# Epigenomics Data Analysis Workshop 2023
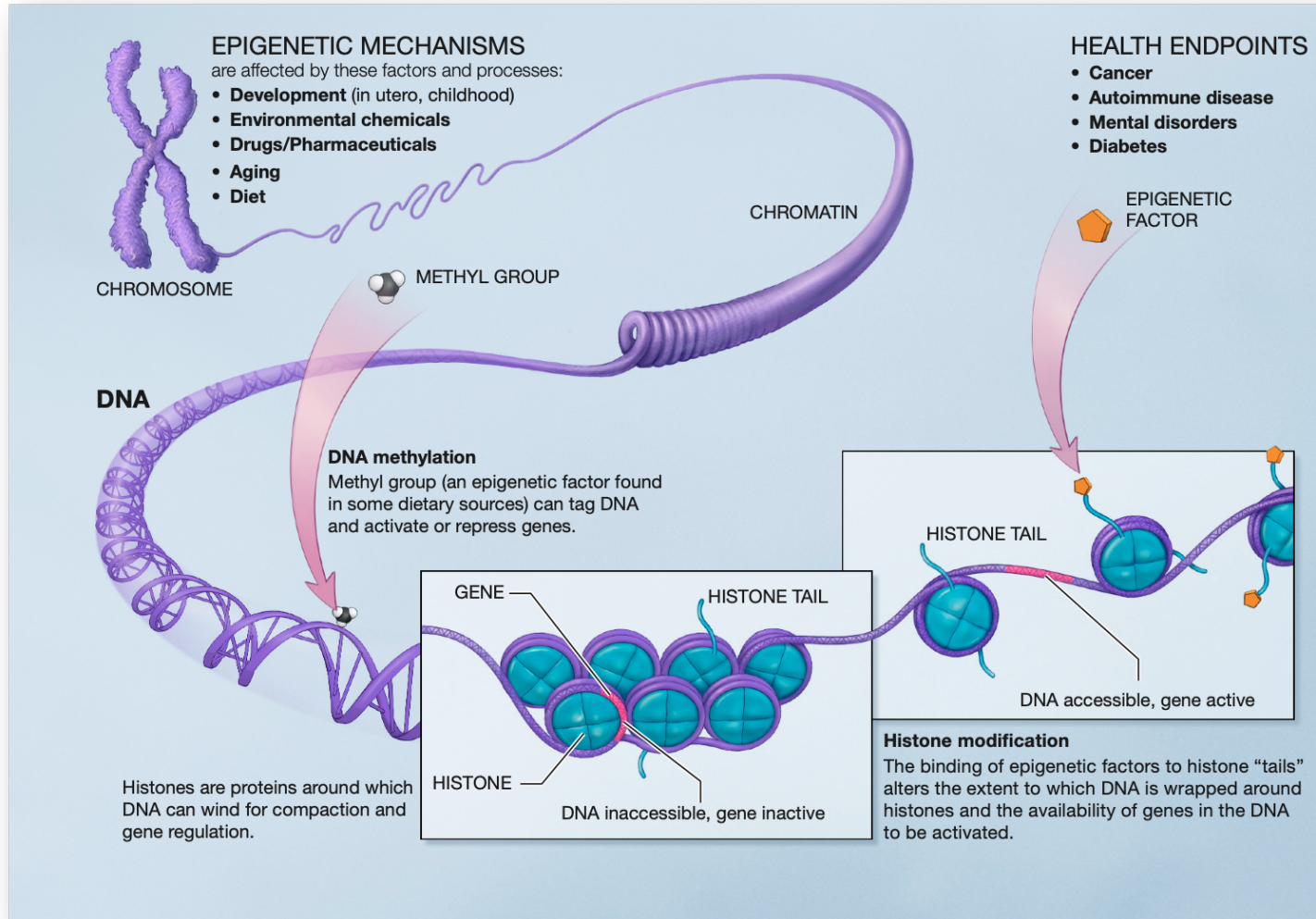
DNA Methylation

# Schedule

- 09:30 - 10:15 Short introduction to DNA methylation + Overview Array exercises

- 10:15 - 10:30 UPPMAX set-up + break

- 10:30 - 12:00 Array exercises

- 10:00 - 13:00 Lunch

- 13:00 - 14:00 DNA Methylation: Methods & Technologies

- 14:00 - 14:15 Break

- 14:15 - 14:30 Overview Exercises Bisulfite Sequencing

- 14:30 - 16:30 Bisulfite Sequencing Exercise
