

Bisulfite-Sequencing Theory and Quality Control

Felix Krueger, Simon Andrews

felix.krueger@babraham.ac.uk

simon.andrews@babraham.ac.uk

v2022-04



Bisulfite-Seq theory and Quality Control

coffee

a.m.

Mapping and QC practical

Visualising and Exploring talk

Lunch

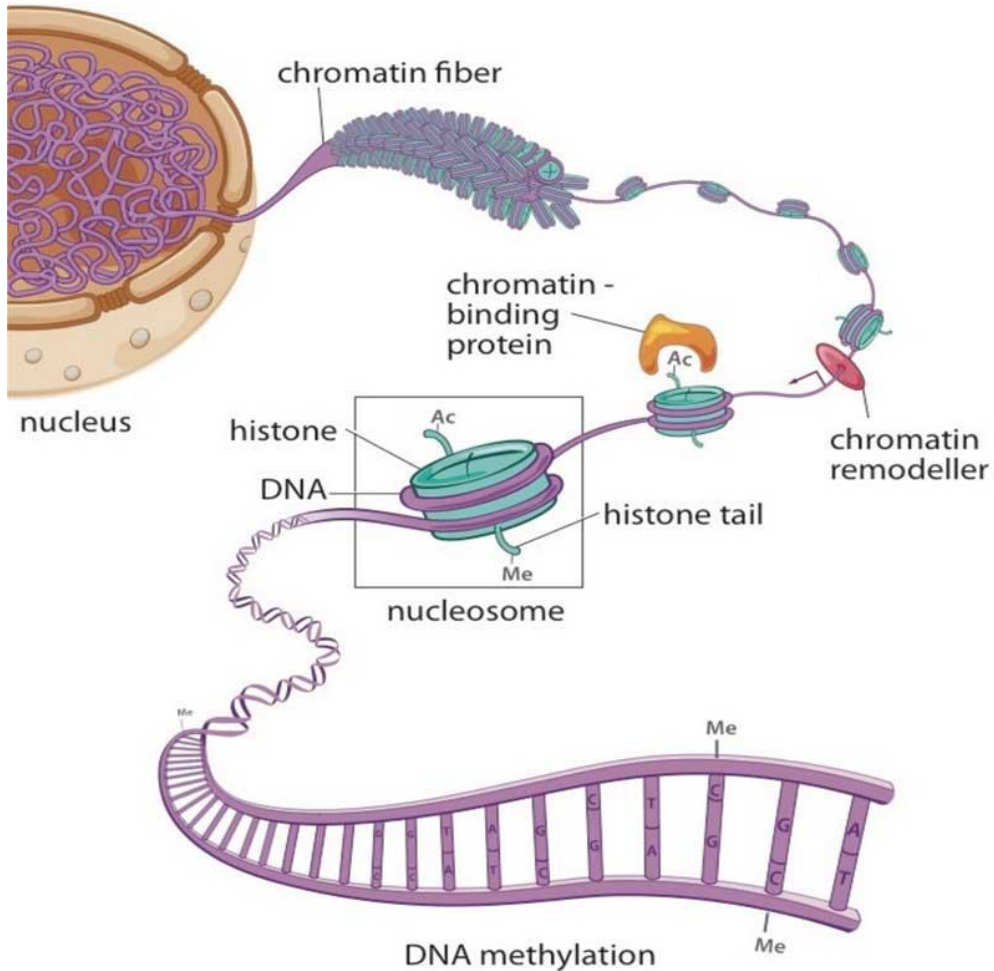
Visualising and Exploring practical

coffee

p.m.

Differential methylation talk & practical

Epigenetics



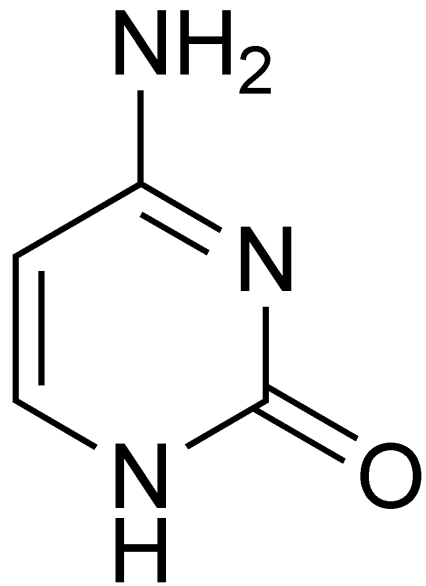
Studies changes in gene expression which are not encoded by the underlying DNA sequence

Chromatin

- histone modification
- non-coding RNAs
- higher order structure (accessibility/compaction)

DNA cytosine methylation

DNA Methylation

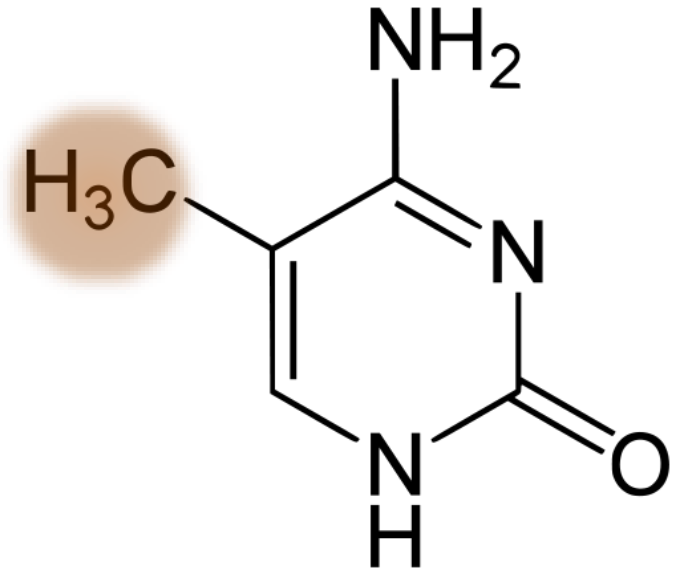


Cytosine

DNA methyl-transferases



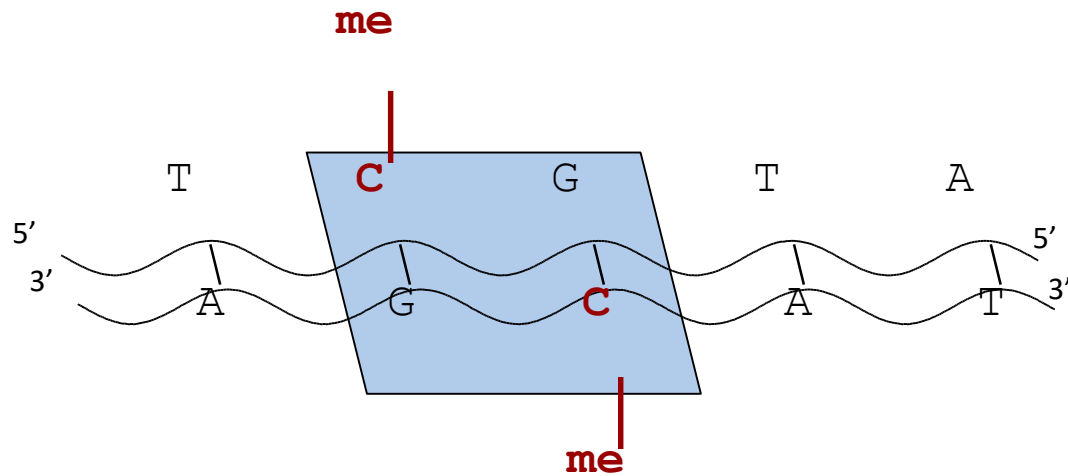
DNA-demethylase(s)
TET enzymes
Passive demethylation



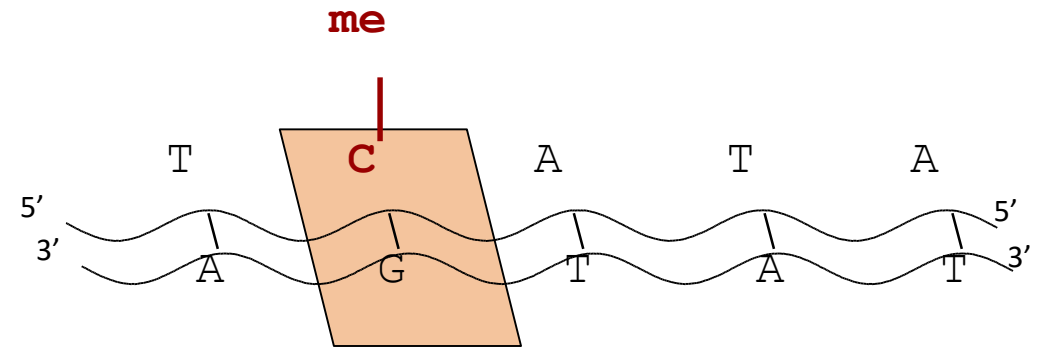
5-methyl Cytosine

Context of DNA methylation

CG context

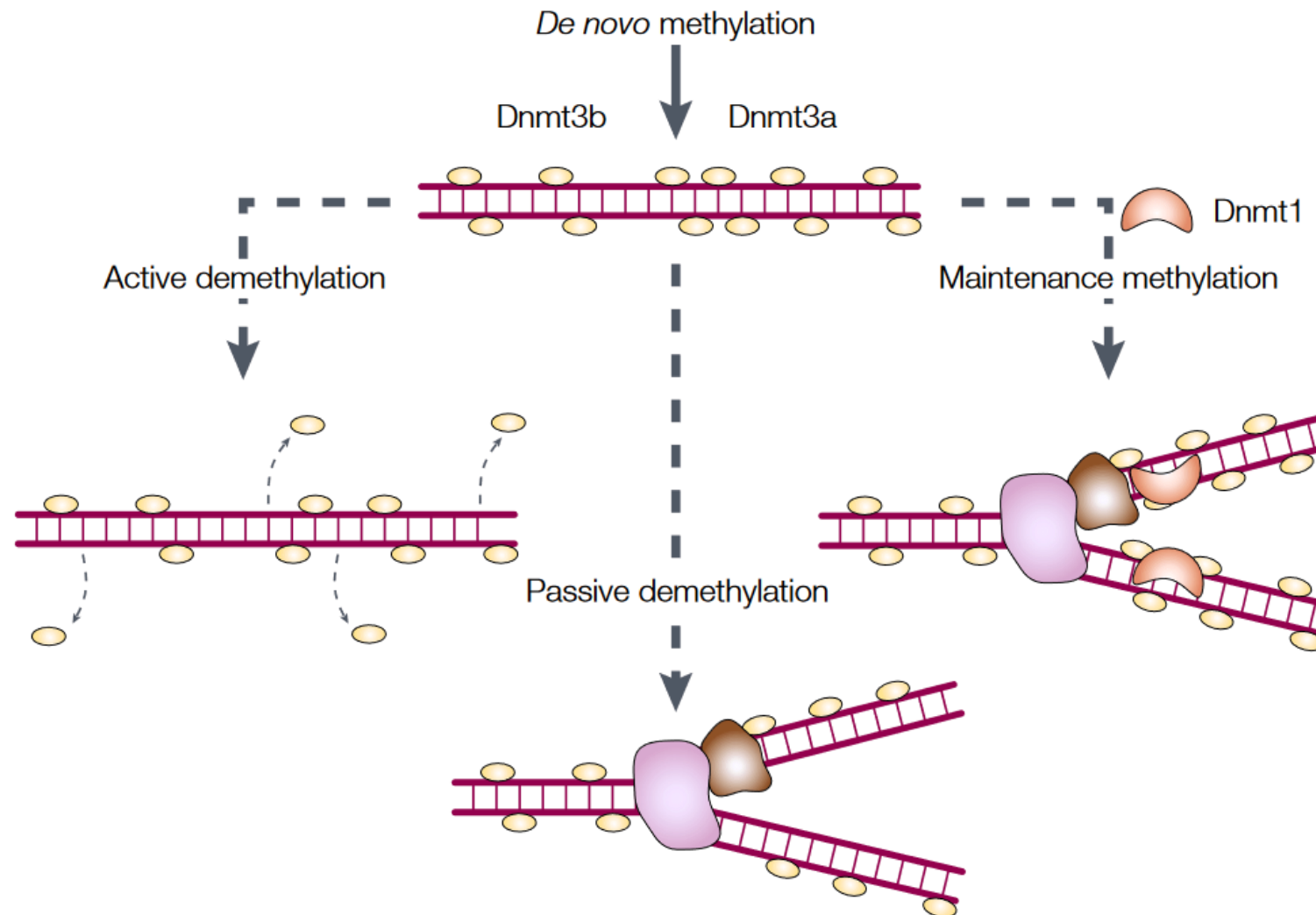


non-CG context

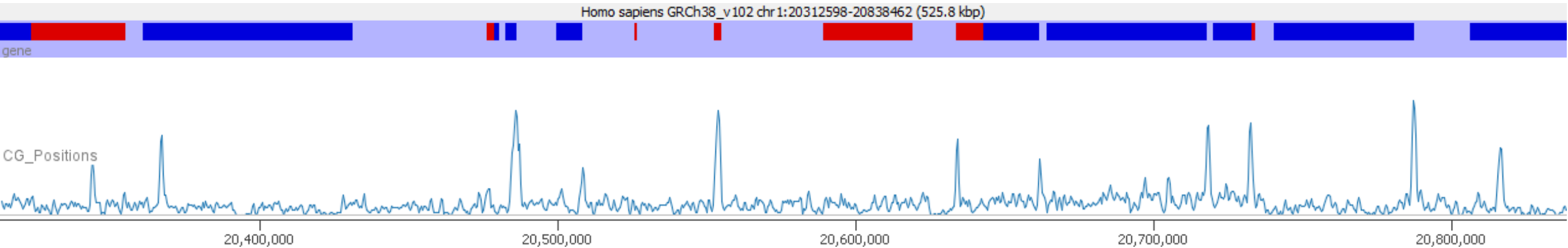


	Mammals	Plants
CG	present	present

DNA methylation is maintained through cell division

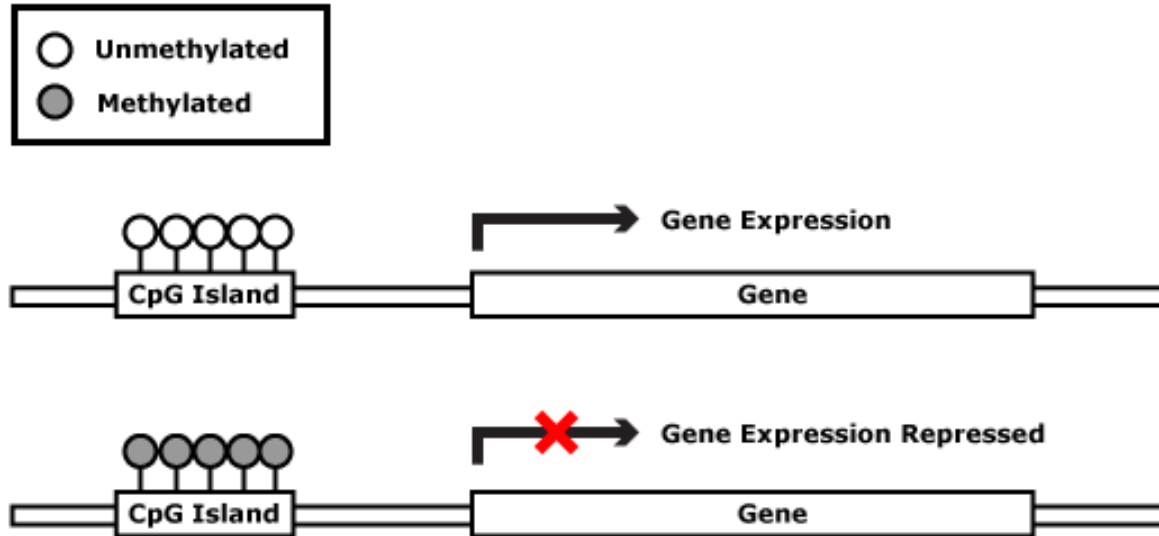


Distribution of CG



- CpG dinucleotides are not evenly distributed
- Most occur within high density regions called CpG Islands
- Most of the genome contains methylated CpGs
- CpGs in CpG islands are largely unmethylated

