

NGS Analysis and Galaxy

Part 1

Kathryn Crouch
kathryn.crouch@glasgow.ac.uk

Introduction



University
of Glasgow

Who am I?

- PhD in comparative immunology 2005
- Worked in industry (big pharma and small biotech)
- MSc bioinformatics 2013
- Core bioinformatician WCIP 2013 - 2021
- Create online resources for public use



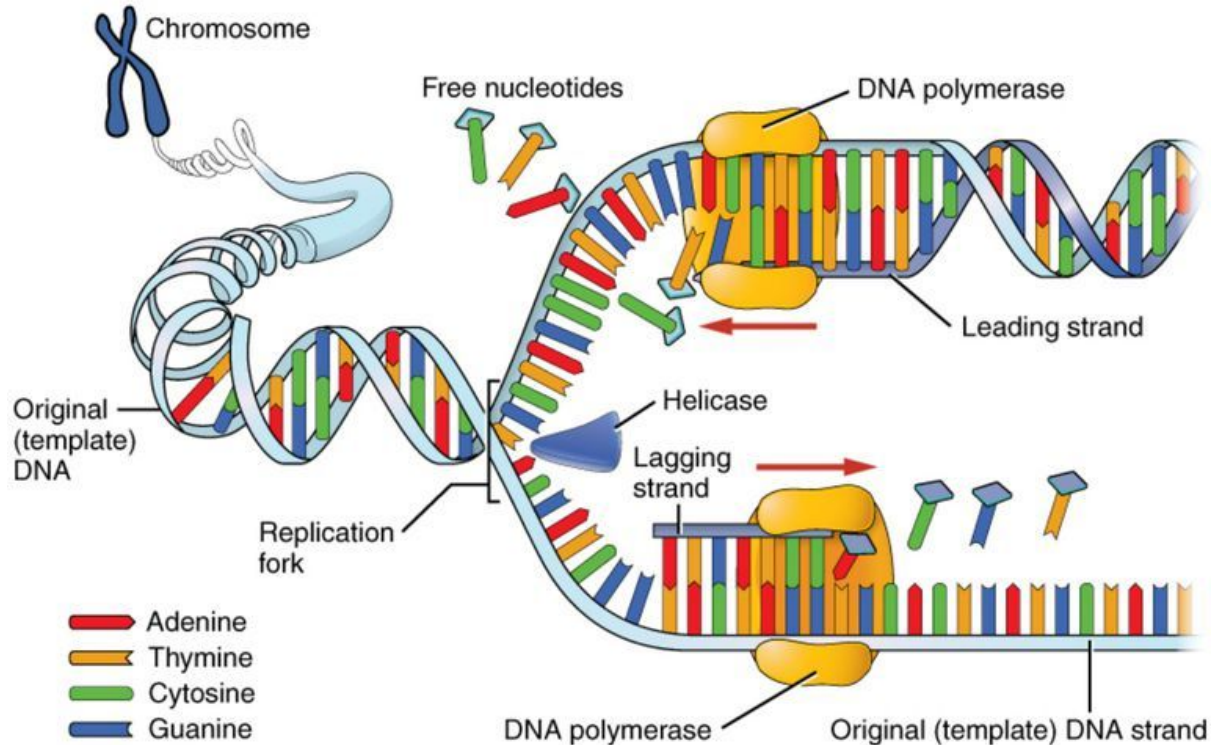
COMPANION

Outline

- Background
 - How does sequencing work?
 - What does the data look like?
- Pre-processing for Analysis
 - QC and trimming
 - Aligning to a reference genome
- RNA-seq and differential expression
- Using Galaxy

How Does Sequencing Work?

- Carry out replication under controlled conditions
- Artificially slow down the reaction to see the order in which bases are incorporated



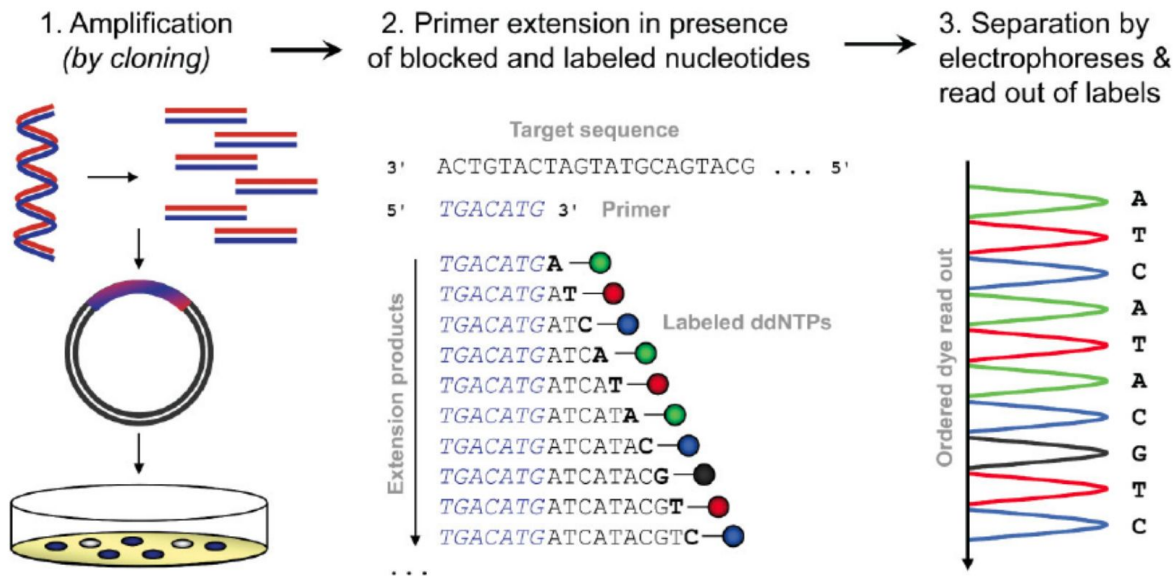
How Does Sequencing Work?

Sanger Sequencing (1975)

Uses modified nucleotides (ddNTPs) that cannot be extended

Each ddNTP is labelled with a different dye so you can see the order in which they are incorporated

Long reads and low error rate, but low throughput



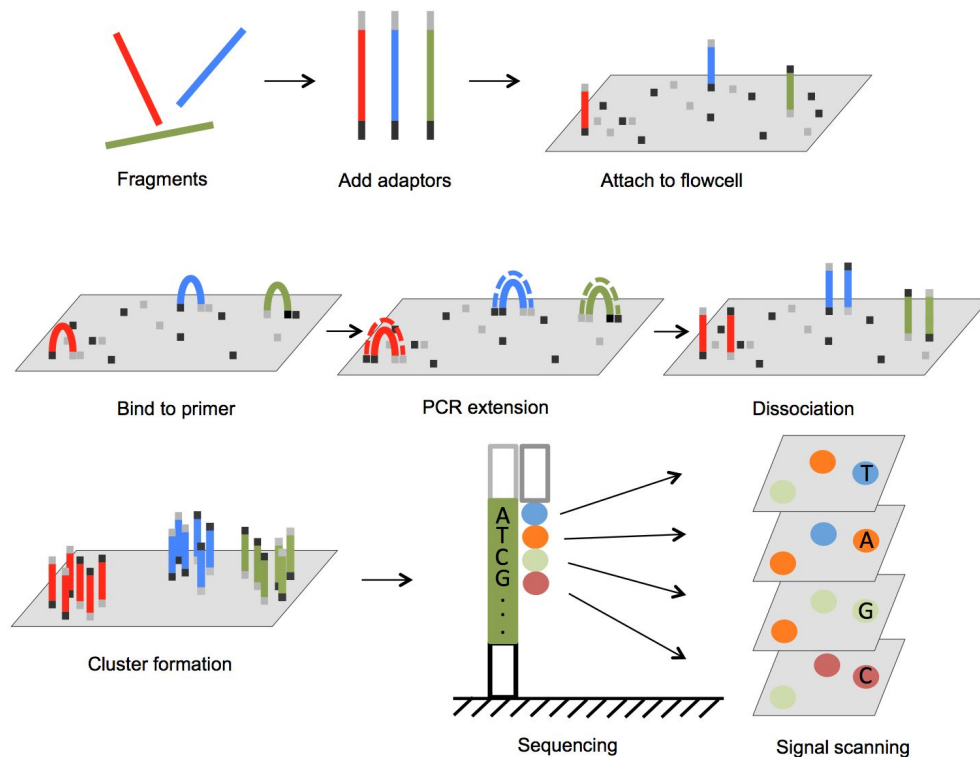
How Does Sequencing Work?