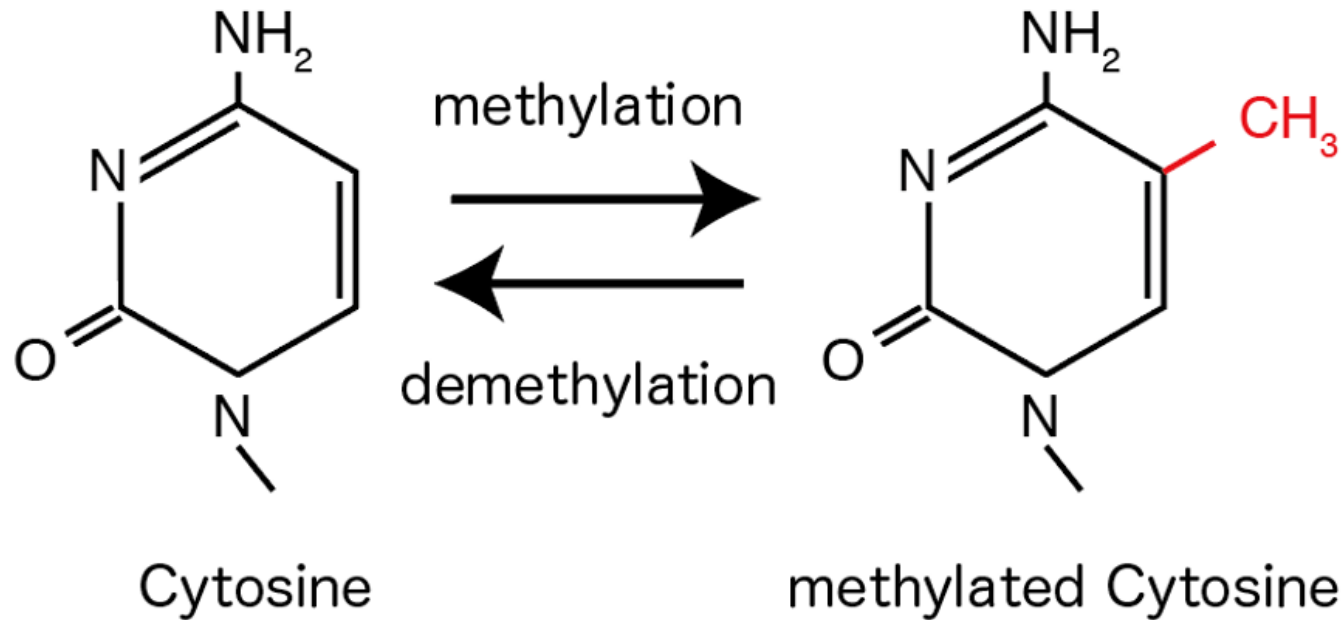
- 16:30 - 17:00 Test Yourself

# Introduction to DNA methylation

# Epigenetics



*source: NIH*

Epigenomics Data analysis 2023: Methylation

# What is DNA methylation?

# What is DNA methylation?

- Mostly found at cytosines followed by guanines

  - 90% in CpG sites

  - Default state is methylated

  - Prone to mutation -> depleted

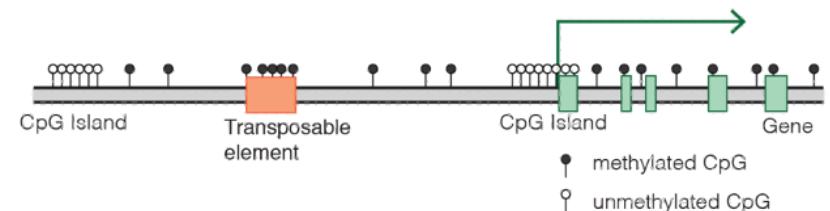# What is DNA methylation?

