

# Answers to QAA

## Part 1

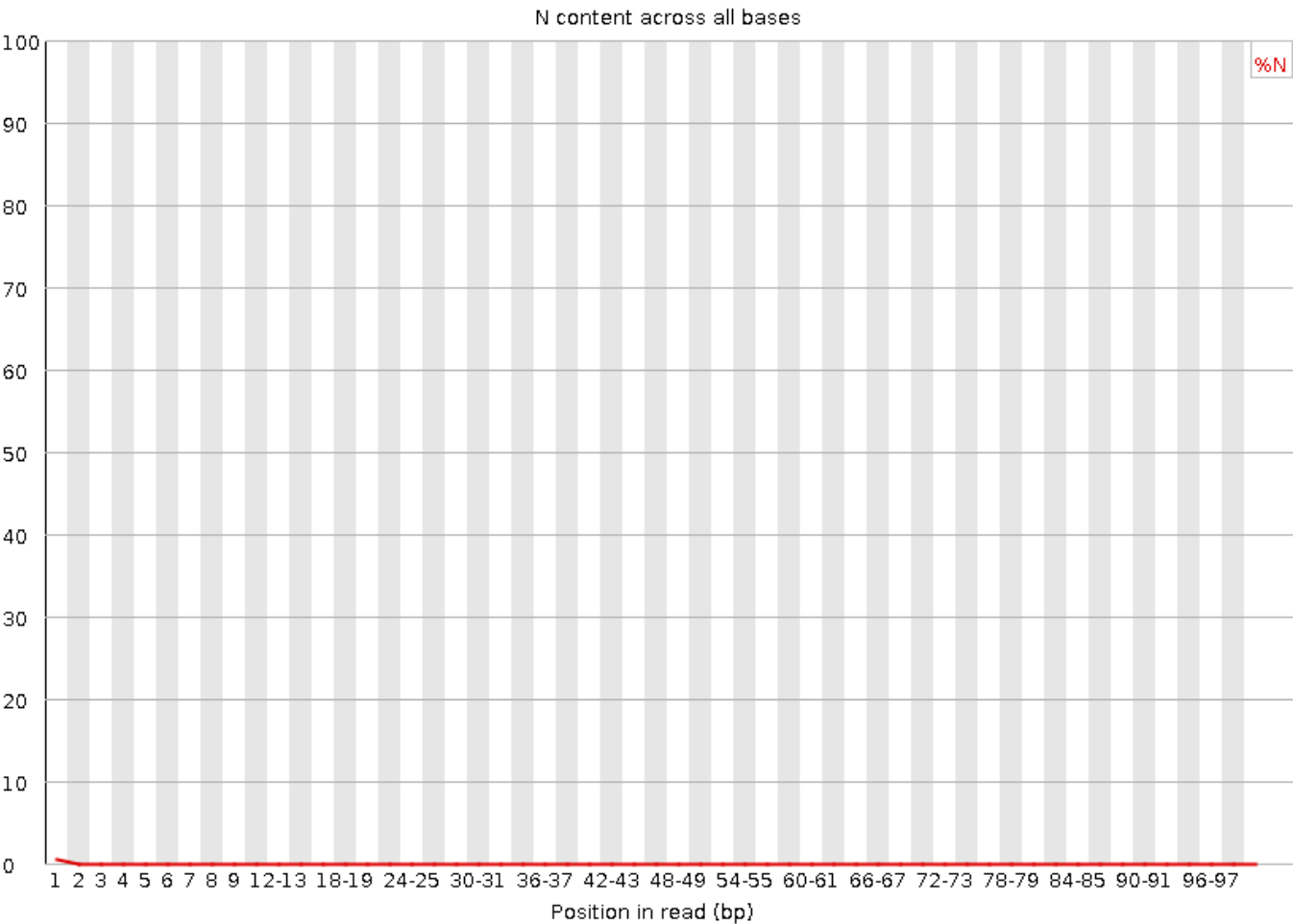
### Question 1

I believe that the data looks pretty consistent across the Fastqc. A little bit of Ns at the beginning of the reads. Also, R1 always has a better quality overall then R2 which makes sense.