Illumina Sequencing (2005)

DNA fragments are
adaptor-ligated and attached to
a flow cell

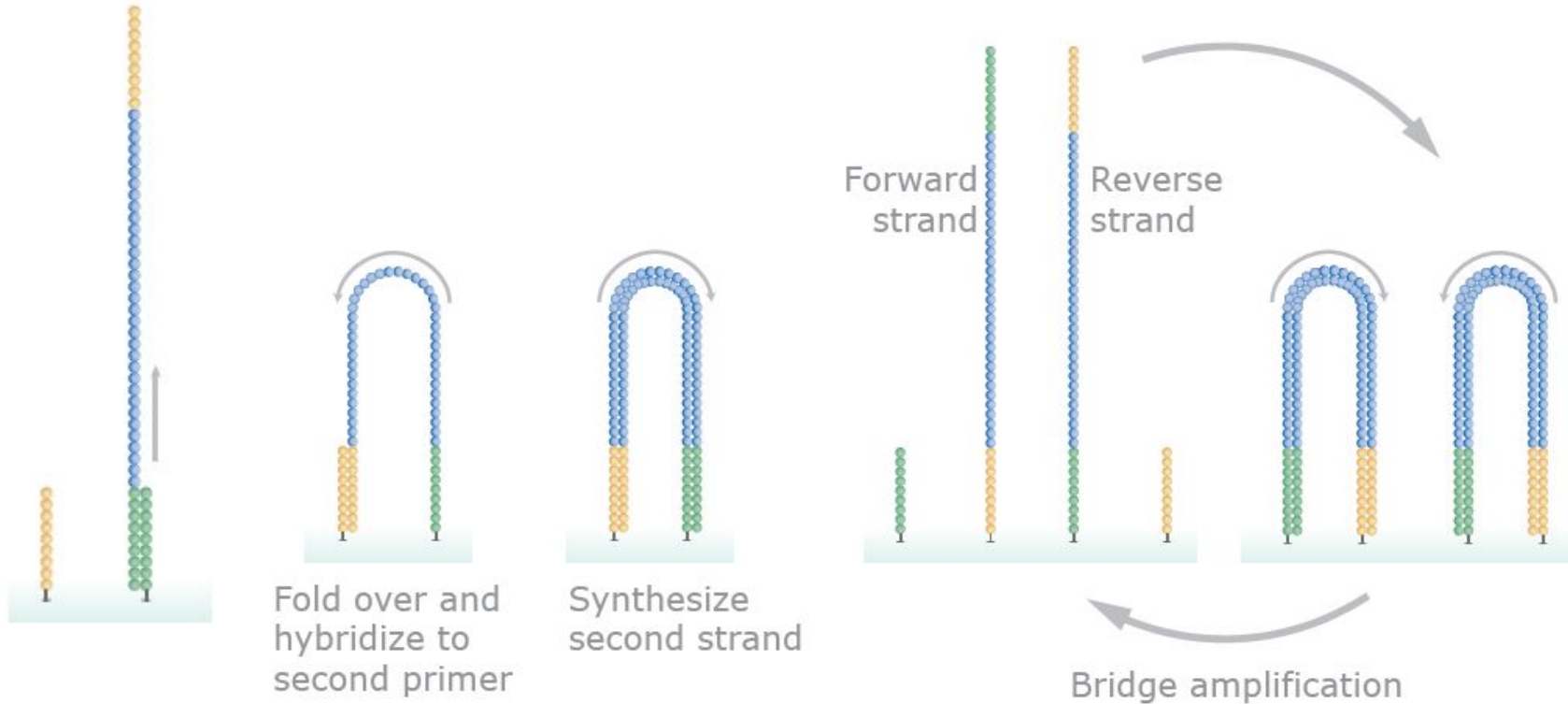
PCR is carried out in situ to form
clusters

Sequencing can be carried out
on millions of clusters
simultaneously

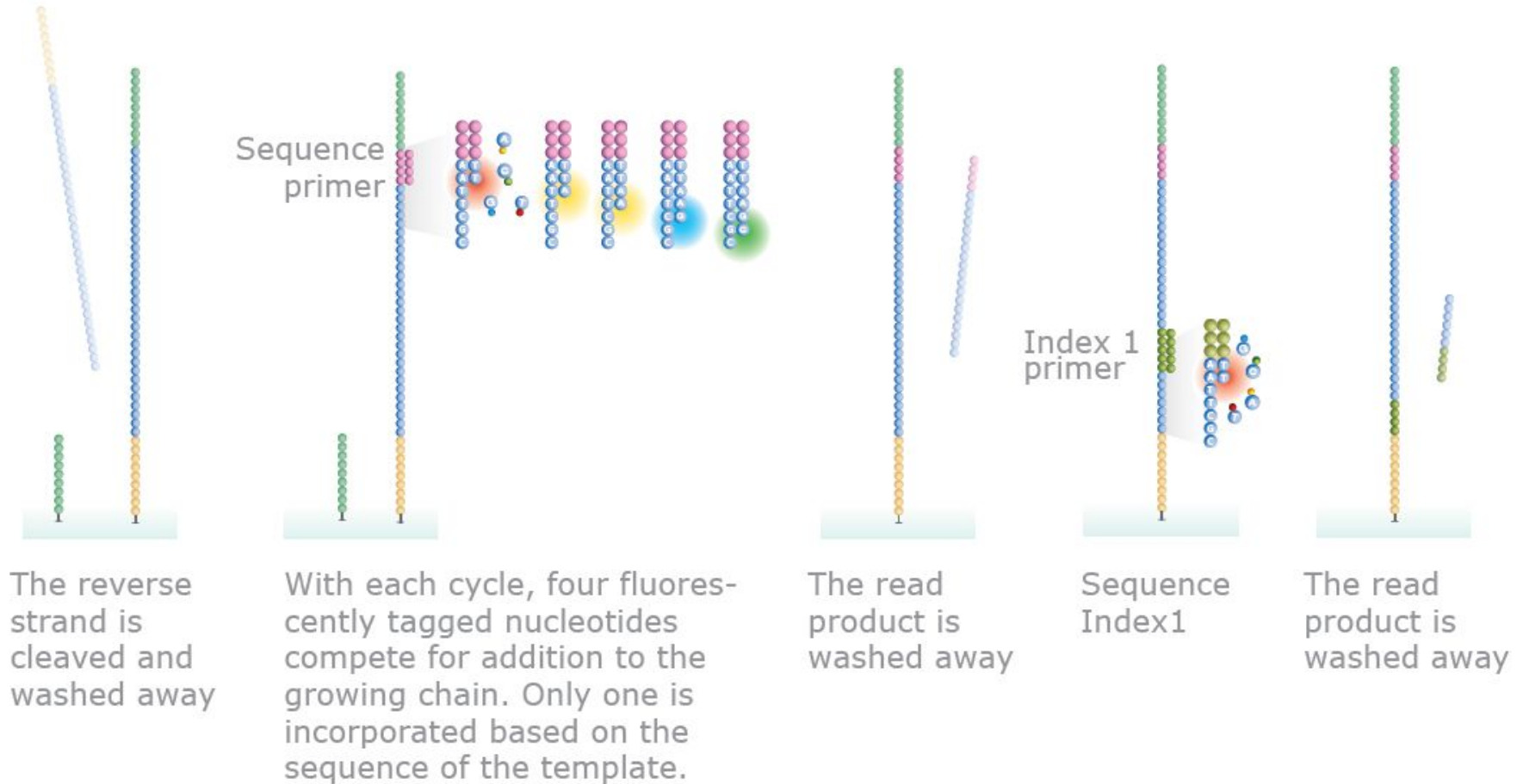


illumina®

How Does Sequencing Work?



How Does Sequencing Work?



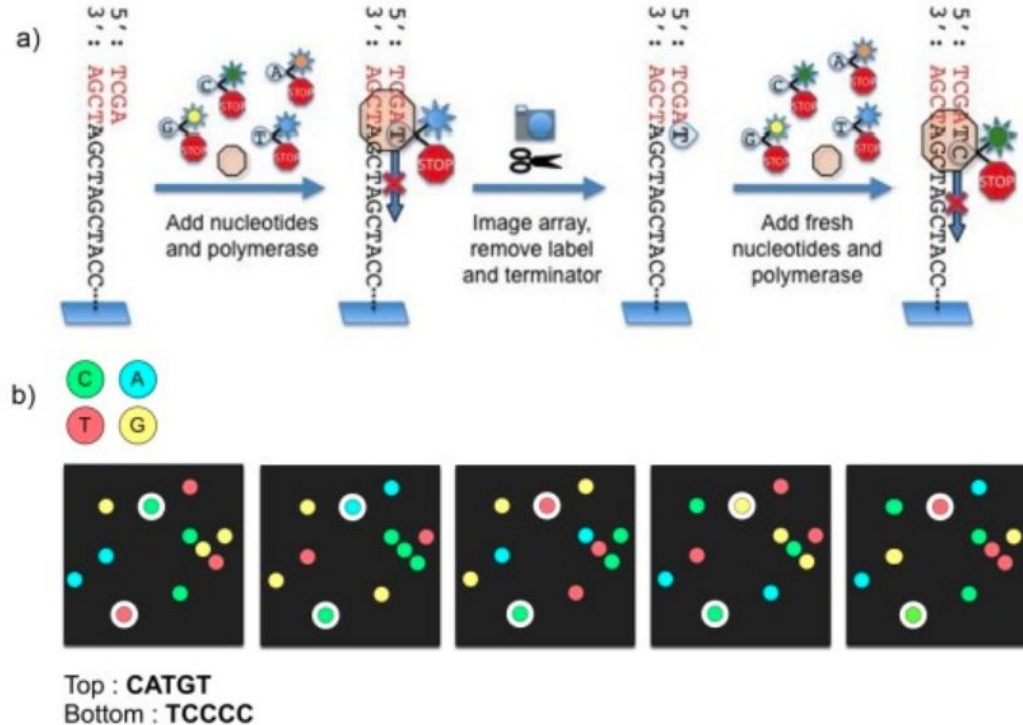
How Does Sequencing Work?

