# Bioinformatics Lessons Schedule

- RNA-seq

- single cell RNA-seq

- RRBS

| Date | Subject |
|------|---------|
| 12-24 | Christmas break |
| 12-31 | Christmas break |
| 01-07 | Process RNA-seq |
| 01-14 | Process RNA-seq, continued |
| 01-21 | Process RNA-seq, continued |
| 01-28 | Analyze RNA-seq |
| 01-28 | Analyze RNA-seq, continued |

# RNA-seq

# Quality Check

# FastQC

- Before going forward, we want to check the quality of the data
  - How much did the sequencer fail?
  - Did we sequence mostly our sample DNA?
- FastQC is a program from the Babraham Institute in the UK that creates an html report on the quality of the sequencing data
  - Has 11 quality control checks that it does

# Basic Statistics

## Good Quality

✅ **Basic Statistics**

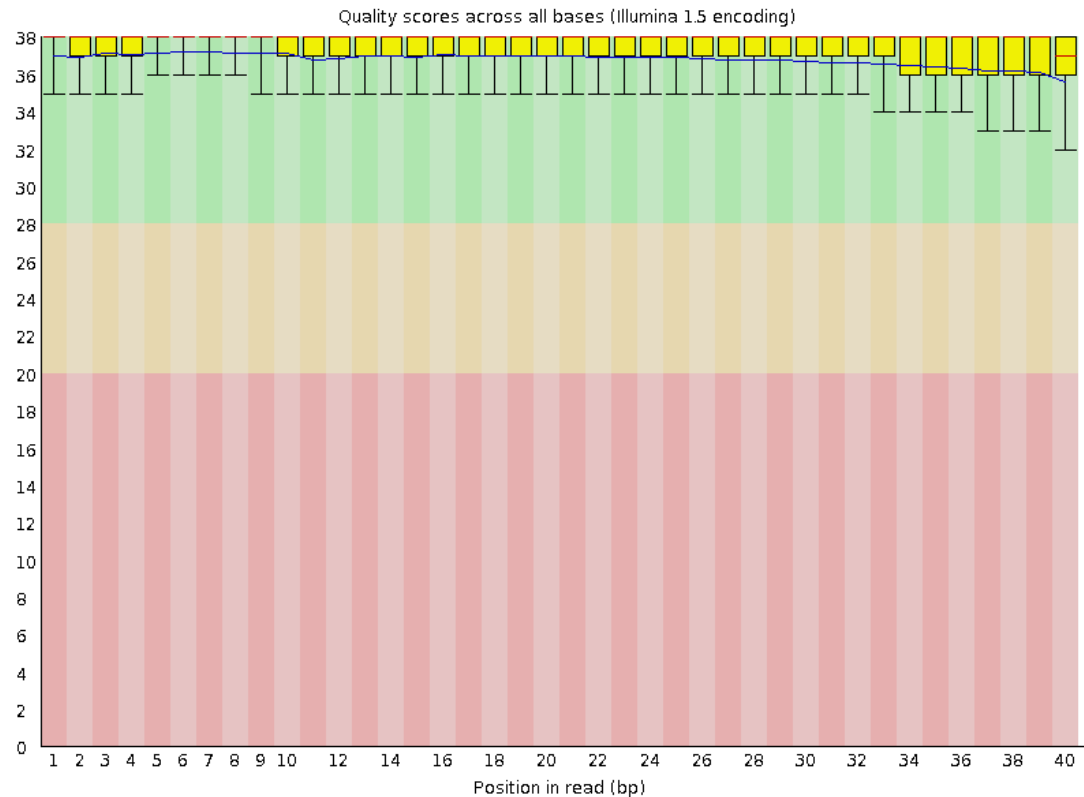| Measure | Value |
|---|---|
| Filename | good_sequence_short.txt |
| File type | Conventional base calls |
| Encoding | Illumina 1.5 |
| Total Sequences | 250000 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 40 |
| %GC | 45 |

