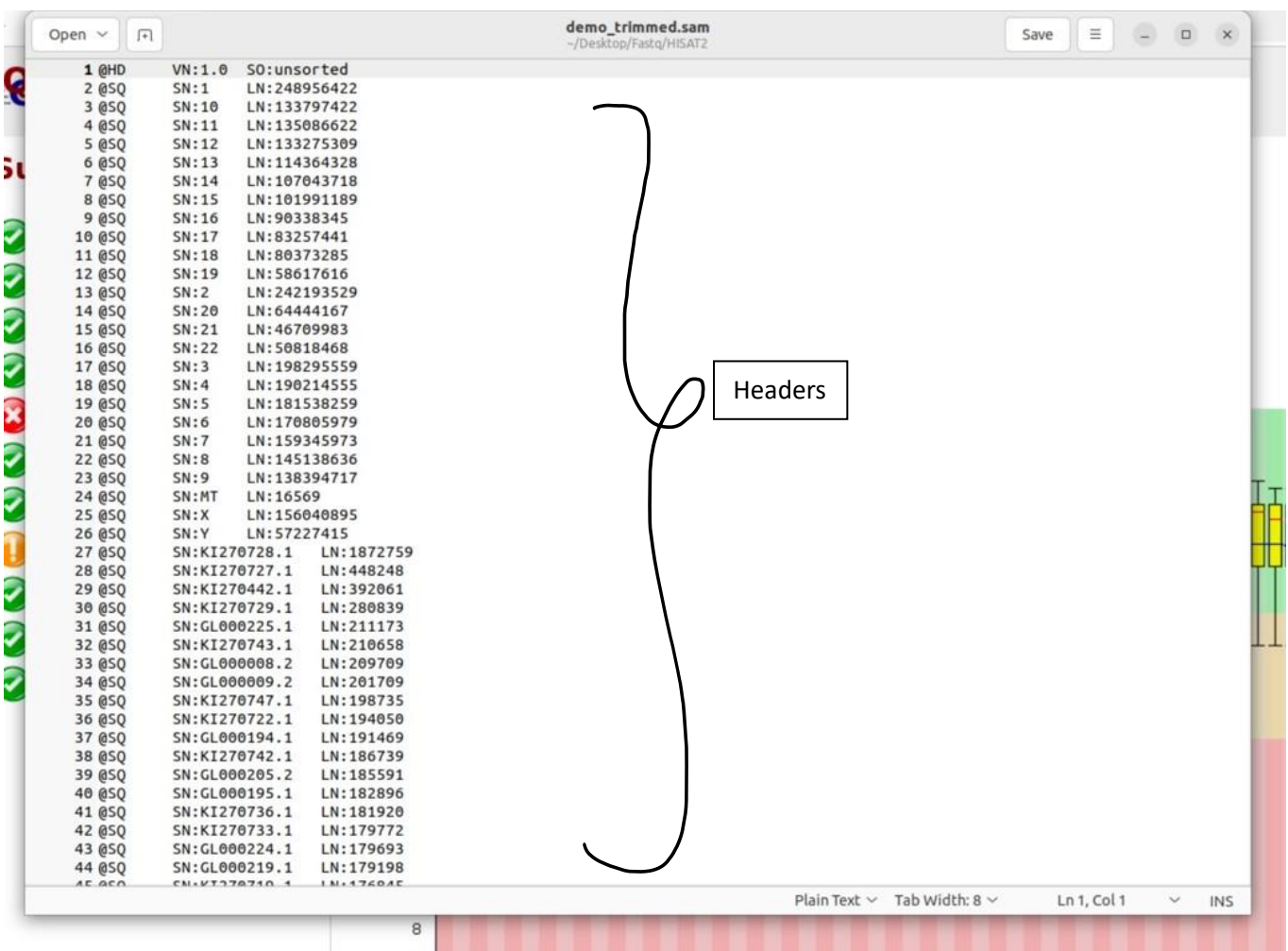
```
raj@raj-VirtualBox: ~/Desktop/Fastq
raj@raj-VirtualBox:~/Desktop$ cd Desktop/
raj@raj-VirtualBox:~/Desktop$ cd Fastq/
raj@raj-VirtualBox:~/Desktop/Fastq$ hisat2 -q --rna-strandness R -x HISAT2/grch38/genome -U data/demo_trimmed.fastq -S HISAT2/demo_trimmed.sam
^C[^X^C(ERR): hisat2-align died with signal 2 (INT)
^Craj@raj-VirtualBox:~/Desktop/Fastq$ open test.sh
raj@raj-VirtualBox:~/Desktop/Fastq$ ./test.sh
./test.sh: line 3: SECONDS: command not found
^C(ERR): hisat2-align died with signal 2 (INT)
./test.sh: line 26: echo0 minutes and 26 seconds elapsed.: command not found
raj@raj-VirtualBox:~/Desktop/Fastq$ ./test.sh
1249800 reads; of these:
  1249800 (100.00%) were unpaired; of these:
    86563 (6.93%) aligned 0 times
    1082380 (86.60%) aligned exactly 1 time
    80857 (6.47%) aligned >1 times
93.07% overall alignment rate
./test.sh: line 26: echo48 minutes and 59 seconds elapsed.: command not found
raj@raj-VirtualBox:~/Desktop/Fastq$
```