- Mostly found at cytosines followed by guanines

  - 90% in CpG sites

  - Default state is methylated

  - Prone to mutation -> depleted

- CpG sites often occurs as clusters: CpG Islands

  - Region with high frequency of CpG

  - Often associated with promoters

  - Unmethylated if gene is expressed





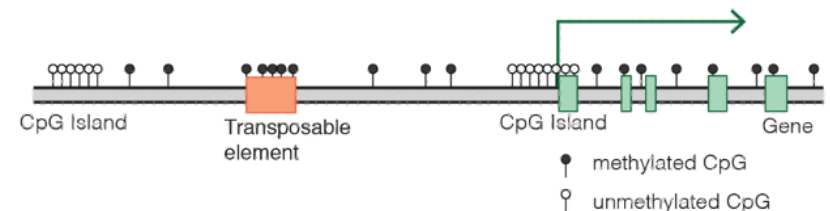Typical mammalian DNA methylation landscape
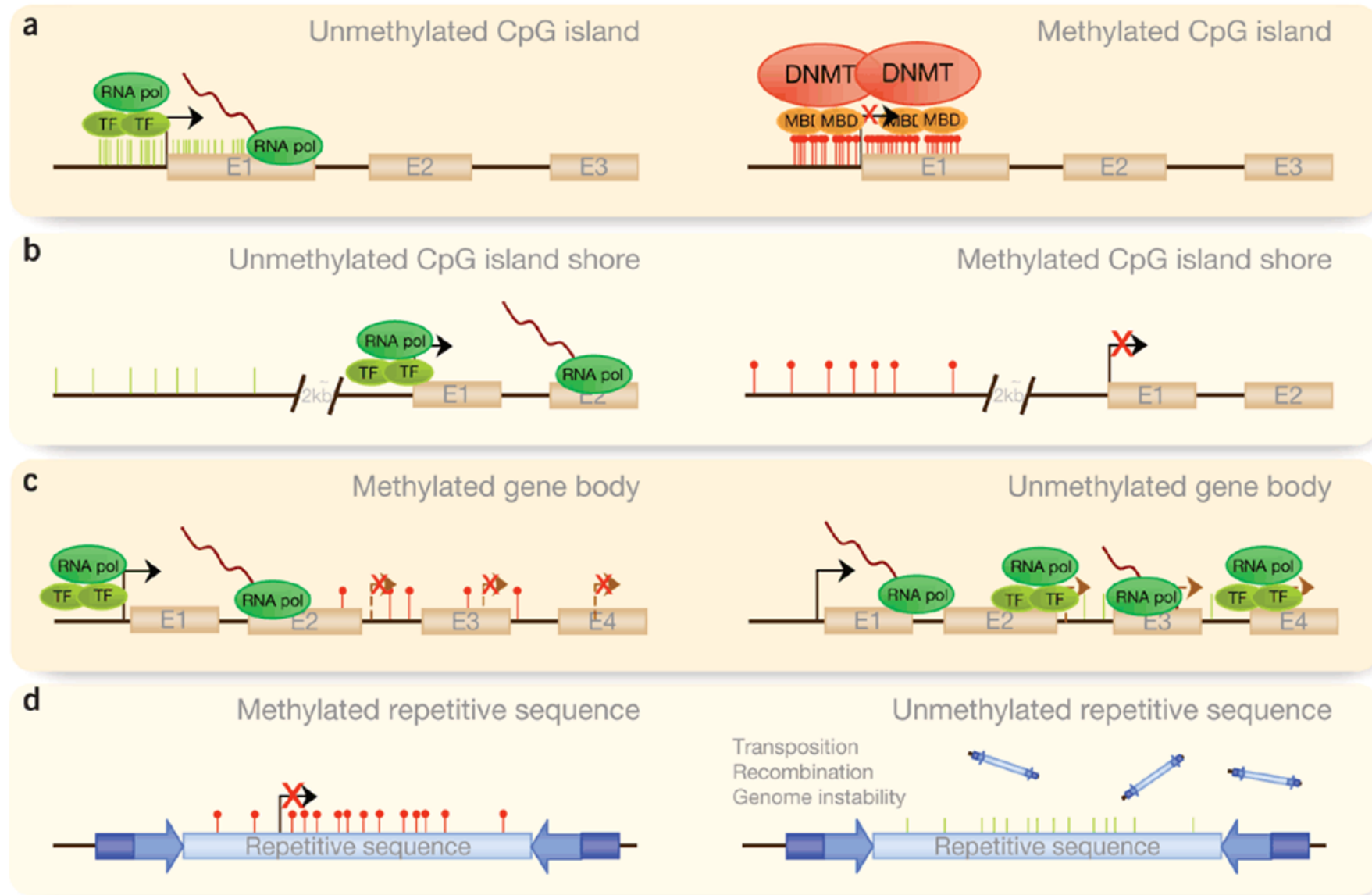
# What is DNA methylation?