## Bad Quality

✅ **Basic Statistics**

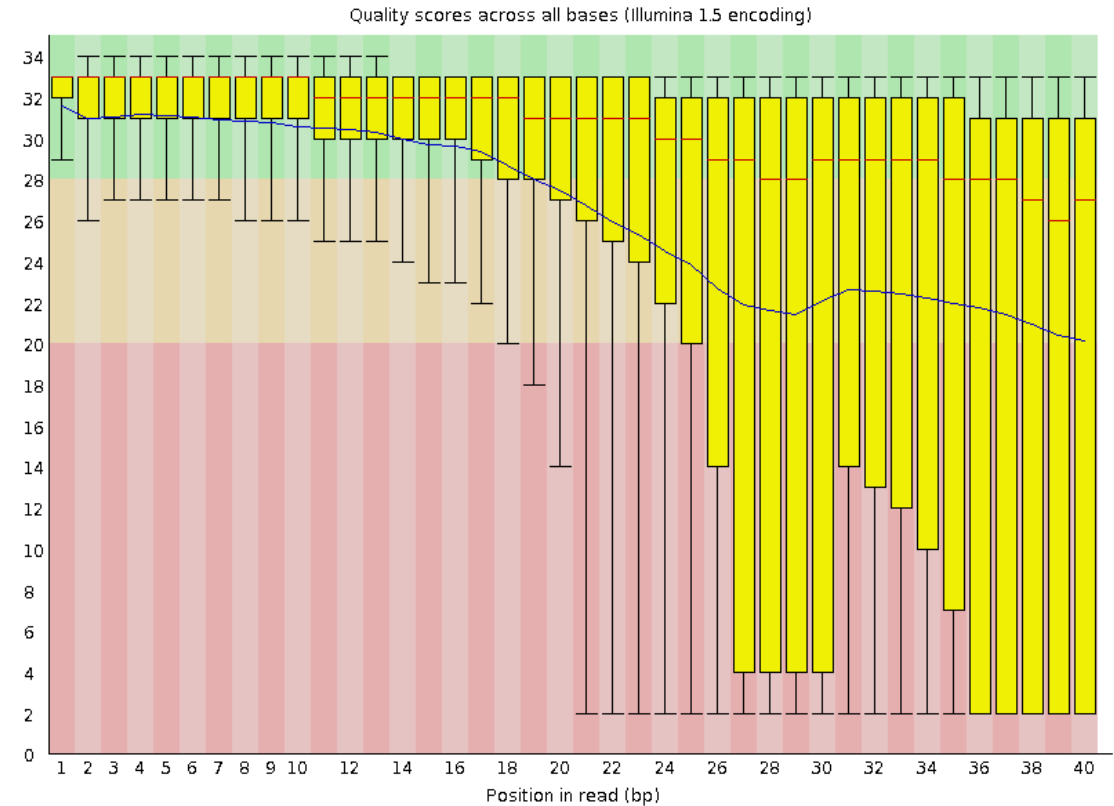| Measure | Value |
|---|---|
| Filename | bad_sequence.txt |
| File type | Conventional base calls |
| Encoding | Illumina 1.5 |
| Total Sequences | 395288 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 40 |
| %GC | 47 |