Illumina Sequencing

Blocked and labelled nucleotides are added

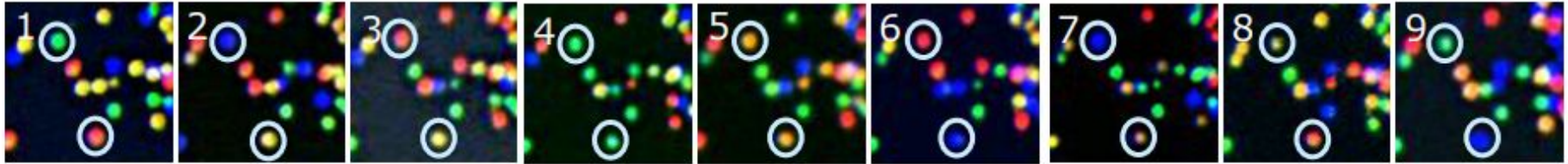
1 nucleotide is incorporated and an image is taken of the array

Label and block are removed and cycle repeats



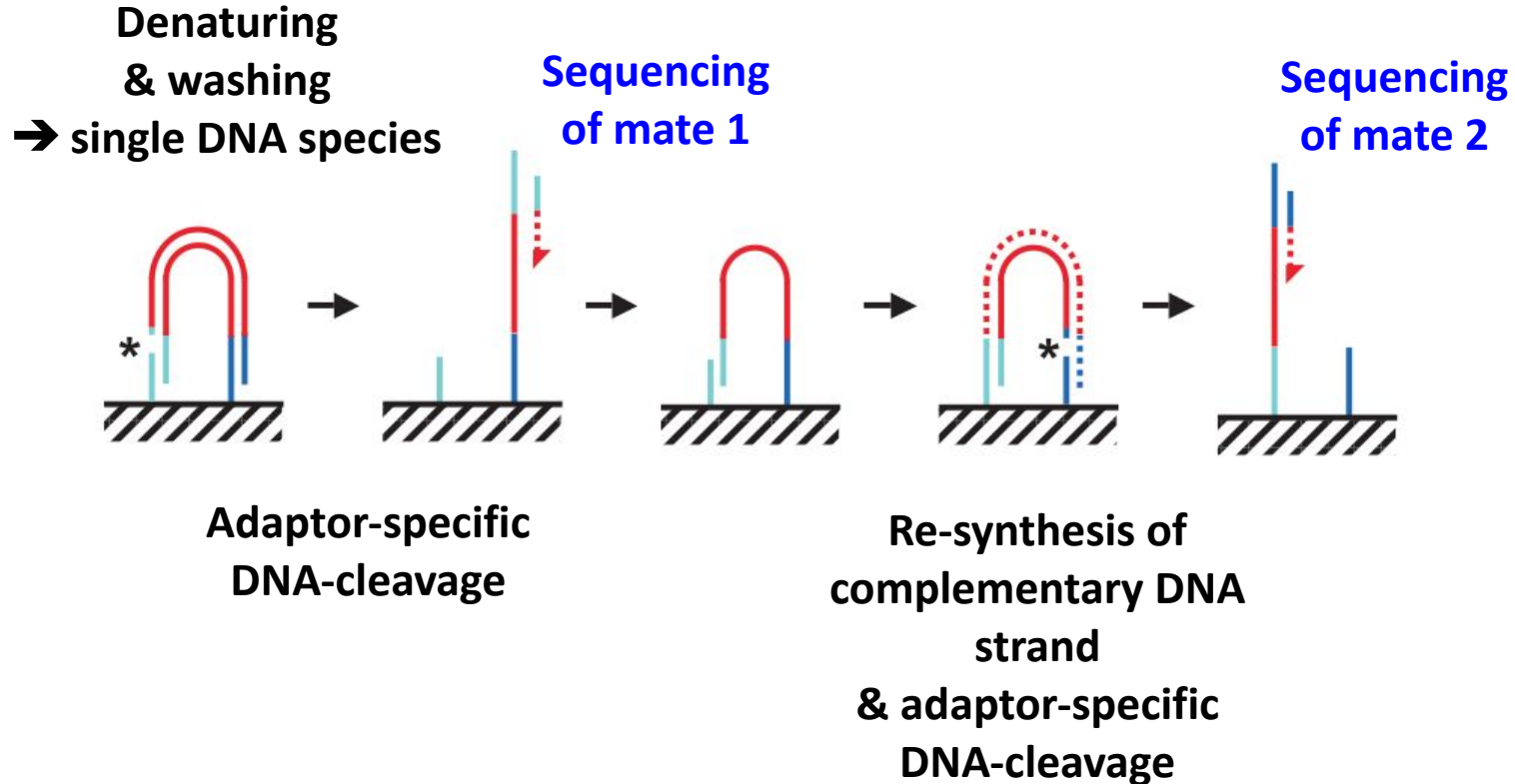
How Does Sequencing Work?

T G C T A C G A T



C A A T A G A C G

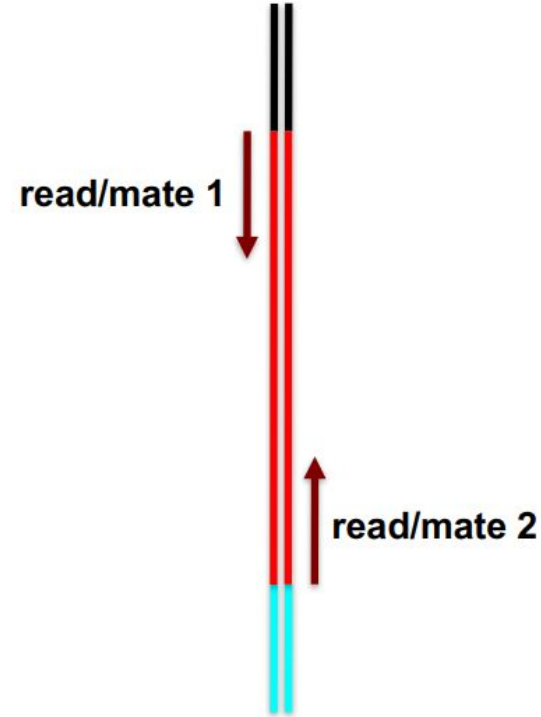
How Does Sequencing Work?



Paired End Illumina Sequencing

A short read is sequenced from
each end of each fragment

Blue and black are adaptors

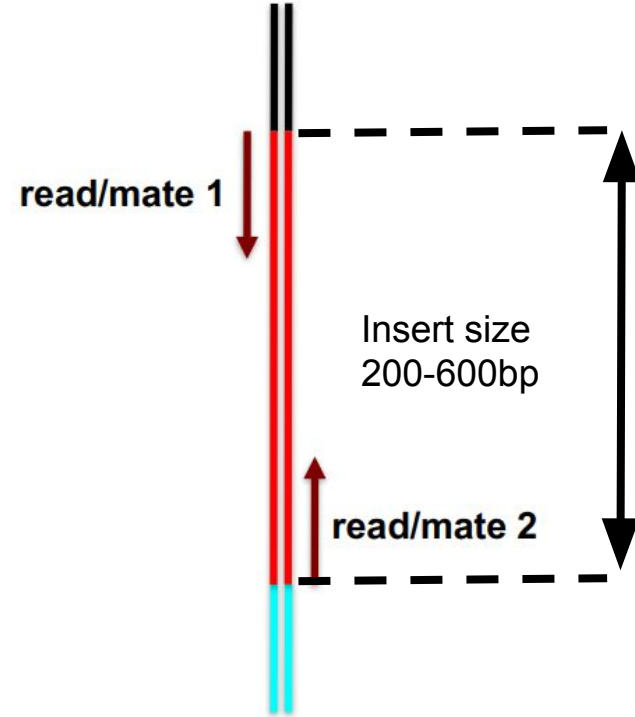


Paired End Sequencing

A short read is sequenced from each end of each fragment

Blue and black are adaptors

Insert or fragment size is the distance between adaptors



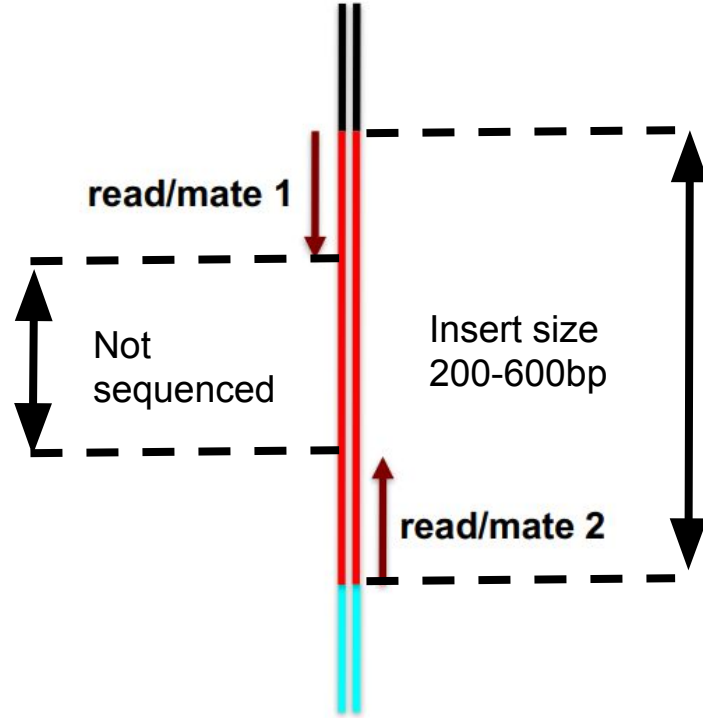
Paired End Sequencing

A short read is sequenced from each end of each fragment

Blue and black are adaptors

Insert or fragment size is the distance between adaptors

The region between the reads is not sequenced - it is covered by other fragments



Paired End Sequencing

A short read is sequenced from each end of each fragment

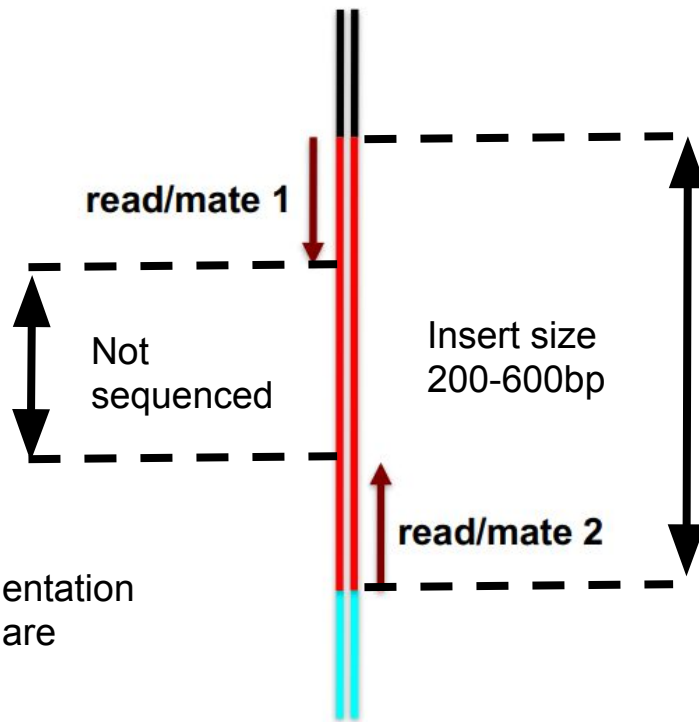
Blue and black are adaptors

Insert or fragment size is the distance between adaptors

The region between the reads is not sequenced - it is covered by other fragments

Reads that map in the correct orientation and the expected distance apart are “concordant” or “proper pairs”

Concordant alignments are prioritised



File Formats: FASTQ

```
@NS500205:127:HW272BGXX:1:11101:7788:1040 1:N:0:TCTCGCGC+AGGCTATA
```

For paired-end reads you will have two files

4 lines per read

- Line 1 is a unique header (this will be shared between the pairs)

File Formats: FASTQ

```
@NS500205:127:HW272BGXX:1:11101:7788:1040 1:N:0:TCTCGCGC+AGGCTATA  
CCTTGNTCCGTCATATTTTTTAGCATTGCAATGACGCTAAGTCCCGATTGACGCGCACGTGCTCACCCGGTTTCC
```

For paired-end reads you will have two files

4 lines per read

- Line 1 is a unique header (this will be shared between the pairs)
- Line 2 is the sequence of the read

File Formats: FASTQ

```
@NS500205:127:HW272BGXX:1:11101:7788:1040 1:N:0:TCTCGCGC+AGGCTATA
CCTTGNTCCGTCATATTTTTTAGCATTGCAATGACGCTAAGTCCCGATTGACGCGCACGTGCTCACCCGGTTTCC
+
AAAAA#EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE
```

For paired-end reads you will have two files

4 lines per read

- Line 1 is a unique header (this will be shared between the pairs)
- Line 2 is the sequence of the read
- Line 4 is the quality for each base
 - Quality is encoded using ASCII
 - <http://www.asciitable.com/>

Phred quality scores are logarithmically linked to error probabilities

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%

Alignment-based Analysis

Quality
Control



Read
Trimming



Alignment

Analysis of NGS Sequencing

Quality
Control



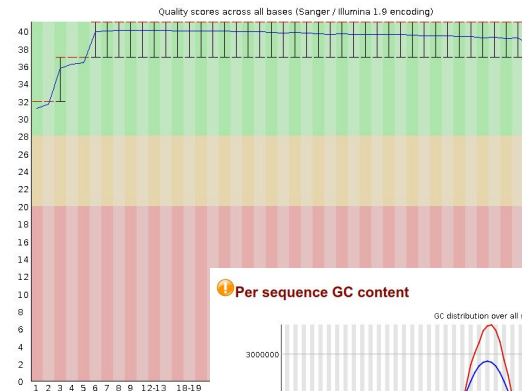
Read
Trimming



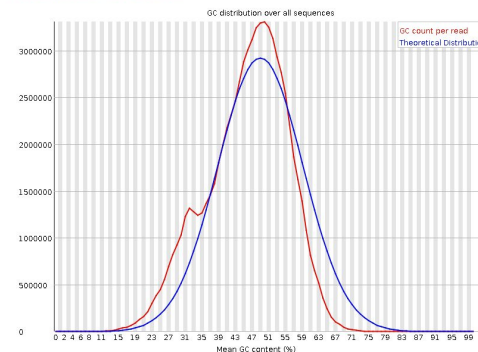
Alignment

- FASTQC
<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
 - Overall sequencing quality
 - GC content
 - N content
 - Read length distribution
 - Over-represented sequences
 - Adaptor content
- Output is an html file that can be opened in a web browser

✓ Per base sequence quality



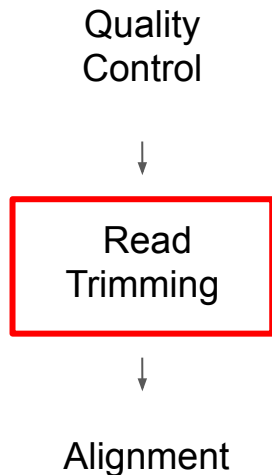
! Per sequence GC content



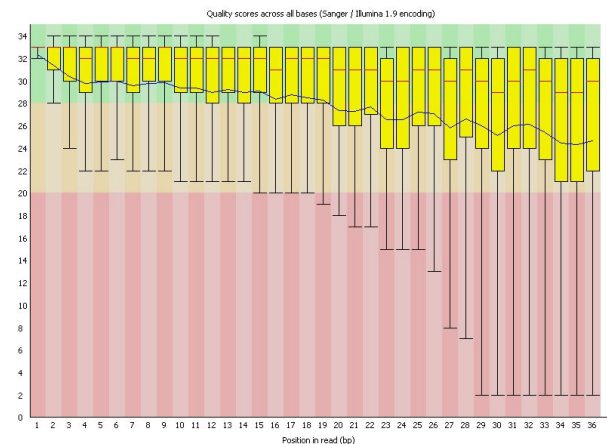
! Overrepresented sequences

Sequence	Count	Percentage	Possible Source
ACAAGTGTGTAAACATTAATTTGCAAGTTTGCAACGCTGTTCTTTAGTGTT	70896	0.12562741276052788	No Hit

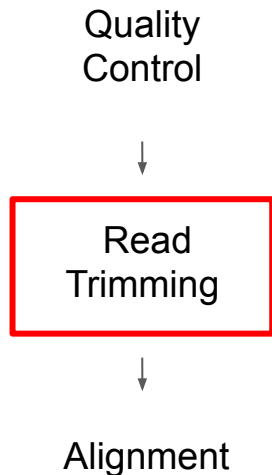
Analysis of NGS Sequencing



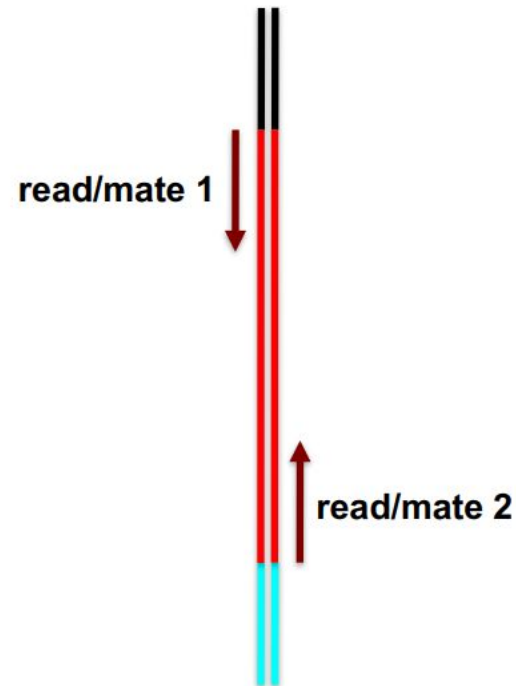
- Trimmomatic
<https://github.com/usadellab/Trimmomatic>
- Sequence quality tends to decrease towards the 3' end of reads - these can affect mapping
- Sickle will:
 - Remove poor quality reads from the 3' end of each read
 - Check for reads that are too short and discard them
 - Check that all reads still have a pair and discard those that don't
- Some trimming tools can also remove adaptors



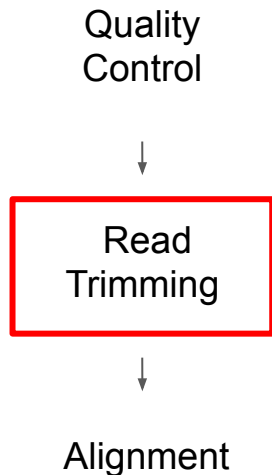
Analysis of NGS Sequencing



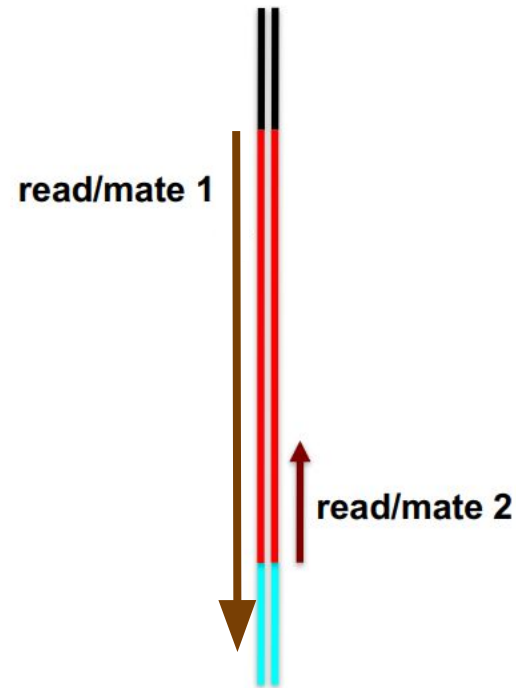
- Trimmomatic
<https://github.com/usadellab/Trimmomatic>
- Sequence quality tends to decrease towards the 3' end of reads - these can affect mapping
- Sickle will:
 - Remove poor quality reads from the 3' end of each read
 - Check for reads that are too short and discard them
 - Check that all reads still have a pair and discard those that don't
- Some trimming tools can also remove adaptors



Analysis of NGS Sequencing



- Trimmomatic
<https://github.com/usadellab/Trimmomatic>
- Sequence quality tends to decrease towards the 3' end of reads - these can affect mapping
- Sickle will:
 - Remove poor quality reads from the 3' end of each read
 - Check for reads that are too short and discard them
 - Check that all reads still have a pair and discard those that don't
- Some trimming tools can also remove adaptors



Analysis of NGS Sequencing

Quality
Control



Read
Trimming



Alignment

What is an alignment?

Two sequences:

ATTGAAAGCTA
GAAATGAAAAGG

How would you align one to the other?

Analysis of NGS Sequencing

Quality
Control



Read
Trimming



Alignment

What is an alignment?

Two sequences:

ATTGAAAGCTA
GAAATGAAAAGG

How would you align one to the other?

```
--ATTGAAA-GCTA
  | | | | | |
GAAATGAAAAGG
```

```
ATTGAAA-GCTA---
  | | | | |
---GAAATGAAAAGG
```

Which one is better??

Analysis of NGS Sequencing

Quality
Control



Read
Trimming



Alignment

What is an alignment?

Two sequences:

ATTGAAAGCTA
GAAATGAAAAGG

How would you align one to the other?

```
--ATTGAAA-GCTA
  | | | | | |
GAAATGAAAAGG--
```

```
ATTGAAA-GCTA---
  | | | | |
---GAAATGAAAAGG
```

Which one is better??

Alignment scoring:

- 1 for a match
- -1 for a mismatch
- -2 for a gap

Now alignment 1 scores -4 and alignment 2 scores -10 so we would choose alignment 1

Alignment algorithms are designed to align data in a reasonable time on a standard computer

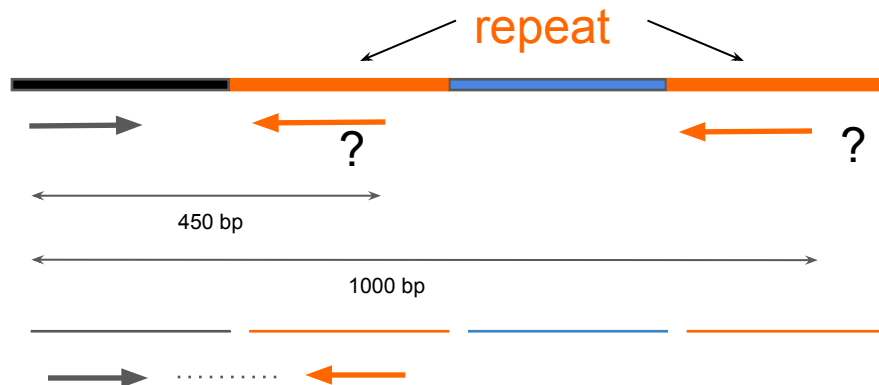
Heuristic - not exhaustive, but “good enough”

Improving Alignments with Paired End Reads

Paired end reads can resolve alignments in repetitive regions

A read that could map in multiple locations due to repetitions in the genome can be located accurately by inference from the position of its pair

In this case, we know that the insert size is ~400 bp, so we can infer that the first alignment is more likely to be correct.