- Mostly found at cytosines followed by guanines

  - 90% in CpG sites

  - Default state is methylated

  - Prone to mutation -> depleted

- CpG sites often occurs as clusters: CpG Islands

  - Region with high frequency of CpG

  - Often associated with promoters

  - Unmethylated if gene is expressed
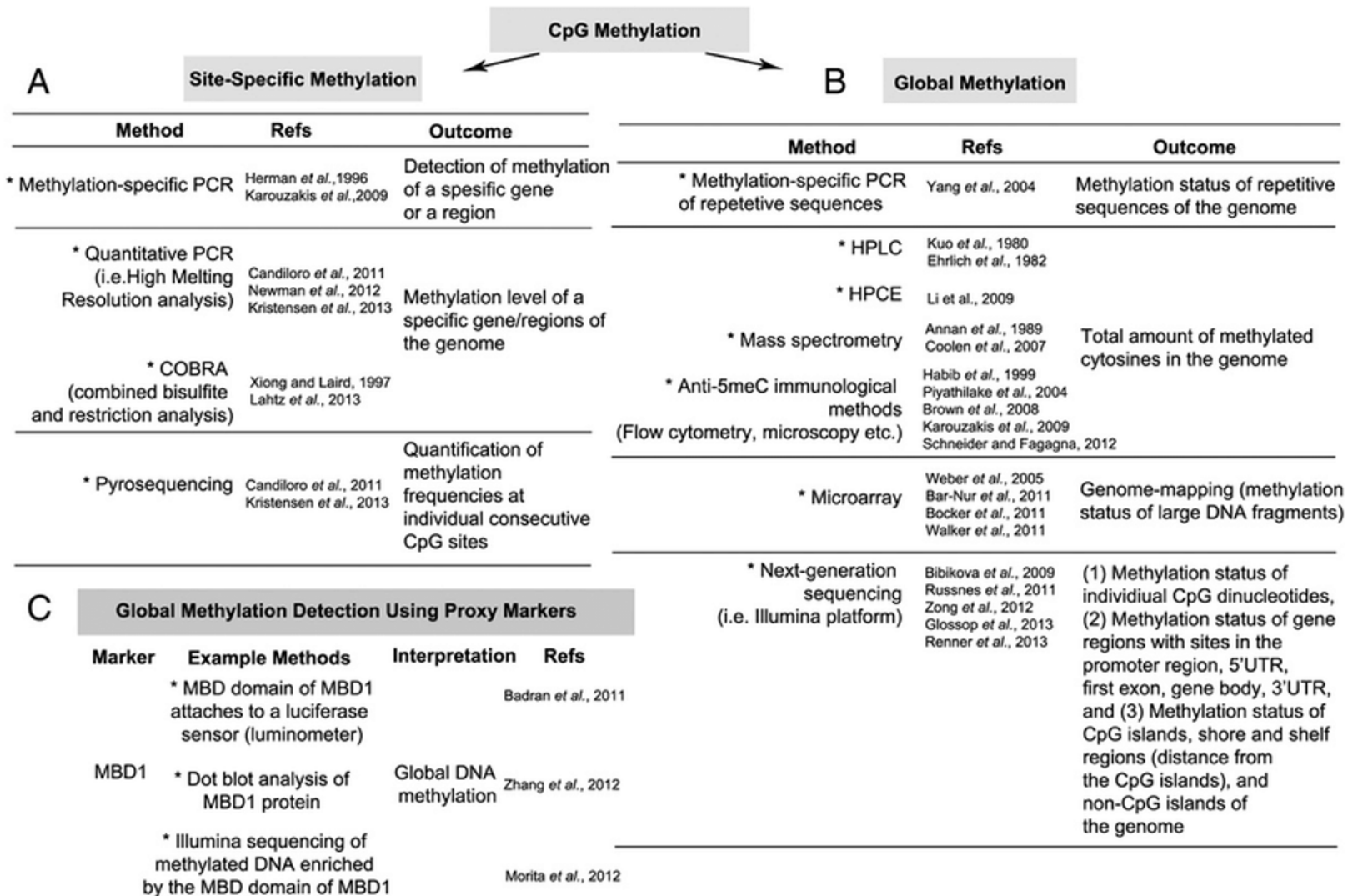
- Role in development, aging, cancer, exercise, …