# Per base sequence quality

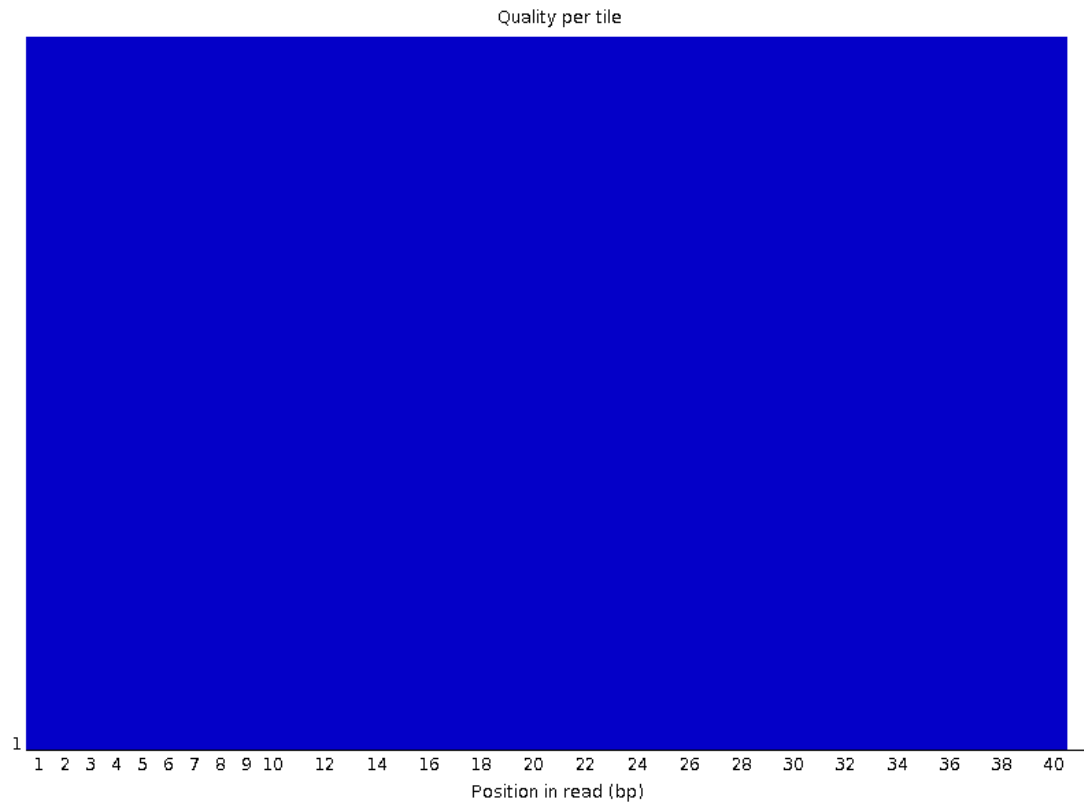Good Quality

Bad Quality

# Per tile sequence quality

## Good Quality



## Bad Quality

# Per sequence quality scores

## Good Quality



## Bad Quality

# Per base sequence content

Good Quality

Bad Quality

# Per sequence GC content

## Good Quality

## Bad Quality

# Per base sequence quality

Good Quality

Bad Quality

# Per base sequence quality

## Good Quality

## Bad Quality

# Sequence Duplication Levels



Good Quality

Bad Quality

# Overrepresented sequences

## Good Quality



## Bad Quality

# Adapter Content

## Good Quality



## Bad Quality

# Sidenote - Loops in Bash

# What's a loop?

# What's a loop?

A loop is a piece of code that repeatedly executes a command for a certain condition until that condition is no longer true.

# What's a loop?

A loop is a piece of code that repeatedly executes a command for a certain condition until that condition is no longer true.

```
for i in *.fastq.gz;
   do echo $i;
done
```

# What's a loop?

A loop is a piece of code that repeatedly executes a command for a certain condition until that condition is no longer true.

```
for i in *.fastq.gz;
   do echo $i;
done
```

condition

# What's a loop?

A loop is a piece of code that repeatedly executes a command for a certain condition until that condition is no longer true.

```
for i in *.fastq.gz;          ← condition
   do echo $i;                ← command
done
```

# What's a loop?

A loop is a piece of code that repeatedly executes a command for a certain condition until that condition is no longer true.

```
for i in *.fastq.gz;          ← condition
   do echo $i;                ← command
done                          ← stop
```

# How do you write a loop in bash?

How do you write a loop in bash?

```
for i in *.fastq.gz;
  do echo $i;
done
```

# How do you write a loop in bash?

starting
a loop

```
for i in *.fastq.gz;
   do echo $i;
done
```

# How do you write a loop in bash?

starting a loop

variable name

```
for i in *.fastq.gz;
    do echo $i;
done
```

# How do you write a loop in bash?

| starting a loop | variable name | | meeting this condition |
|---|---|---|---|

```
for i in *.fastq.gz;
   do echo $i;
done
```

# How do you write a loop in bash?

starting a loop

variable name

meeting this condition

```
for i in *.fastq.gz;
   do echo $i;
done
```

execute

# How do you write a loop in bash?

starting a loop

variable name

meeting this condition

execute

command

```
for i in *.fastq.gz;
    do echo $i;
done
```

# How do you write a loop in bash?

starting a loop

variable name

meeting this condition

```
for i in *.fastq.gz;
   do echo $i;
done
```

execute

command

variable

# How do you write a loop in bash?

starting a loop | variable name | meeting this condition

```
for i in *.fastq.gz;
    do echo $i;
done
```

execute | command | variable

stop

# How do you write a loop in bash?

starting a loop

variable name

meeting this condition

```
for i in *.fastq.gz;
  do echo $i;
done
```

execute

command

variable

stop

Now go try this loop on the server

# Run FastQC

1. Go to the RNA-seq data directory

2. Make a directory to put the FastQC reports into, `fastqc`

3. Run fastqc on the samples

```
for i in *.fastq.gz; do fastqc $i -o fastqc/;
done
```

# Trim Bad Quality Sequences

# What is trimming and why do it?

# What is trimming and why do it?

- Trimming removes sequencing adapters, bad quality sequences, and/or other biased sequence information

# What is trimming and why do it?

- Trimming removes sequencing adapters, bad quality sequences, and/or other biased sequence information