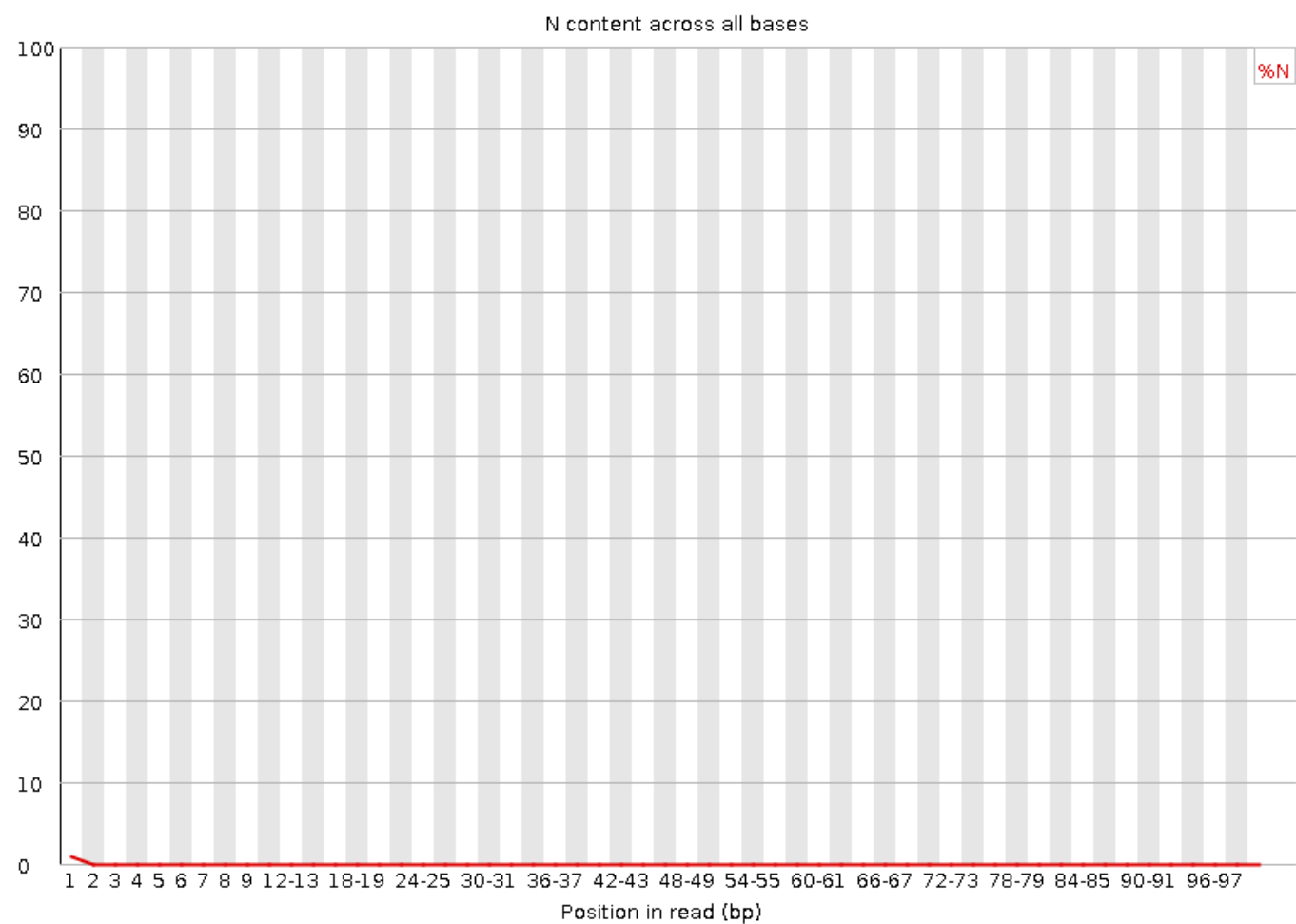
### Fastqc Data

#### Per Base N Content

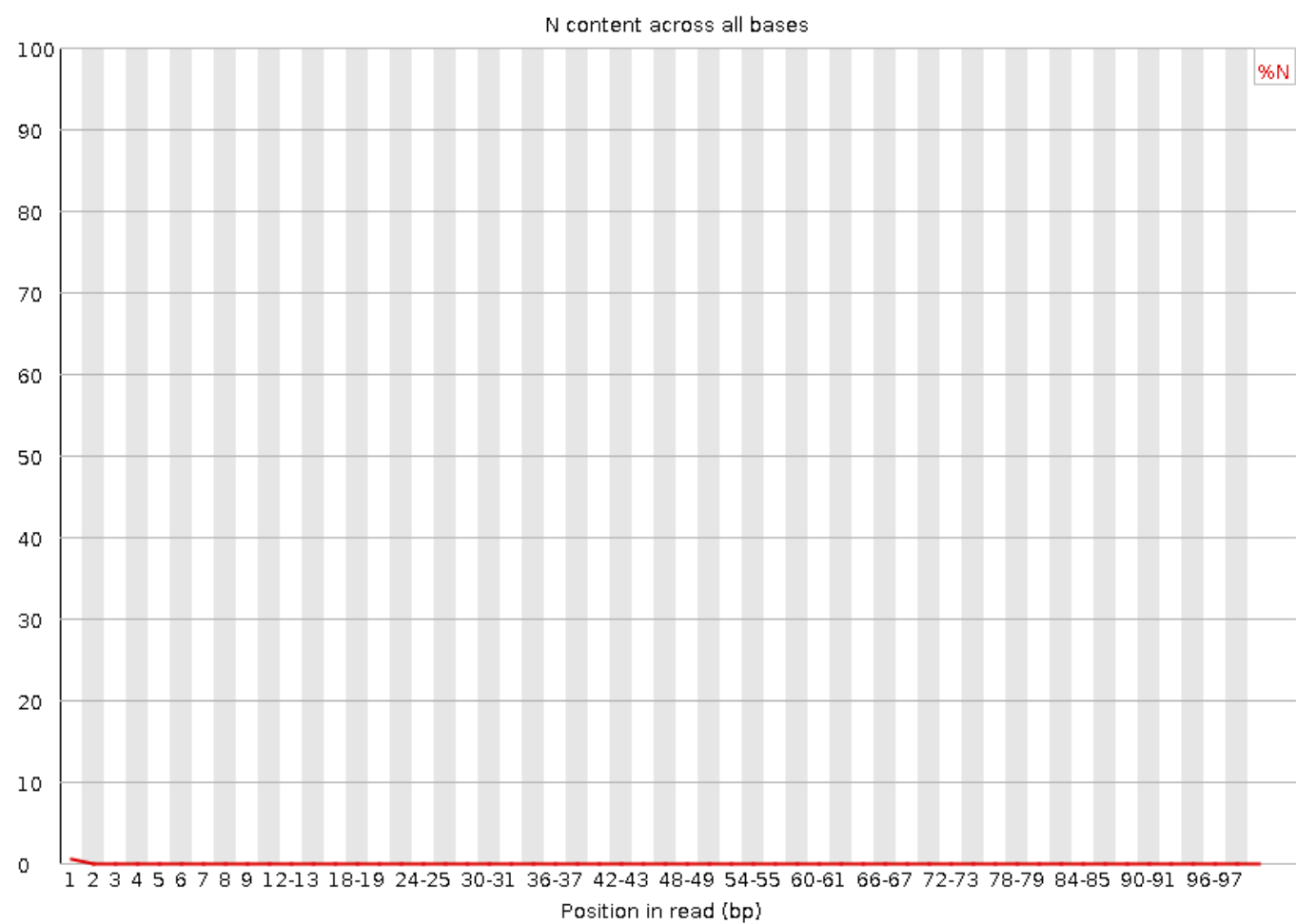
19\_3F\_R1



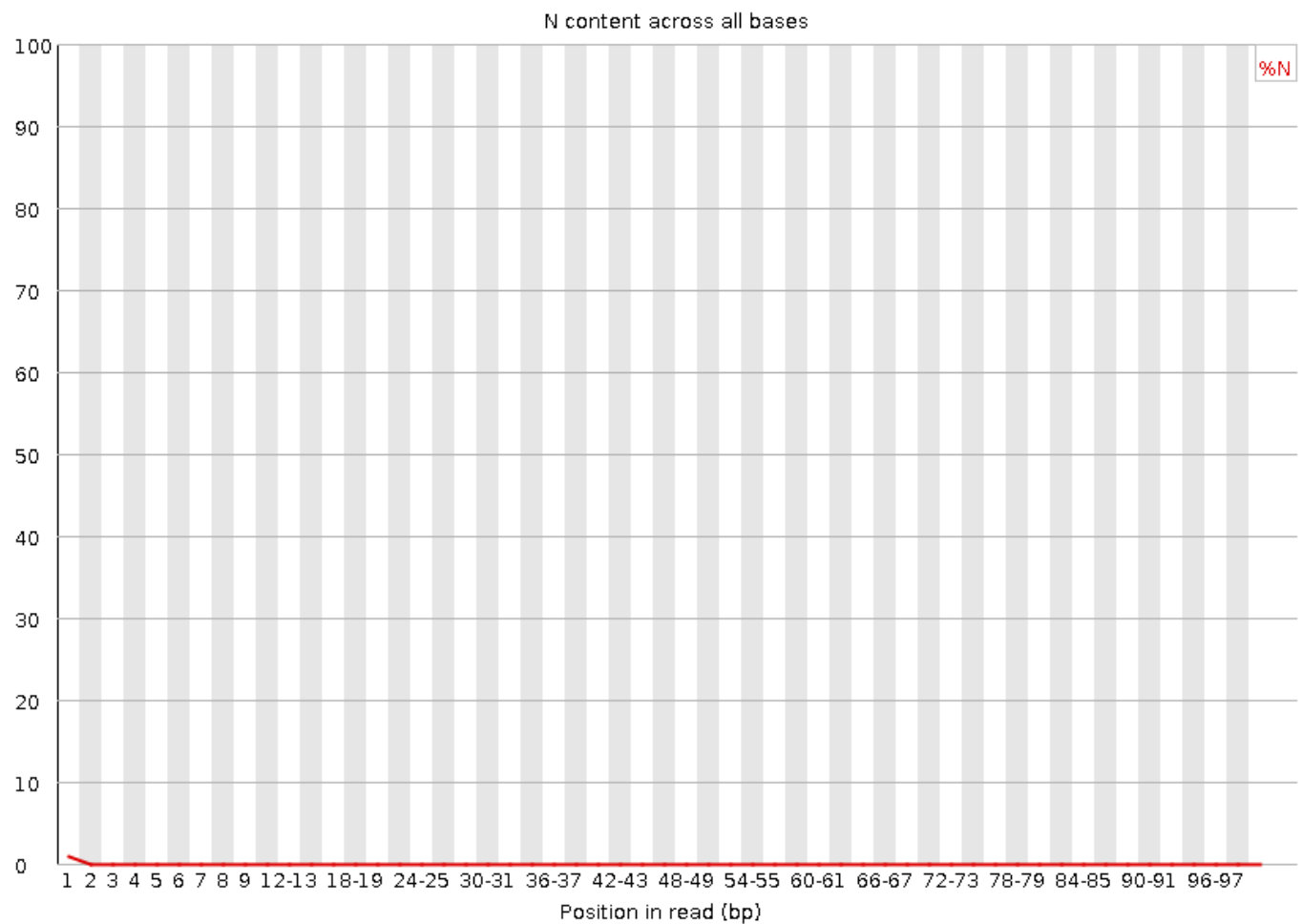
19\_3F\_R2



2\_2B\_R1



2\_2B\_R2



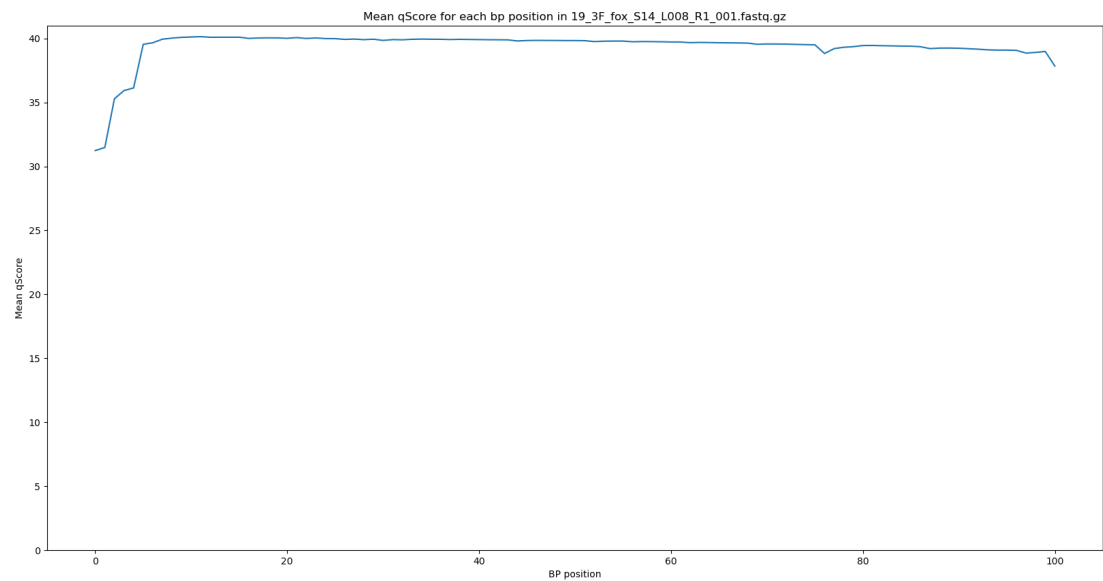
## Question 2

Fastqc runs faster, probably because Fastqc is written in Java. Also, the