After this we mapped the trimmed file to the human genome data, and using HISAT2 we outputted the trimmed fastq file as a SAM file and using HISAT2 and samtools we outputted the trimmed fastq file as a BAM file. From the above terminal snippet, you can see that it took 48 minutes and 59 seconds for this mapping and conversion.

- **Snippet of the files converted:**



Output of the SAM file:



Open

demo_trimmed.sam

Save

~\Desktop\Fastq\HISAT2

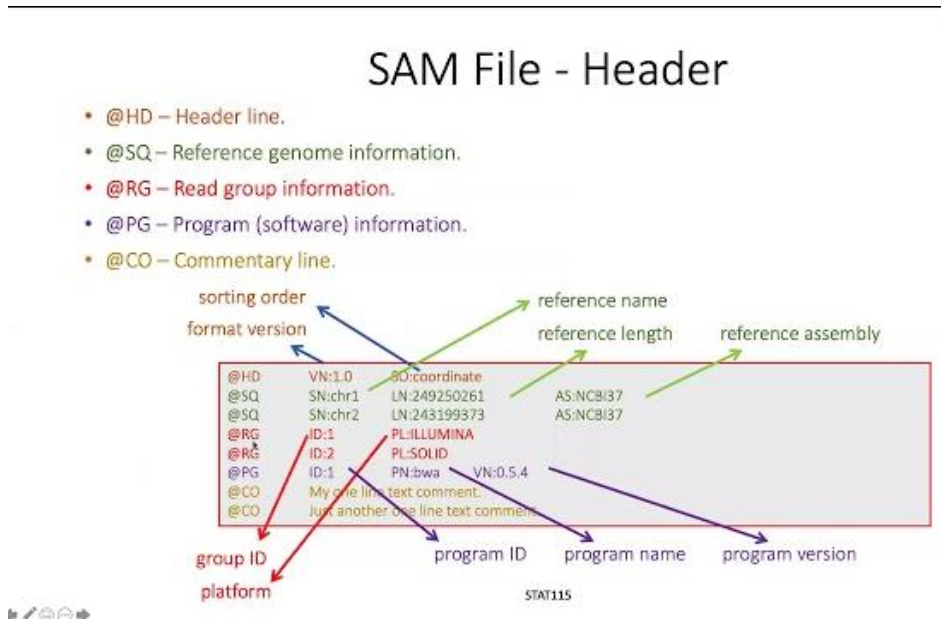
12

```

182 @SQ      SN:KI270376.1  LN:1136
183 @SQ      SN:KI270337.1  LN:1121
184 @SQ      SN:KI270335.1  LN:1048
185 @SQ      SN:KI270378.1  LN:1048
186 @SQ      SN:KI270379.1  LN:1045
187 @SQ      SN:KI270329.1  LN:1040
188 @SQ      SN:KI270419.1  LN:1029
189 @SQ      SN:KI270336.1  LN:1026
190 @SQ      SN:KI270312.1  LN:998
191 @SQ      SN:KI270539.1  LN:993
192 @SQ      SN:KI270385.1  LN:990
193 @SQ      SN:KI270423.1  LN:981
194 @SQ      SN:KI270392.1  LN:971
195 @SQ      SN:KI270394.1  LN:970
196 @PG      ID:hisat2      PN:hisat2      VN:2.2.1
HISAT2/demo_trimmed.sam -U data/demo_trimmed.fastq"
197 SRR960459.1 4 15 34981220 60 100M * 0 0
NAGAACTGGCGCGCAATGGGCTGACCGCTTCTCTGCTGCTTACGGTATCGCCGCTCCCGATTGCGCAGCGCATCGCCTTCTATCGCCTTCTTGACGAGTT
#1=DDFFHHHGHIIJJJJJJGGGAFGBHHEHGFBBFDEDECDDA==CB@BDDDD?;B-<CBDD>BBBDDDB5<@DDDCDD@-9ACDDDB7B<? YT:Z:UU
198 SRR960459.2 0 15 34981220 60 100M * 0 0
CTCCTTACTATGCAGGACAGCAAAATGCGCTGAAACCATTCAGTTTAAATGACAGGCTTTTCAATAAAAAATGCATTTTAAATAATACAGGCTTTAAAAATA @CCFFFFFHGHJJJBHIJEEIJJJIGGIIJJIIIGIAHGIIG?
HDIEIGIJEHGHIDHGHGEIIIJIGGEIJJHHHHHFCDF@CACCCCA@CDD AS:i:0 XN:i:0 XM:i:0 XO:i:0 XG:i:0 NM:i:0 MD:Z:100 YT:Z:UU XS:A:- NH:l:1
199 SRR960459.3 0 6 34282136 60 75M1I24M * 0 0
CCTTGGACCACTGATTAGTGACACATTTTCCACACAAAGGGAGGTTAAAGACAAAGTTTAACTGAGAAACTTTTTTTTGGCAAGATAATAAAG ?@DDEFFHHGGIIG<BCHFHGH@HEHEHGEHIIIEHGH=?