Regulation by DNA methylation

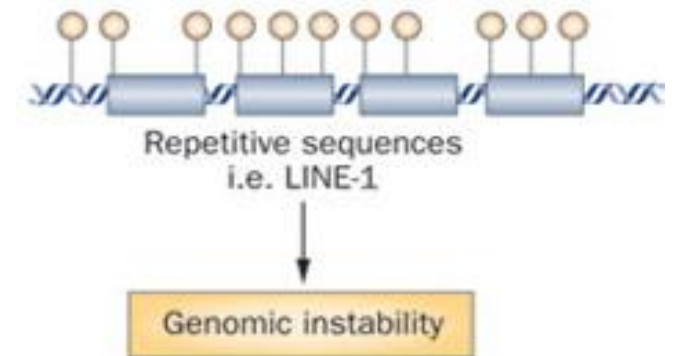
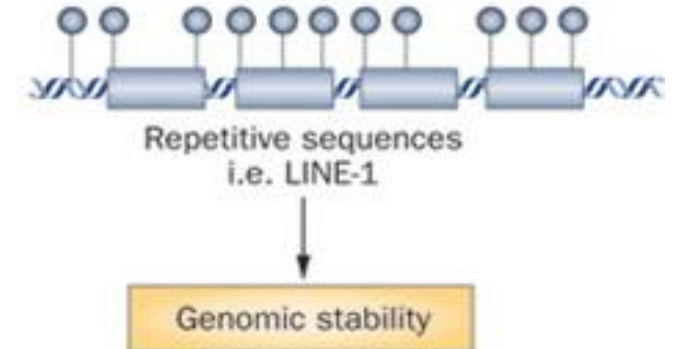


Silencing of gene expression

Tissue differentiation and embryonic development

Faults in correct DNA methylation may result in

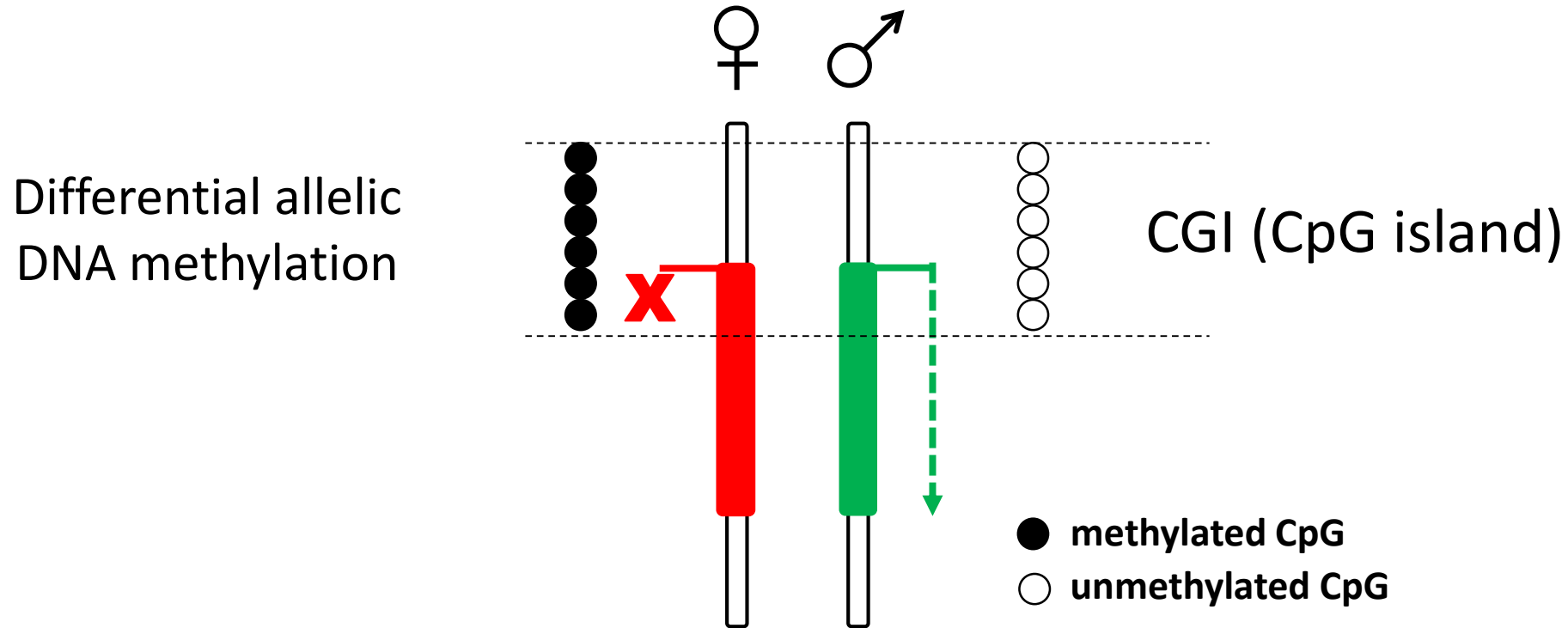
- early development failure
- epigenetic syndromes
- cancer



Repeat activity

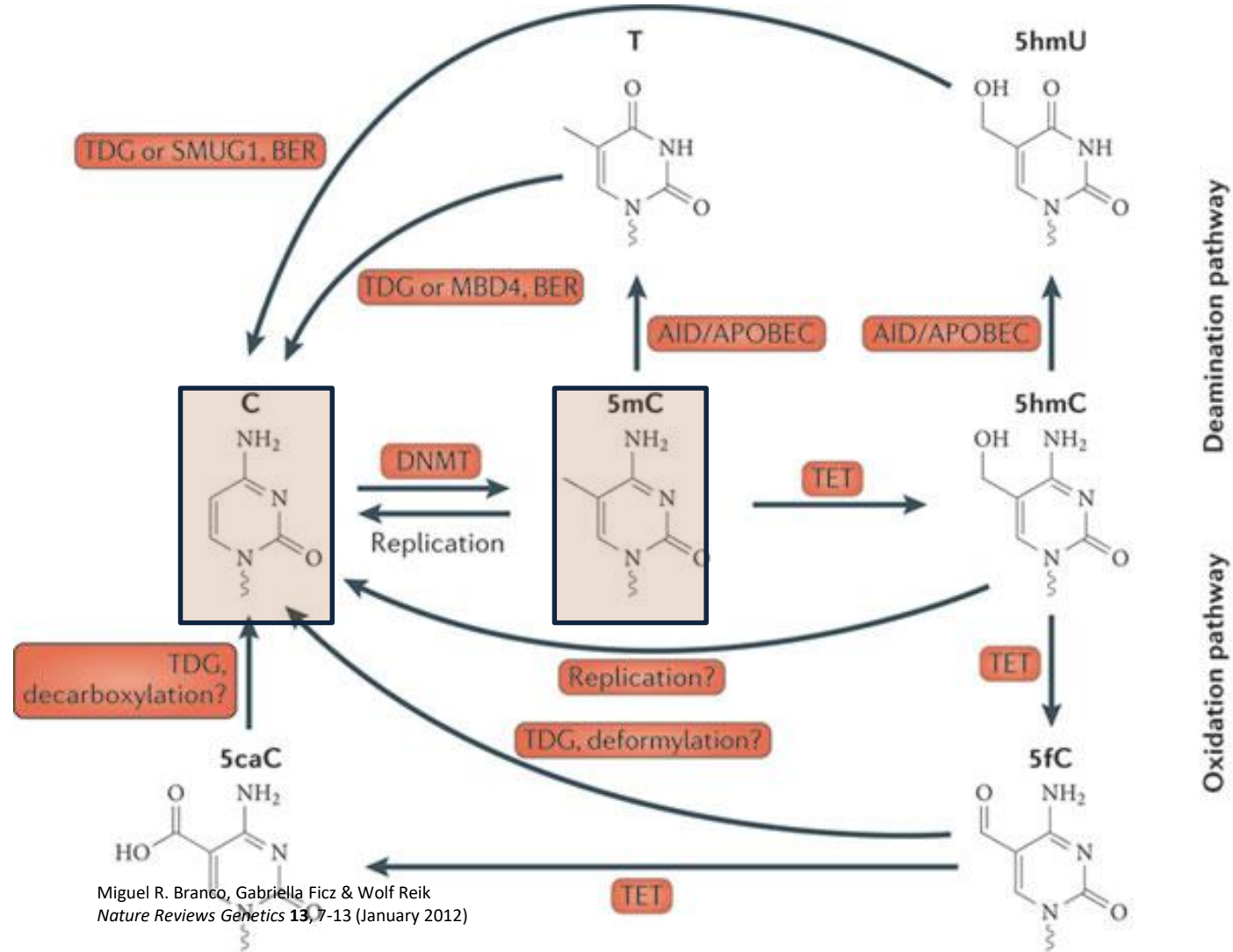
Genomic stability

Genomic Imprinting: mono-allelic expression



Imprinted Genes: Mono-allelic expression with parent-of-origin specificity. Have key roles in energy metabolism, placenta functions.

Cytosine Modifications



Miguel R. Branco, Gabriella Ficz & Wolf Reik
Nature Reviews Genetics **13**, 7-13 (January 2012)

Measuring DNA methylation by Bisulfite-sequencing

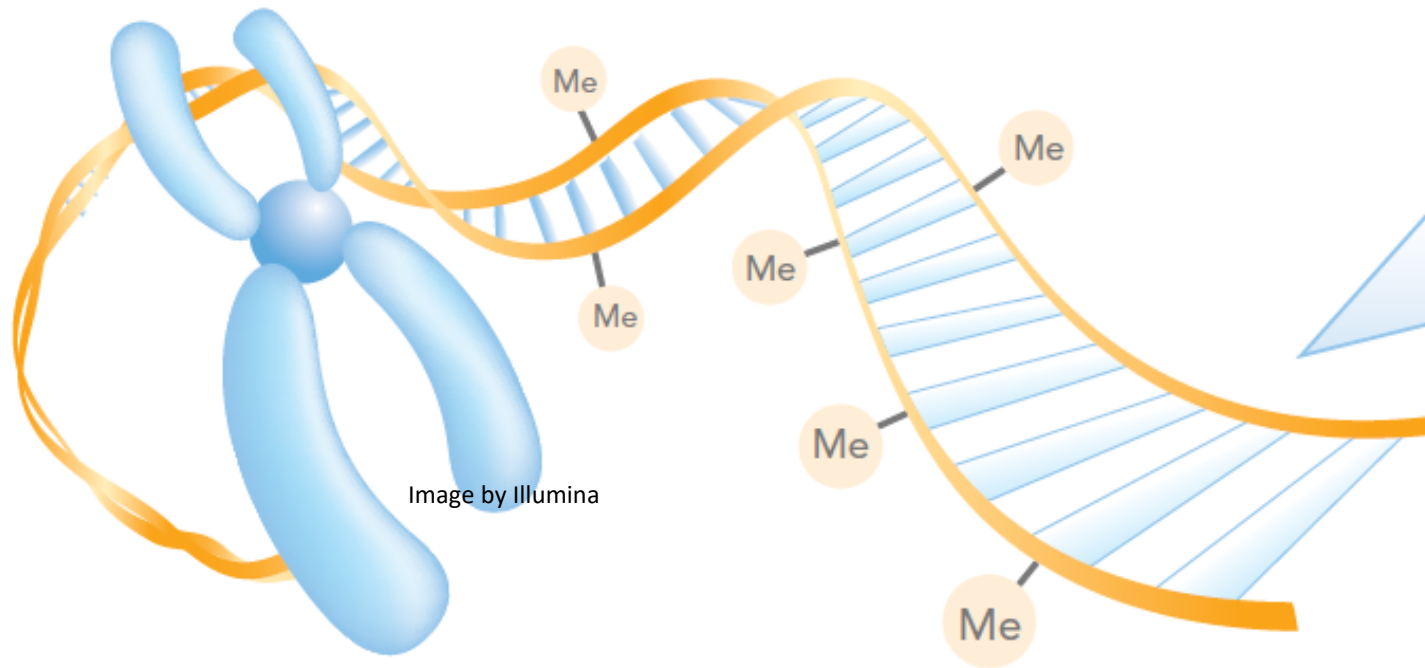
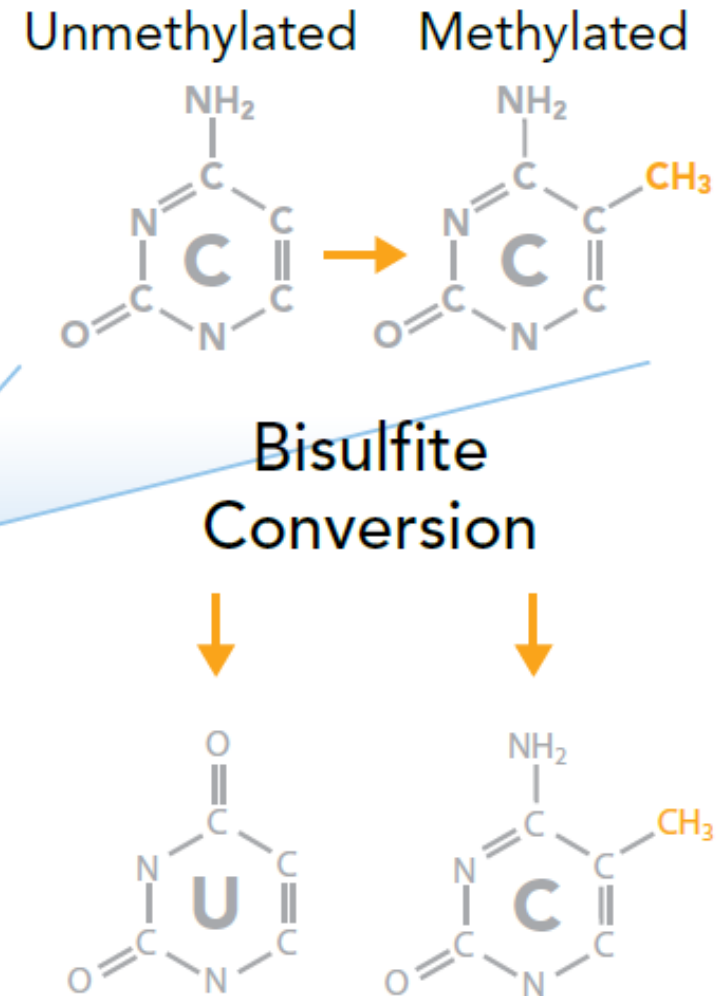
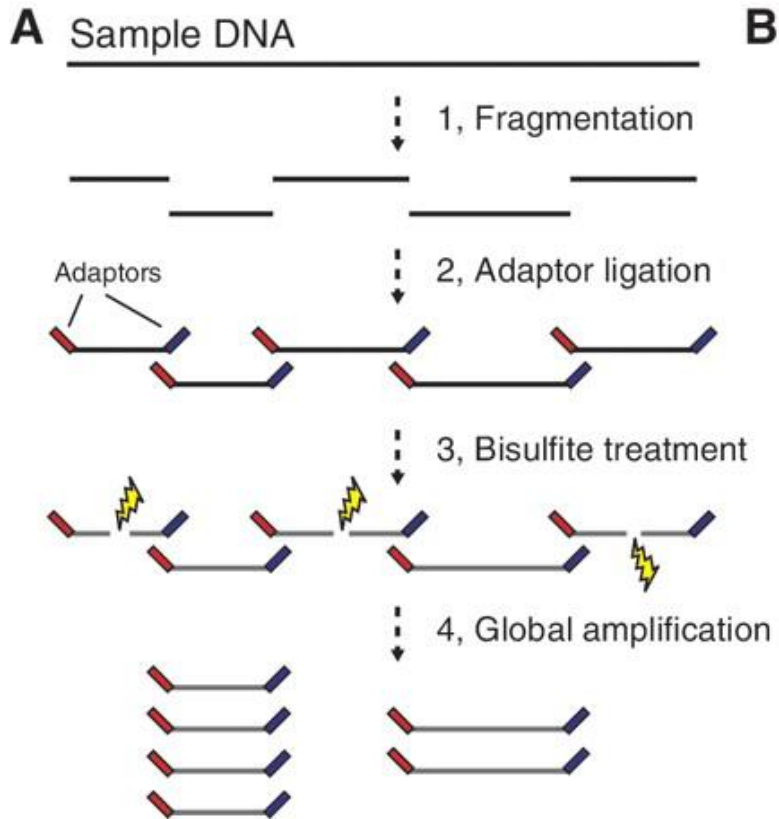