data seems to look pretty similar to the Fastqc graphs.

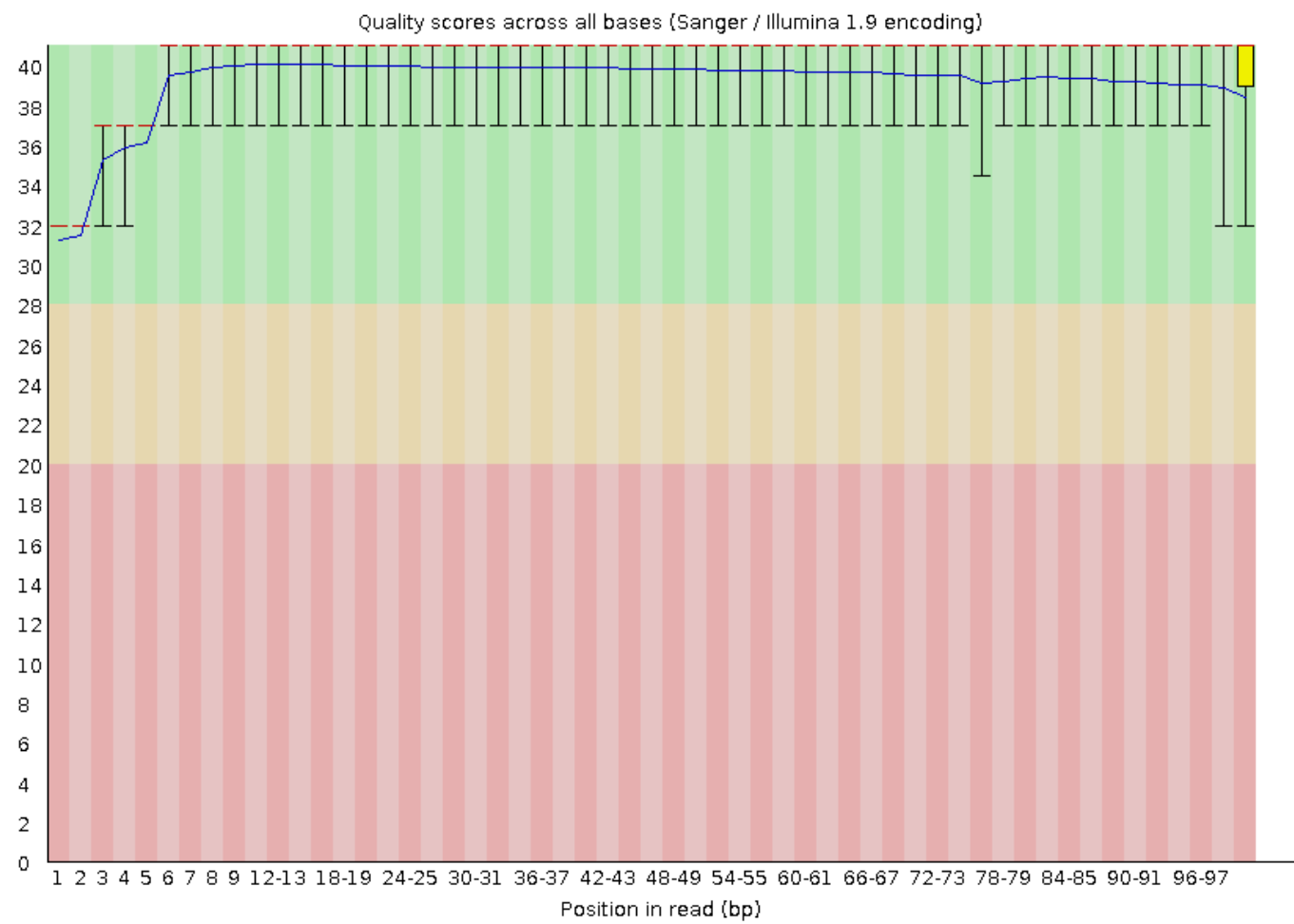
## Question 3

The overall state of the data is pretty good. I feel like my graphs and the Fastqc graphs are very similar. They have the same pattern, however, since I did not graph the error bars I don't know if their error is similar.

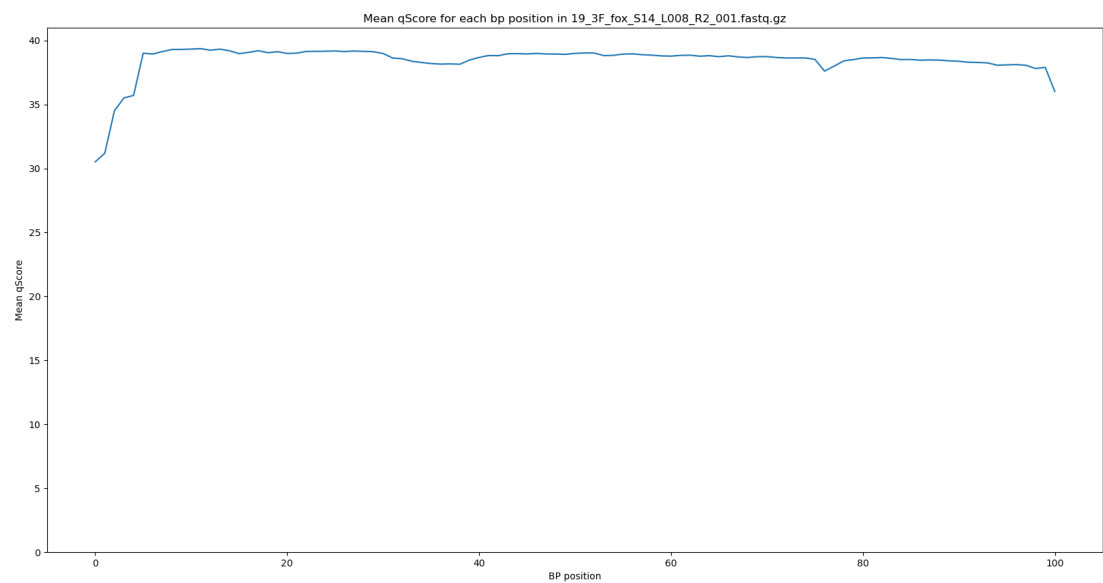
19\_3F\_fox\_S14\_L008 R1



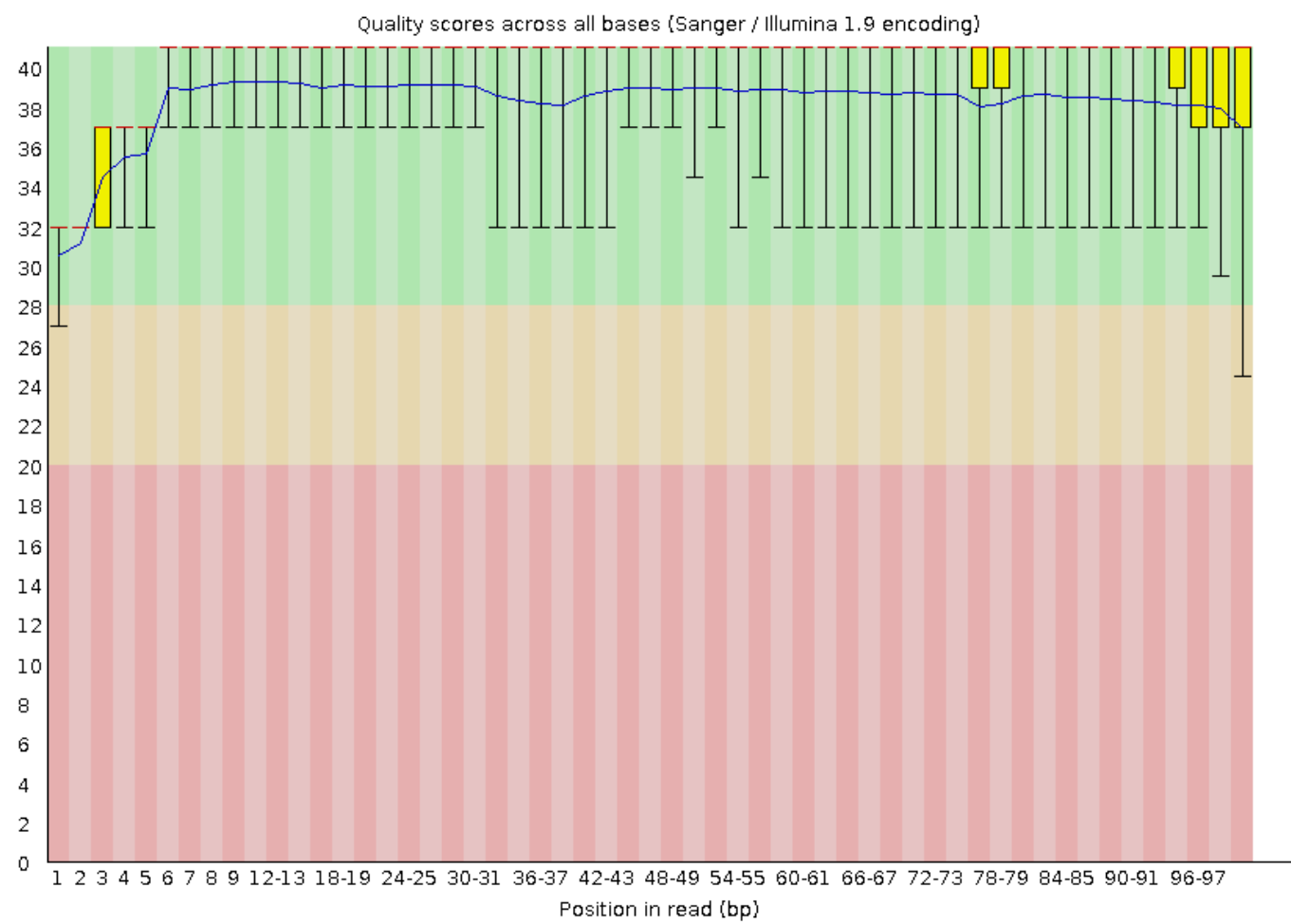
Fastqc R1



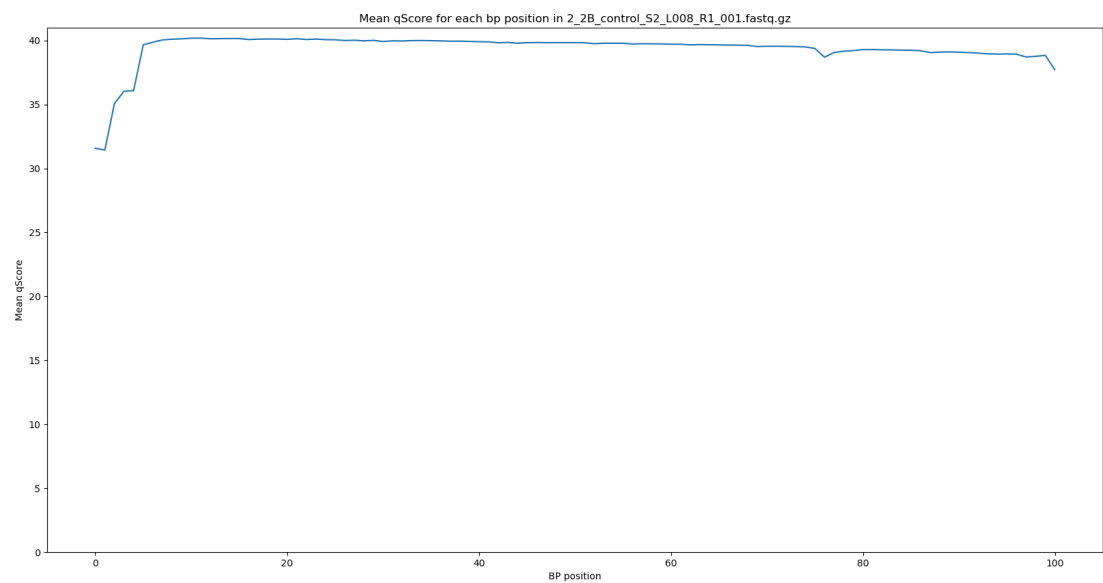
R2



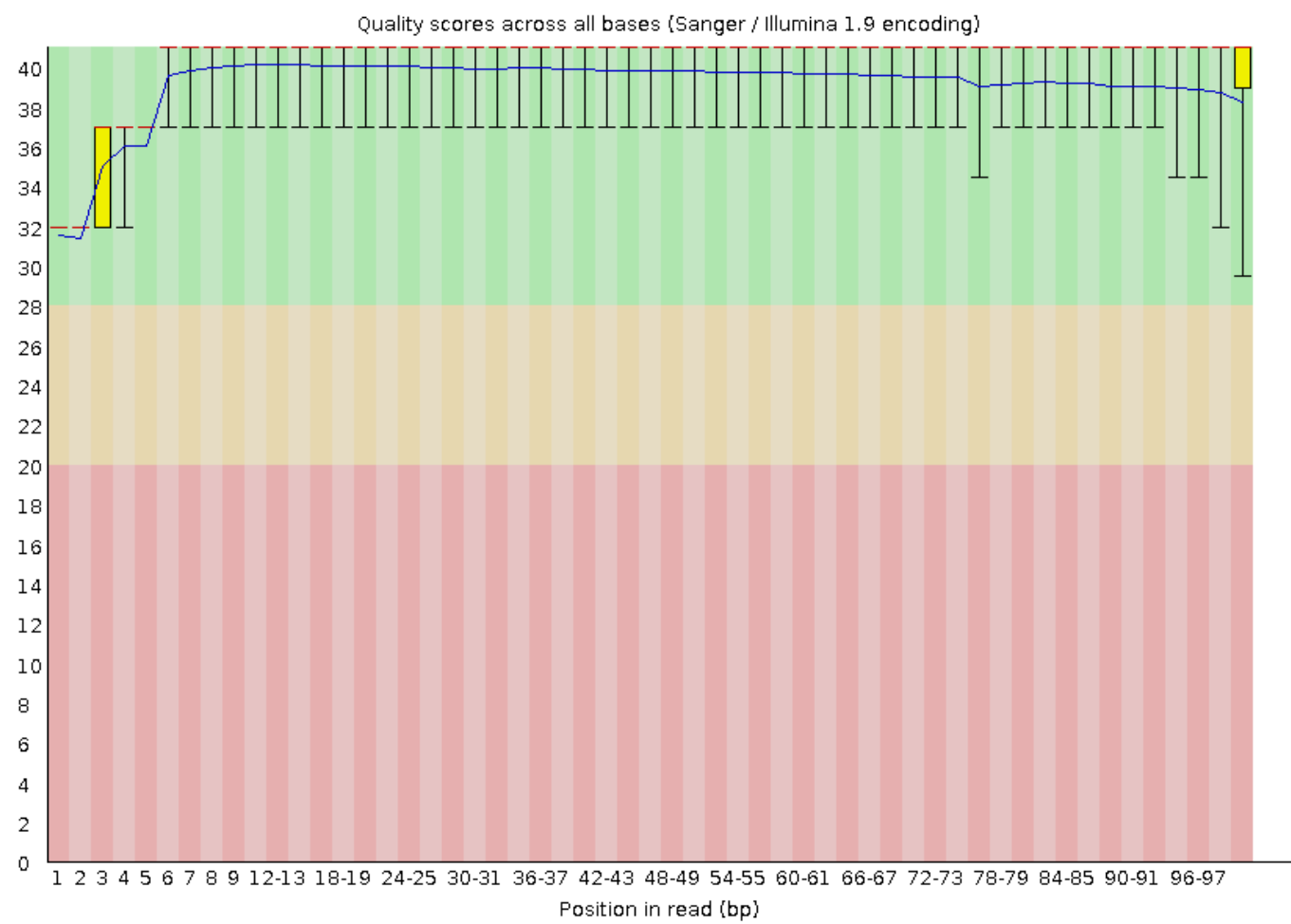
Fastqc R2



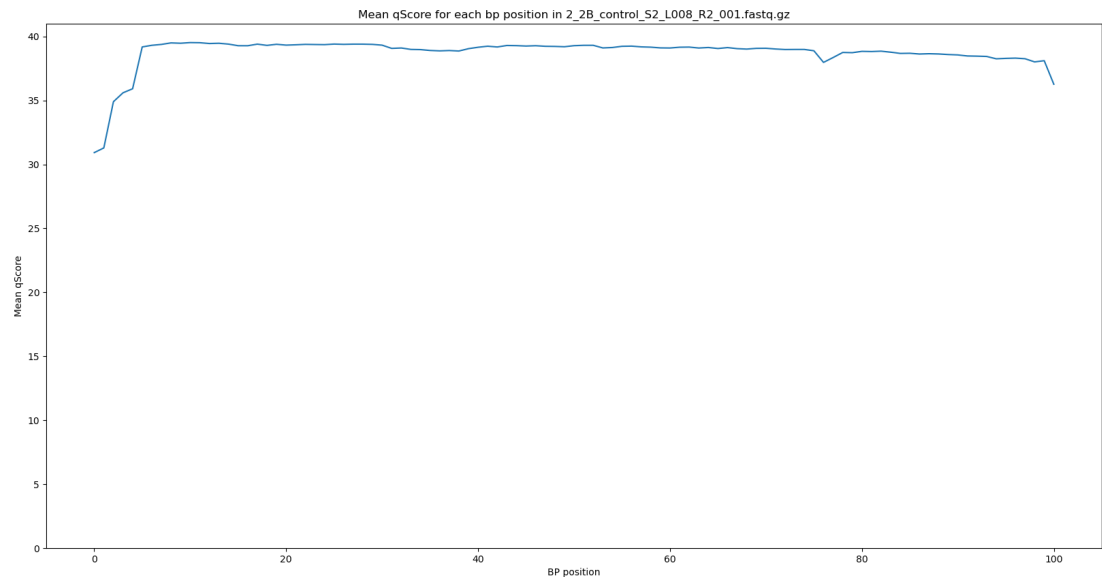
2\_2B\_control\_S2\_L008 R1



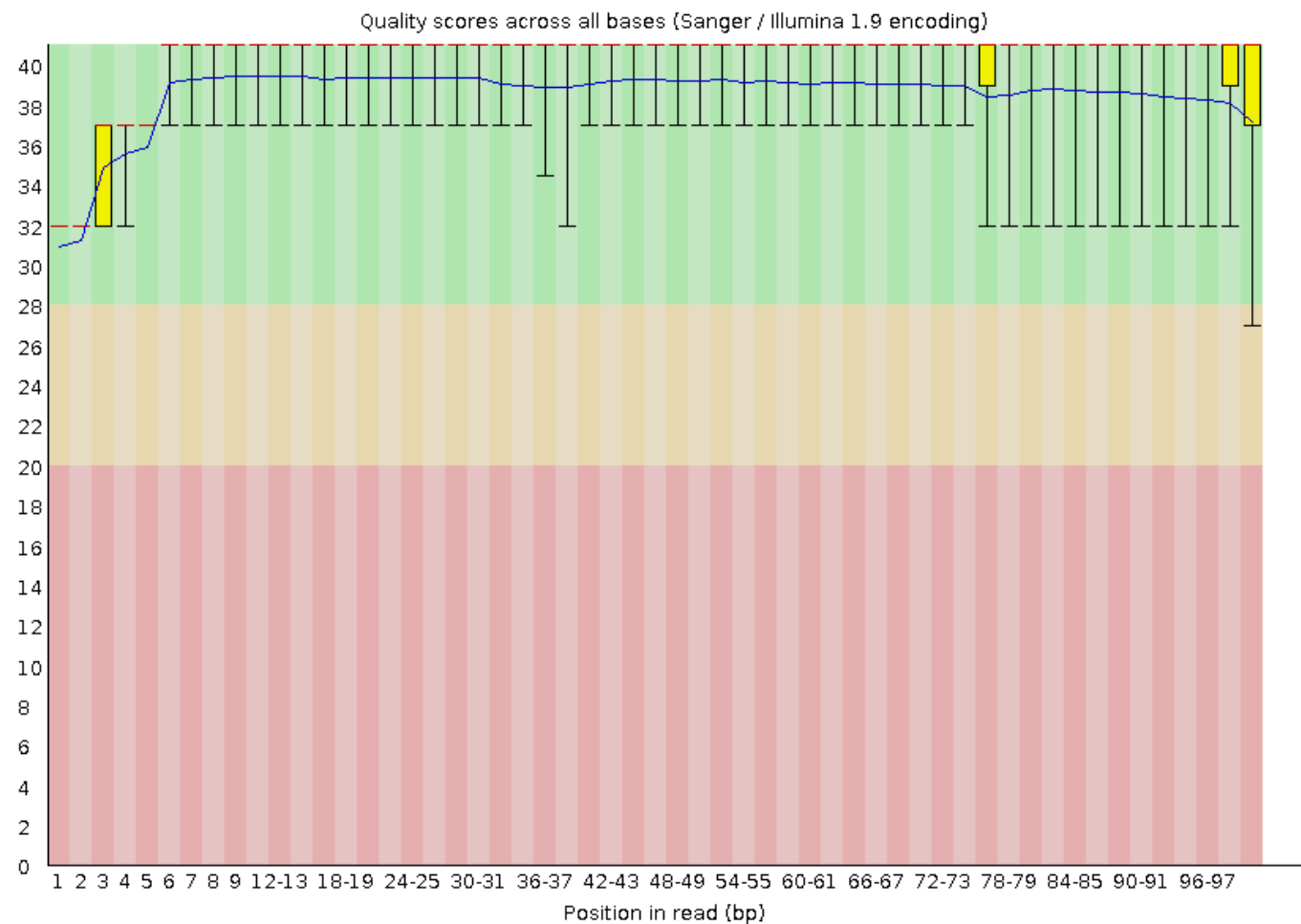
Fastqc R1



R2



Fastqc R2



Part 2

Question 5

You can see the adapters outlined in red if you run the command. Unfortunately it doesn't translate well to a markdown format.

```
zcat
/projects/bgmp/shared/2017_sequencing/demultiplexed/2_2B_control_S2_L008_R1
_001.fastq.gz | sed -n '2~4p' | grep --color=always
'AGATCGGAAGAGCACACGTCTGAACTCCAGTCA' | awk '{print $0}' | head -n 20
ANGCCCCAAACCCAAACAACACACACACACACACACACACACACACACAAGATCGGAAGAGCACACG
TCTGAACTCCAGTCACCGATCGATAT
TGACTAGTGAAGTACCGGCTCTAGGCCATTAATGCCCAGGAGTGTGATGGCATTAGATCGGAAGAGCACAC
GTCTGAACTCCAGTCACCGATCGATA
GGATGATCAGCCATCCTTGATCAGCTTCTGATCTGCTGACGGGAGTTGGCATTGGCGATTTTCATTAGATCGGA
AGAGCACACGTCTGAACTCCAGTCAC
GTCCGGATAGCCGTGGCCTGTCACGTTGCGGAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCGATCGATAT
CTCGTATGCCGTCTTCTGCTTGA
TGCTTAAAGTCAGTTCGGACACGCCAGCCTTTATCATAAGCCAATGTGTGCCGGTCTTTCTGGAAGAGAGATCGG
AAGAGCACACGTCTGAACTCCAGTCA
CCGCCATGTGGTTGCTGGGATTTGAACTCCGGACCTTCGGAAGAGCAGTCGGGTGCTCTTACCCACTAGATCGGA
AGAGCACACGTCTGAACTCCAGTCAC
CGGCACTTCTGTCTCTGTTTCAGTTAATGGCCAGCCACAGGGGCACGTCCACGGAGATCGGAAGAGCACACGTCT
GAACTCCAGTCACCGATCGATATCT
GTAGAGTGAATGATCCCCCTGTGCTTGGTGACAGTATGGTGCATTCCATCCTTTTCTGAGAGATCGGAAGAGCA
CACGTCTGAACTCCAGTCACCGATCG
GCACCGATTTAACAACAGTTTTTCGAAAATTCACAGTTACTGTTGGCTTTTCTGTAGTGGAGATCGGAAGAGCACA
CGTCTGAACTCCAGTCACCGATCGAT
CGTGATCATTTCAAAATCATTCCCCTTCCGCTACTGTGTTGCGAGCGGTGAGATCGGAAGAGCACACGTCTGA
ACTCCAGTCACCGATCGATATCTCGT
CCCCGAGGAGCTCCTCACCCACAGCTTCTTCCAGCTTTATTGGTGTCTGATGGCCTTGGGAGATCGGAAGAGCA
CACGTCTGAACTCCAGTCACCGATCG
TGCTTGTGACAATGATGCCACGGCATGCTGGGTGACATTGTAGACTCTTCCGGTTTTGCAGATCGGAAGAGCA
CACGTCTGAACTCCAGTCACCGATCG
TGGGGTTGGAGTTTCCCTCAGCTTACACCATTTGTTTGGGGCAAGCAGATCTGAGAGTTCCAGATCGGAAGAGCA
CACGTCTGAACTCCAGTCACCGATCG
CACCAACTTACGAGCCACCTCTTCATACTTCTATCGGCCTCTTCTGCAATGTGCAGATCGGAAGAGCACACGTCT
GAACTCCAGTCACCGATCGATATCT
ATTCTTGTGAGTAGCCAGTTTGTGACAGTTCCAGTAGTGAAGTACACAGATCGGAAGAGCACACGTCTGAAC
TCCAGTCACCGATCGATATCTCGTAT
CGCTCCTCCTTGCACTGTTTCTGCTGTTCCACACCAGGACCCAGTCTGAGCTGTCATCTGTGAGATCGGAAGA
GCACACGTCTGAACTCCAGTCACCGA
CCTGCACTATAGGCACCAGACCTCATCTATCCAGACCTGCCTGTTAGATCGGAAGAGCACACGTCTGAACTCCA
GTCACCGATCGATATCTCGTATGCCG
GCAGCATATATGTTGGATTTTTAGGAAAGACCAATTCACAGCCCTCATGTGGGCTATAATTTTTAGATCGGAAGA
GCACACGTCTGAACTCCAGTCACCGA
CTGGGCTGGAACCCTGCGGTCTACTTGAGCAGGTTCTGCAGCATGGCCGTAGCATAGGAGATCGGAAGAGCACAC
GTCTGAACTCCAGTCACCGATCGATA
CCCCACTTTGTTCTCCACAGGCTCTGTGCTCTCCTGCCATCGCCGCTCTGCCTCTCCCCAGAGATCGGAAGAGC
ACACGTCTGAACTCCAGTCACCGATC
(fastqc) [apowers4@n226:/projects/bgmp/apowers4/Bi623/PS/ps3-QAA]
└─ zcat
/projects/bgmp/shared/2017_sequencing/demultiplexed/2_2B_control_S2_L008_R2
_001.fastq.gz | sed -n '2~4p' | grep --color=always
'AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT' | awk '{print $0}' | head -n 20
NGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTTGTTGGGTTGGGGCTTAGATCGGAAGAGCGTCGT
GTAGGGAAAGAGTGTATCGATCGGTG
NATGCCATCAACATCCTGGCGGGCATTAAATGGCCTAGAGGCCGGTCAGTCACTAGTCAAGATCGGAAGAGCGTCG
TGTAGGGAAAGAGTGTATCGATCGGT
```



```

NATGAAATCGCCAATGCCAACTCCCGTCAGCAGATCAGGAAGCTGATCAAGGATGGGCTGATCATCCAGATCGGA
AGAGCGTCGTGTAGGGAAAGAGTGTATCGATCGGTG
NCGCAACGTGACAGGCCACGGCTATCCGGACAGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGTATCGATCGGTG
TAGATCTCGGTGGTCGCCGTATCATT
NTCTTCCAGAAAGACCGGCACACATTGGCTTATGATAAAGGCTGGCGTGTCCGAACTGACTTTAAGCAAGATCGG
AAGAGCGTCGTGTAGGGAAAGAGTGT
NGTGGGTAAGAGCACCCGACTGCTCTTCCGAAGGTCCGGAGTTCAAATCCCAGCAACCACATGGCGGAGATCGGA
AGAGCGTCGTGTAGGGAAAGAGTGTATCGATCGGTGTA
NCGTGGACGTGCCCTGTGGCTGGCCATTAACCTGAAACAGAGACAGAAGTGCCGAGATCGGAAGAGCGTCGTGT
AGGGAAAGAGTGTATCGATCGGTGTA
CTCAGAAAAGGATGGAATGCACCATACTGTACCAAGCACAGGGGGATCATTACACTCTACAGATCGGAAGAGCG
TCGTGTAGGGAAAGAGTGTATCGATC
CCACTACAGAAAAGCCAACAGTAACTGTGAATTTTCGAAAAGTGTGTTAAATCGGTGCAGATCGGAAGAGCGTC
GTGTAGGGAAAGAGTGTATCGATCGG
CGACCGCTCGCAACACAGTAGCGGGAACGGGAATGATTTTGAAATGATCACGAGATCGGAAGAGCGTCGTGTAGG
GAAAGAGTGTATCGATCGGTGTAGAT
CCCAAGGCCATCAGACACCAATAAAGCTGGAAGAAGCTGTGGGGTGAGGAGCTCCTCGGGGAGATCGGAAGAGCG
TCGTGTAGGGAAAGAGTGTATCGATC
GACACCACAGCGACCTCAGAGAACAAAGAGCGGCTTCAACTTTGGAACCCTAGACACAAAGAGTGTGAGATCGGAA
GAGCGTCGTGTAGGGAAAGAGTGTAT
GTAAGAACGTGATGCCAAAAGAGGAGACGCCTGCTGAGGATGAAAGTGAAAAGATCGGAAGAGCGTCGTGTAGGG
AAAGAGTGTATCGATCGGTGTAGATC
GCAAAACCGGAAGAGTCTACAATGTCACCCAGCATGCCGTGGGCATCATTGTCAACAAGCAAGATCGGAAGAGCG
TCGTGTAGGGAAAGAGTGTATCGATC
GGAAGTCTCAGATCTGCTTGCCCCAAACAAATGGTGTAAGCTGAGGGAACTCCAACCCCAAGATCGGAAGAGCG
TCGTGTAGGGAAAGAGTGTATCGATC
GCACATTGCAGAAGAGGCCGATAGGAAGTATGAAGAGGTGGCTCGTAAGTTGGTGAGATCGGAAGAGCGTCGTGT
AGGGAAAGAGTGTATCGATCGGTGTA
GGCCCATTCATCATCTGCTTGTCTGCACTTCCACAGCCTTGCCACTGTCACTTTCATCACTGTAGATCGGAAG
AGCGTCGTGTAGGGAAAGAGTGTATC
GTGTGAATCAGTCACTACTGGAAGTGCACAACTGGCTACTGACAAGAATAGATCGGAAGAGCGTCGTGTAGGGA
AAGAGTGTATCGATCGGTGTAGATCT
CACAGATGACAGCTGCAGGACTGGGTCTGGTGTGGAACAGACGAAACAGTGCAAGGAGGAGCGAGATCGGAAGA
GCGTCGTGTAGGGAAAGAGTGTATCG
AACAGGCAGGTCTGGATAGGATGAGGTCTGGTGCCTATAGTGCAGGAGATCGGAAGAGCGTCGTGTAGGGAAAGA
GTGTATCGATCGGTGTAGATCTCGGT

```

cutadapt commands run:

```

#19_3F
cutadapt -a AGATCGGAAGAGCACACGTCTGAACTCCAGTCA -A
AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT -o trimmed_19_3F_1R.fastq.gz -p
trimmed_19_3F_2R.fastq.gz 19_3F_fox_S14_L008_R1_001_fastqc.zip
19_3F_fox_S14_L008_R2_001_fastqc.zip

#2_2B
cutadapt -a AGATCGGAAGAGCACACGTCTGAACTCCAGTCA -A
AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT -o trimmed_2_2B_1R.fastq.gz -p
trimmed_2_2B_2R.fastq.gz 2_2B_control_S2_L008_R1_001_fastqc.zip
2_2B_control_S2_L008_R2_001_fastqc.zip