- Why is that important?
  - Helps prevent incorrect base calls by removing poor quality information
  - Increases speed and accuracy of alignment by removing artificial sequences and low quality sequences

# What is trimming and why do it?

- Trimming removes sequencing adapters, bad quality sequences, and/or other biased sequence information
- Why is that important?
  - Helps prevent incorrect base calls by removing poor quality information
  - Increases speed and accuracy of alignment by removing artificial sequences and low quality sequences

- Trimming does two complementary things:
  1. Removes any sequence information that comes from library preparation or sequencing
  2. Removes low quality bases / low quality reads

# Trim Sequences

1. Go back up to the rnaseq directory

2. Make a folder to put the analysis results in, `analysis`

3. Make a folder inside the analysis folder to put the trimmed reads in, `analysis/01_trim`

# Trim Sequences

```
for i in *R1.fastq.gz;
    do trim_galore
        --paired
        --fastqc
        --illumina
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

# Trim Sequences

```
for i in *R1.fastq.gz;
    do trim_galore
        --paired
        --fastqc
        --illumina
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

loop condition

# Trim Sequences

```
for i in *R1.fastq.gz;                          loop condition
    do trim_galore                              call the program
            --paired
            --fastqc
            --illumina
            --output analysis/01_trim/
            --retain_unpaired
            -q 30
            $i
            ${i/R1/R2};
done
```

# Trim Sequences

```
for i in *R1.fastq.gz;                          ← loop condition
    do trim_galore                              ← call the program
        --paired                                ← reads are paired-end
        --fastqc
        --illumina
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

# Trim Sequences

```
for i in *R1.fastq.gz;                          ← loop condition
    do trim_galore                              ← call the program
        --paired                                ← reads are paired-end
        --fastqc                                ← run FastQC again after trimming
        --illumina
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

# Trim Sequences

```
for i in *R1.fastq.gz;                    ← loop condition
    do trim_galore                        ← call the program
        --paired                          ← reads are paired-end
        --fastqc                          ← run FastQC again after trimming
        --illumina                        ← trim Illumina adapters
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

# Trim Sequences

```
for i in *R1.fastq.gz;                    ← loop condition
    do trim_galore                        ← call the program
        --paired                          ← reads are paired-end
        --fastqc                          ← run FastQC again after trimming
        --illumina                        ← trim Illumina adapters
        --output analysis/01_trim/        ← output goes here
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

# Trim Sequences

```
for i in *R1.fastq.gz;
    do trim_galore
        --paired
        --fastqc
        --illumina
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

loop condition

call the program

reads are paired-end

run FastQC again after trimming

trim Illumina adapters

output goes here

keep reads where one mate fails trimming but the other doesn't

# Trim Sequences

```
for i in *R1.fastq.gz;
    do trim_galore
        --paired
        --fastqc
        --illumina
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

loop condition

call the program

reads are paired-end

run FastQC again after trimming

trim Illumina adapters

output goes here

keep reads where one mate fails trimming but the other doesn't

read files

# Trim Sequences

```
for i in *R1.fastq.gz;
    do trim_galore
        --paired
        --fastqc
        --illumina
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

loop condition

call the program

reads are paired-end

run FastQC again after trimming

trim Illumina adapters

output goes here

keep reads where one mate fails trimming but the other doesn't

Keep bases at this quality or above

# Trim Sequences

```
for i in *R1.fastq.gz;
    do trim_galore
        --paired
        --fastqc
        --illumina
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

loop condition

call the program

reads are paired-end

run FastQC again after trimming

trim Illumina adapters

output goes here

keep reads where one mate fails trimming but the other doesn't

Keep bases at this quality or above

By default bases quality less than 20 will be trimmed and if the read falls below 20 bp, it will be discarded; we set the minimum quality to be 30

# Trim Sequences

```
for i in *R1.fastq.gz;
    do trim_galore
        --paired
        --fastqc
        --illumina
        --output analysis/01_trim/
        --retain_unpaired
        -q 30
        $i
        ${i/R1/R2};
done
```

By default bases quality less than 20 will be trimmed and if the read falls below 20 bp, it will be discarded; we set the minimum quality to be 30

loop condition

call the program

reads are paired-end

run FastQC again after trimming

trim Illumina adapters

output goes here

keep reads where one mate fails trimming but the other doesn't

Keep bases at this quality or above

read files

# Trim Command

```
for i in *R1.fastq.gz; do trim_galore
--paired --fastqc --illumina --output analysis/01_trim/
--retain_unpaired -q 30 $i ${i/R1/R2}; done
```