GBBDEEHGEGIEFF=CEIGDADCHHCFDECCBBB@B:4<C@C@:C> AS:i:-8 ZS:l:-18 XN:i:0 XM:i:0 XO:i:1 XG:i:1 NM:l:1 MD:Z:99 YT:Z:UU XS:A:- NH:l:1
200 SRR960459.4 0 7 6023889 60 75M * 0 0
CTGCCAATCCATGGCAGACCTTTCTGGGATTCAAAACCAATTCATCAGATCGCTGCCTCTGAGGAGGTACA ??@FDDDBDBCFDE3B@FCDGGGEE+<F*) :E*?<D3? )003B>FB?E*?8B7EGB.@37=4@CFE;
77B;? AS:l:0 XN:l:0 XM:l:0 XO:l:0 XG:l:0 NM:l:0 MD:Z:75 YT:Z:UU XS:A:- NH:l:1
201 SRR960459.5 0 12 56669698 60 4596M * 0 0
NTCCGCTCCAGCGTTCAAGCAATCCCGTTTCAGCCTCAGGAGTAGTGGGACACAGGTGCGCGCCACCATGCCAGCTAATTTTTTTTGTATTTTAG #1=DDD?
DAHMHHEHAGFCEDCFDGGGCHJJJJFBEE9B0BDG=B;FHHGDEC27=ADFB@BBBAAAC:CCBDD<CCDC@B<B?CCDEDA? AS:i:-4 XN:i:0 XM:i:0 XO:i:0 XG:i:0 NM:l:0 MD:Z:
96 YT:Z:UU XS:A:- NH:l:1
202 SRR960459.6 16 16 791514 60 100M * 0 0
AGAACTGGAGCGTTGACGTAACTCGGGGGAAGTCGTTGTCCTGAAGGGCTCTGCGCTGGCCAGGATGATCTGGTGGAGCGTGGTGGTTCTGN @CCCA@><895(4+4(@@A8:A88B>B?<A:(883BBA;55?;
7)3HEAA9D@544;70(80??*B@HDFDCBD=F?<C<HC<CF@=DFBFHDDDB?1# AS:l:-1 ZS:l:-1 XN:i:0 XM:l:1 XO:i:0 XG:i:0 NM:l:1 MD:Z:99G0 YT:Z:UU XS:A:+
NH:l:1
203 SRR960459.7 16 21 8260828 1 80M * 0 0
TGAGCTCTCGCTGGCCTTGAAATCCGGGGGAGAGGTGTAATCTCGCGCGGGCGGTACCCATATCCGCAGCAGGTC C@:2)0099BA?8(<3:>>>@5/9<582:5(;?;96;.(6>??645/
AQHFHDBADGHCIGI@HHFHBEEDD@=8 AS:l:0 ZS:l:0 XN:l:0 XM:l:0 XO:l:0 XG:l:0 NM:l:0 MD:Z:80 YT:Z:UU XS:A:+ NH:l:5
204 SRR960459.7 272 21 8216599 1 80M * 0 0
TGAGCTCTCGCTGGCCTTGAAATCCGGGGGAGAGGTGTAATCTCGCGCGGGCGGTACCCATATCCGCAGCAGGTC C@:2)0099BA?8(<3:>>>@5/9<582:5(;?;96;.(6>??645/
AQHFHDBADGHCIGI@HHFHBEEDD@=8 AS:l:0 ZS:l:0 XN:l:0 XM:l:0 XO:l:0 XG:l:0 NM:l:0 MD:Z:80 YT:Z:UU XS:A:+ NH:l:5
205 SRR960459.7 272 21 8443868 1 80M * 0 0
TGAGCTCTCGCTGGCCTTGAAATCCGGGGGAGAGGTGTAATCTCGCGCGGGCGGTACCCATATCCGCAGCAGGTC C@:2)0099BA?8(<3:>>>@5/9<582:5(;?;96;.(6>??645/
AQHFHDBADGHCIGI@HHFHBEEDD@=8 AS:l:0 ZS:l:0 XN:l:0 XM:l:0 XO:l:0 XG:l:0 NM:l:0 MD:Z:80 YT:Z:UU XS:A:+ NH:l:5

```

Header understanding of the SAM file:



Steps Followed:

