

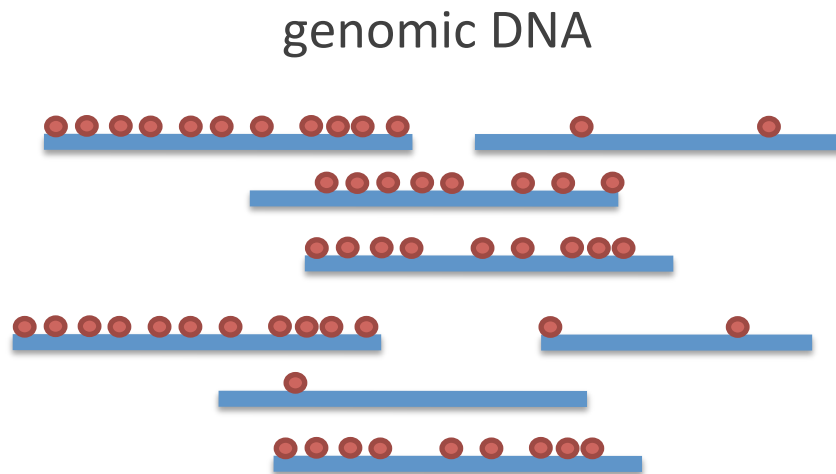
# Measuring DNA methylation using high-throughput sequencing

Mackenzie Gavery  
11/10/16

DNA methylation –  
how to we measure it?

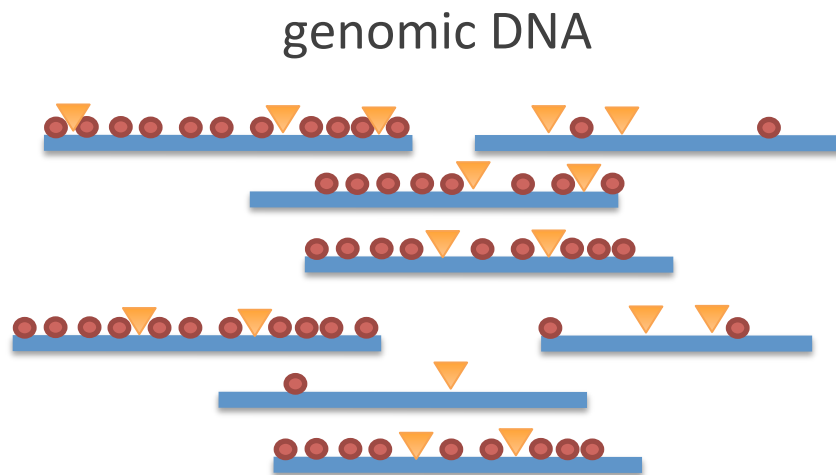
# DNA methylation – how to we measure it?

## Reduced Representation Bisulfite Sequencing (RRBS)



# DNA methylation – how to we measure it?

## Reduced Representation Bisulfite Sequencing (RRBS)

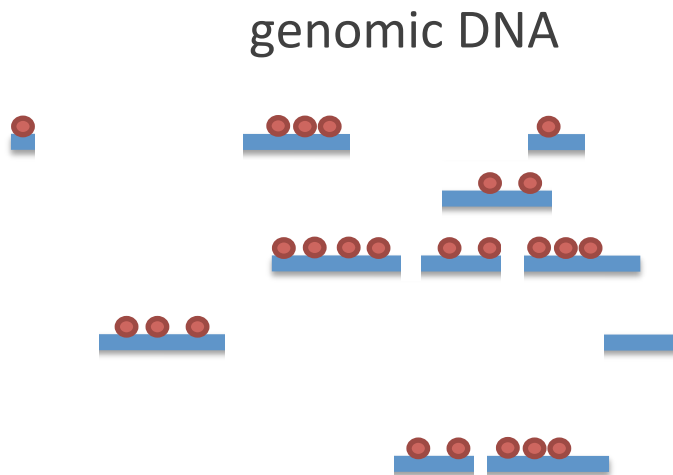


▼ =CCGG



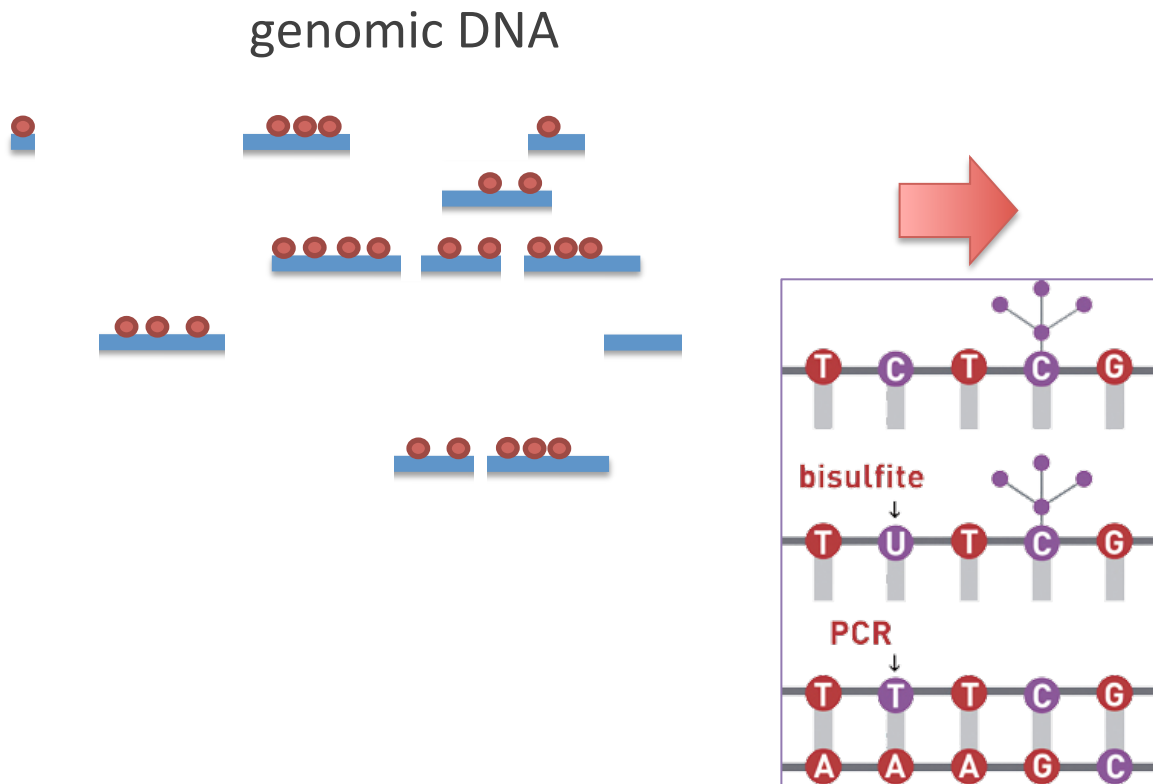
# DNA methylation – how to we measure it?

## Reduced Representation Bisulfite Sequencing (RRBS)



# DNA methylation – how to we measure it?

## Reduced Representation Bisulfite Sequencing (RRBS)



# DNA methylation analysis

- Library prep
- Sequencing
- Bioinformatics (DNA methylation)
  - QC/Trimming
  - Mapping
  - Extracting Methylation Data
  - Interpretation (e.g. differential methylation)

# Library Prep

Fragment DNA



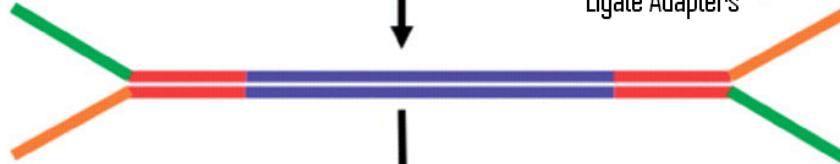
End Repair (Blunt ends)



Add 3' A Tail



Ligate Adapters



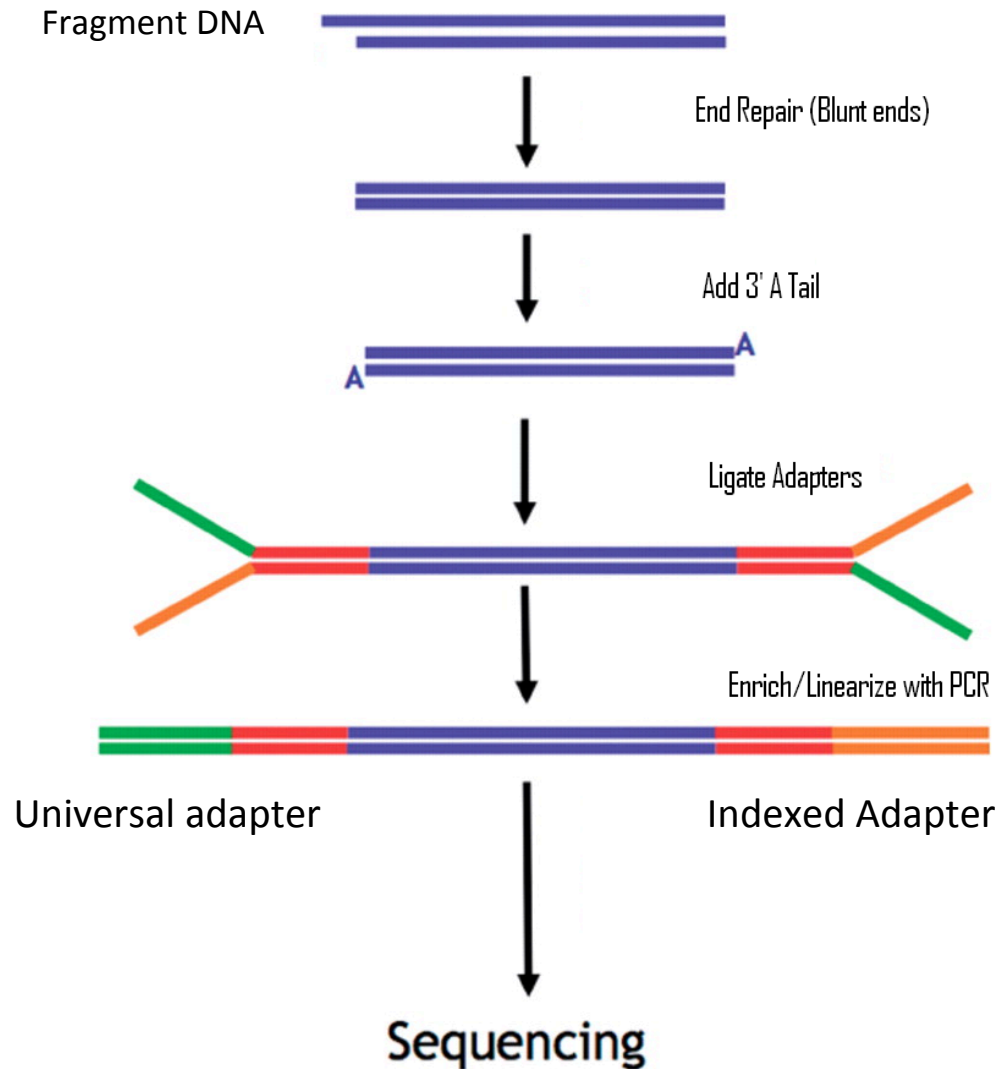
Enrich/Linearize with PCR



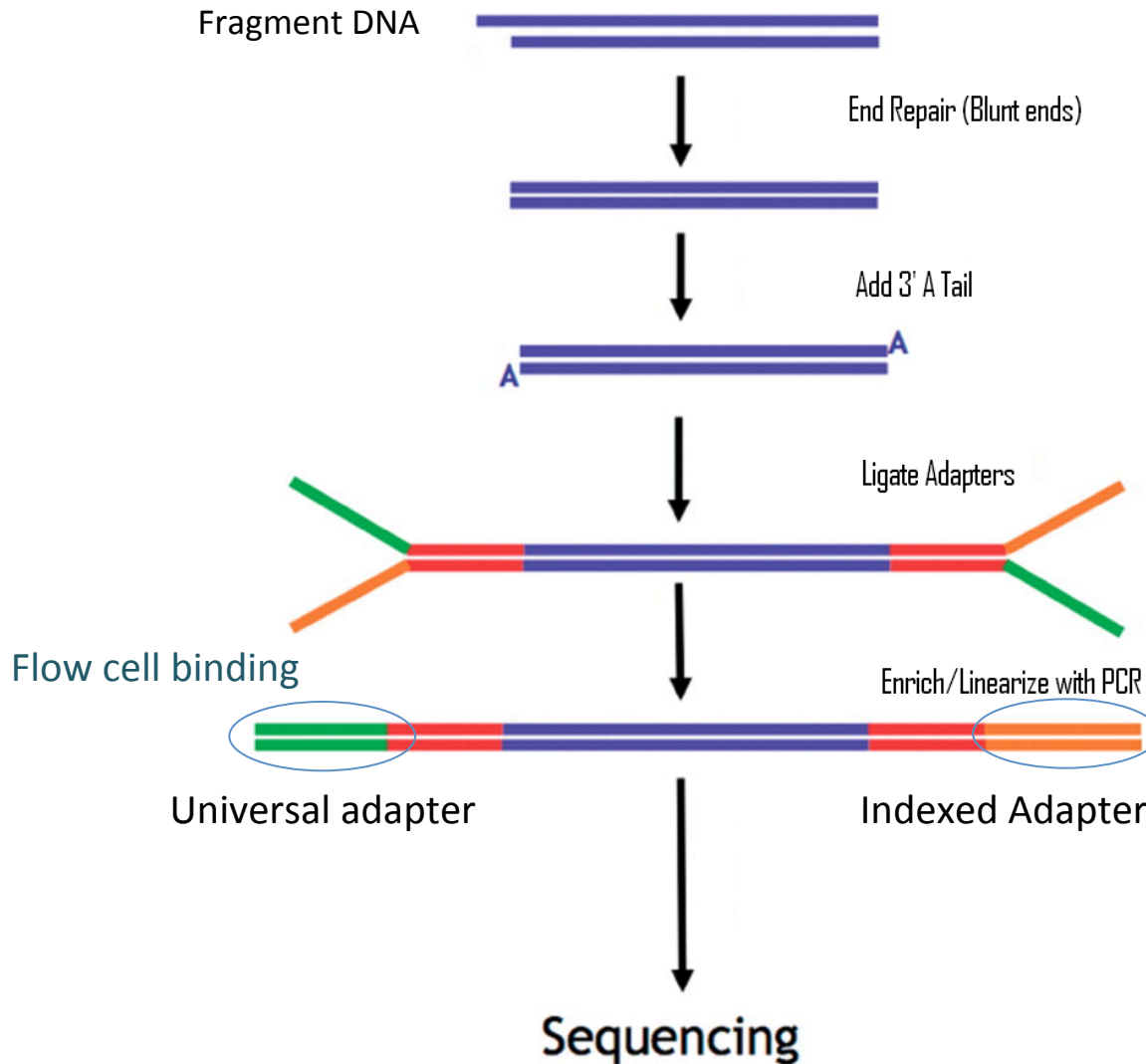
Sequencing



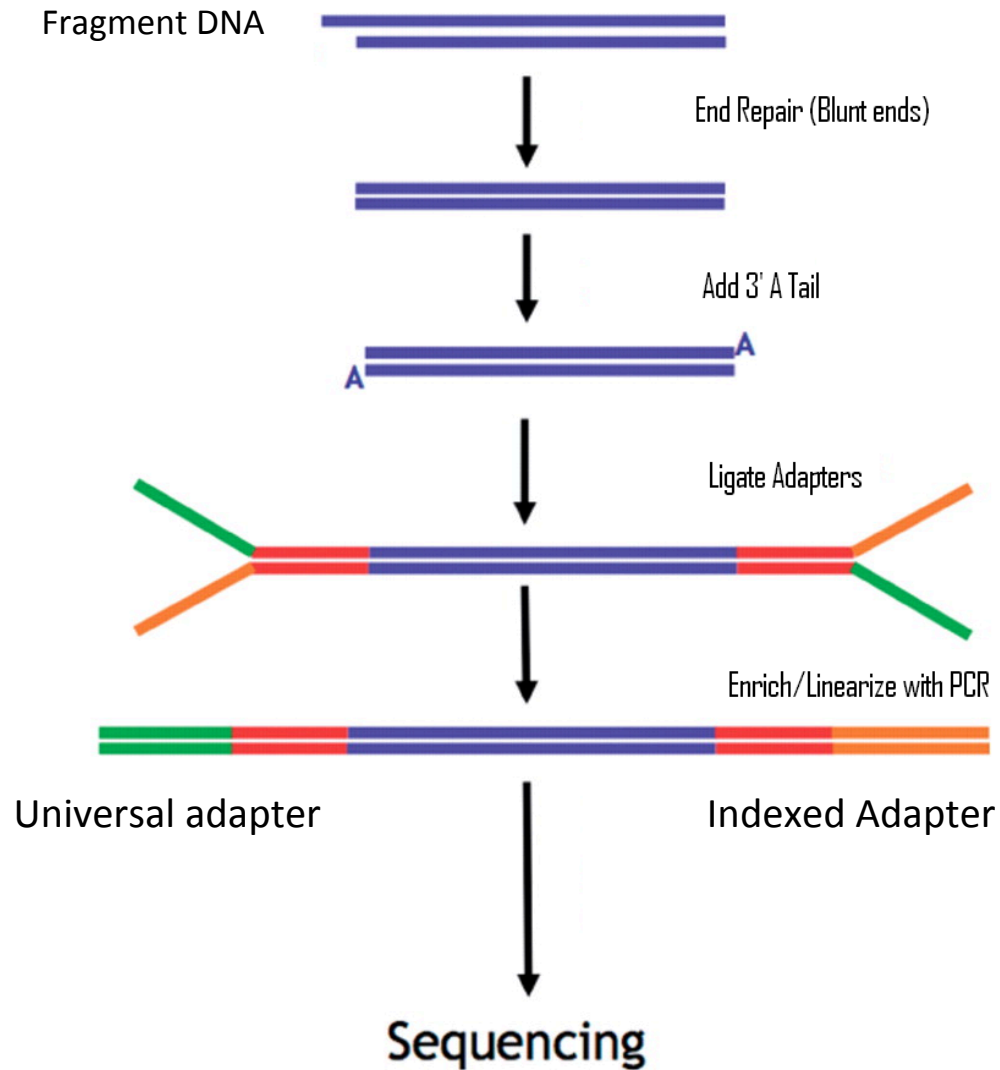
# Library Prep



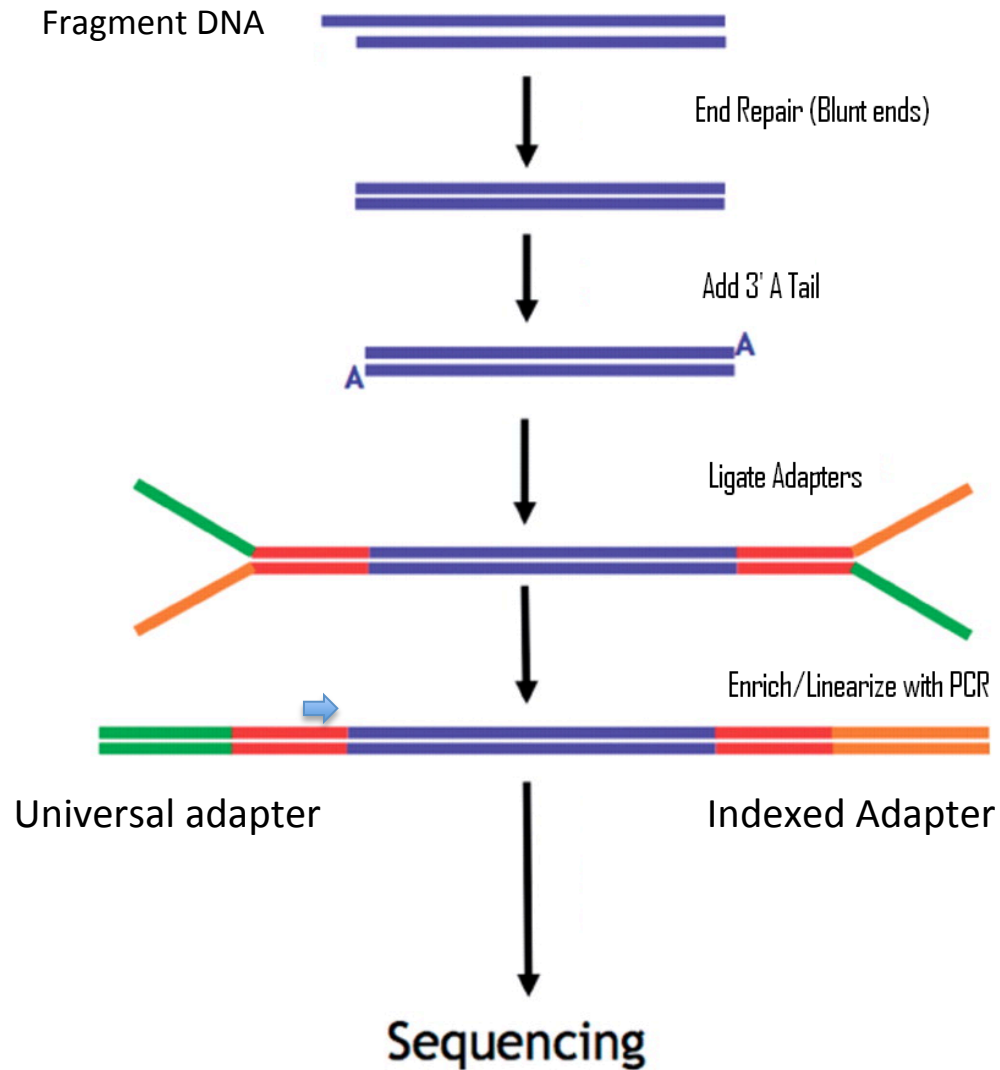
# Library Prep



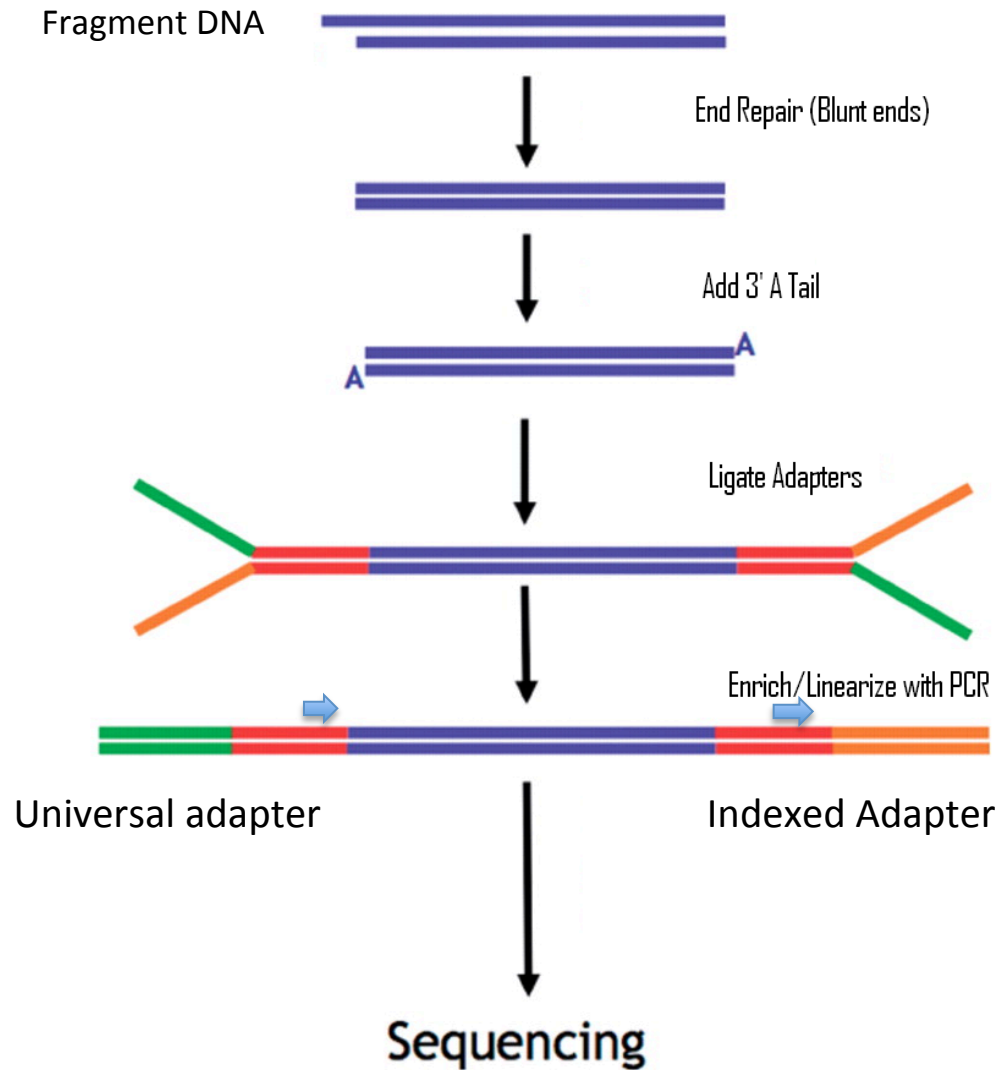
# Library Prep



# Library Prep



# Library Prep



# Library Prep

Fragment DNA



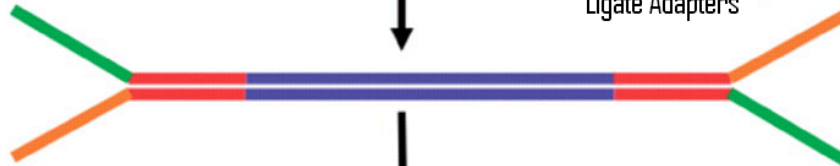
End Repair (Blunt ends)



Add 3' A Tail



Ligate Adapters

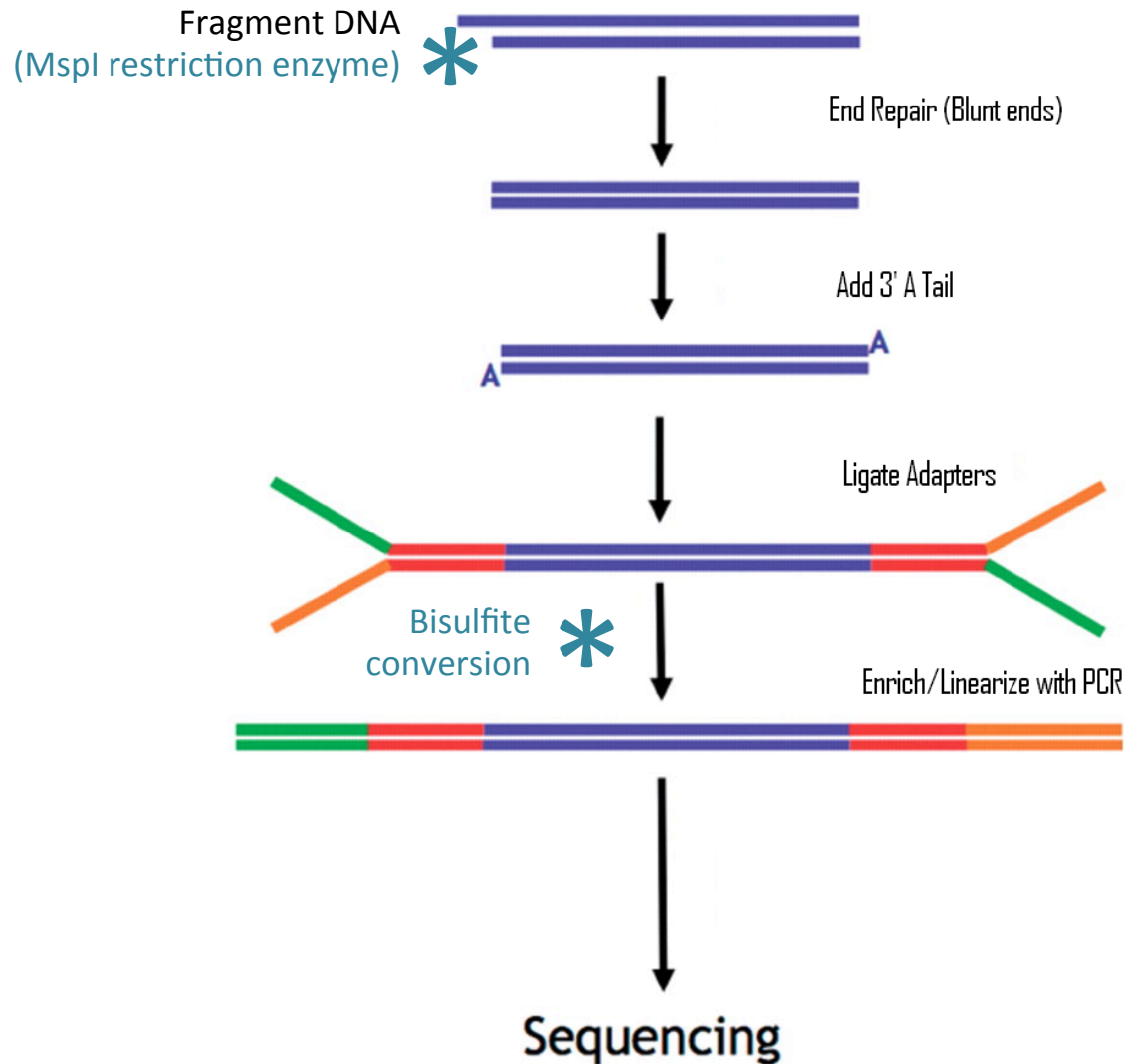


Enrich/Linearize with PCR



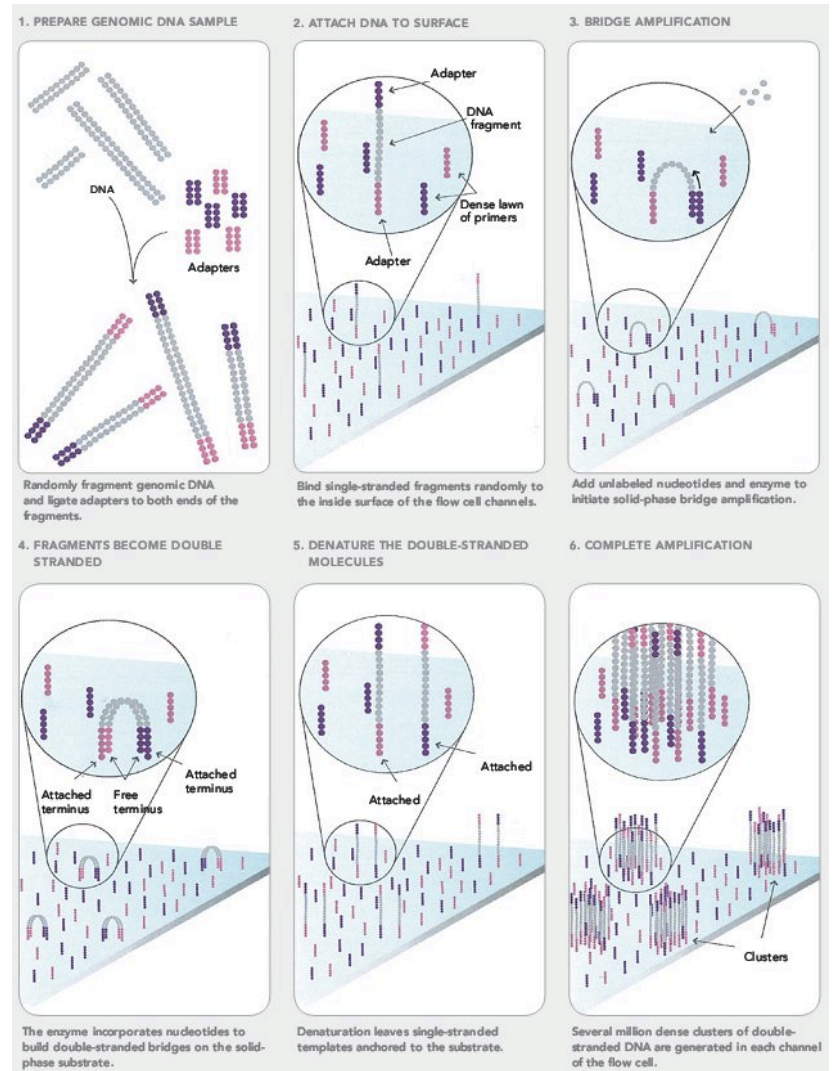
Sequencing

# Library Prep



# Sequencing

- Indexed libraries are pooled together
- Illumina HiSeq  
(single read 100bp)
- ~200 million reads





# Fastq Files

```
@SN747:551:C99B9ACXX:7:1110:1855:1151 1:N:0:CAGATC
CGGGTGATGTAGATGGTGGTGGGGTGGTTGGATCGATTGTGGGGAGTTGGGAAGGTGGTGTAAATTTGTAGTTGGCGTACGATTTGAGATCGGAAGAGCACA
+
@@C?A@?4CFFFFG@EEGG@AGHA?CAD=FHDFB;;FHF=FEBHFD?C;>CAC==CAB8<<ACCDDCEDCDD:A@09<9<@B?B>@<8@:(087&+2<@A<
```

# Fastq Files

```
@SN747:551:C99B9ACXX:7:1110:1855:1151 1:N:0:CAGATC
CGGGTGATGTAGATGGTGGTGGGGTGGTTGGATCGATTGTGGGGAGTTGGGAAGGTGGTGTAAATTTGTAGTTGGCGTACGATTTGAGATCGGAAGAGCACA
+
@@C?A@?4CFFFFG@EEGG@AGHA?CAD=FHDFB;;FHF=FEBHFD?C;>CAC==CAB8<<ACCDCEDECDD:A@09<9<@B?B>@<8@:(087&+2<@A<
```

**How do you analyze 200 million reads!?**

# Bioinformatics (DNA methylation)

Quality Control & Trimming



Mapping RRBS reads to the genome

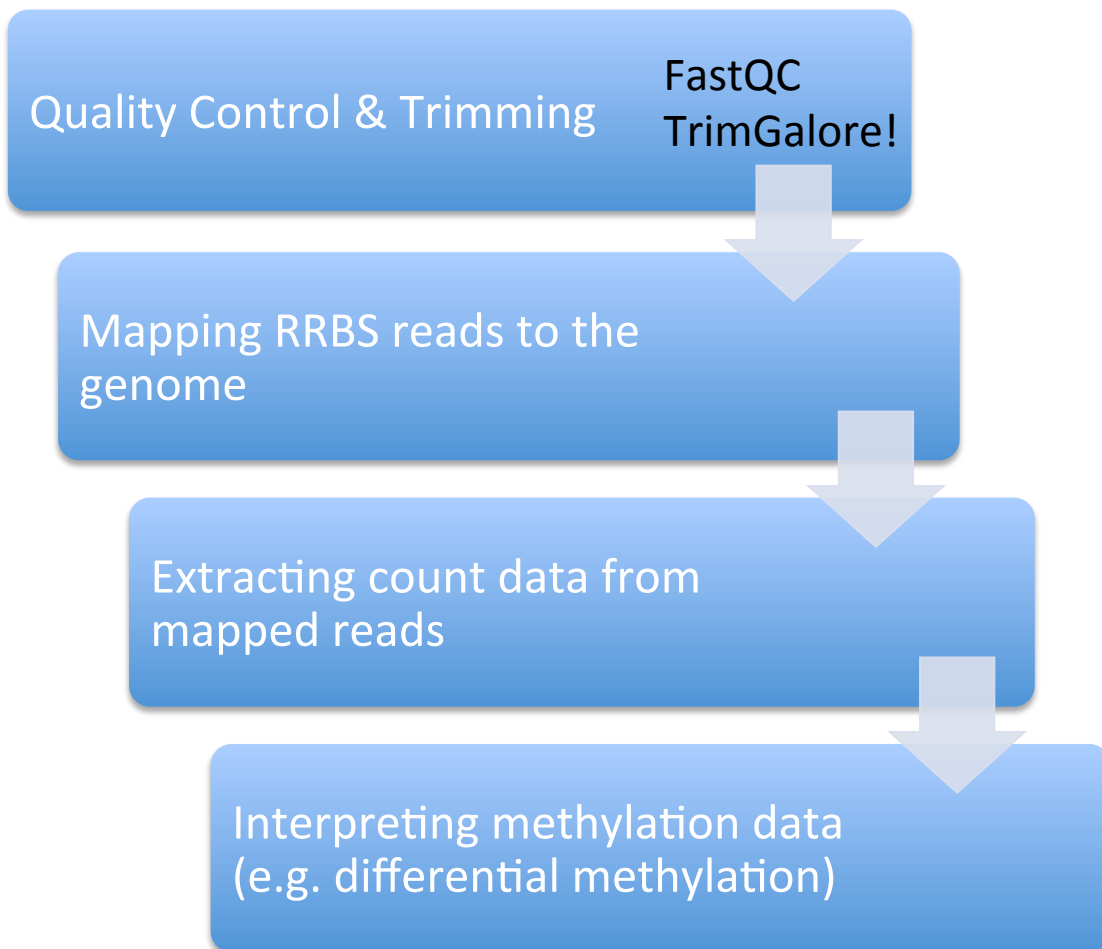


Extracting count data from mapped reads



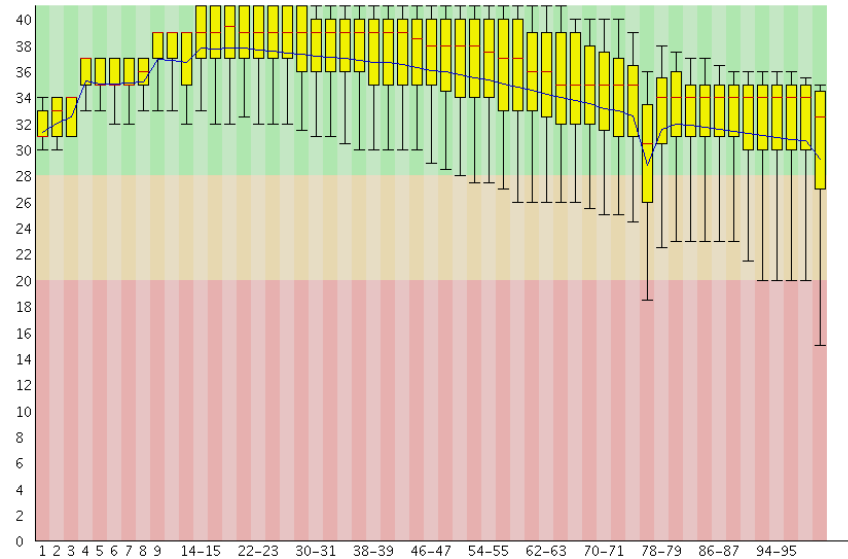
Interpreting methylation data  
(e.g. differential methylation)

# Bioinformatics (DNA methylation)



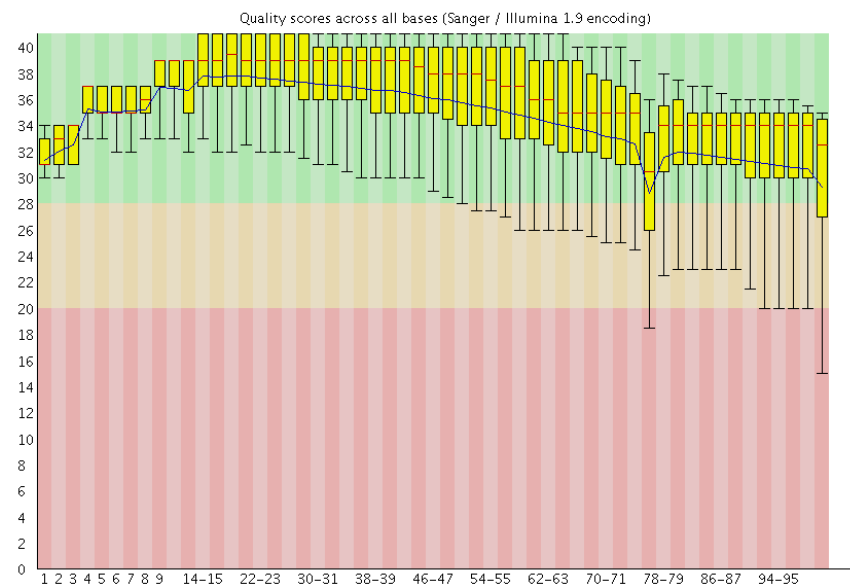
✔ Per base sequence quality

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



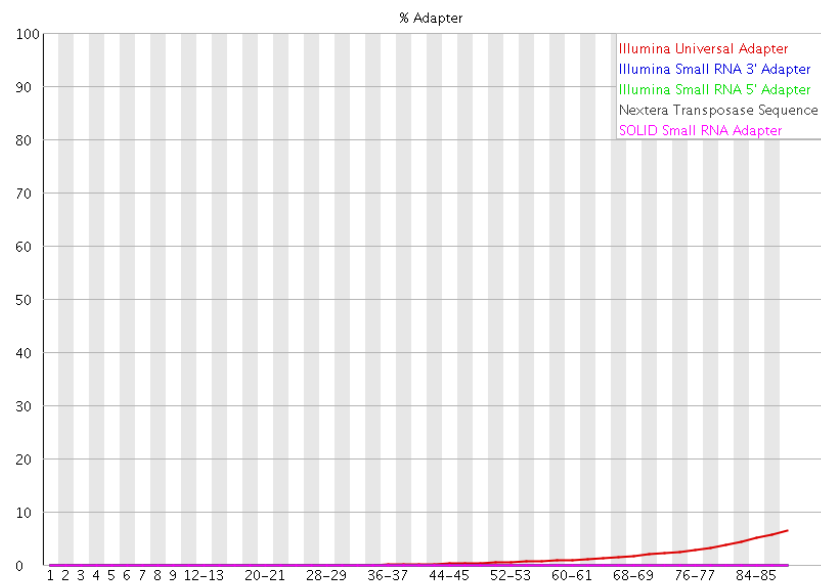
Position in read (bp)

## ✓ Per base sequence quality



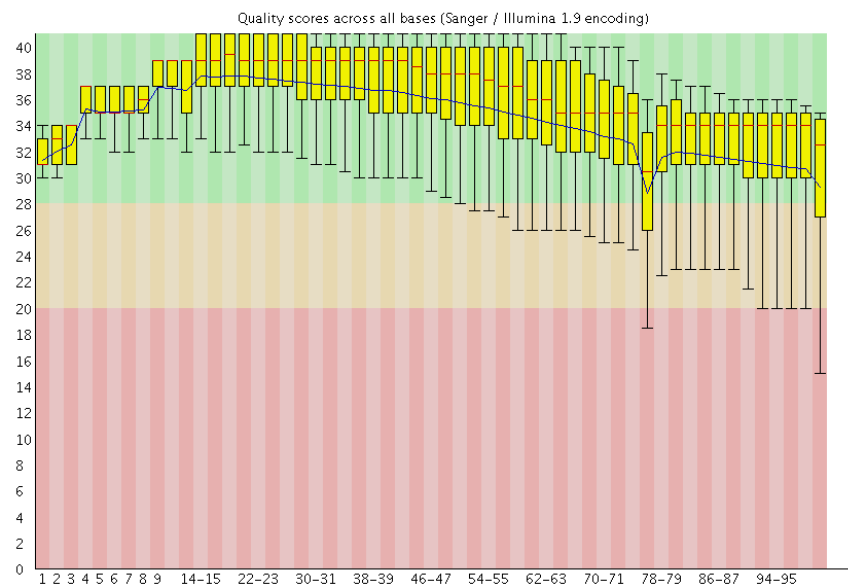
Position in read (bp)

## 🚩 Adapter Content



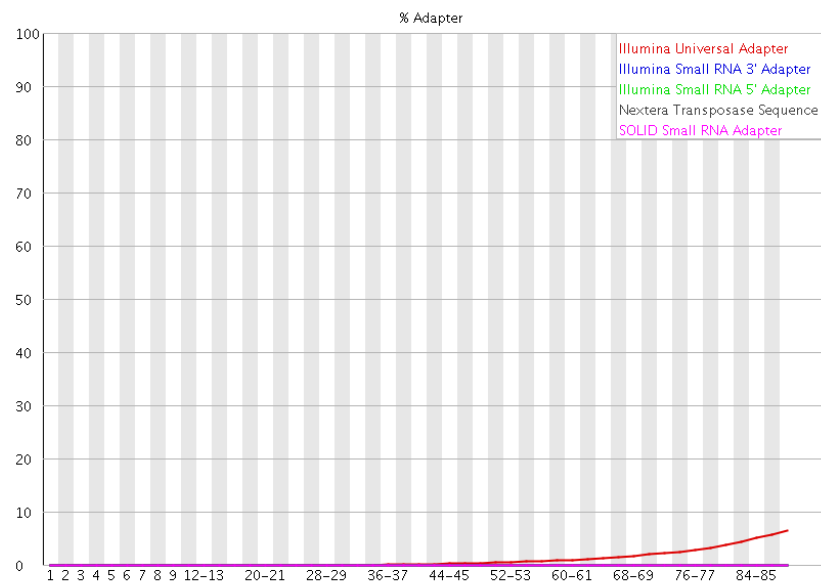
Position in read (bp)

## ✔ Per base sequence quality



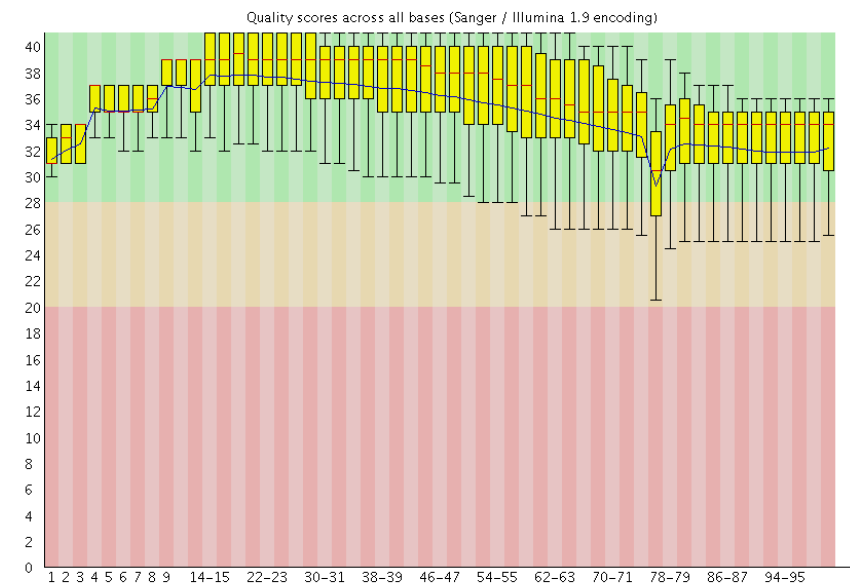
Position in read (bp)

## ⚠ Adapter Content

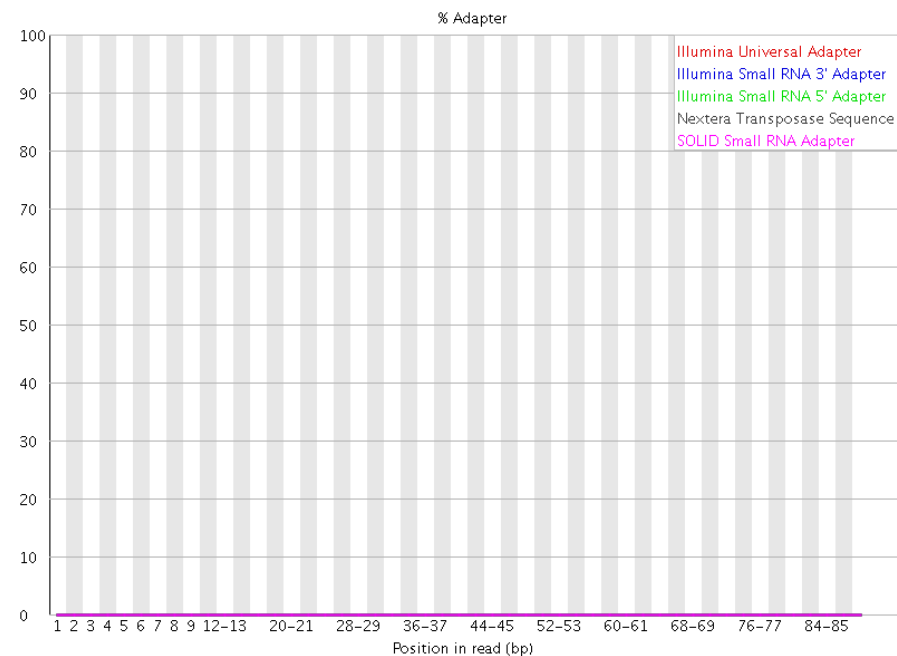


Position in read (bp)

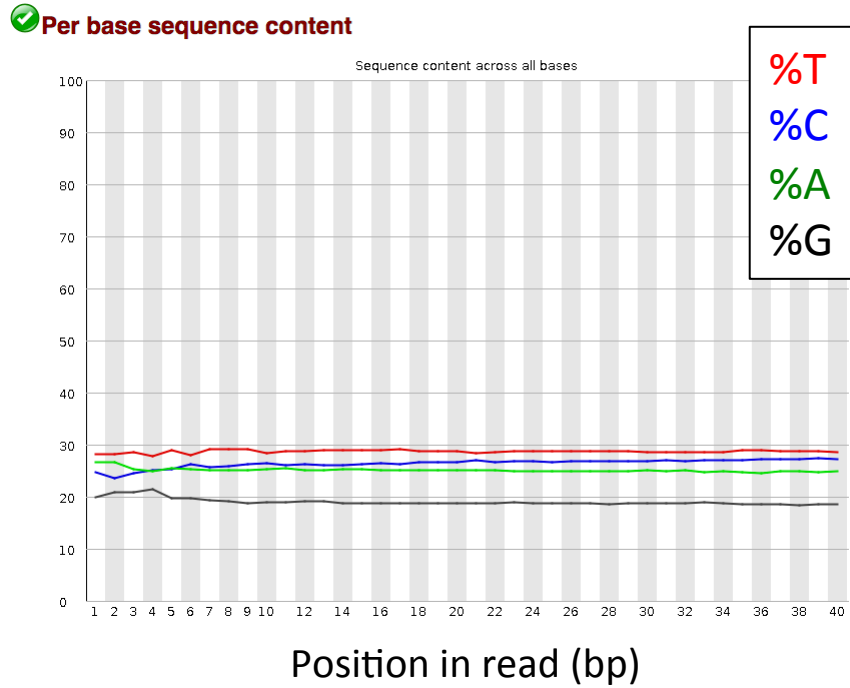
## ✔ Per base sequence quality



## ✔ Adapter Content



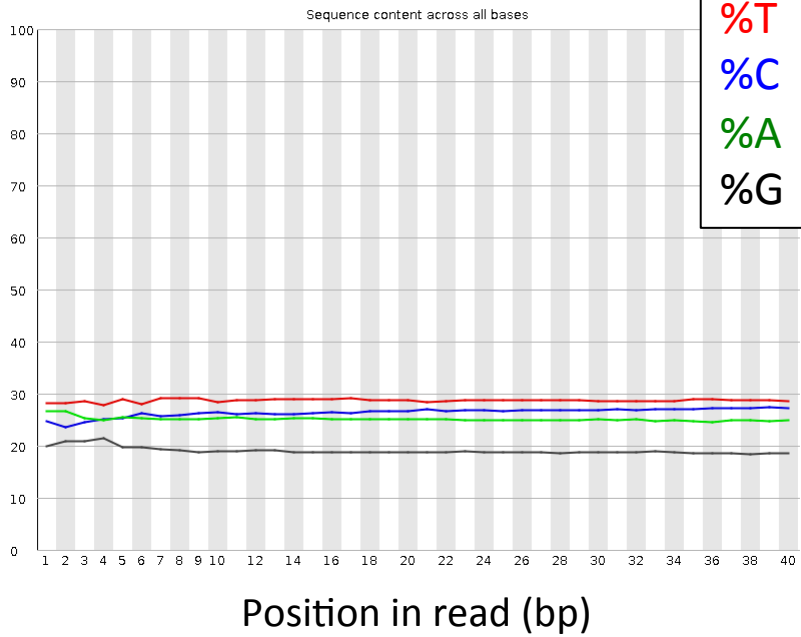
# QC



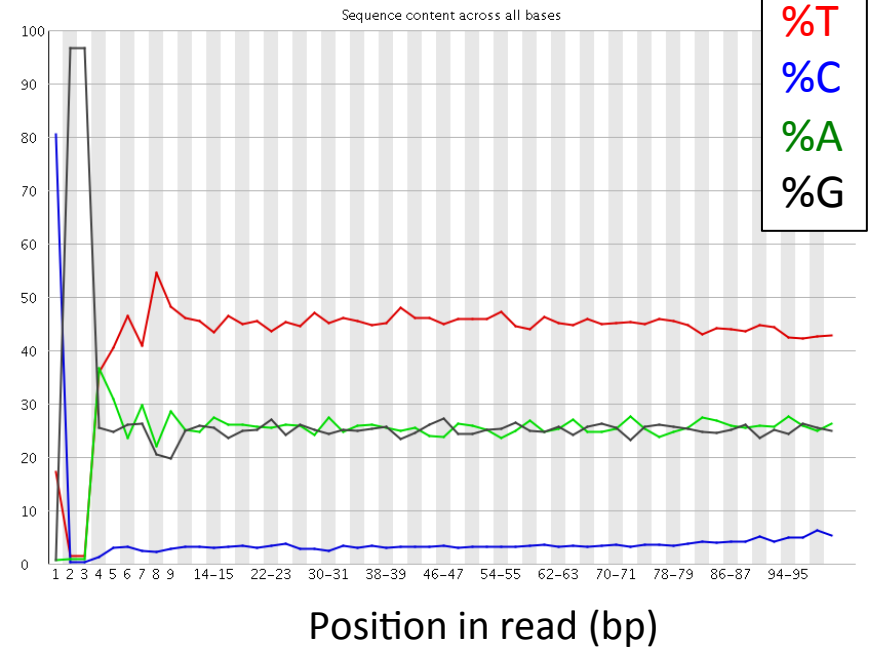


# QC

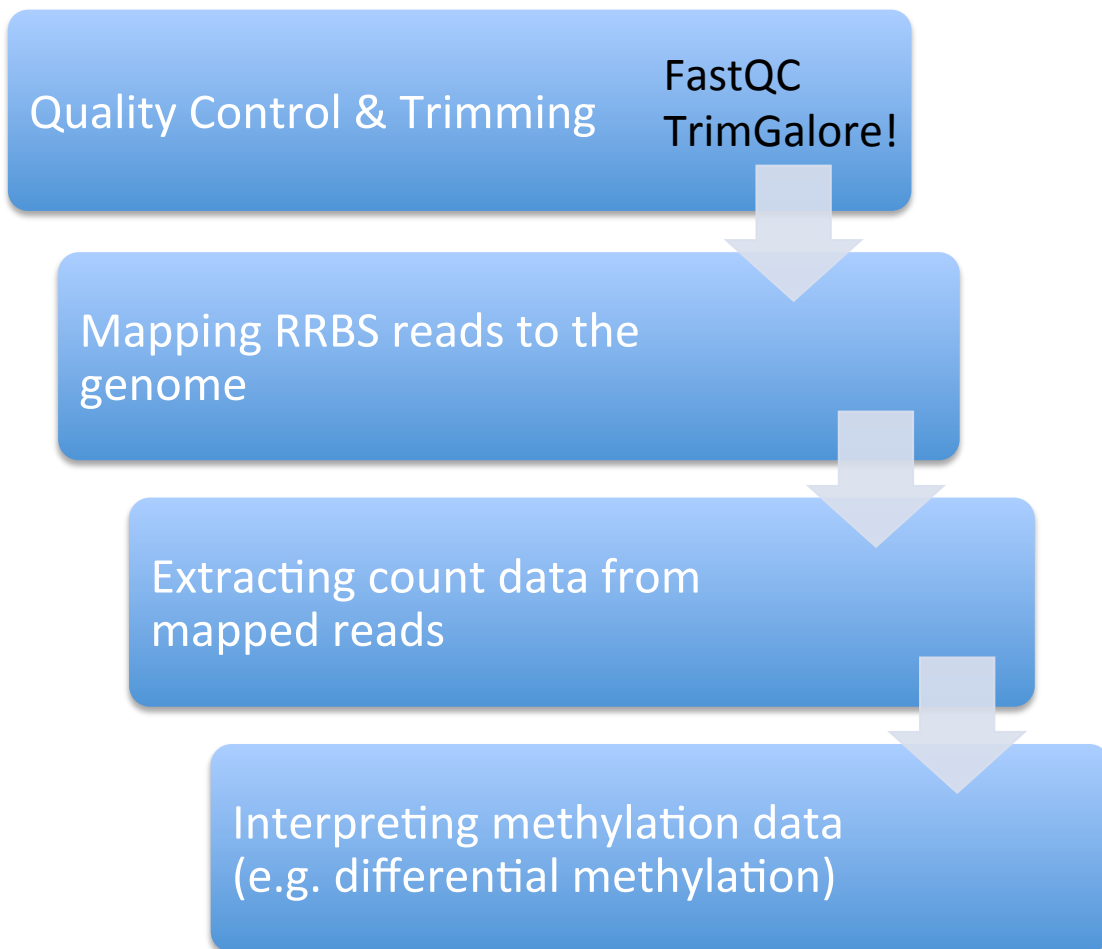
## ✔ Per base sequence content



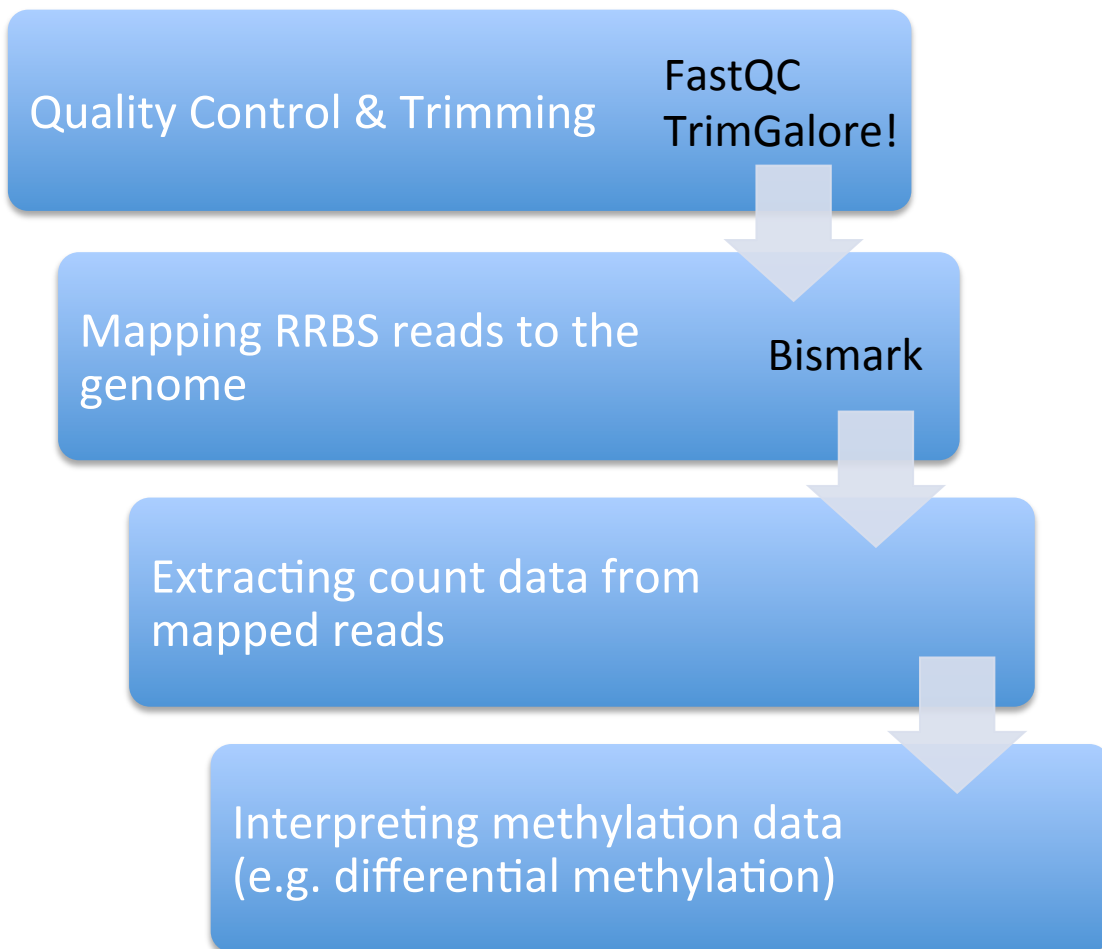
## ✖ Per base sequence content



# Bioinformatics (DNA methylation)



# Bioinformatics (DNA methylation)



## *O. mykiss* scaffold 13

```
SN747:551:C99B9ACXX:8:2315:19986:99910_1:N:0:ACAGTG 16 scaffold_570 338476 4 73M * 0 0
ATCACAAAACGCGCTAACCAAAATTACCAAATACACGATATTTCTCCAACACATTAACGCGACTAACTTCCG
EDDDDDDDDFGHGIEHCJJJJIGHFFJJJJIFDHGJJJJIGGJJJJIGHGJJJJJJHDFHHFDFFFCCC NM:i:18 MD:Z:
0G0G4G1G12G0G2G3G0G1G4G1G8G7G0G4G2G0G6XM:Z:x.....x.h.Z.....hh..h...xh.h...Zx.h.....z.....hh.Z.Zx..xh.....Z XR:Z:CT
XG:Z:GA

SN747:551:C99B9ACXX:8:2315:19801:99933_1:N:0:ACAGTG 0 scaffold_2861 36547 42 97M * 0 0
CGGAAGGTGTTATAGAGGGTAGTGCGTACGGTTAGTATATTATTGGGGTAAATTTTTTGATATTTAGGGGAAGGTTTAAAAAATTGTGAAAGATT
@@@DDDDBB:CFA22AGE6+<:FFH8FFE@<FFFEA0BFFCECBD4BFB;;@CF>FCEFDCCEDB@@@B?B==?<>@BB>@>@BBAA@>>>:3A@
NM:i:17 MD:Z:9C2C18C0C0C4C2C1C10C2C0C3C3C9C0C0C16C1
XM:Z:Z.....h..x.....Z...Z..hhx...h..h.x.....h..hx...h...x.....hhh.....h. XR:Z:CT XG:Z:CT

SN747:551:C99B9ACXX:8:2315:20190:99881_1:N:0:ACAGTG 16 scaffold_614 418590 0 101M * 0 0
CAAATACAAAAAATCAAAACAAATAAATACCCACCCAAAACCTCAAACCATAAAACAAACAAAAAACATCATCTTATAAATAATCAAATAACCAACCACCA
>EEDA?5BBBDA>1595BB@EB@;@5;(,89,, "FFFE@3D=52FBFFIIFFFFF?9G:FFGD?F>@FFFCE9HCAEA?CCIIFF<C2)DB@B;A@?: NM:i:24
MD:Z:
2G0G1G3G7G15A6A5G2G8G1G0G0G4G9G3G0G2G0G0G3G2G2G2G0XM:Z:..xh.h...h.....h.....h..z.....h.zxh....z.....h...hh..zxh
.....x..z..z XR:Z:CT XG:Z:GA

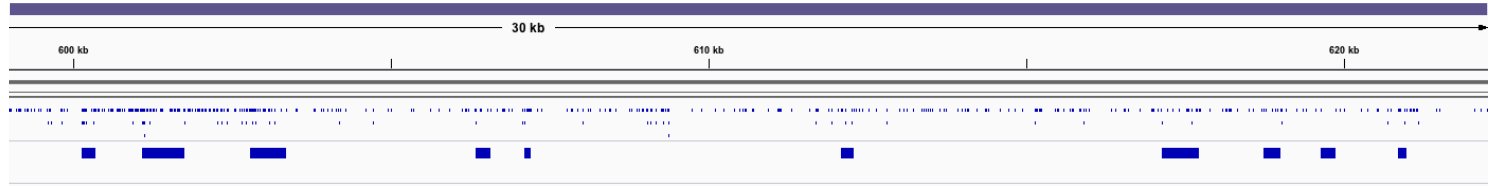
SN747:551:C99B9ACXX:8:2315:20074:99951_1:N:0:ACAGTG 16 scaffold_1684 42005 3 98M * 0 0
CATTAAACGACGAATTAATCTTCGTACCCGAACCTCGTAAAAAAAACCTCACGCTCTCCGTCGCCCCCTCTCCGCATCACTACCATCTCTCTACCG
CDDDBDB@BBCDDEEDDDDDDB?<C??8DDDD@DBDDDDDDDBCC@>D@<7DDDDDDDDDBDBBIGIGIGIJJIJGHFHGCIJHFADDBGFFDBDBB?
NM:i:15 MD:Z:4G1G5G0G2G14G5G0G3G0G1G3G2T7T33G3
XM:Z:....h.h.Z..Zxh..h.....Z.....Z.h...hh..h...z.Z.....Z..H.....Z.....x..Z XR:Z:CT XG:Z:GA
```

“Looking” at data using IGV

# “Looking” at data using IGV

*O. mykiss* scaffold 13

CG motif  
RRBS fragments



# “Looking” at data using IGV

*O. mykiss* scaffold 13

CG motif  
RRBS fragments

mapped  
reads



# “Looking” at data using IGV

*O. mykiss* scaffold 13

CG motif  
RRBS fragments

mapped  
reads





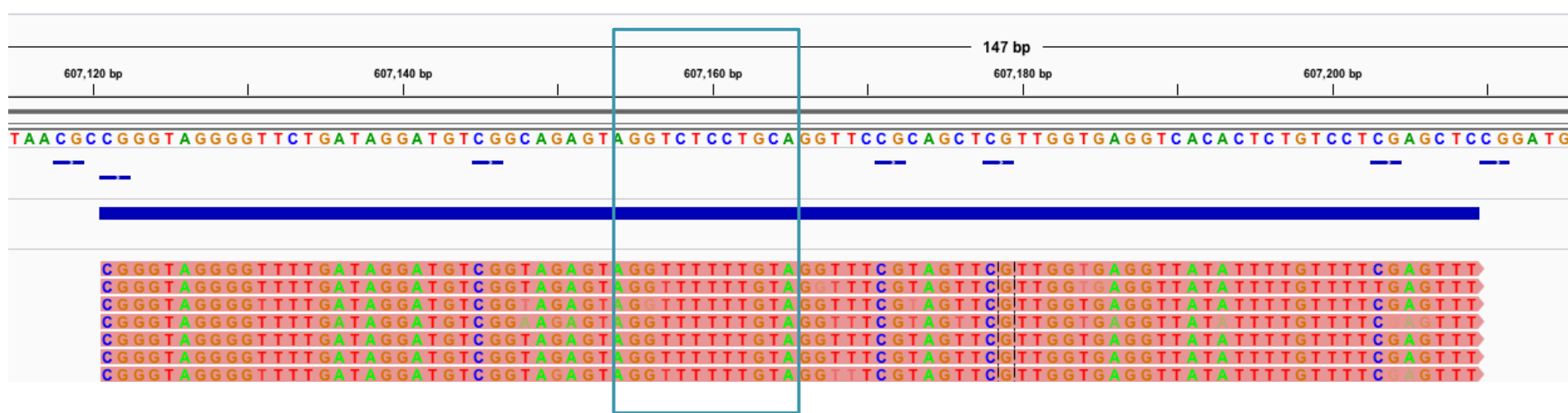
# “Looking” at data using IGV

*O. mykiss* scaffold 13



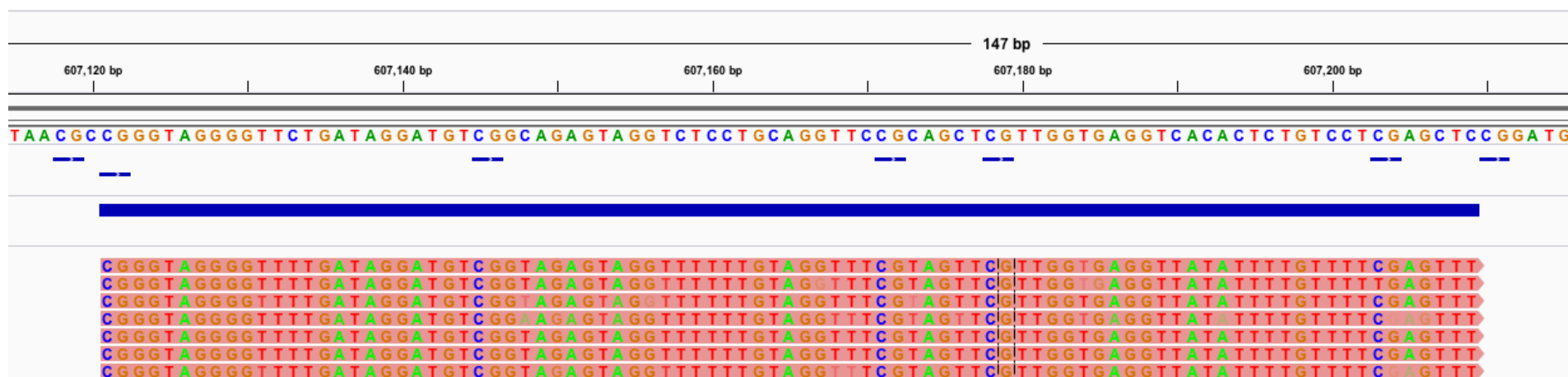
# “Looking” at data using IGV

*O. mykiss* scaffold 13



# “Looking” at data using IGV

*O. mykiss* scaffold 13



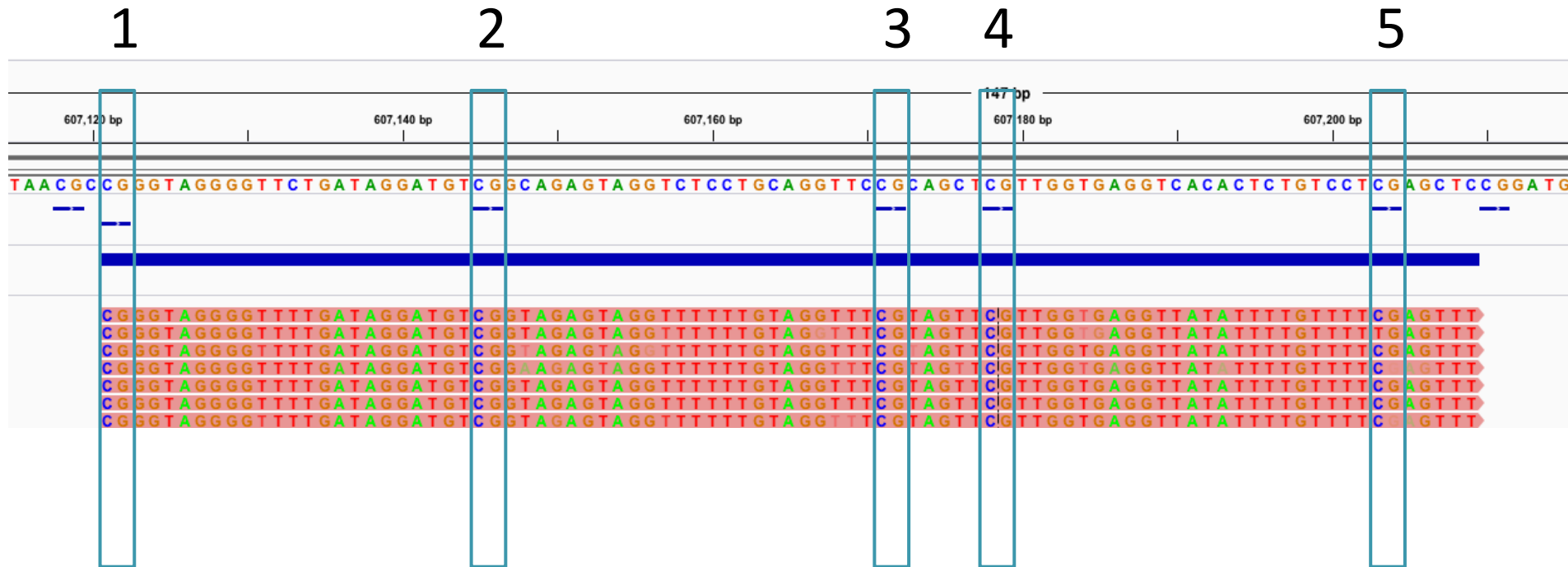
# “Looking” at data using IGV

*O. mykiss* scaffold 13



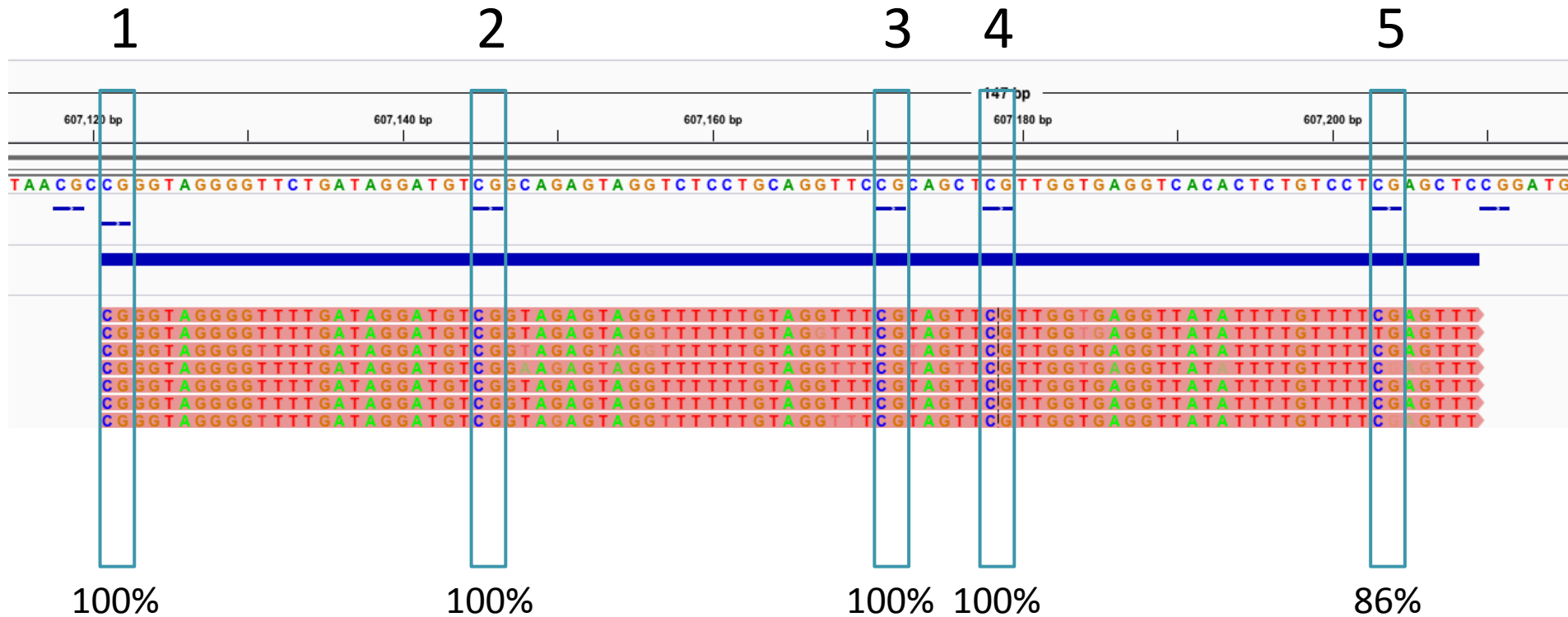
# “Looking” at data using IGV

*O. mykiss* scaffold 13

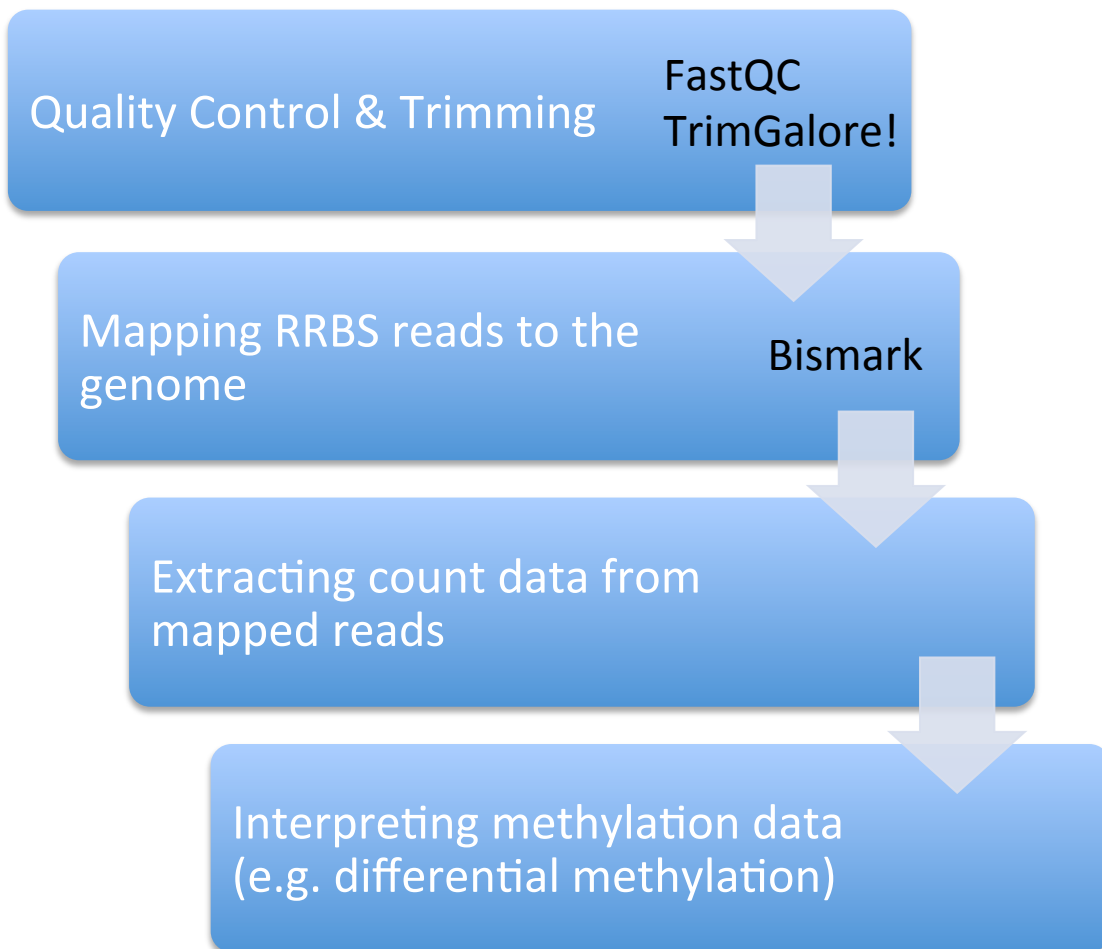


# “Looking” at data using IGV

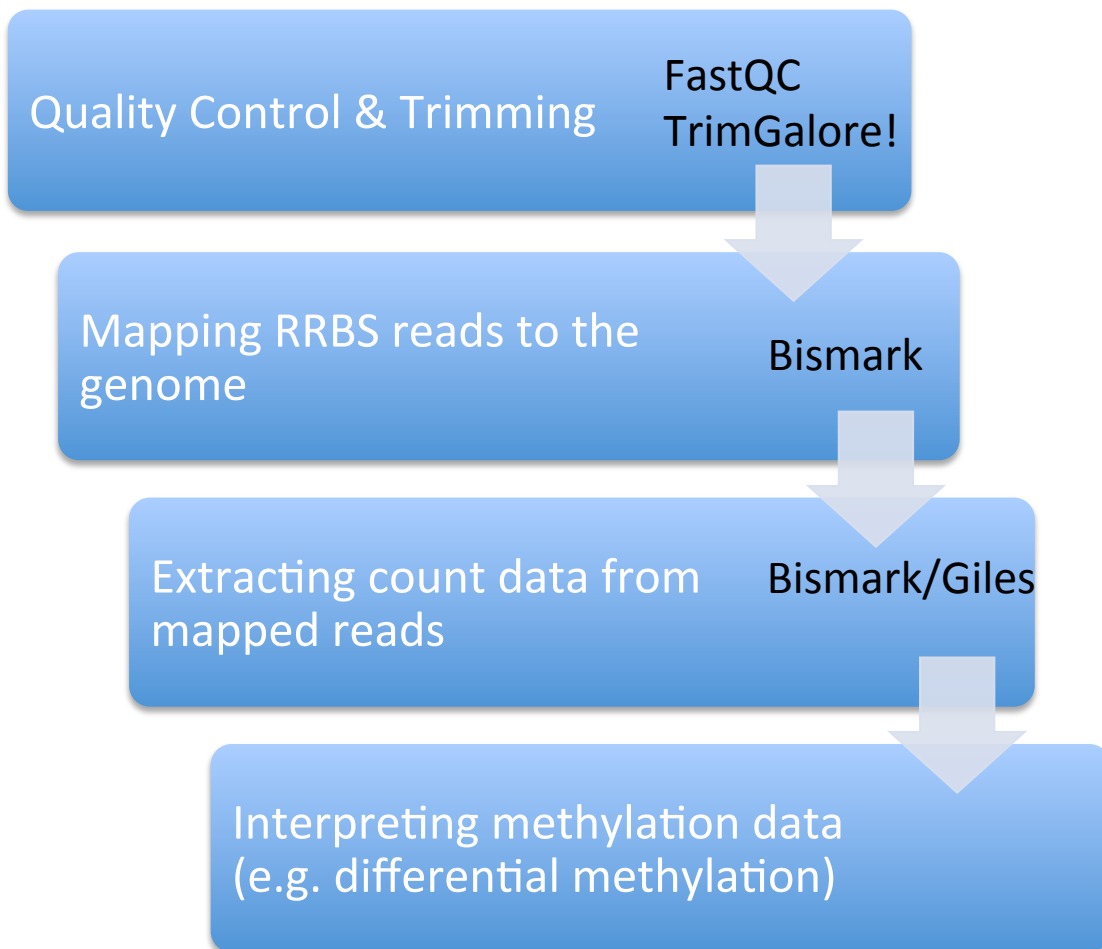
*O. mykiss* scaffold 13



# Bioinformatics (DNA methylation)

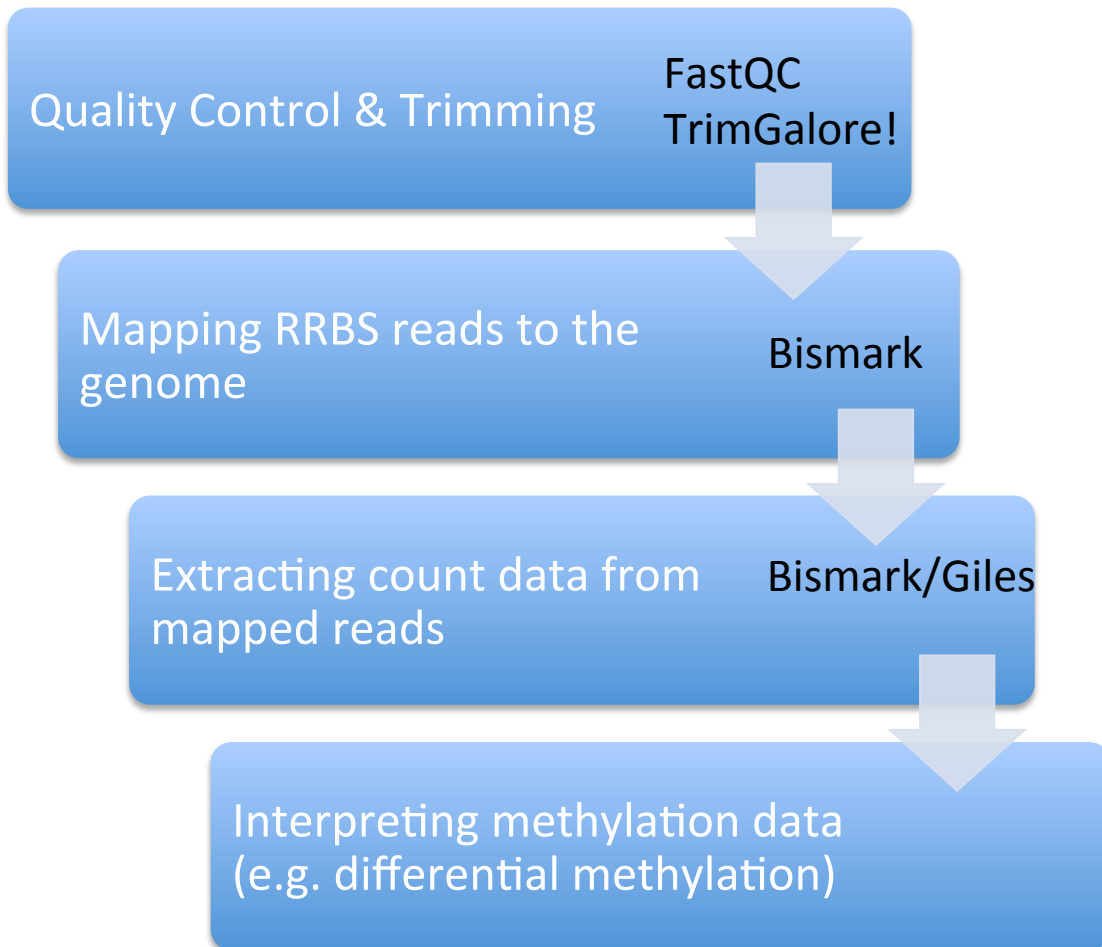


# Bioinformatics (DNA methylation)





# Bioinformatics (DNA methylation)



# Bioinformatics (DNA methylation)



Quality Control & Trimming

FastQC  
TrimGalore!

Mapping RRBS reads to the  
genome

Bismark

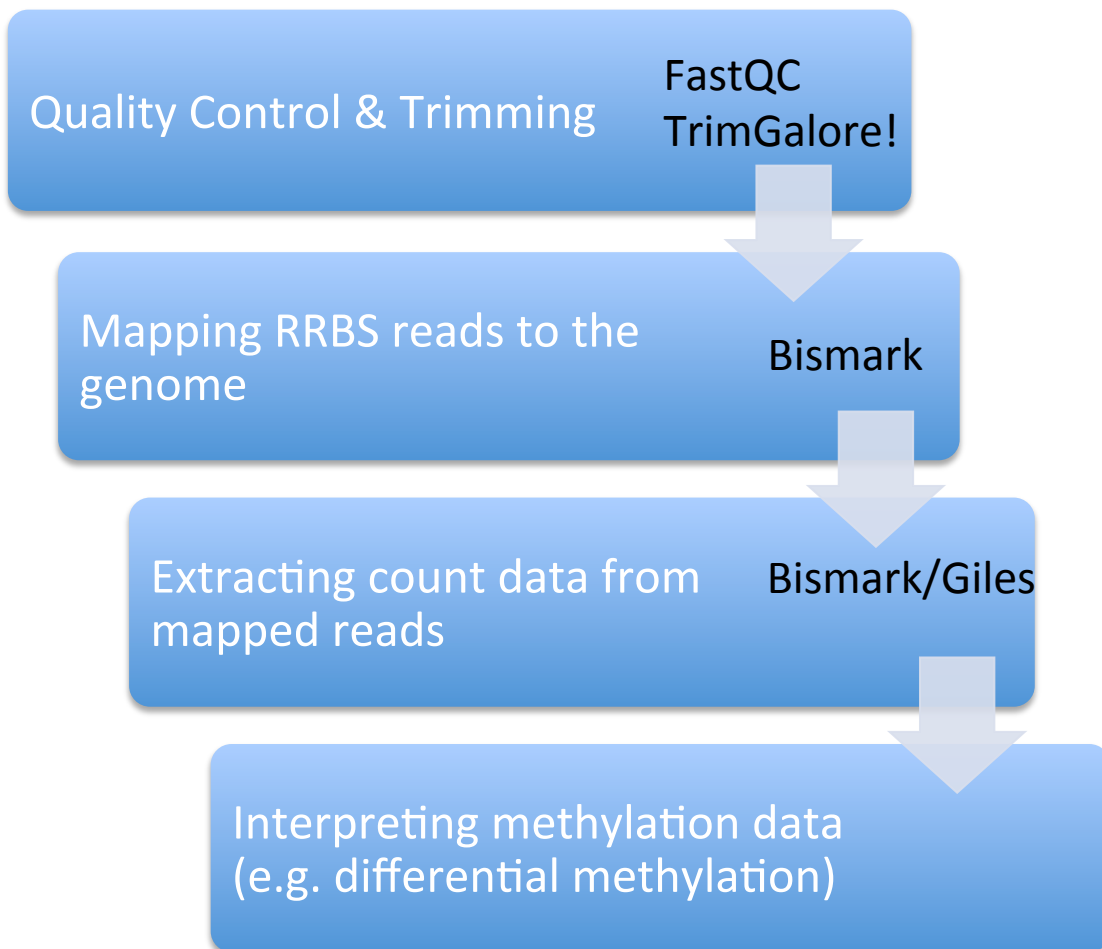
Extracting count data from  
mapped reads

Bismark/Giles

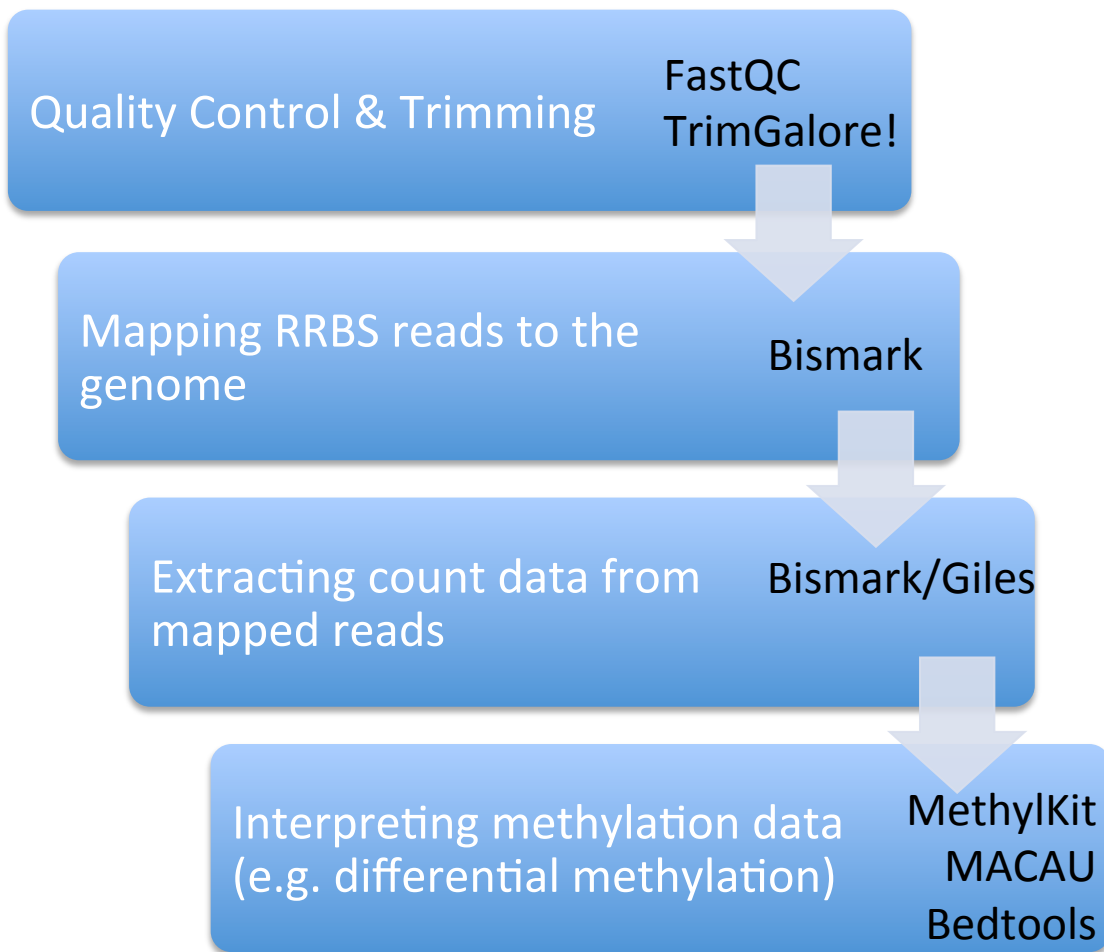
Interpreting methylation data  
(e.g. differential methylation)

chr	base	coverage	freqC	freqT
scaffold_1	5941	20	1.00	0.00
scaffold_1	5973	20	0.95	0.05
scaffold_1	5982	20	0.95	0.05
scaffold_1	5994	20	0.95	0.05
scaffold_1	5998	8	0.88	0.12
scaffold_1	6012	20	1.00	0.00
scaffold_1	6101	2	1.00	0.00
scaffold_1	6103	2	1.00	0.00
scaffold_1	6278	9	0.89	0.11
scaffold_1	6285	42	0.98	0.02

# Bioinformatics (DNA methylation)



# Bioinformatics (DNA methylation)

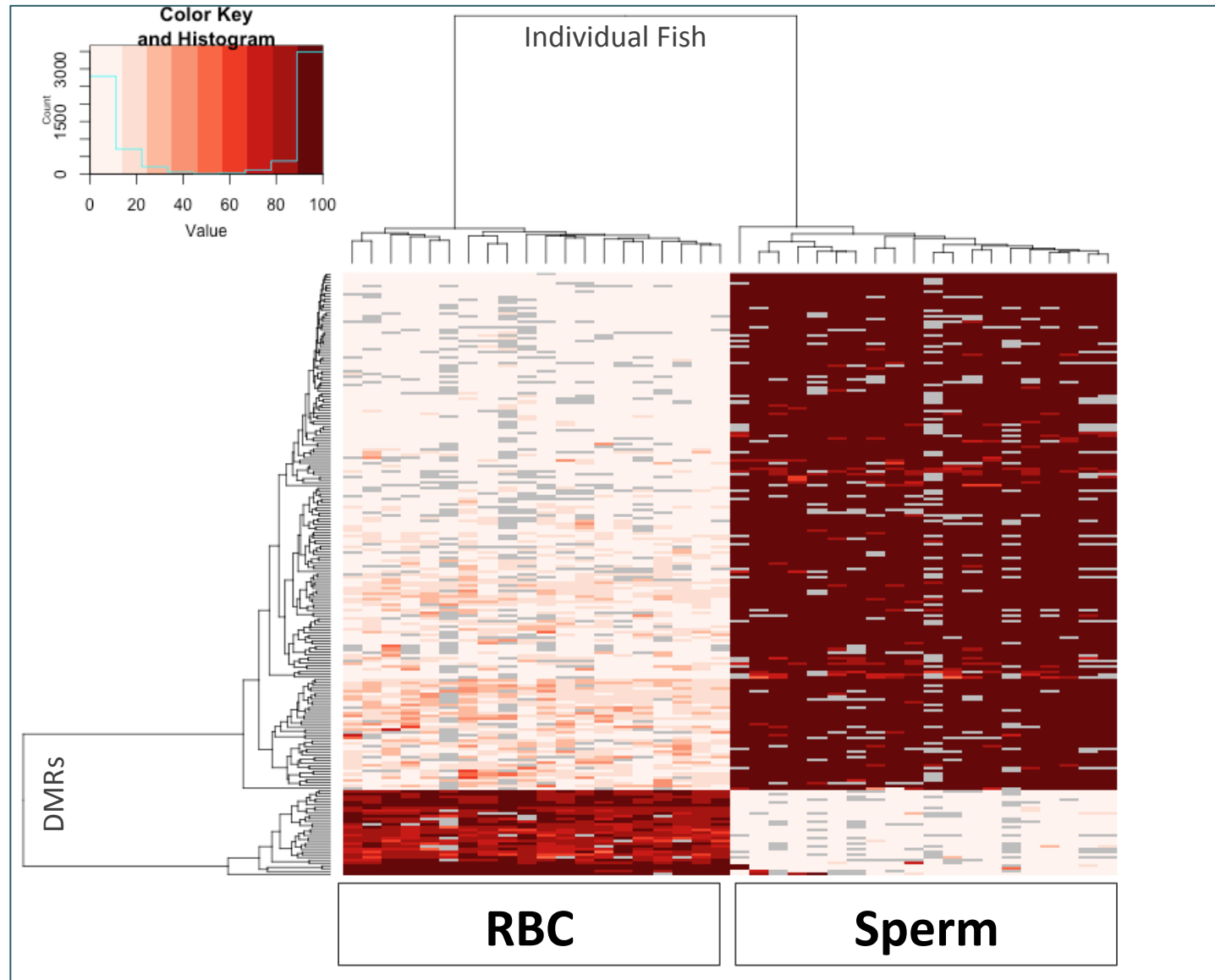


# Differential Methylation Analysis

# Differential Methylation Analysis

- MethylKit
  - R package
  - Logistic regression
  - Single variable (e.g. treatment v. control)
- MACAU
  - Beta-binomial mixed model
  - Multiple variables including relatedness

# Differential Methylation Analysis



# Descriptive Methylome Data

0bp

200,000bp

CG

genes

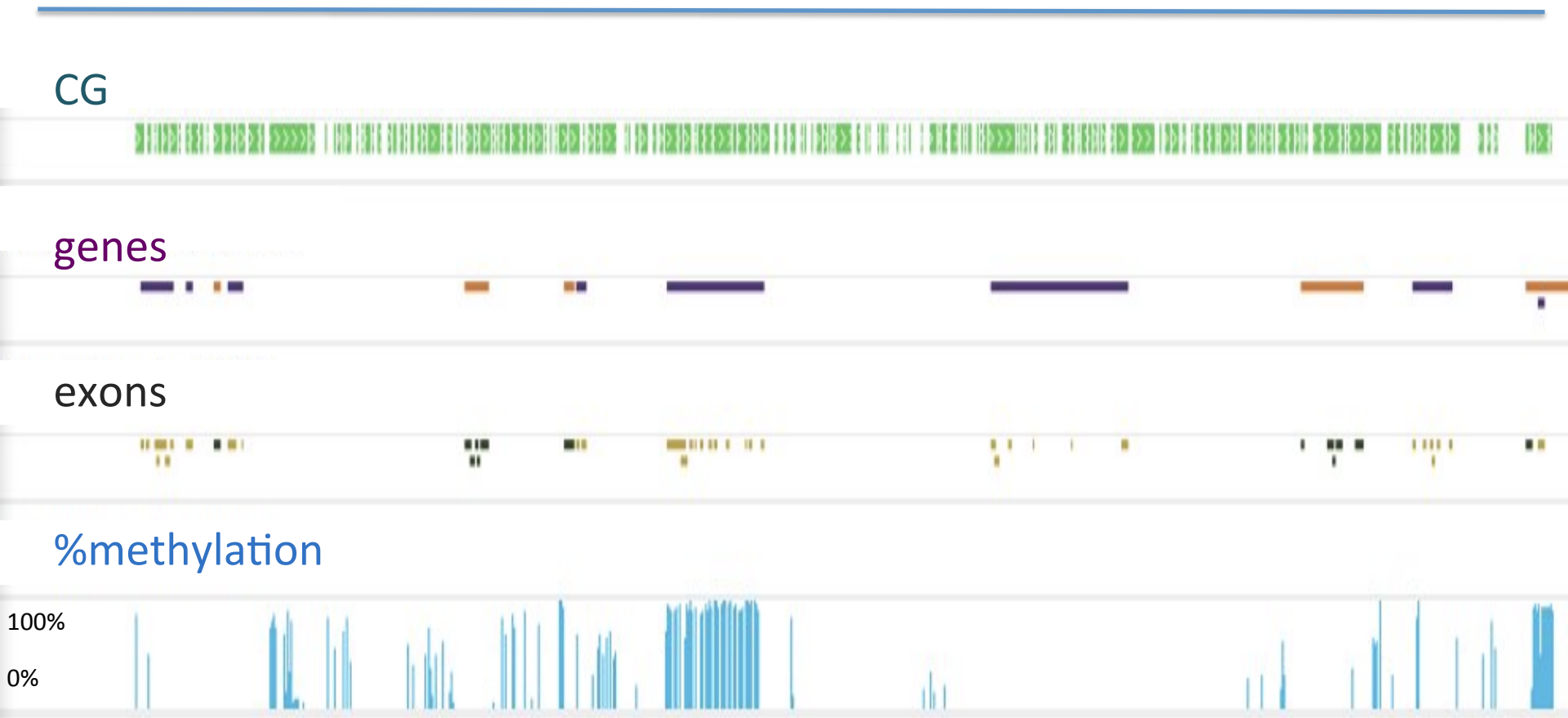
exons

%methylation

100%

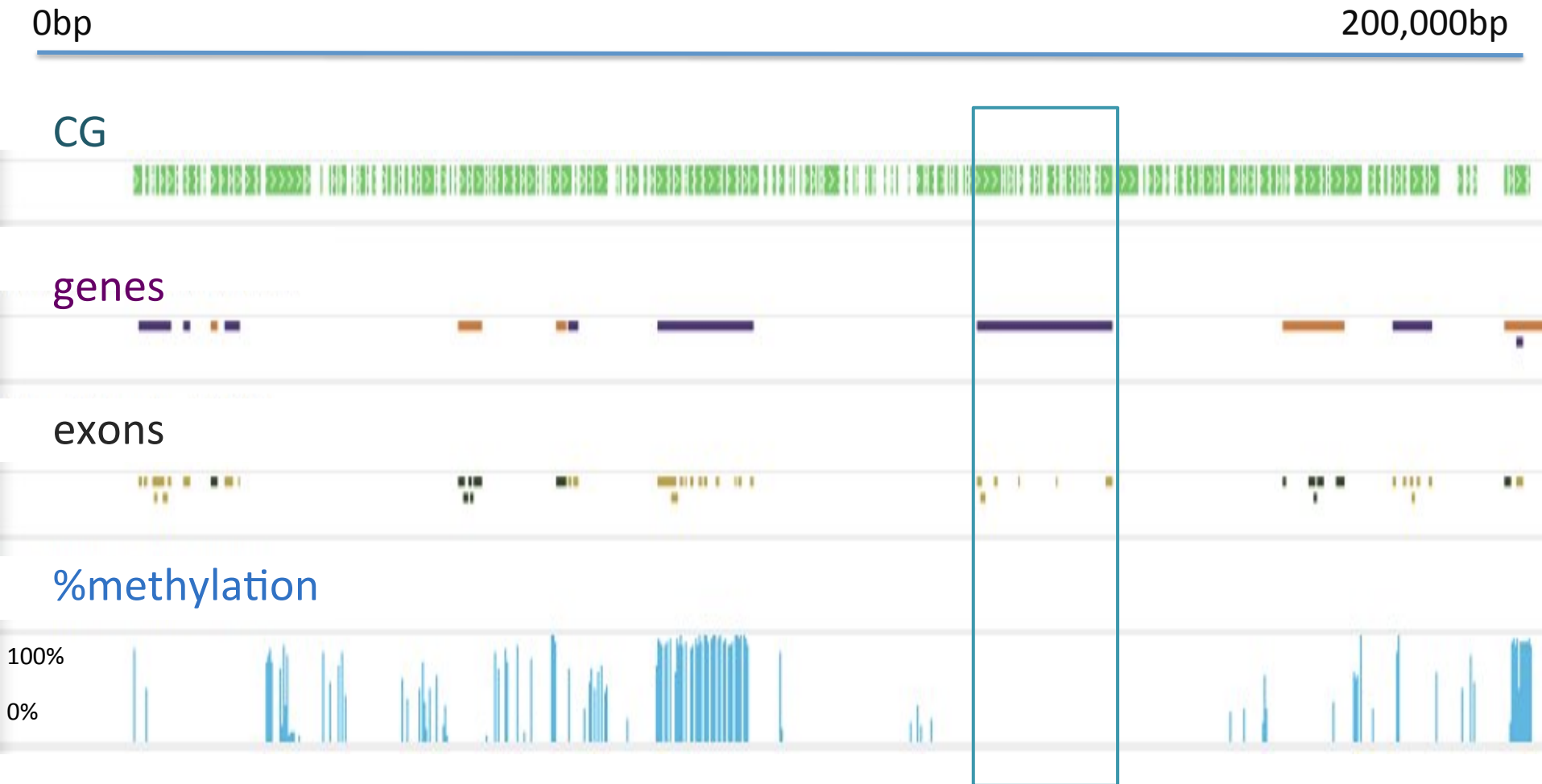
0%

(Galaxy Trackster)





# Descriptive Methylome Data



(Galaxy Trackster)

# Summary

- It can be hard to work with data that you can't **see**
  - Moving data takes a lot of time and space
  - Formatting data is a constant challenge
  - Software packages are great, but it's always important to understand what they are doing
  - Bisulfite sequencing has unique challenges