Image by Illumina

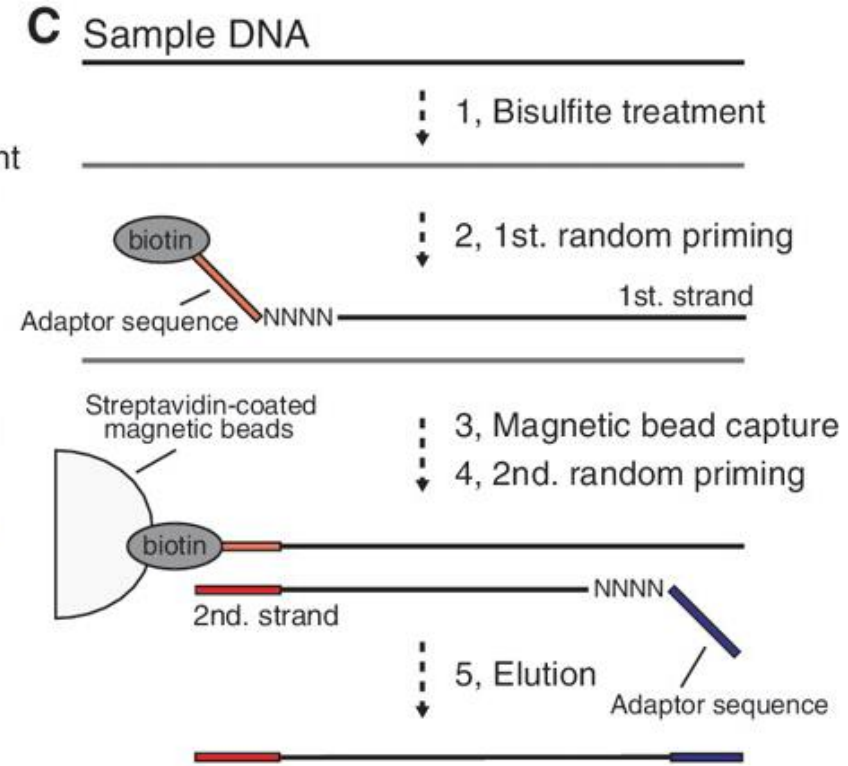
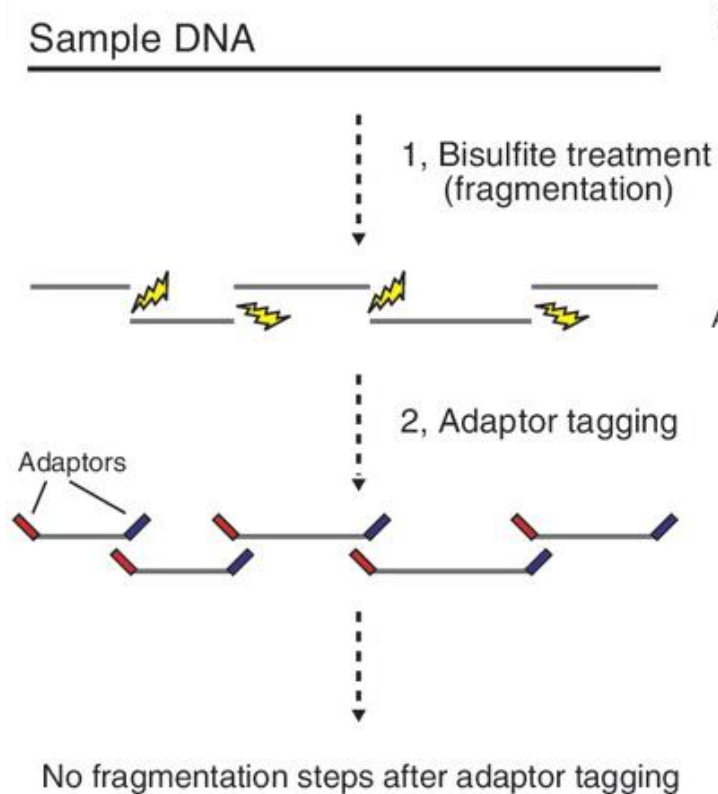


Bisulphite Library Preparation

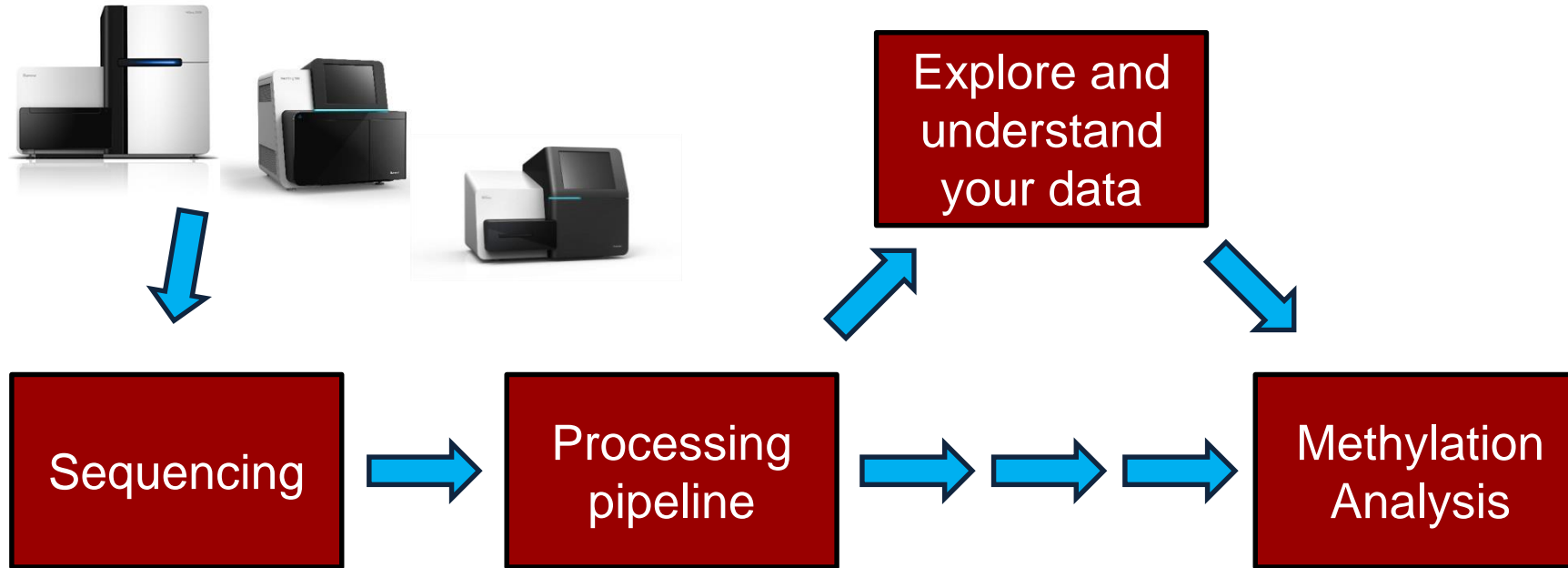
Whole Genome Bisulphite Sequencing (WGBS)



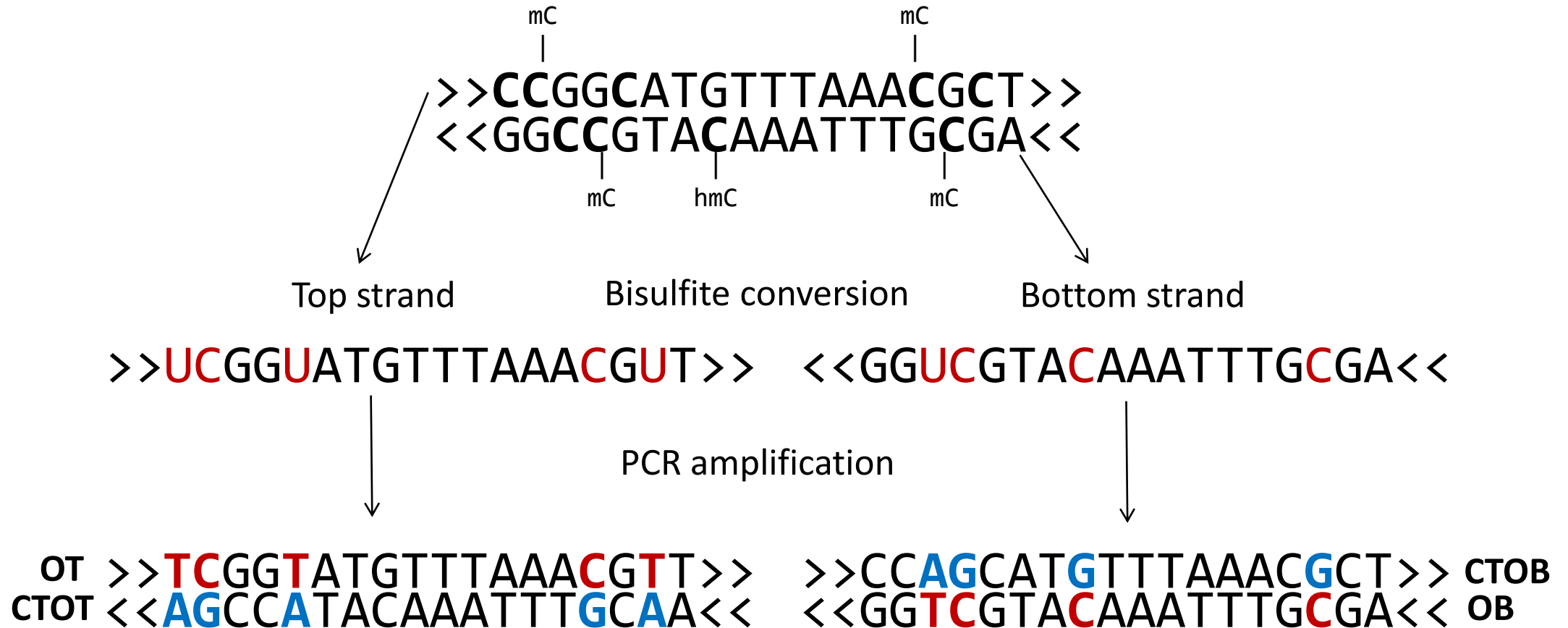
Post Bisulphite Adapter Tagging (PBAT)



BS-Seq Analysis Workflow



Bisulfite conversion of a genomic locus

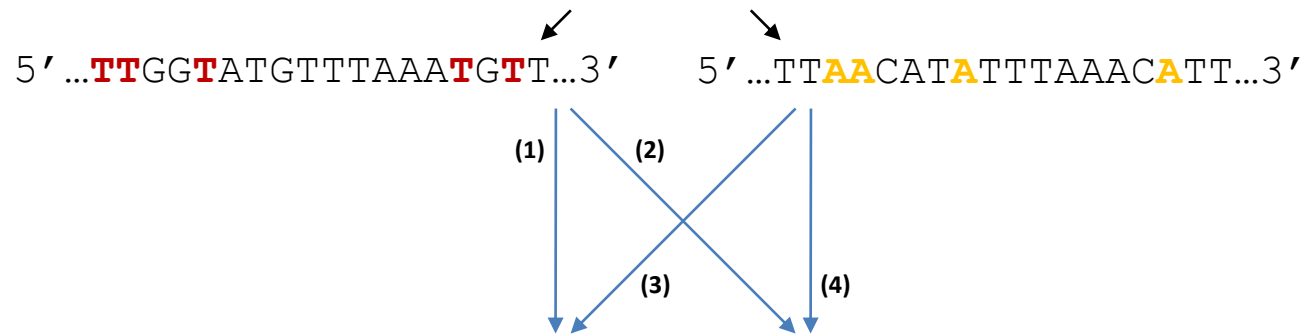


- 2 different PCR products and 4 possible different sequence strands from one genomic locus
- each of these 4 sequence strands can theoretically exist in any possible conversion state

3-letter alignment of Bisulfite-Seq reads

sequence of interest

TTGGCATGTTTAAACGTT



...TTGGTATGTTTAAATGTT...

...CCAACATATTTAAACACT...

...AACCATACAAATTTACAA...

...GGTTGTATAAATTTGTGA...

forward strand C -> T converted genome

forward strand G -> A converted genome
(equals reverse strand C -> T conversion)

(1) (2) (3) (4)

5' ...CCGGCATGTTTAAACGCT...3'

read sequence TTGGCATGTTTAAACGTTA

genomic sequence CCGGCATGTTTAAACGCTA

methylation call xz . . **H** **Z** . h . .

Fully bisulfite convert read
(as both forward and reverse strand)

Align to bisulfite converted genomes

Read all 4 alignment outputs and extract
the unmodified genomic sequence if the
sequence could be mapped uniquely

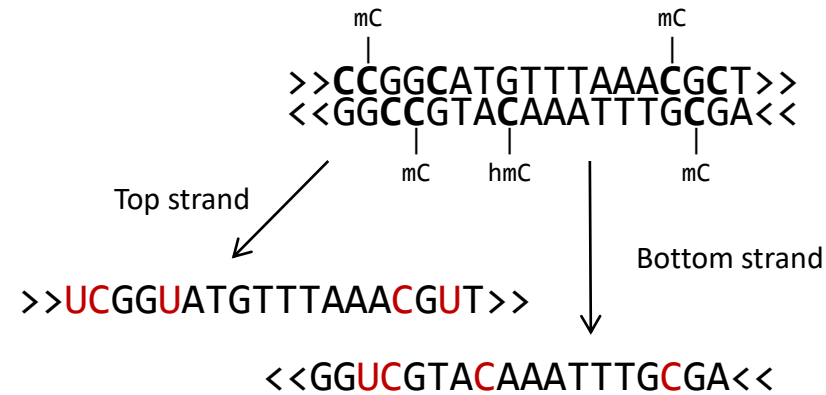
Methylation Call

h unmethylated C in CHH context
H methylated C in CHH context
x unmethylated C in CHG context
X methylated C in CHG context
z unmethylated C in CpG context
Z methylated C in CpG context



Bismark

Common sequencing protocols



1) Directional libraries

(vast majority of kits, also EpiGnome/Truseq)

OT >>TCGGTATGTTTAAACGTT>>
 <<GGTCGTACAAATTTGCGA<< OB

2) PBAT libraries

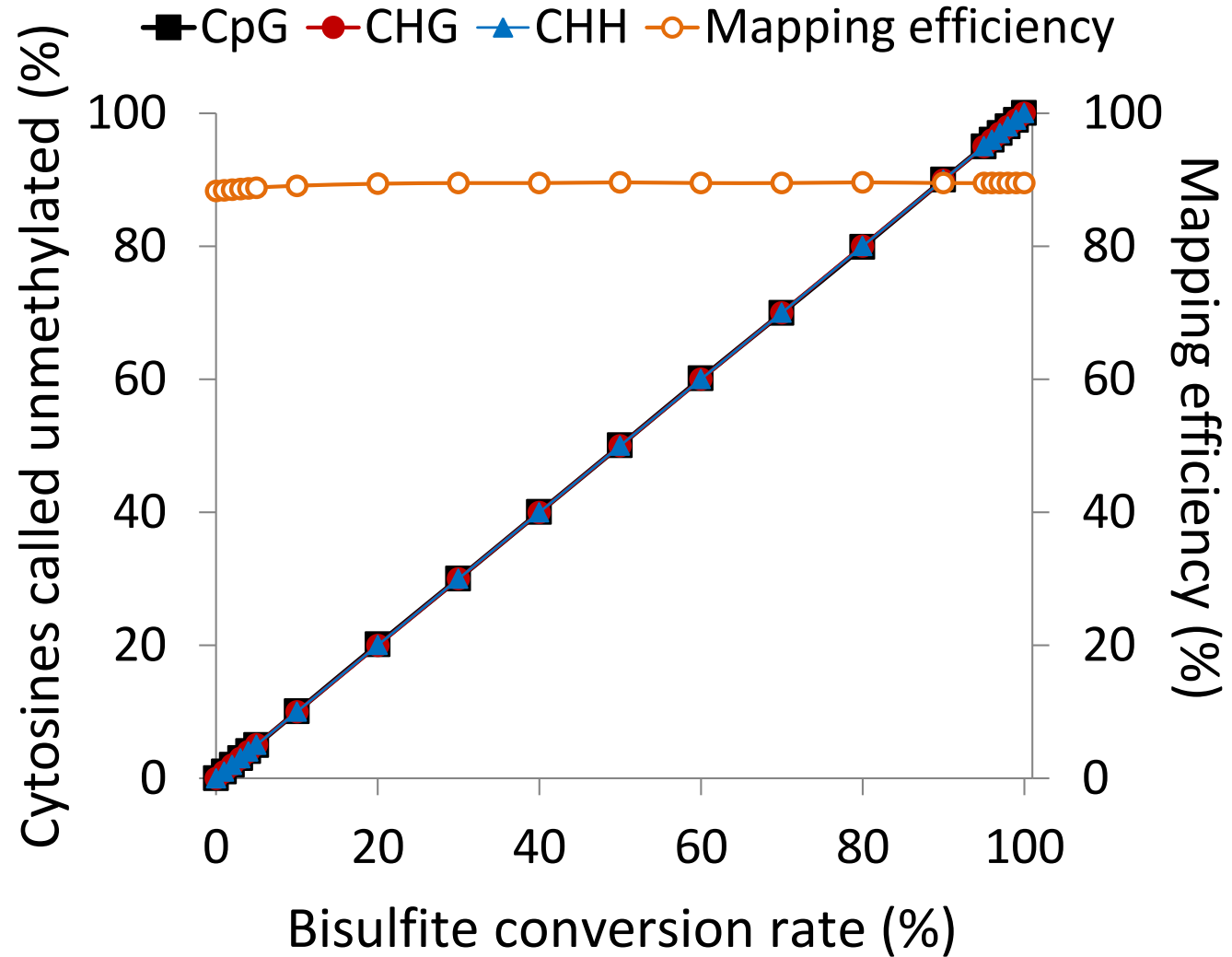
CTOT <<AGCCATACAAATTTGCAA<<
 >>CCAGCATGTTTAAACGCT>> CTOB

3) Non-directional libraries

(e.g. single-cell BS-Seq, Zymo Pico Methyl-Seq)

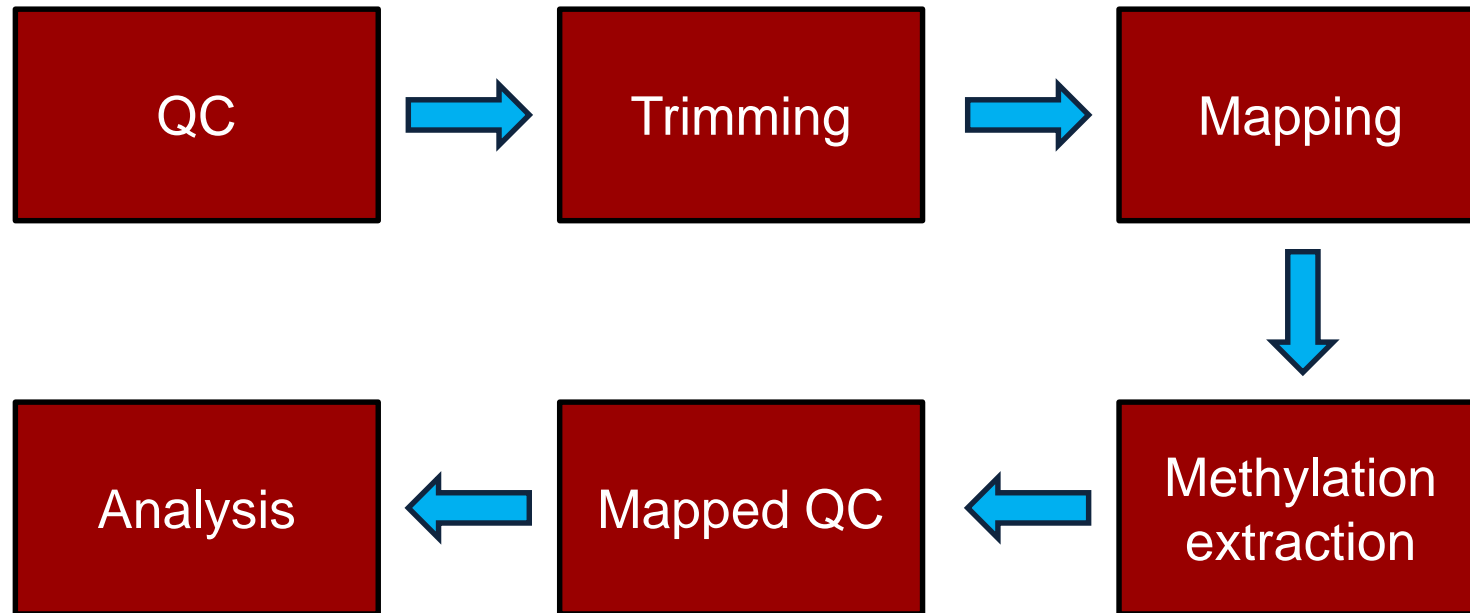
OT >>TCGGTATGTTTAAACGTT>>
 CTOT <<AGCCATACAAATTTGCAA<<
 >>CCAGCATGTTTAAACGCT>> CTOB
 <<GGTCGTACAAATTTGCGA<< OB

Validation





BS-Seq Analysis Workflow



Part I: Initial QC - What does QC tell you about your library?

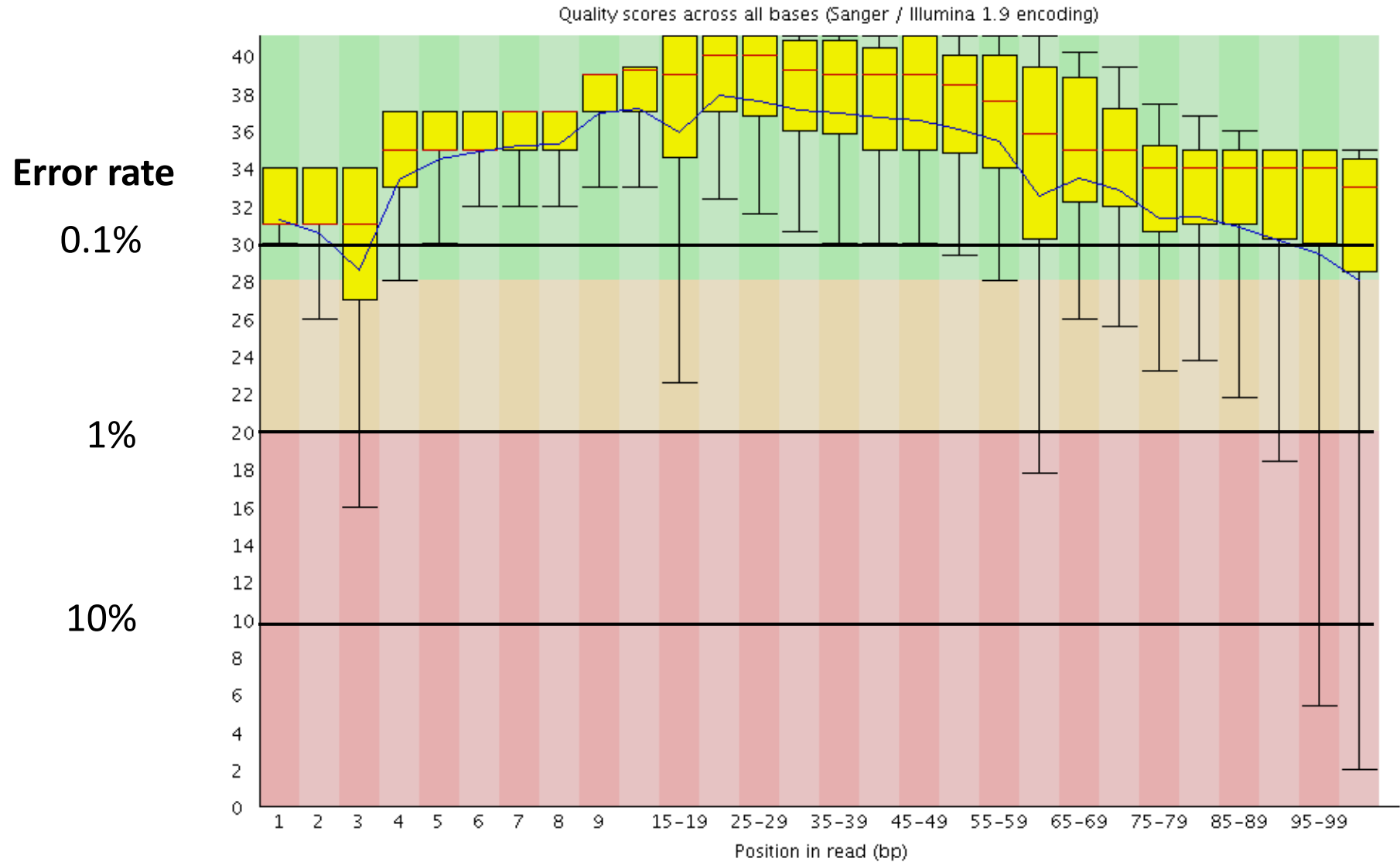
- # of sequences
- Basecall qualities
- Base composition
- Potential contaminants
- Expected duplication rate



Basic Statistics

Measure	Value
Filename	s_4_1_sequence.txt
File type	Conventional base calls
Encoding	Illumina 1.5
Total Sequences	35290120
Sequence length	40
%GC	46

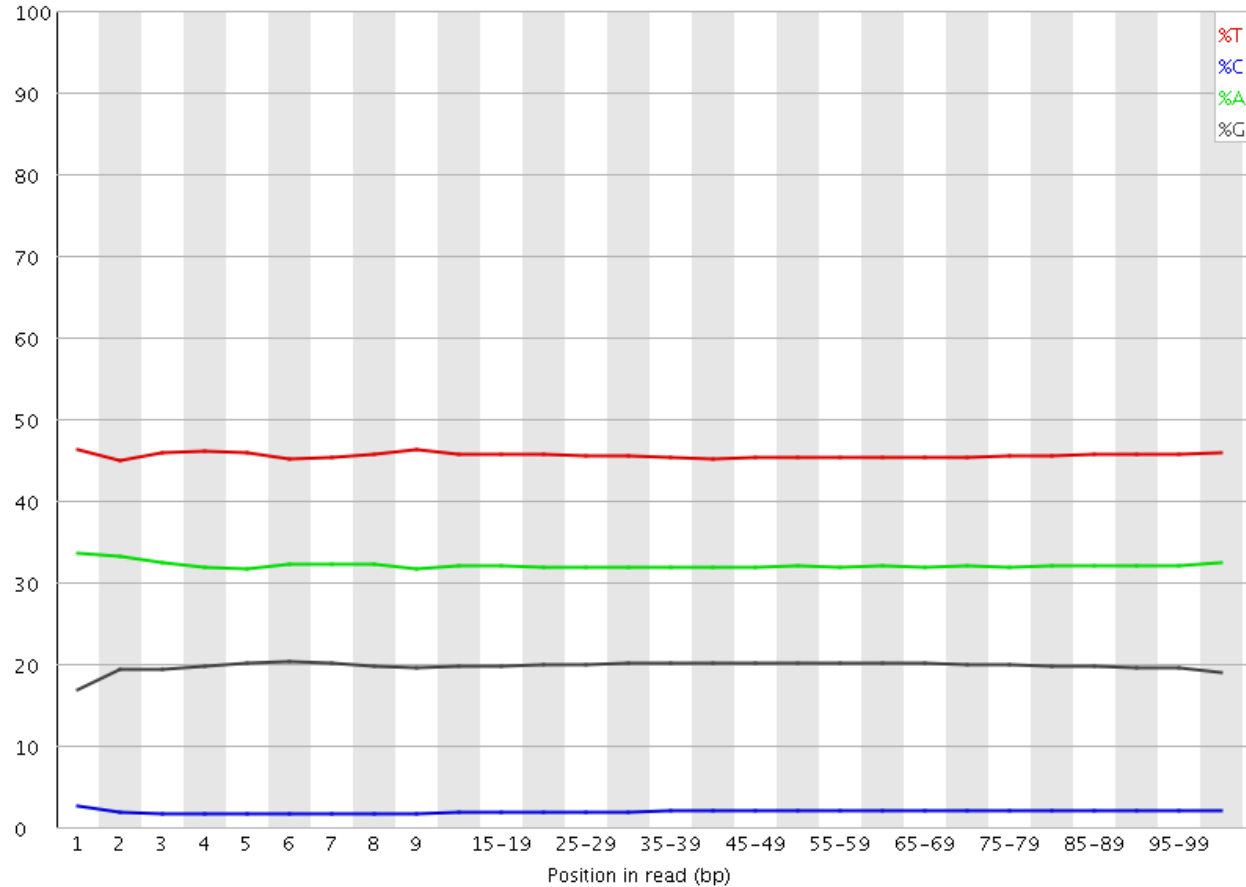
QC Raw data: Sequence Quality



QC: Base Composition

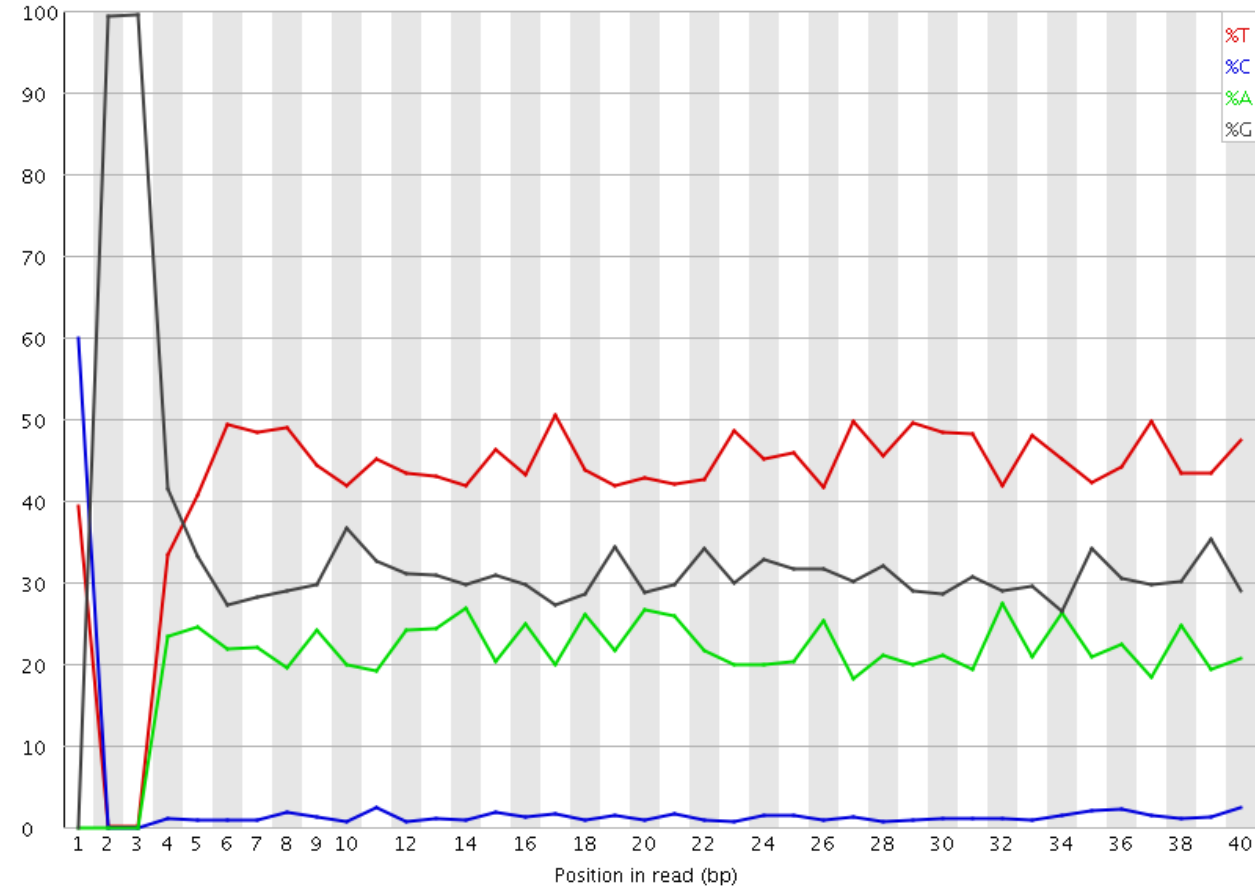
WGSBS

Sequence content across all bases

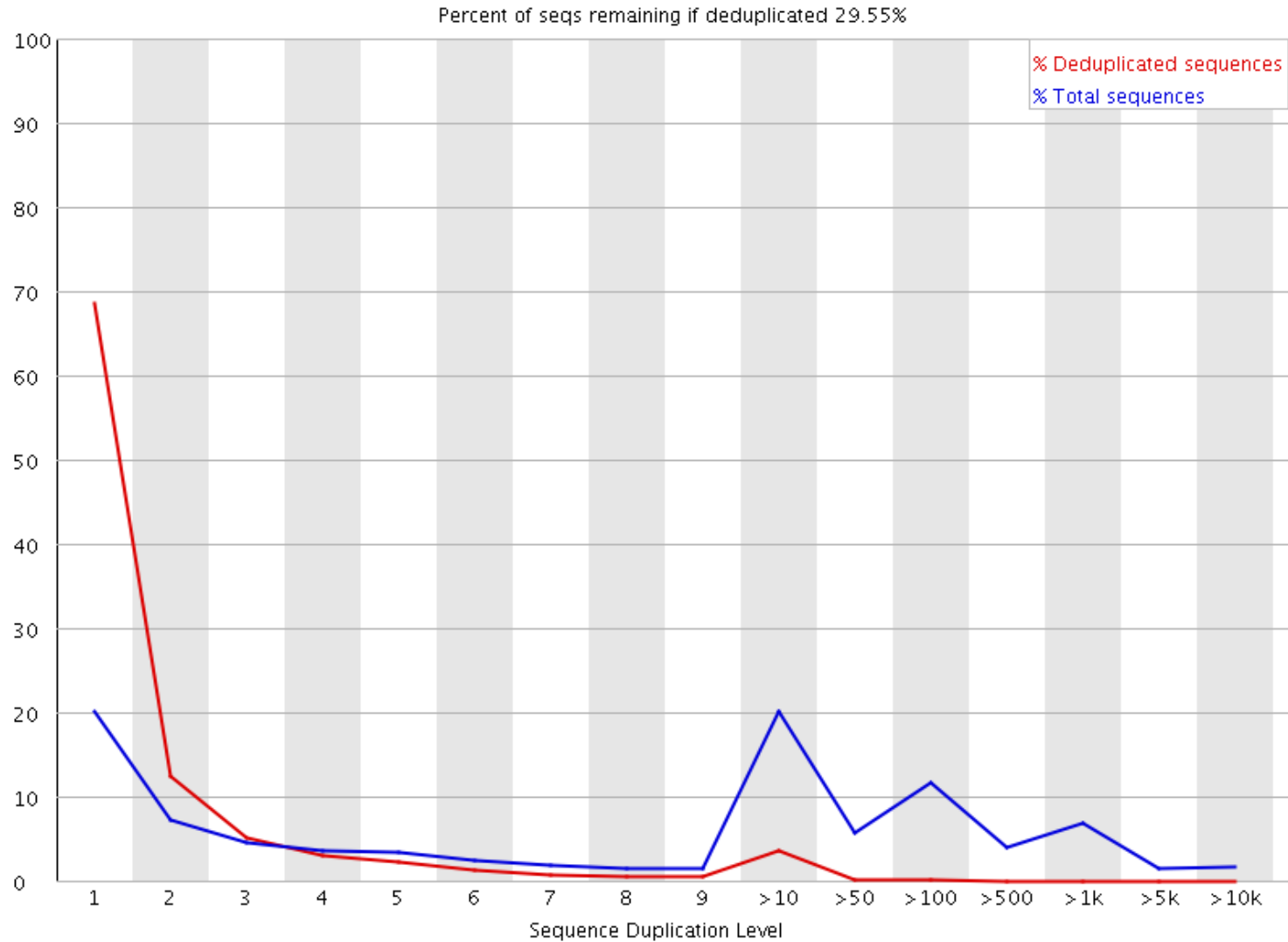


RRBS

Sequence content across all bases



QC: Duplication rate

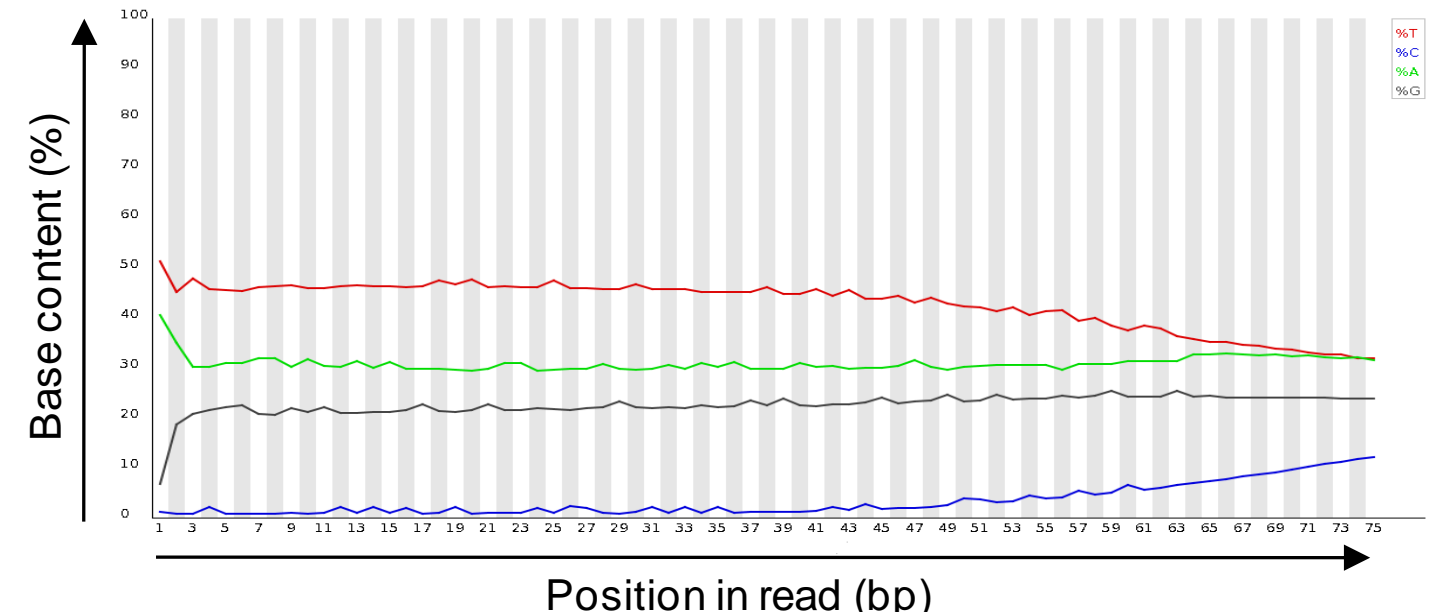
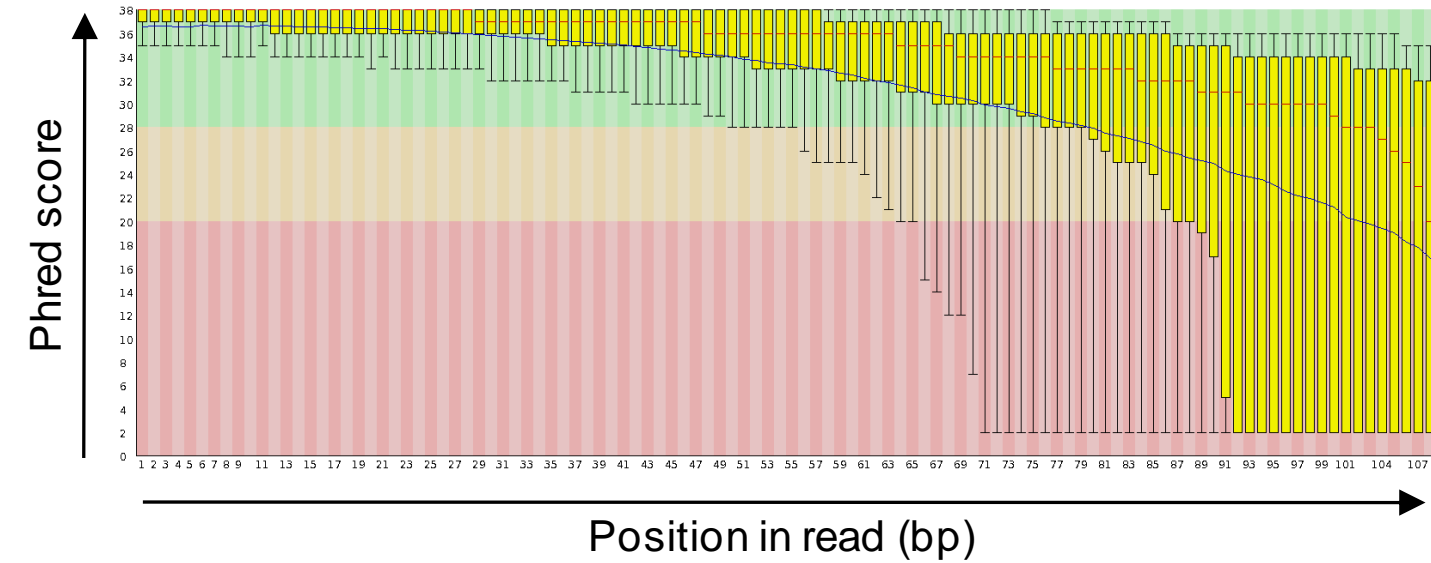


QC: Overrepresented sequences

Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GAAGAGCGGTTTCAGCAGGAATGCCGAGACCGATCTCGTAT	6254891	23.52739098691508	Illumina Paired End PCR Primer 2 (100% over 40bp)
GATCGGAAGAGCGGTTTCAGCAGGAATGCCGAGACCGATCT	1956005	7.357393503317777	Illumina Paired End PCR Primer 2 (100% over 40bp)
GAAGAGCGGTTTCAGCAGGAATGCCGAGATCGGAAGAGCGG	774763	2.9142237687587667	Illumina Paired End PCR Primer 2 (96% over 31bp)
GAAGAGCGGTTTCAGCAGGAATGCCGAGGATCGGAAGAGCG	140148	0.5271581538405985	Illumina Paired End Adapter 2 (100% over 27bp)
AAGAGCGGTTTCAGCAGGAATGCCGAGATCGGAAGAGCGGT	105720	0.3976593317352233	Illumina Paired End PCR Primer 2 (96% over 30bp)
NAAGAGCGGTTTCAGCAGGAATGCCGAGACCGATCTCGTAT	98639	0.37102458213233724	Illumina Paired End PCR Primer 2 (97% over 40bp)
AAGAGCGGTTTCAGCAGGAATGCCGAGACCGATCTCGTATG	82413	0.30999147281777295	Illumina Paired End PCR Primer 2 (100% over 40bp)
GATCGGAAGAGCGGTTTCAGCAGGAATGCCGAGATCGGAAG	53872	0.20263624214188372	Illumina Paired End PCR Primer 2 (97% over 36bp)
NNAGAGCGGTTTCAGCAGGAATGCCGAGACCGATCTCGTAT	36541	0.137446742725471	Illumina Paired End PCR Primer 2 (100% over 38bp)
ATCGGAAGAGCGGTTTCAGCAGGAATGCCGAGACCGATCTC	35781	0.13458804908076072	Illumina Paired End PCR Primer 2 (100% over 40bp)
CGGAAGAGCGGTTTCAGCAGGAATGCCGAGACCGATCTCGT	33905	0.1275315895051338	Illumina Paired End PCR Primer 2 (100% over 40bp)
NATCGGAAGAGCGGTTTCAGCAGGAATGCCGAGACCGATCT	30564	0.1149646217854272	Illumina Paired End PCR Primer 2 (97% over 40bp)
GAAGAGCGGTTTCAGCAGGAATGCCGAGACGGATCTCGTAT	28274	0.10635092646123442	Illumina Paired End PCR Primer 2 (97% over 40bp)
CAAACAACCTTCTAAAACAAACAAAAACACAAAACCACTAA	27952	0.10513974310123876	No Hit

Common problems in BS-Seq

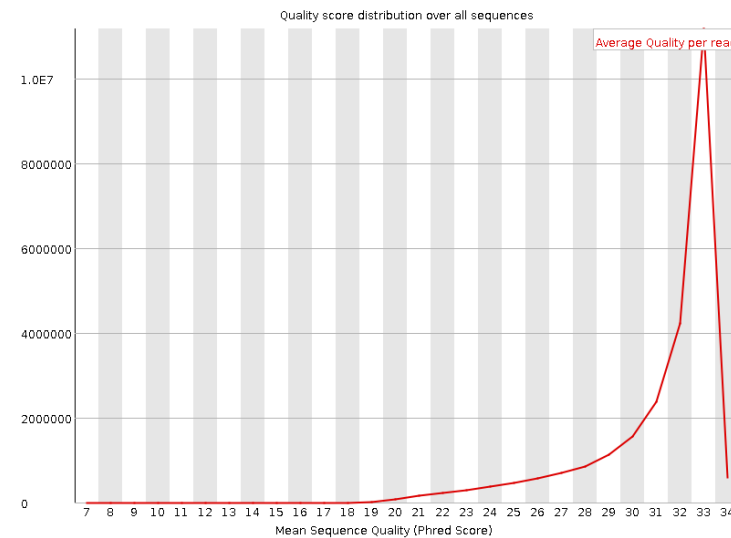
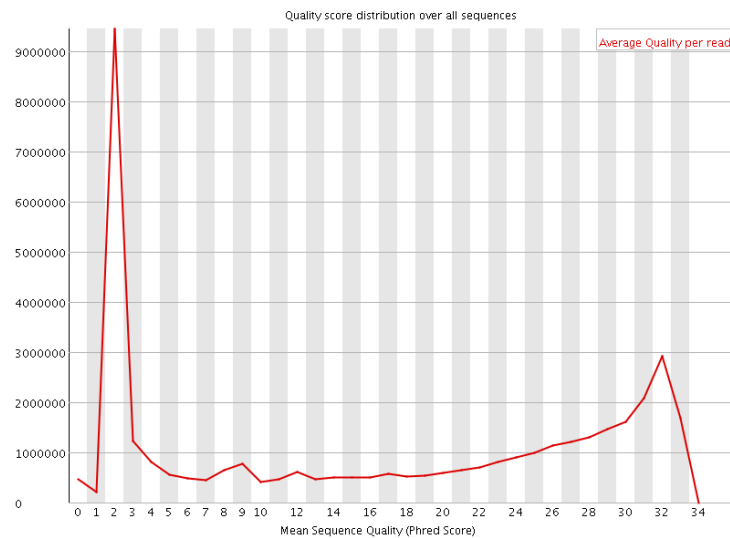
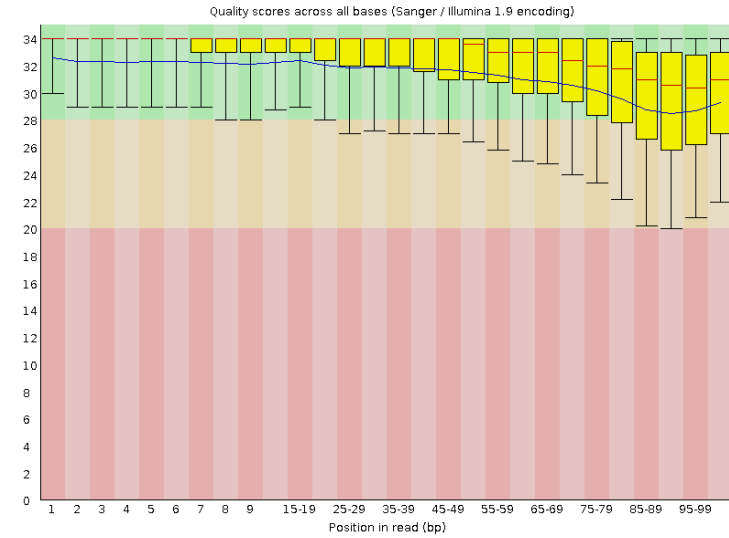
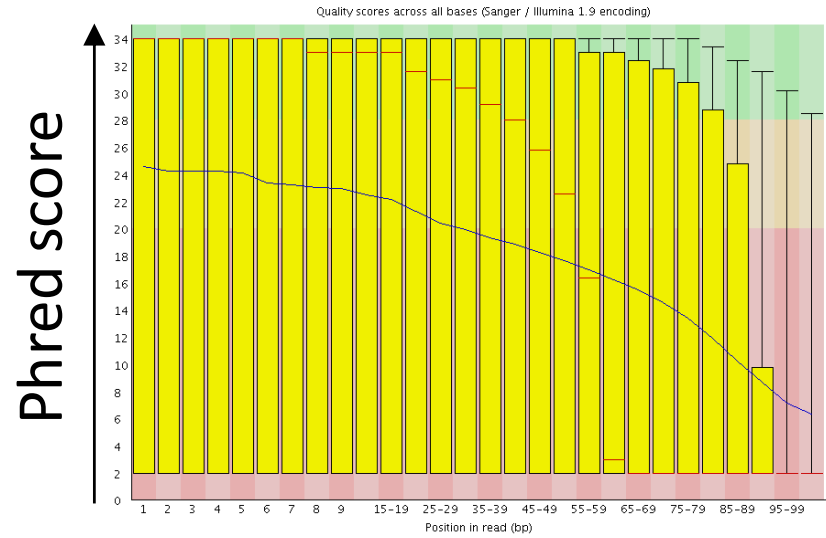


Not observed
in 'normal' libraries,
e.g. ChIP or RNA-Seq

Removing poor quality basecalls

before trimming

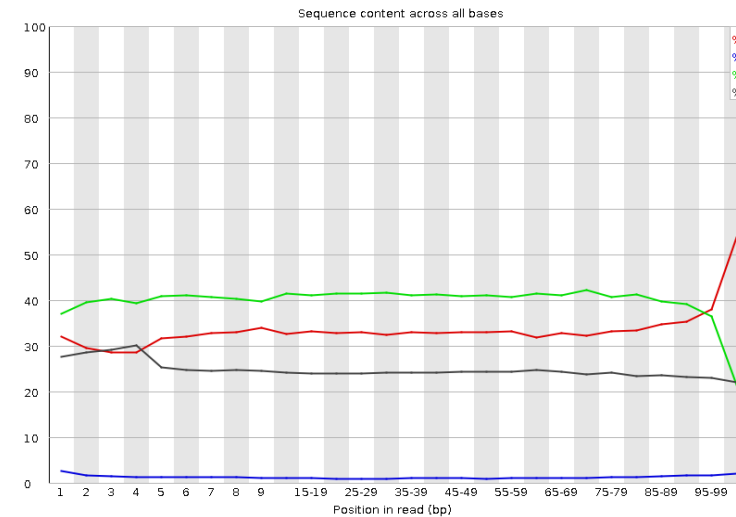
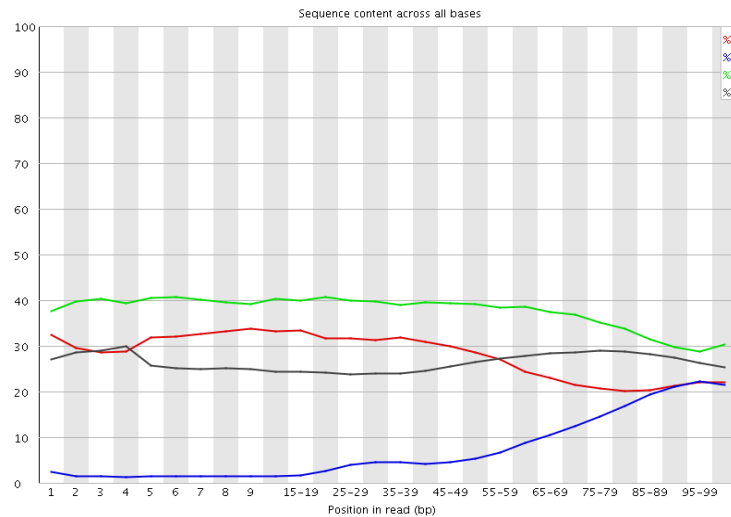
after trimming



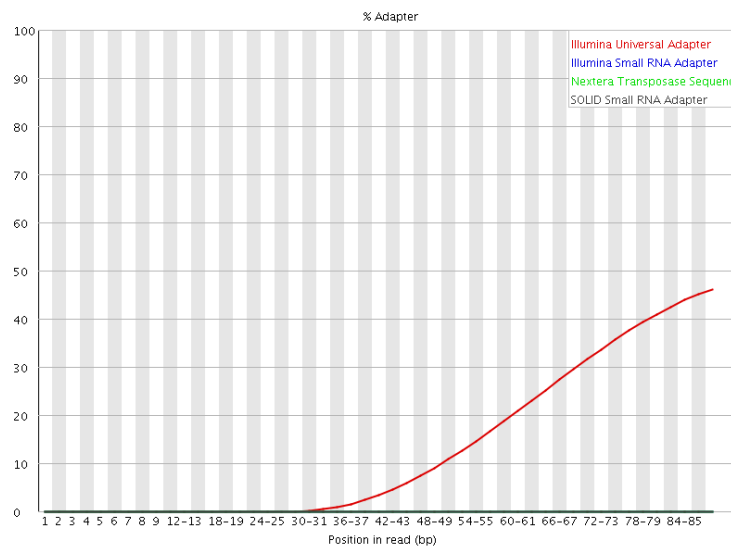
Removing adapter contamination

before trimming

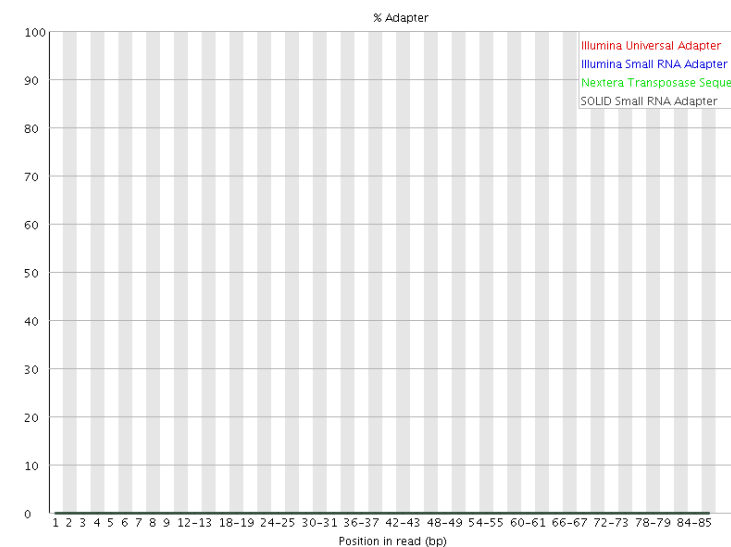
after trimming



✘ Adapter Content



✔ Adapter Content



Summary Adapter/Quality Trimming

Important to trim because failure to do so might result in:

- Low mapping efficiency
- Mis-alignments
- Errors in methylation calls since adapters are methylated
- Basecall errors tend toward 50% (C:mC)

Part II: Sequence alignment – Bismark primary alignment output (BAM file)

chromosome position

Read 1

```
HISEQ2000-06:366:C3G4NACXX:3:1101:1316:2067_1:N:0:            99        16            71322125            255        100M        =
71322232            207
NTTATTTAGTTTTTTAGGGTTTGTGTGTAGGAGTGTGGGAATTATGTTTTTTATGGTTGATATTTATTTAAAAGTGAGTATAAATTATATATATTTTTTTT
#1=DDDDDAAFFHIIIA:<FGHCCEFGHD?CFFBBBGEHHGHIII<FEHIIIII==DE??EHHFHEEEEEEEEC>;>66;@CDEEEDCEEEEEEEEDDDCBB
NM:i:14 XX:Z:G8C2C7C21C13C6CC1C17CC3C4CC4
XM:Z:.....h..h.....x.....h.....x.....hh.h.....hh...h...hh....
XR:Z:CT XG:Z:CT XA:Z:1
```

sequence
quality

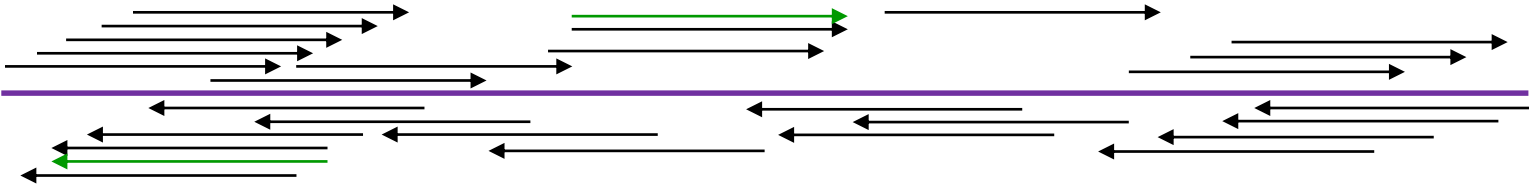
```
HISEQ2000-06:366:C3G4NACXX:3:1101:1316:2067_1:N:0:            147        16            71322232            255        100M        =
71322125            -207
GGTTATTTTATTTAGGGTTATTGTTTTAGAGTTTTATTGTTGTGAACAGATATATGATTAAGGTAATTTTTATAAGGATAATATTTAATTGGAGTTGGTT
CCCEEECADCFFFFHHGHGHIIGIHFIJJJIJHFGHGGGEHIJIIJGIGFJJJJJJJJJIGJJJJGJJJJIIIIJJIJIIJJJJJIJHHHHHHFFFFCC
NM:i:21 XX:Z:2G2CC1C1C1C11C11C2C10C1C4CC4C2C1C3C5C2C12C3C1
XM:Z:.....hh.h.h.x.....h.....x..x.....X..h.h...hh...h..h..h...h...h.....x...h.
XR:Z:GA XG:Z:CT XB:Z:1
```

methylation call

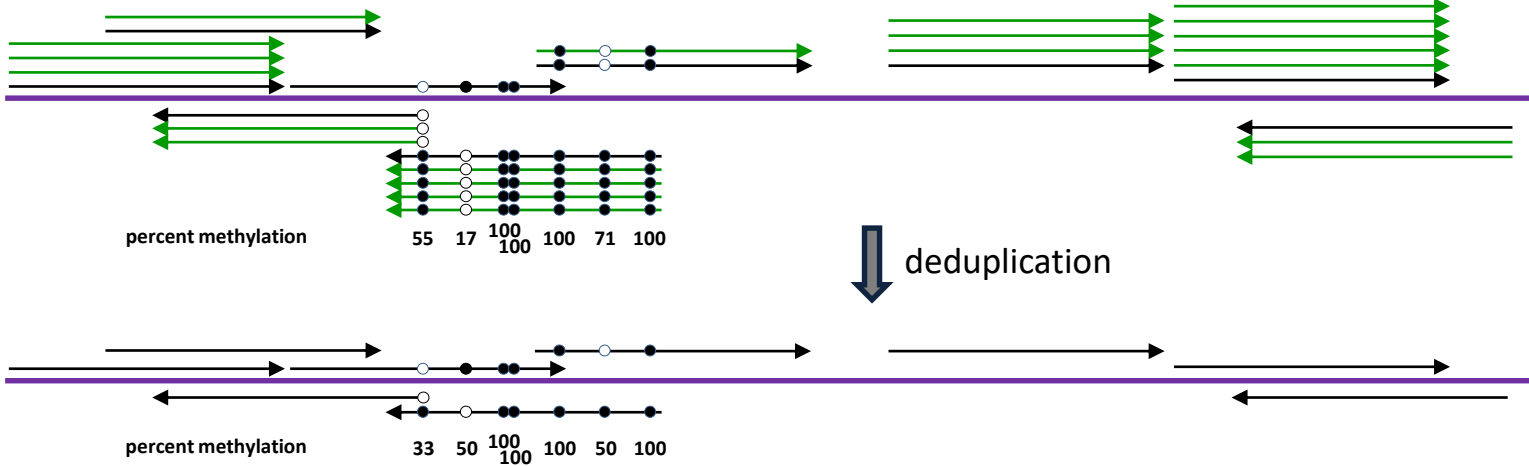
Read 2

Sequence duplication

Complex/diverse library:



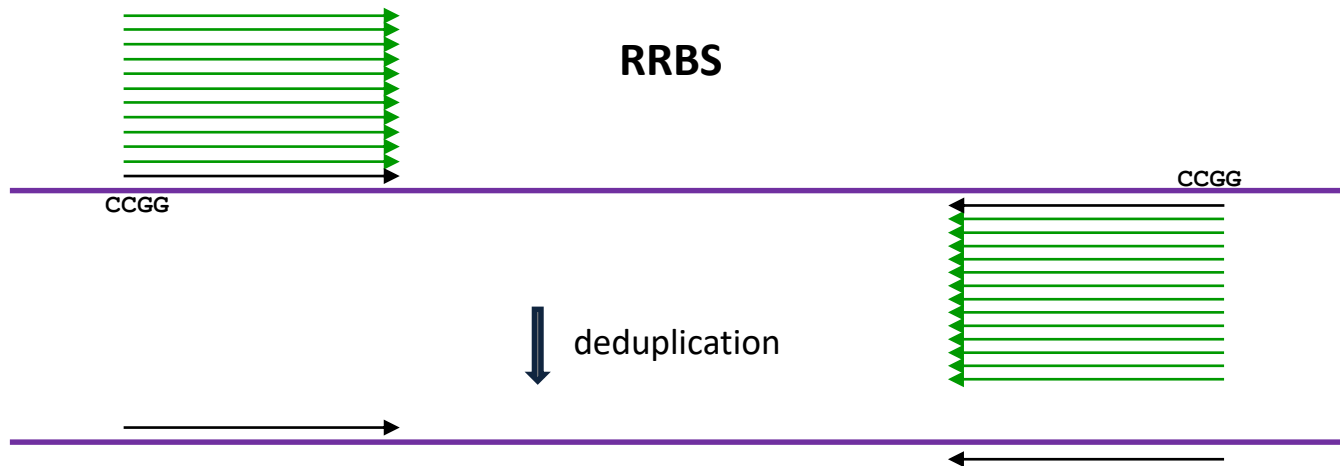
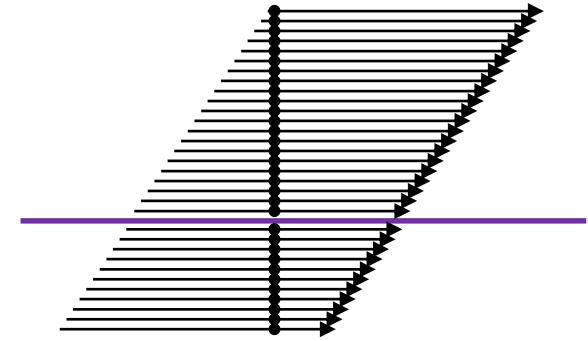
Duplicated library:



Deduplication - considerations

Advisable for large genomes and moderate coverage

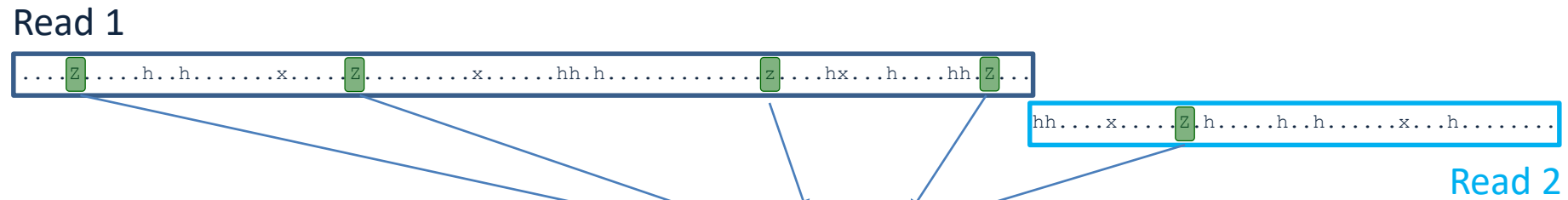
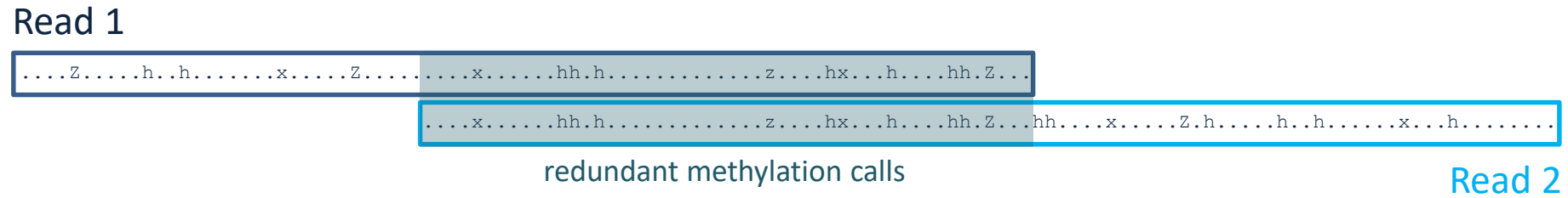
- Unlikely to sequence several genuine copies of the same fragment by chance
- Could limit coverage in high coverage studies, but would need to be very deep



NOT advisable for RRBS or target enrichment methods

- Non-random start positions (restriction sites)
- Higher local density means random collisions are likely

Methylation extraction



CpG methylation output

```

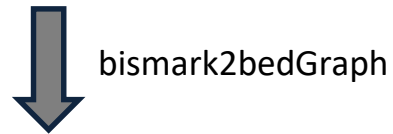
Bismark methylation extractor version v0.10.1
HS9_11915:8:2311:4022:38651#13/1      +      1      3029229 Z
HS9_11915:8:1208:13025:95413#13/1    +      1      3079409 Z
HS9_11915:8:1301:11752:81850#13/1    -      1      3104640 z
HS9_11915:8:2112:15483:84166#13/1    +      1      3104862 Z
HS9_11915:8:2110:8777:33683#13/1     -      1      3104862 z
HS9_11915:8:2208:16561:25806#13/1    +      1      3104862 Z
HS9_11915:8:2308:15290:100335#13/1   -      1      3124392 z
HS9_11915:8:2308:15290:100335#13/1   +      1      3124416 Z
HS9_11915:8:2212:13818:79056#13/1    +      1      3124416 Z
HS9_11915:8:2105:9522:91783#13/1     +      1      3124392 Z
HS9_11915:8:2105:9522:91783#13/1     +      1      3124416 Z
  
```

read ID meth state chr pos context

Methylation extraction

```
Bismark methylation extractor version v0.10.1
HS9_11915:8:2311:4022:38651#13/1      +      1      3029229 Z
HS9_11915:8:1208:13025:95413#13/1    +      1      3079409 Z
HS9_11915:8:1301:11752:81850#13/1    -      1      3104640 z
HS9_11915:8:2112:15483:84166#13/1    +      1      3104862 Z
HS9_11915:8:2110:8777:33683#13/1     -      1      3104862 z
HS9_11915:8:2208:16561:25806#13/1    +      1      3104862 Z
```

CpG methylation output



```
1      5705370 5705370 100      1      0
1      5706335 5706335 60       3      2
1      5706336 5706336 100      3      0
1      5706453 5706453 75       3      1
1      5706454 5706454 0        0      2
1      5706845 5706845 71.4285714285714      5      2
1      5706846 5706846 66.6666666666667     2      1
1      5707925 5707925 0        0      1
1      5707926 5707926 66.6666666666667     2      1
1      5709177 5709177 100      2      0
1      5709178 5709178 0        0      1
1      5710030 5710030 66.6666666666667     4      2
```

bedGraph/coverage output

chr	pos	methylation percentage	meth	unmeth
-----	-----	------------------------	------	--------

Methylation extraction

1	10525	10525	66.6666666666667	2	1
1	10542	10542	100 3 0		
1	10563	10563	66.6666666666667	2	1
1	10571	10571	100 3 0		
1	10577	10577	66.6666666666667	2	1
1	10579	10579	100 3 0		
1	10589	10589	50 2 2		
1	10609	10609	0 0 1		
1	10617	10617	0 0 1		
1	10620	10620	0 0 1		

coverage output



coverage2cytosine

1	10525	+	2	1	CG	CGC
1	10526	-	0	0	CG	CGG
1	10542	+	3	0	CG	CGA
1	10543	-	0	0	CG	CGG
1	10563	+	2	1	CG	CGC
1	10564	-	0	0	CG	CGT
1	10571	+	3	0	CG	CGC
1	10572	-	0	0	CG	CGG
1	10577	+	2	1	CG	CGC
1	10578	-	0	0	CG	CGA
1	10579	+	3	0	CG	CGG
1	10580	-	0	0	CG	CGC
1	10589	+	2	2	CG	CGG

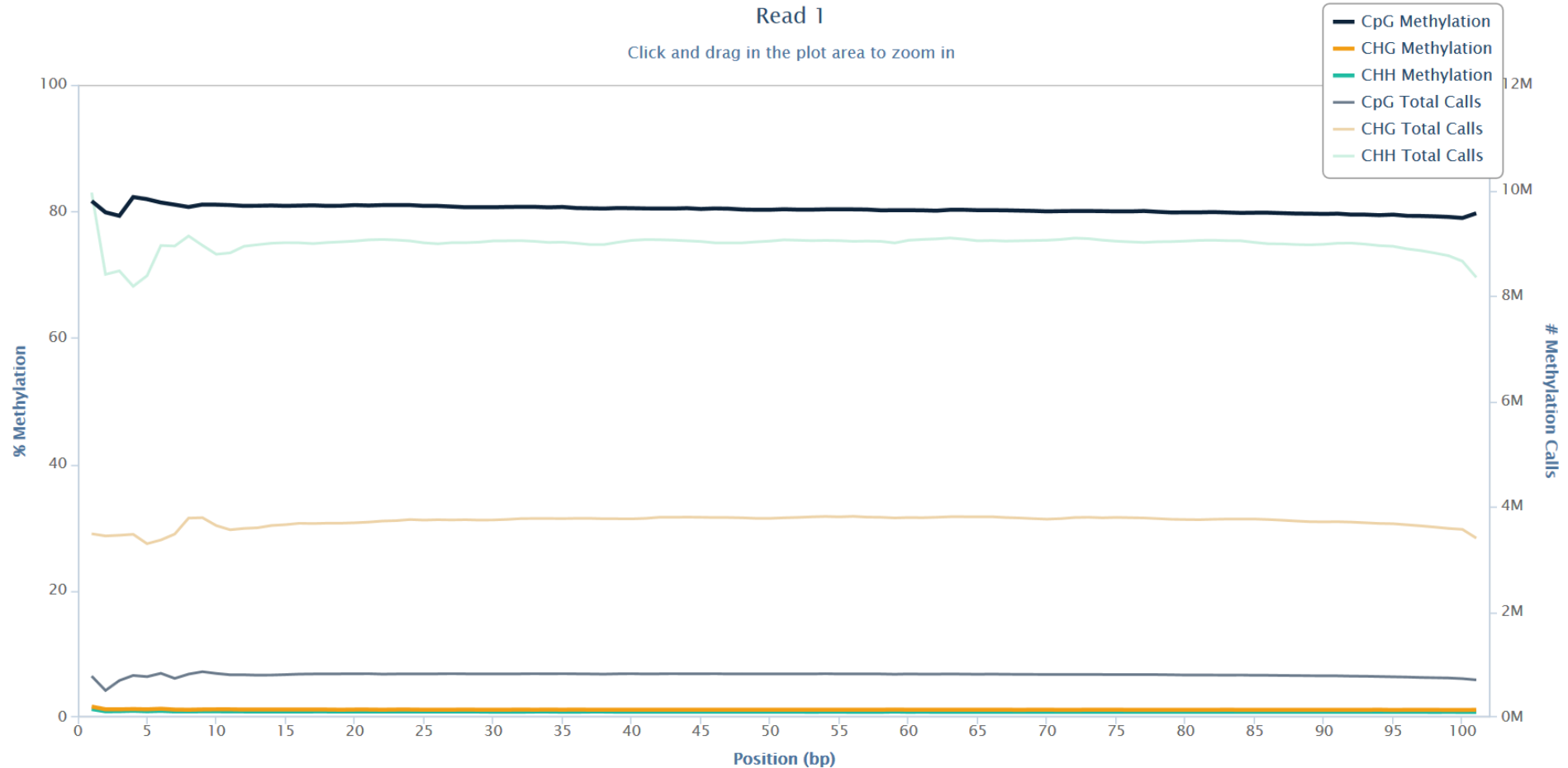
Optional: merge into CpG dinucleotide entities

Genome wide CpG report

chr pos strand meth unmeth di-nuc tri-nuc

Part III: Mapped QC - Methylation bias

M-Bias Plot



good opportunity to look at conversion efficiency

Specialist applications

(e)RRBS

WGBS

PBAT

NOMe-seq

target
enrichment

amplicon

single-cell

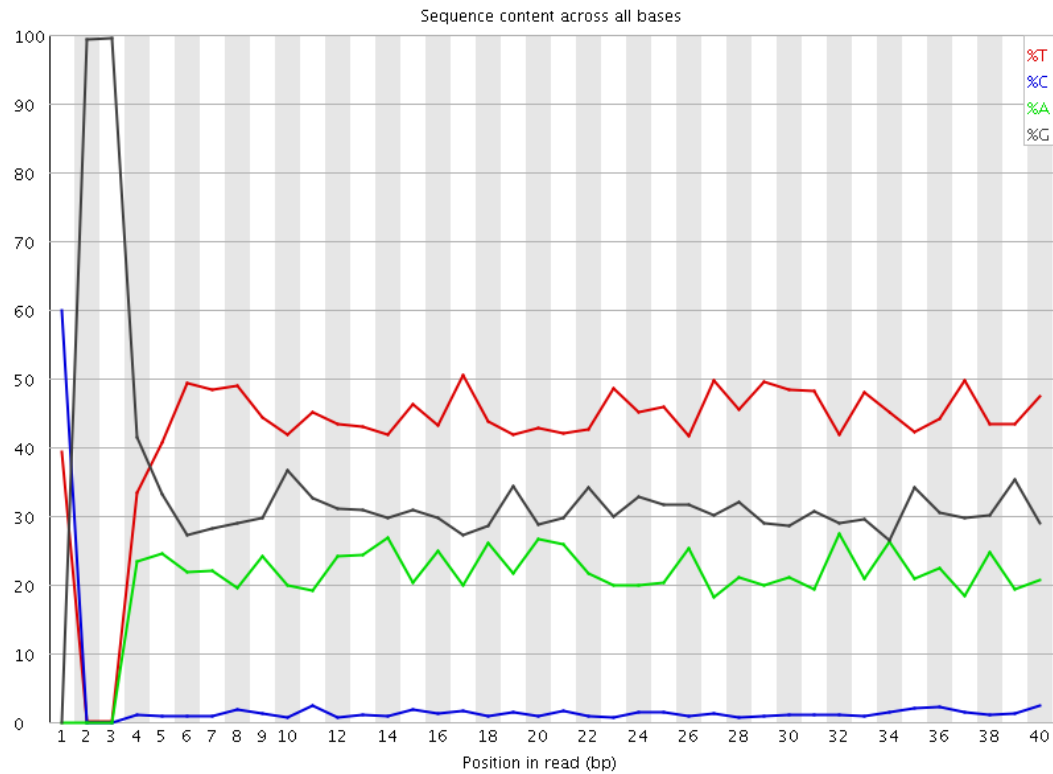
NMT-seq

non-directional

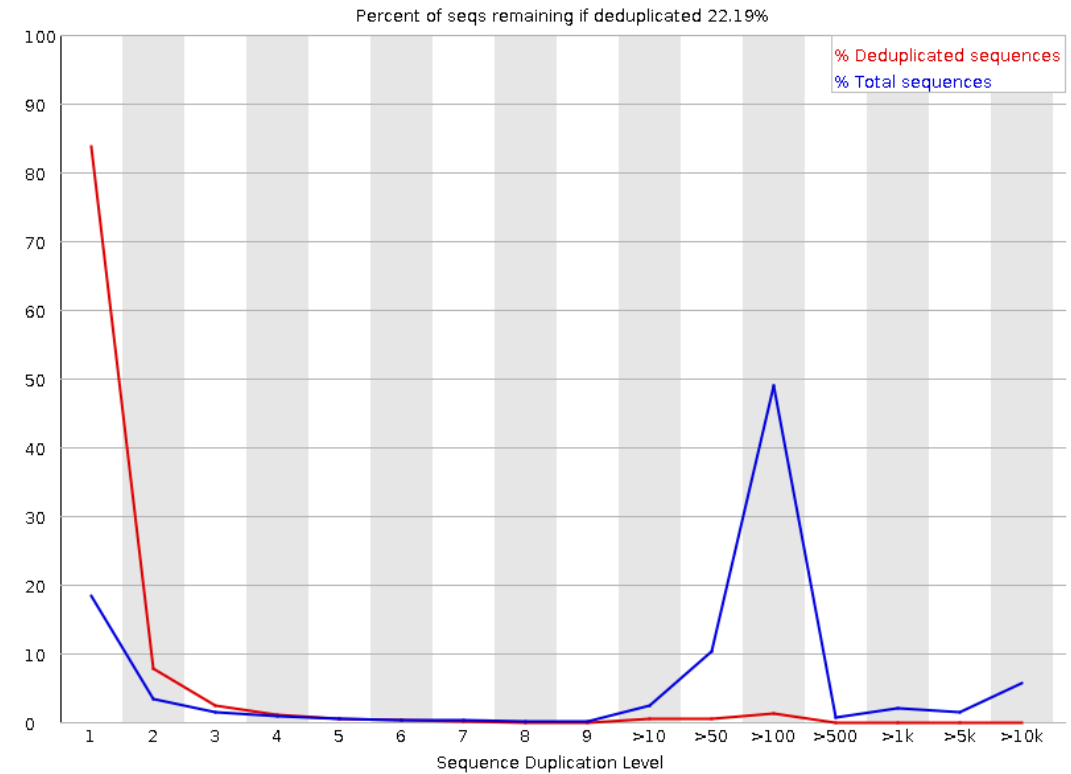
+ different library kit protocols

Reduced representation BS-Seq (RRBS)

Sequence composition bias



High duplication rate

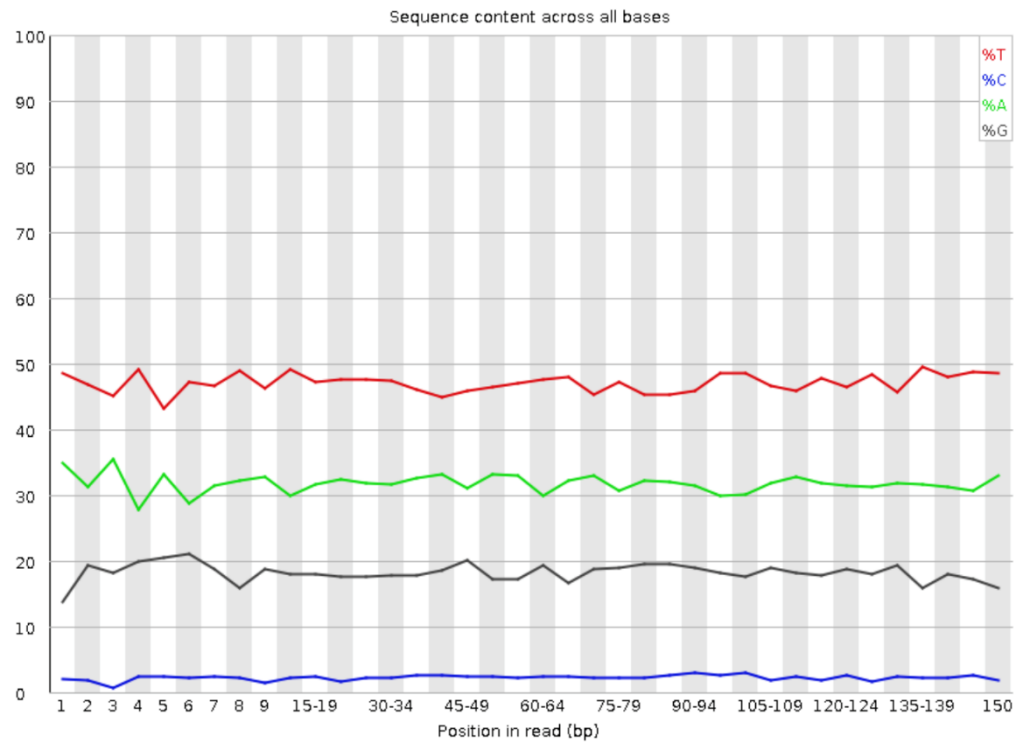


Accel Swift kit

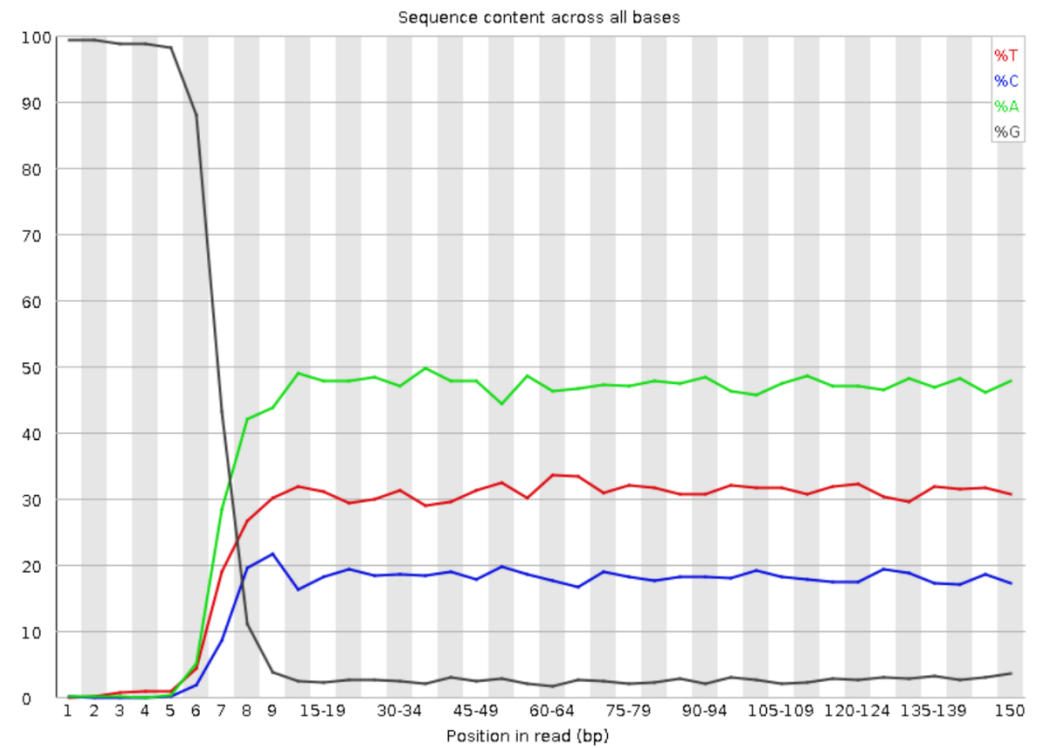
Read 1

Read 2

Per base sequence content



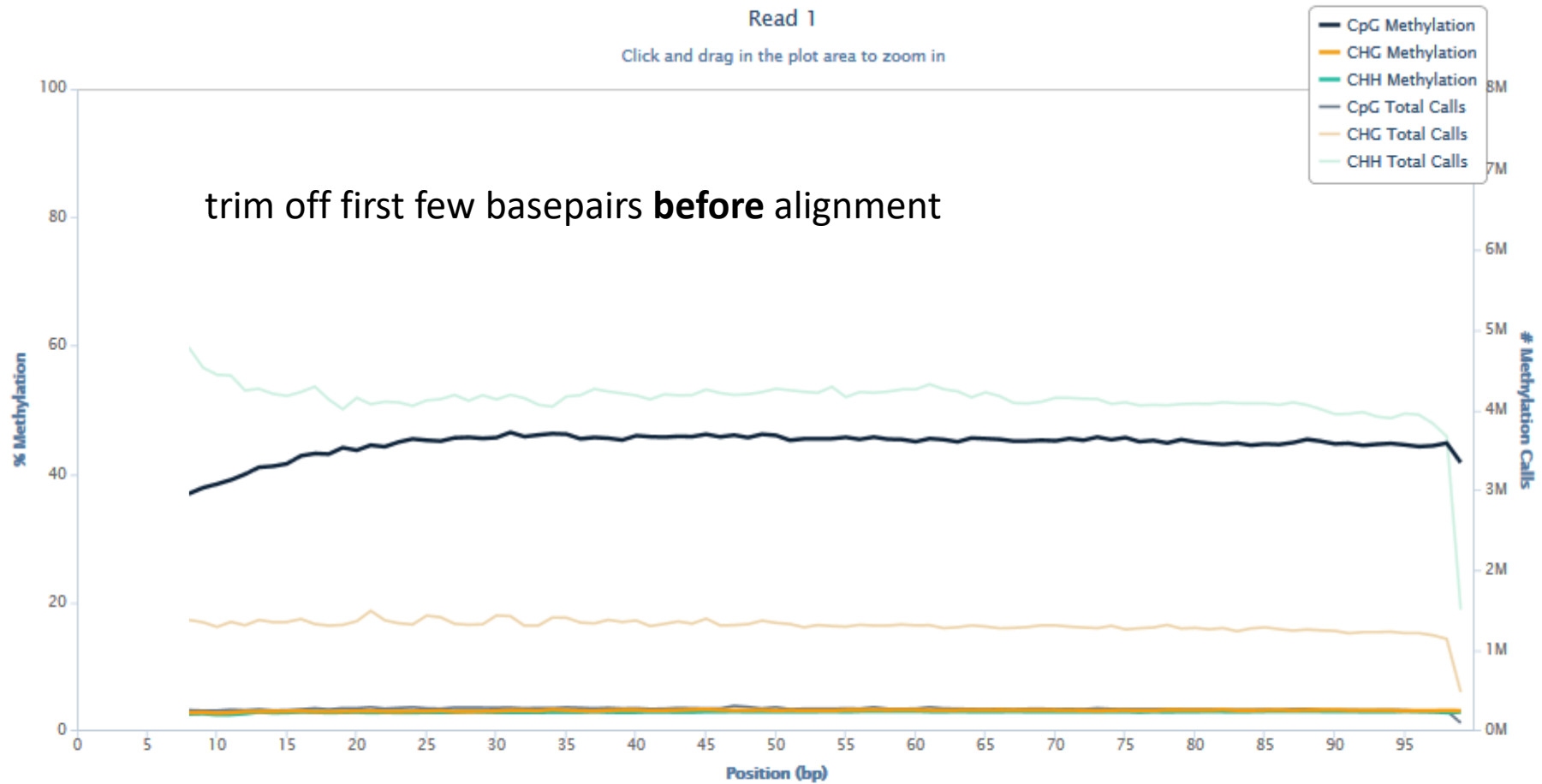
Per base sequence content



Use appropriate trimming (`trim_galore --clip_r1 10 --clip_r2 15`)

PBAT-Seq

M-Bias Plot



Bismark Bisulfite Mapper

A tool to map bisulfite converted sequence reads and determine cytosine methylation states

[View on GitHub](#)

Bismark Bisulfite Mapper



User Guide - v0.23.0

30 September, 2020

This User Guide outlines the Bismark suite of tools and gives more details for each individual step. For troubleshooting some of the more commonly experienced problems in sequencing in general and bisulfite-sequencing in particular please browse through the sequencing section at QCFail.com.

Technique	5' Trimming	3' Trimming	Mapping	Deduplication	Extraction
BS-Seq	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<code>--ignore_r2 2</code>
RRBS	<code>--rrbs (R2 only)</code>	<code>--rrbs (R1 only)</code>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
RRBS (NuGEN Ovation)	special processing	special processing	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<code>--ignore_r2 2</code>
PBAT	6N / 9N	(6N / 9N)	<code>--pbat</code>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
single-cell (scBS-Seq)	6N	(6N)	<code>--non_directional;</code> single-end mode	<input checked="" type="checkbox"/>	<input type="checkbox"/>
TruSeq (EpiGnome)	8 bp	(8 bp)	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Accel-NGS (Swift)	R1: 10, R2:15bp	(10 bp)	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Zymo Pico-Methyl	10 bp	(10 bp)	<code>--non_directional</code>	<input checked="" type="checkbox"/>	<input type="checkbox"/>

<http://felixkrueger.github.io/Bismark/Docs/>

Genomic sequence not in the genome assembly creates mapping artefacts

Probably the single biggest problem with the mapping of reads to a reference sequence is dealing with reads which come from parts of the genome which aren't in the assembly. These reads can cause significant amounts of noise in analyses performed on genomic data.

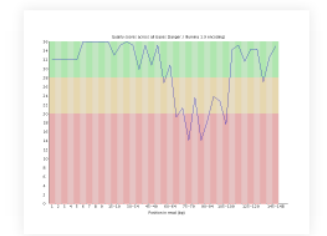
March 21, 2016 | [Simon Andrews](#) | [All Technologies](#), [All Applications](#)



Illumina 2 colour chemistry can overcall high confidence G bases

With the introduction of the NextSeq system Illumina changed the way their image data was acquired so that instead of capturing 4 images per cycle they needed only 2. This speeds up image acquisition significantly but also introduces a problem where high quality calls for G bases can be made where there is actually no signal on the flowcell.

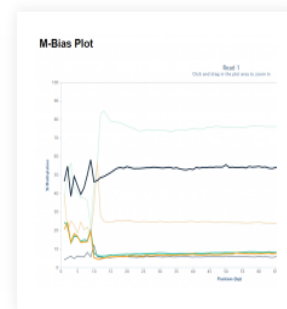
May 4, 2016 | [Simon Andrews](#) | [NextSeq](#), [All Applications](#), [Cutadapt](#), [FastQC](#)



Mispriming in PBAT libraries causes methylation bias and poor mapping efficiencies

Random priming in PBAT libraries introduces drastic biases in the base composition and methylation levels especially at the 5' end of all reads. As a result, affected bases should be removed from the libraries before the alignment step.

March 11, 2016 | [Felix Krueger](#) | [Illumina](#), [Methylation](#), [PBAT](#), [BamQC](#), [Bismark](#), [FastQC](#), [Trim Galore!](#)



Library end-repair reaction introduces methylation biases in paired-end (PE) Bisulfite-Seq applications

Library construction of standard directional BS-Seq samples often consist of several steps including sonication, end-repair, A-tailing and adapter ligation. Since the end-repair step typically uses unmethylated cytosines for the fill-in reaction the filled-in bases will generally appear unmethylated after bisulfite conversion irrespective of their true genomic methylation state.

February 12, 2016 | [Felix Krueger](#) | [Illumina](#), [BS-Seq](#), [Methylation](#), [Bismark](#), [Data Processing](#)



Bismark workflow

Pre Alignment

FastQC

Initial quality control

Trim Galore

Adapter/quality trimming using Cutadapt; handles RRBS and paired-end reads; Trim Galore and RRBS User guide

Alignment

Bismark

Output BAM

Post Alignment

Deduplication

optional

Methylation extractor

Output individual cytosine methylation calls; optionally bedGraph or genome-wide cytosine report

M-bias analysis

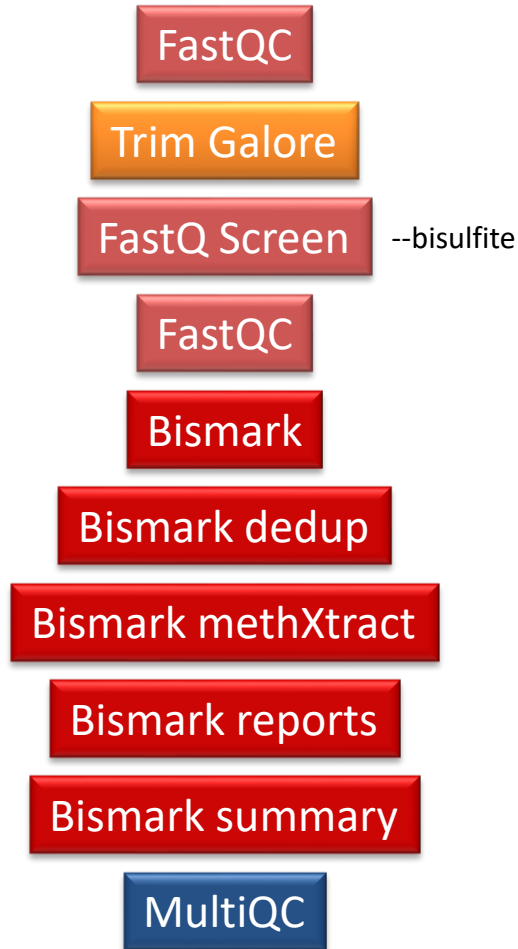
bismark2report

Graphical HTML report generation

Example: http://www.bioinformatics.babraham.ac.uk/projects/bismark/PE_report.html

protocol: *Quality Control, trimming and alignment of Bisulfite-Seq data*

nf_bisulfite_WGBS



Bismark workflow using a workflow manager

nextflow

```
[ce/4a2468] process > FASTQC (lane7561_TTGGTATG_i5F_del_GFP_35_L001_R3) [100%] 96 of 96, cached: 96 ✓
[2c/b83638] process > TRIM_GALORE (lane7561_TTGGTATG_i5F_del_GFP_35_L001_R1) [ 94%] 90 of 96, cached: 90
[be/b14097] process > FASTQ_SCREEN (lane7561_TCCTCAAT_i5F_del_GFP_29_L001_R1) [  0%] 0 of 90
[94/8a2957] process > FASTQC2 (lane7561_TCCAGTCG_i5F_del_GFP_41_L001_R1) [ 46%] 41 of 90, cached: 41
[01/56550d] process > BISMARK (lane7561_GCCAATGT_i5F_del_GFP_6_L001_R1) [  1%] 1 of 90, cached: 1
[-         ] process > BISMARK_DEDUPLICATION [  0%] 0 of 1
[-         ] process > BISMARK_METHYLATION_EXTRACTOR -
[-         ] process > BISMARK2REPORT -
[-         ] process > BISMARK2SUMMARY -
[-         ] process > MULTIQC -
```

Useful links

- **FastQC** www.bioinformatics.babraham.ac.uk/projects/fastqc/
- **Trim Galore** <https://github.com/FelixKrueger/TrimGalore>
- **Cutadapt** <https://code.google.com/p/cutadapt/>
- **Bismark** <https://github.com/FelixKrueger/Bismark>
- **Bowtie 2** <http://bowtie-bio.sourceforge.net/bowtie2/>
- **SeqMonk** www.bioinformatics.babraham.ac.uk/projects/seqmonk/

 <https://sequencing.qcfail.com/>



Sierra: A web-based LIMS system for small sequencing facilities

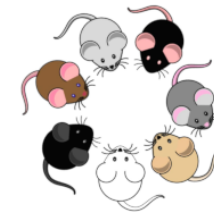


reStrainingOrder



SeqMonk: Genome browser, quantitation and data analysis

Trim Galore! Quality and adapter trimming for (RRBS) sequencing libraries



FastQ Screen: organism and contamination detection

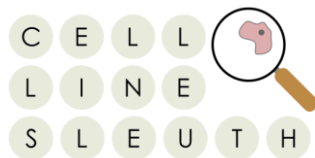


Bismark: Bisulfite-sequencing alignments and methylation calls



Hi-C mapping

ASAP: Allele-specific alignments



FastQC: quality control for high throughput sequencing



GcA_TTAcG_TaA_TGcCcT_A_TGcCcT_A_TGcCcT_A_TGGACT
TGGACT

DNA methylation is reset during reprogramming