Typical mammalian DNA methylation landscape
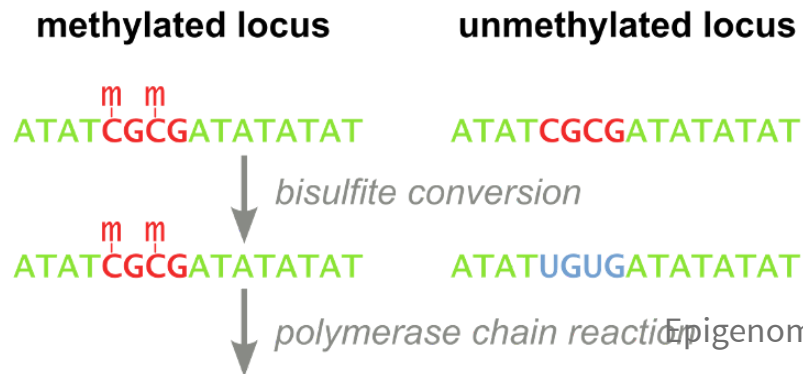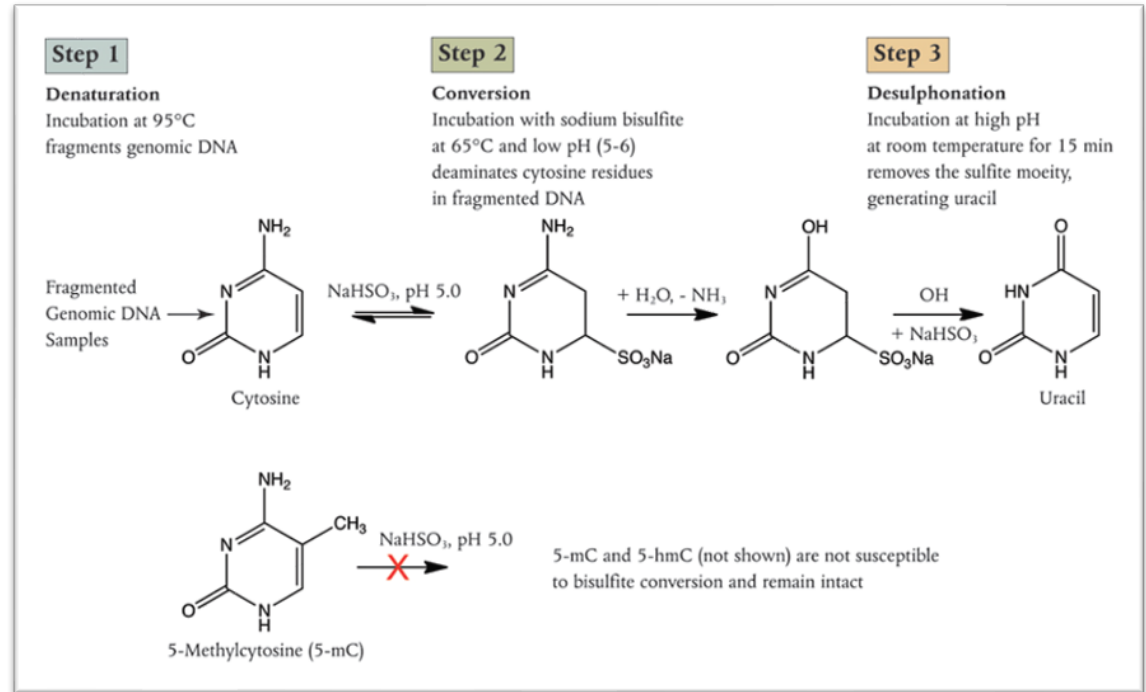
# Effects of Methylation

# Detection of DNA methylation



Celik et al. (2014), Journal of Immunological Methods

# Bisulfite Conversion

- Bisulfite conversion crucial for both arrays and sequencing

- C -> U (->T)

- mC -> mC (-> C)

- methylation-specific PCR, high resolution melting curve analysis, micro-array based approaches and next generation sequencing