Aligning RNA-seq Data

Quality
Control



Read
Trimming

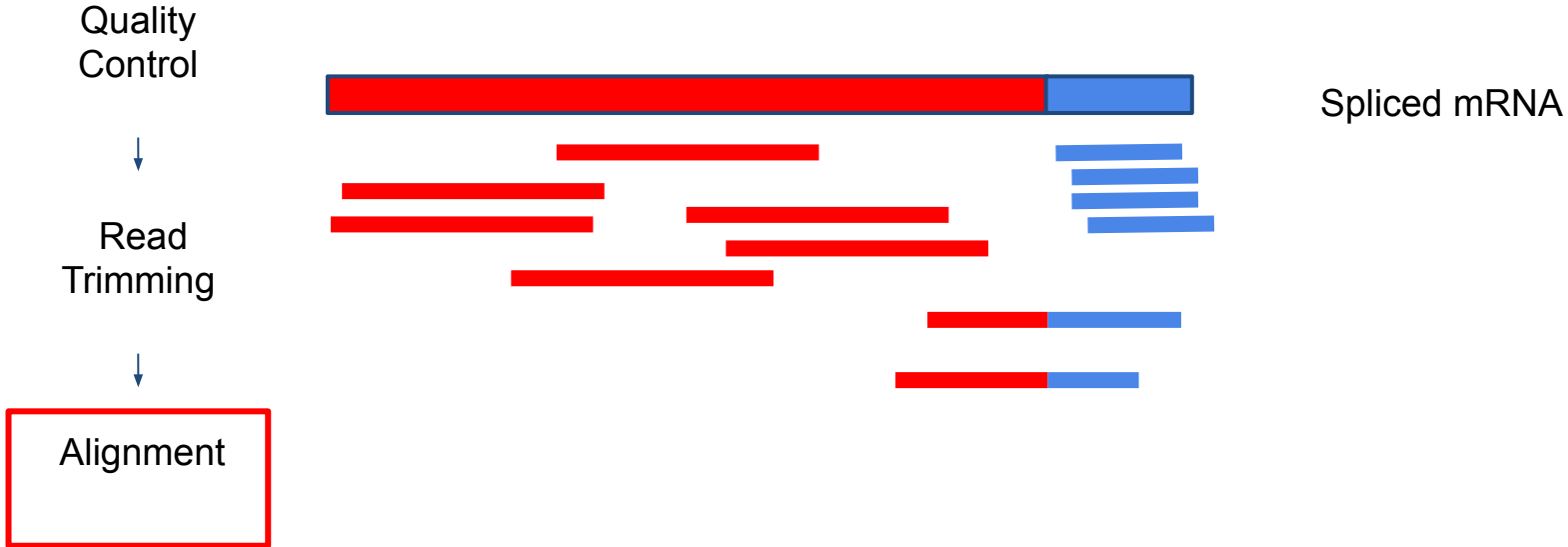


Alignment

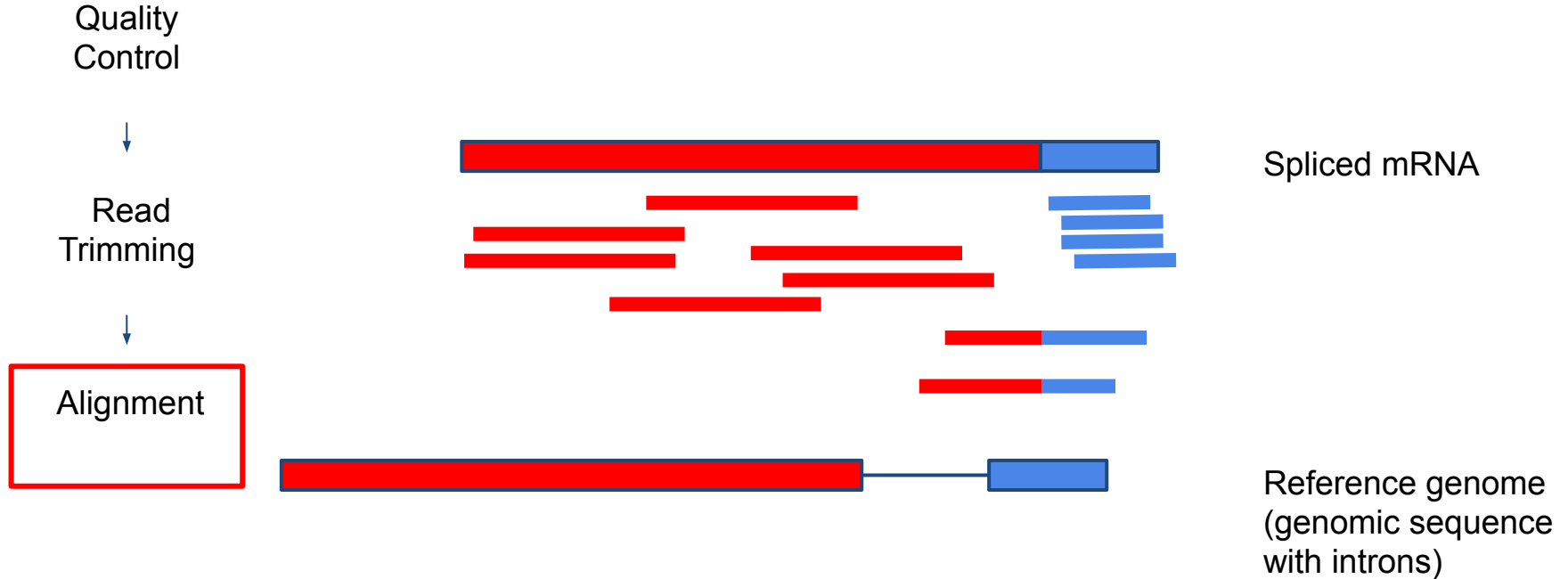


Spliced mRNA

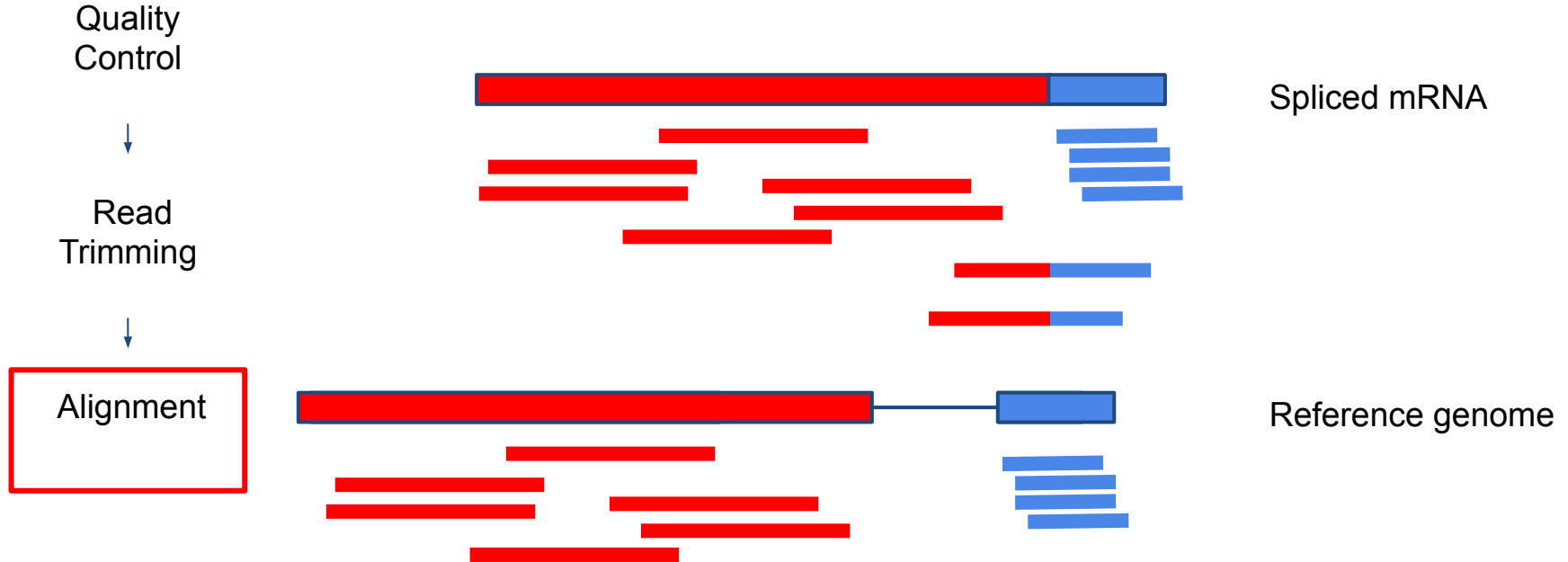
Aligning RNA-seq Data



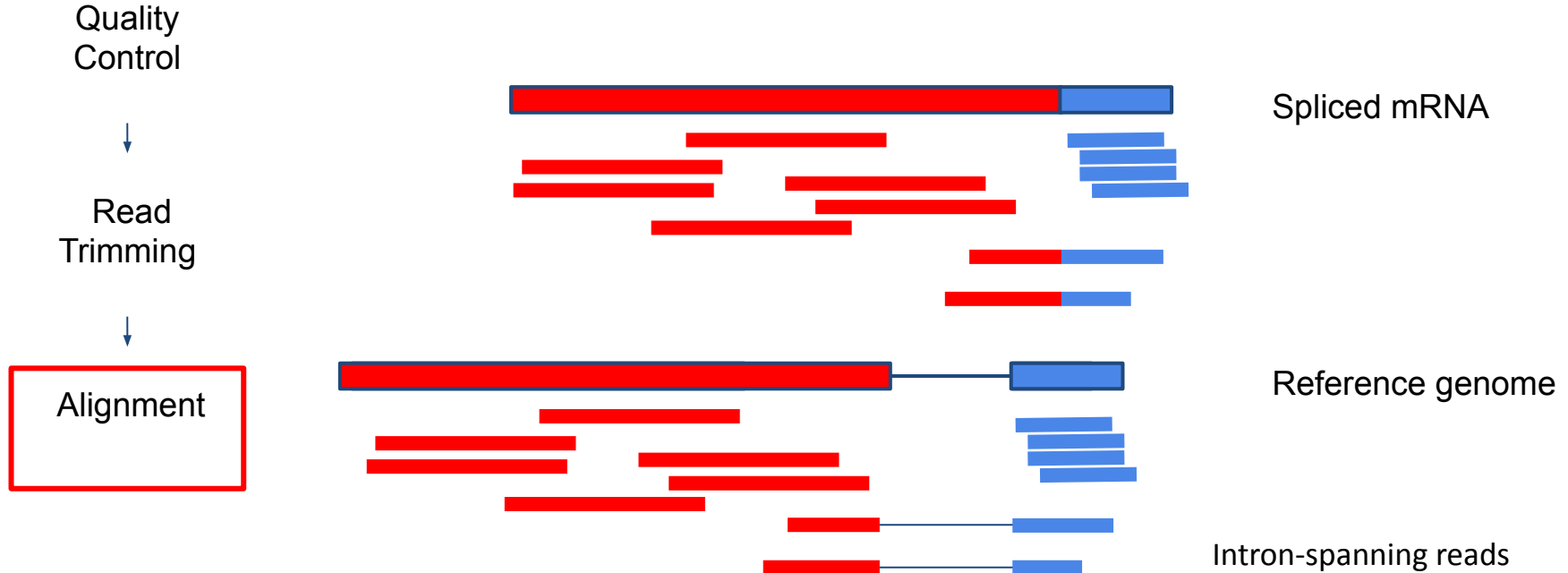
Aligning RNA-seq Data



Aligning RNA-seq Data



Aligning RNA-seq Data



Alignment Tools

Quality
Control



Read
Trimming



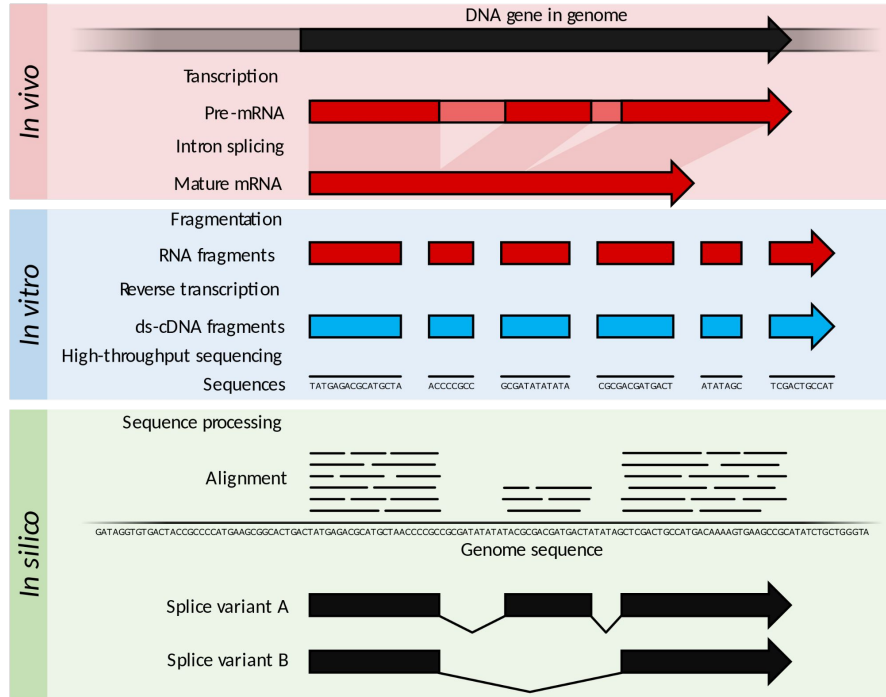
Alignment

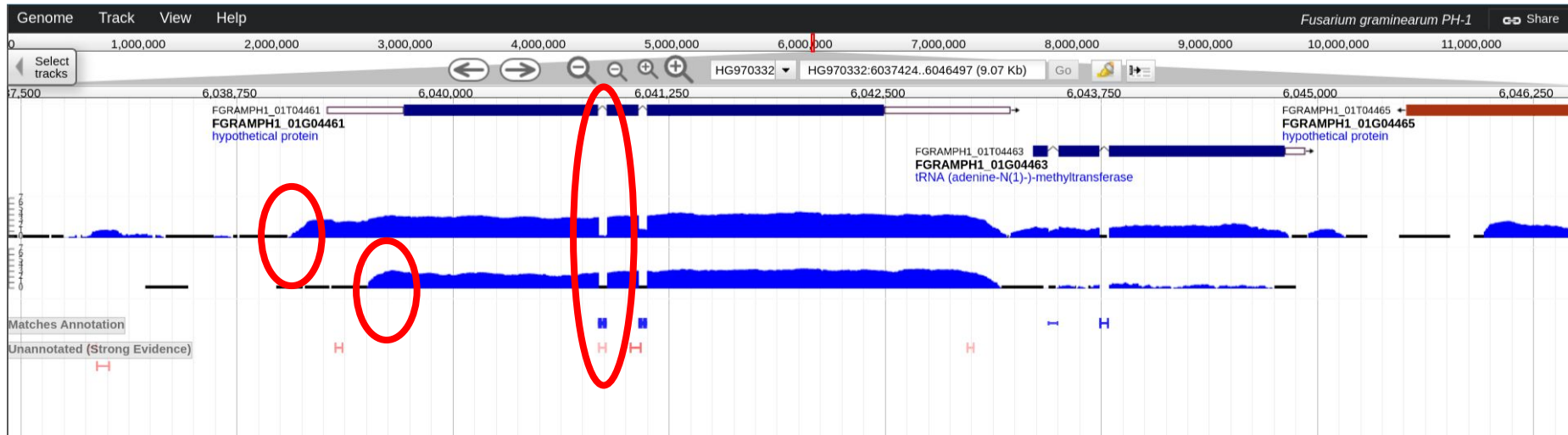
Mapping Tools for RNA-seq data

- Must be capable of aligning intron-spanning reads
- Hisat2
 - Fast, sacrifices sensitivity
 - <http://daehwankimlab.github.io/hisat2/>
- STAR
 - Very sensitive, but slow
 - <https://github.com/alexdobin/STAR>

Transcript Sequencing (RNA-seq)

Transcriptome sequencing

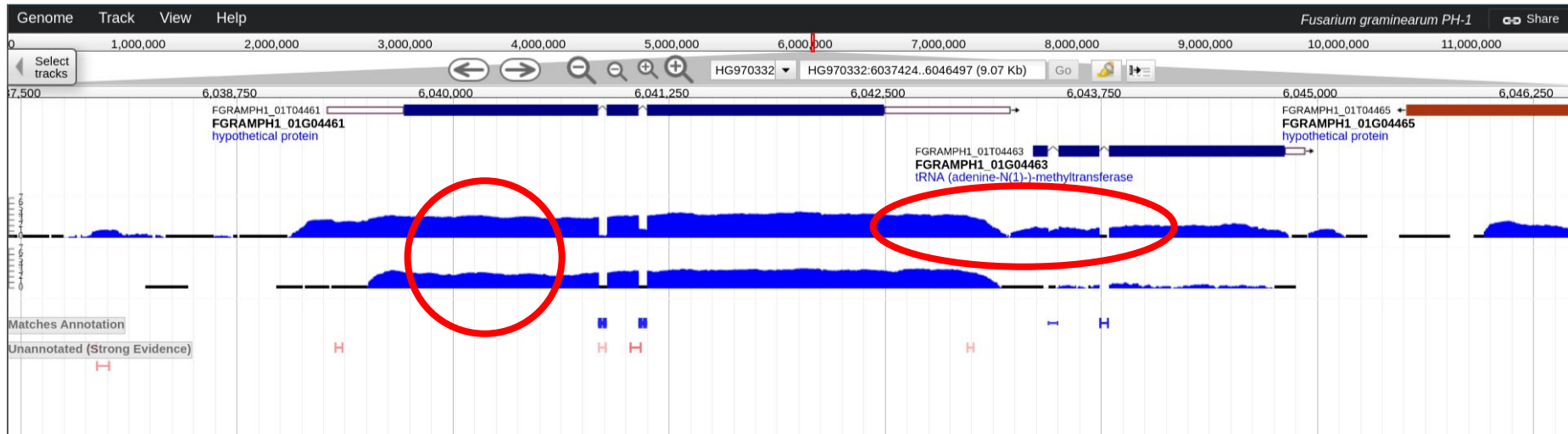




Gene Model Prediction

Alignment of RNA-seq reads to a genomic reference can help us to predict and confirm gene model structure

- Introns can be predicted based on coverage and on individual reads that cross splice junctions
- UTRs can be predicted based on coverage
- Differential splicing can also be predicted from coverage

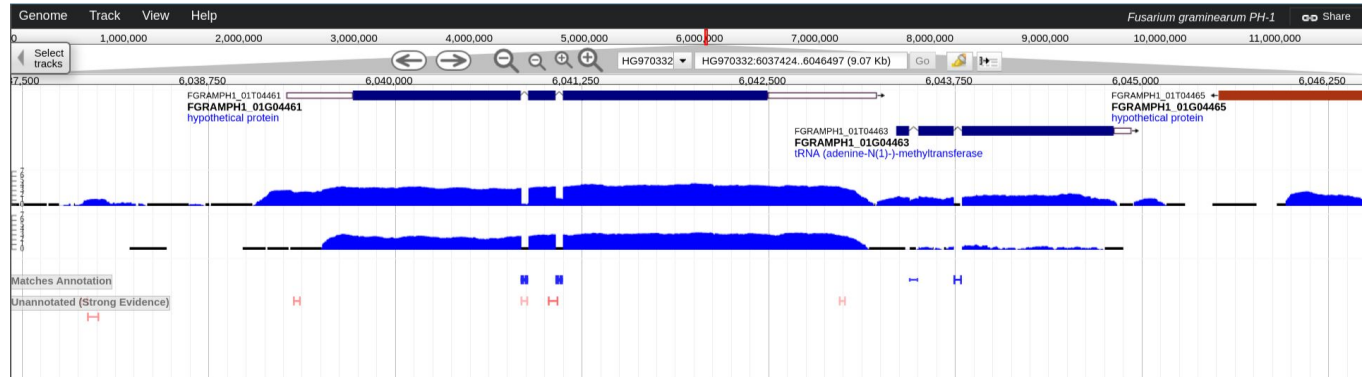
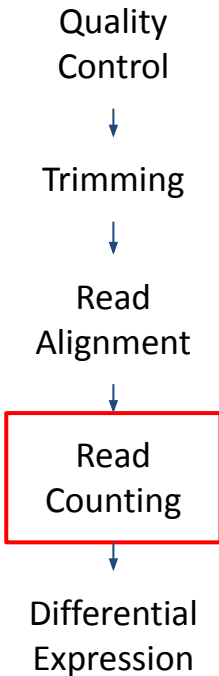


Differential Expression

Depth of coverage can help us learn about transcript abundance

- Differential transcript abundance can be observed both within and between samples

Quantifying Expression



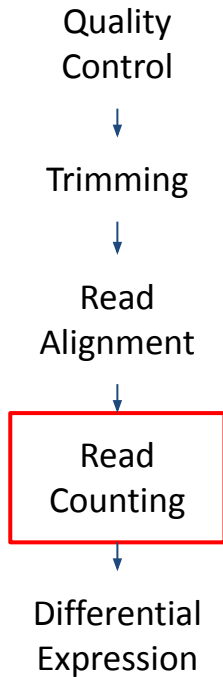
Quantifying Expression

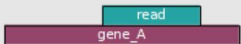
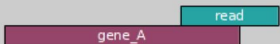




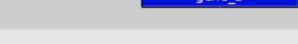

We've seen that we can see expression differences in a genome browser

Looking at plots like this is great for one gene, but it is too much to look at every gene individually and is not statistically robust

To examine transcript expression globally and perform robust statistics, we must count how many reads map to each gene.

Quantifying Expression



	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous (both genes with --nonunique all)	gene_A	gene_A
	ambiguous (both genes with --nonunique all)		
	alignment_not_unique (both genes with --nonunique all)		

Read Counting Tools:

htseq-count:

https://htseq.readthedocs.io/en/release_0.11.1/count.html

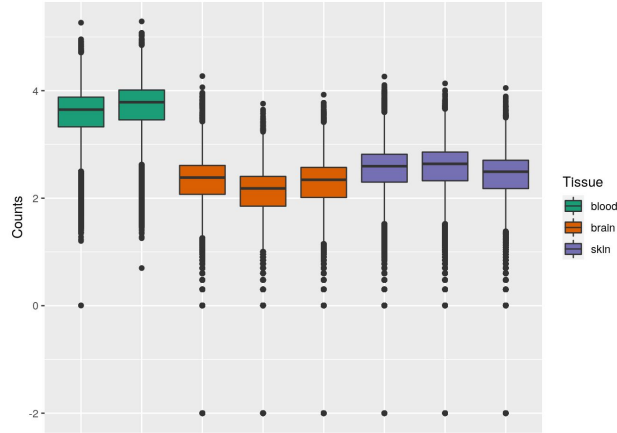
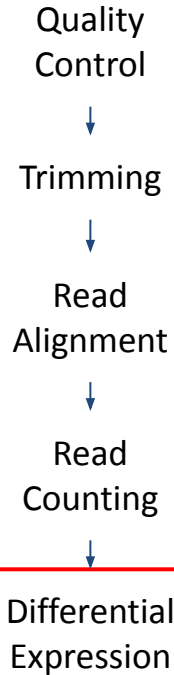
FeatureCounts:

<https://subread.sourceforge.net/featureCounts.html>

Kallisto:

<https://pachterlab.github.io/kallisto/>

Normalisation

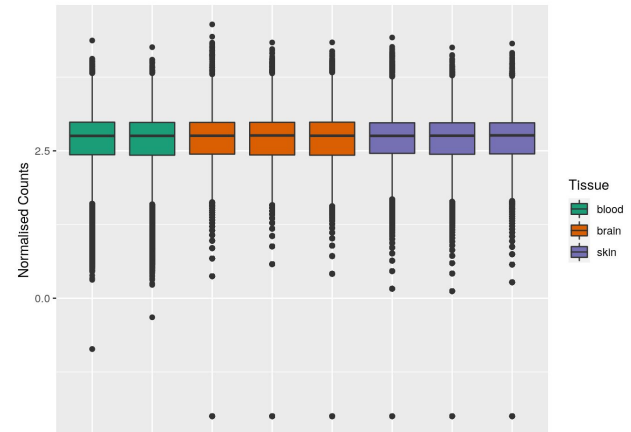


Raw count data

Each box represents is one sample and shows the distribution of read counts for each gene

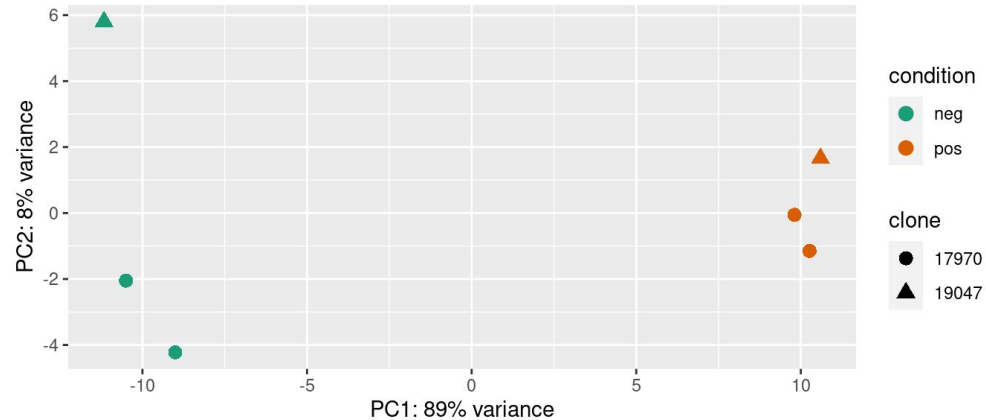
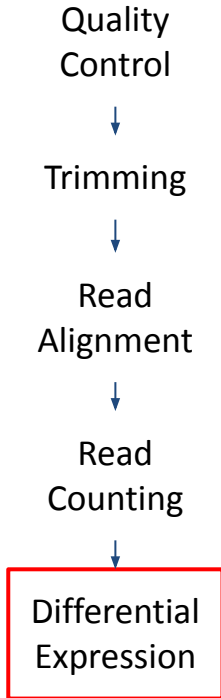
Normalised count data

After normalising, count distributions are aligned so individual genes can be compared



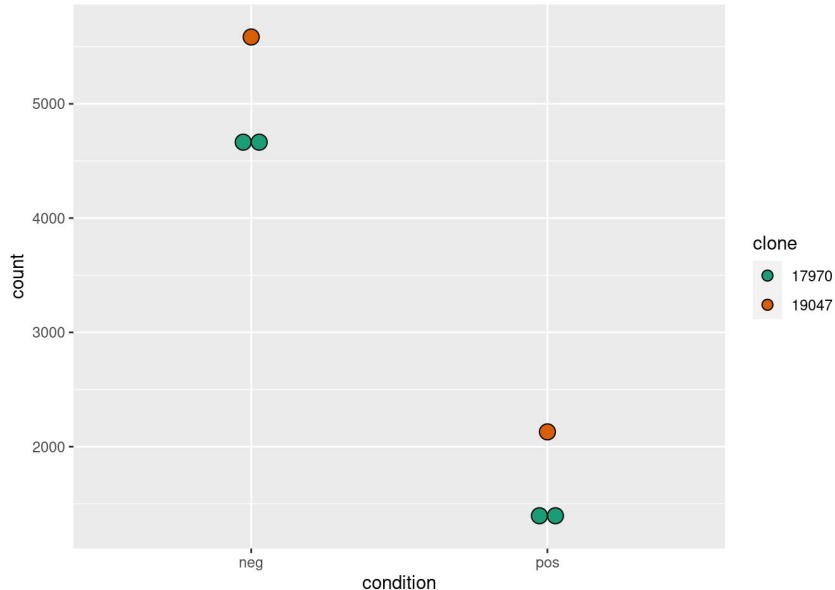
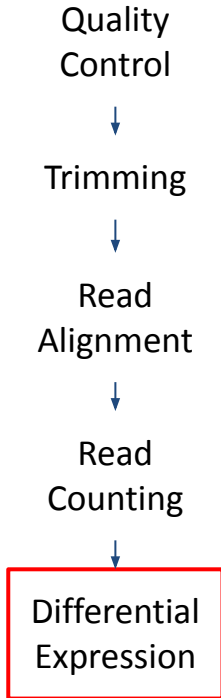
With Quantitative Data We Can...

- Explore our data



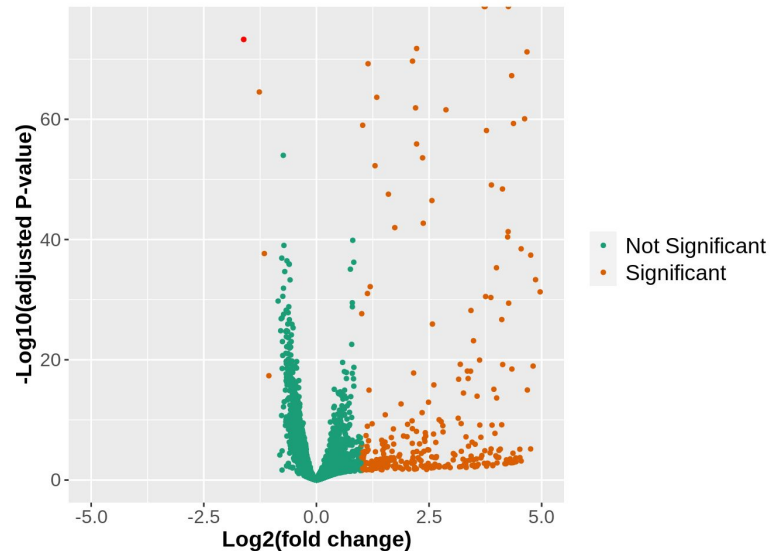
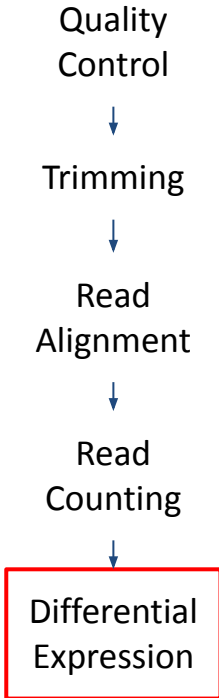
With Quantitative Data We Can...

- Explore our dataset
- Look at expression for individual genes



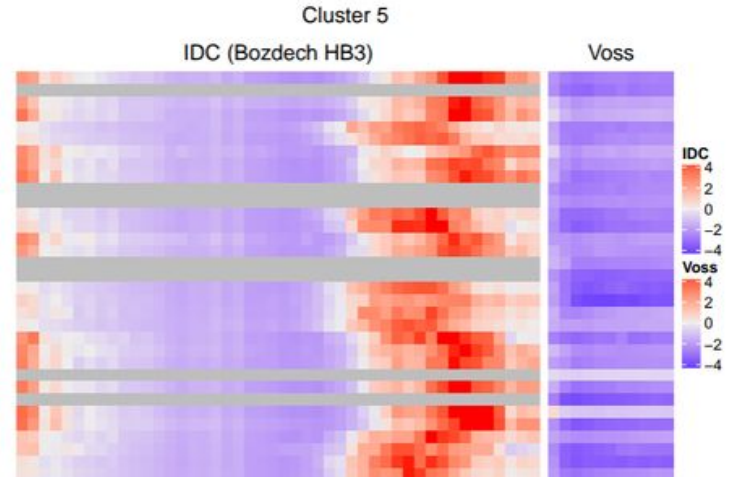
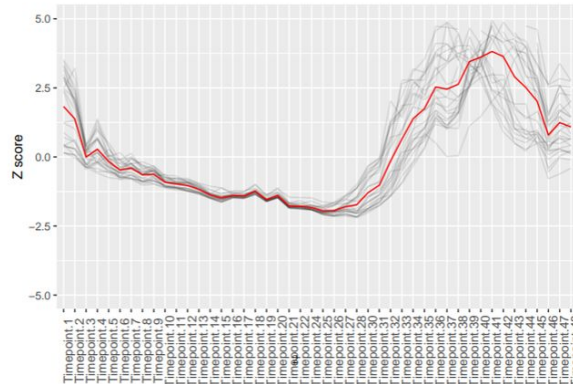
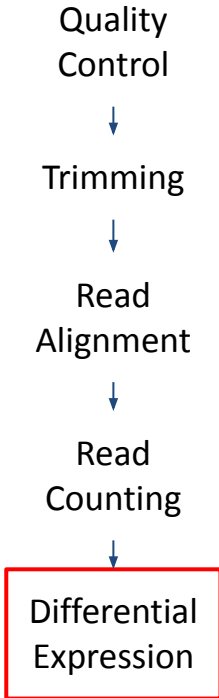
With Quantitative Data We Can...

- Explore our dataset
- Look at expression for individual genes
- Do pairwise statistical tests (differential expression)



With Quantitative Data We Can...

- Explore our dataset
- Look at expression for individual genes
- Do pairwise statistical tests (differential expression)
- Do advanced analysis (clustering, coexpression, etc.)



This session

1. Set up workflow for RNA-seq
2. Set up workflow for SNP calling (we will talk about this later)
3. Let them run overnight - we'll look at the output tomorrow
4. Once they are running you can log off

Galaxy