```

## Question 6

Trimmomatic commands run

```
# 2_2B
trimmomatic PE output_2_2B/trimmed_2_2B_1R.fastq.gz
output_2_2B/trimmed_2_2B_2R.fastq.gz
trimmomatic_2_2B/trimmed_2_2B_1P.fastq.gz
trimmomatic_2_2B/trimmed_2_2B_1u.fastq.gz
trimmomatic_2_2B/trimmed_2_2B_2P.fastq.gz
trimmomatic_2_2B/trimmed_2_2B_2u.fastq.gz LEADING:3 TRAILING:3 MINLEN:35
SLIDINGWINDOW:5:15

# 19_3F
trimmomatic PE output_19_3F/trimmed_19_3F_1R.fastq.gz
output_19_3F/trimmed_19_3F_2R.fastq.gz
trimmomatic_19_3F/trimmed_19_3F_1P.fastq.gz
trimmomatic_19_3F/trimmed_19_3F_1u.fastq.gz
trimmomatic_19_3F/trimmed_19_3F_2P.fastq.gz
trimmomatic_19_3F/trimmed_19_3F_2u.fastq.gz LEADING:3 TRAILING:3 MINLEN:35
SLIDINGWINDOW:5:15
```

## Part 3

### Question 10

For 19\_3F reads: Sequences Mapped: 31075050 Sequences Unmapped: 1569764

For 2\_2B reads: Sequences Mapped: 58690 Sequences Unmapped: 11508686

### Question 12

I propose that these data are not strand-specific, because \_\_no\_feature in 19\_3F\_unstranded.txt = 1256405 out of 16322407 which is 7.69% \_\_no\_feature in 19\_3F\_sam\_stranded.txt = 14086570 out of 16322405 which is 86.30%

More then 3/4 of the data has no features if it is given the stranded comand which would lead me to believe that my data is unstranded.

You can't really draw any conclusions from the 2\_2B reads as they do not seem to align to the Mouse genome.

## Slurm Scripts

### Genome index

```
#!/bin/bash
#SBATCH --account=bgmp
#SBATCH --partition=bgmp
#SBATCH --job-name=star_genomeindex
#SBATCH --output=star-%j-sbatch.out
#SBATCH --nodes=1
#SBATCH --ntasks-per-node=1
#SBATCH --cpus-per-task=8
#SBATCH --time=2:00:00

# Activate the environment
conda activate fastqc

#Unzip the files
gunzip /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/mouse_genome/Mus_musculus.GRCm39.dna.primary_assembly.fa.gz
gunzip /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/mouse_genome/Mus_musculus.GRCm39.104.gtf.gz

# Run the actual database indexer STAR
/usr/bin/time -v STAR --runThreadN 8 \
--runMode genomeGenerate \
--genomeDir /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/mouse_genome/genome_mous_star_2.7.1a \
--genomeFastaFiles /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/mouse_genome/Mus_musculus.GRCm39.dna.primary_assembly.fa \
--sjdbGTFfile /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/mouse_genome/Mus_musculus.GRCm39.104.gtf \

# Rezip the files
zip /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/mouse_genome/Mus_musculus.GRCm39.dna.primary_assembly.fa
zip /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/mouse_genome/Mus_musculus.GRCm39.104.gtf
```

## Star Align

### Star Align 2\_2B

```
#!/bin/bash
#SBATCH --account=bgmp
#SBATCH --partition=bgmp
#SBATCH --job-name=star_align2
#SBATCH --output=staralign2-%j-sbatch.out
#SBATCH --nodes=1
#SBATCH --ntasks-per-node=1
#SBATCH --cpus-per-task=8
#SBATCH --time=10:00:00

# Activate the environment
```

```
conda activate fastqc

# Run the actual database indexer STAR
/usr/bin/time -v STAR --runThreadN 8 --runMode alignReads \
--outFilterMultimapNmax 3 \
--outSAMunmapped Within KeepPairs \
--alignIntronMax 1000000 --alignMatesGapMax 1000000 \
--readFilesCommand zcat \
--readFilesIn /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/trimmomatic_2_2B/trimmed_2_2B_1P.fastq.gz
/projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/trimmomatic_2_2B/trimmed_2_2B_2P.fastq.gz \
--genomeDir /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/mouse_genome/genome_mous_star_2.7.1a \
--outFileNamePrefix aligned_star_sam_2_2B
```

## Star Align 19\_3F

```
#!/bin/bash
#SBATCH --account=bgmp
#SBATCH --partition=bgmp
#SBATCH --job-name=star_align19
#SBATCH --output=staralign19-%j-sbatch.out
#SBATCH --nodes=1
#SBATCH --ntasks-per-node=1
#SBATCH --cpus-per-task=8
#SBATCH --time=10:00:00

# Activate the environment
conda activate fastqc

# Run the actual database indexer STAR
/usr/bin/time -v STAR --runThreadN 8 --runMode alignReads \
--outFilterMultimapNmax 3 \
--outSAMunmapped Within KeepPairs \
--alignIntronMax 1000000 --alignMatesGapMax 1000000 \
--readFilesCommand zcat \
--readFilesIn /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/trimmomatic_19_3F/trimmed_19_3F_1P.fastq.gz
/projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/trimmomatic_19_3F/trimmed_19_3F_2P.fastq.gz \
--genomeDir /projects/bgmp/apowers4/Bi623/PS/ps3-
QAA/mouse_genome/genome_mous_star_2.7.1a \
--outFileNamePrefix aligned_star_sam_19_3F
```

## HTSeq-count

### 2\_2B

```
# Unstranded
#!/bin/bash
#SBATCH --account=bgmp
#SBATCH --partition=bgmp
#SBATCH --job-name=count_2B_unstranded
#SBATCH --time=2:00:00

htseq-count --stranded=no aligned_star_sam_2_2BAligned.out.sam
Mus_musculus.GRCm39.104.gtf > htseq_outputs/2_2B_unstranded.txt

# Stranded
#!/bin/bash
#SBATCH --account=bgmp
#SBATCH --partition=bgmp
#SBATCH --job-name=count_2B_stranded
#SBATCH --time=2:00:00

htseq-count --stranded=yes aligned_star_sam_2_2BAligned.out.sam
Mus_musculus.GRCm39.104.gtf > htseq_outputs/2_2B_stranded.txt
```

## 19\_3F

```
# Unstranded
#!/bin/bash
#SBATCH --account=bgmp
#SBATCH --partition=bgmp
#SBATCH --job-name=count_19_unstranded
#SBATCH --time=2:00:00

htseq-count --stranded=no aligned_star_sam_19_3FAligned.out.sam
Mus_musculus.GRCm39.104.gtf > htseq_outputs/19_3F_unstranded.txt

# Stranded
#!/bin/bash
#SBATCH --account=bgmp
#SBATCH --partition=bgmp
#SBATCH --job-name=count_19_stranded
#SBATCH --time=2:00:00

htseq-count --stranded=yes aligned_star_sam_19_3FAligned.out.sam
Mus_musculus.GRCm39.104.gtf > htseq_outputs/19_3F_sam_stranded.txt
```