**methylated locus**

**unmethylated locus**

m m
ATATCGCGATATATAT

ATATCGCGATATATAT

*bisulfite conversion*

m m
ATATCGCGATATATAT

ATATUGUGATATATAT

*polymerase chain reaction*

# Illumina Methylation Arrays

**GoldenGate**

1500 CpGs,
cancer focused

**Infinium
HumanMethylation450**

480K CpGs, 99% RefSeq genes

| 2007 | 2008 | 2011 | 2015 |
|------|------|------|------|

**MethylationEPIC**

850K CpGs, >90% 450 +
additional regulatory
regions

HumanMethylation450 array content.

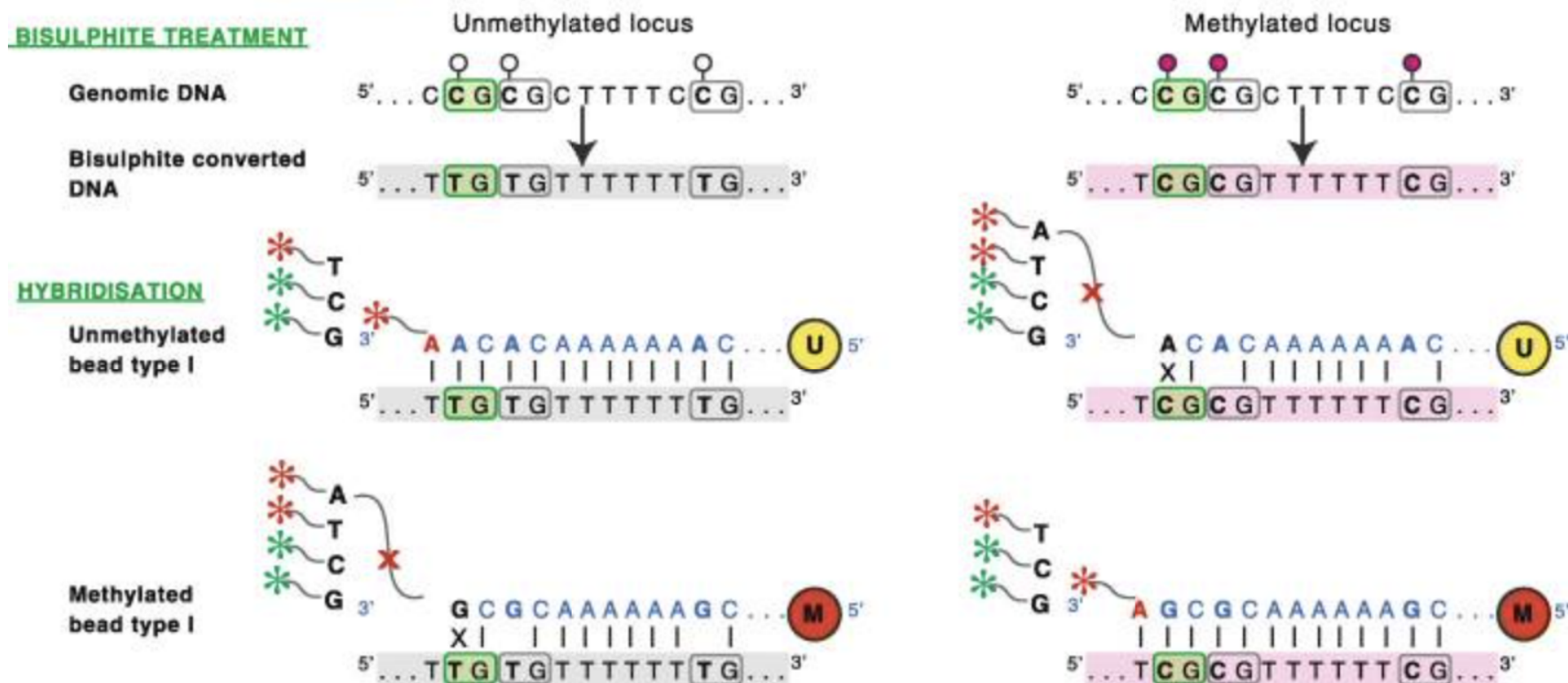| Feature type | Included on array |
|--------------|-------------------|
| Total number of sites | 485,577 |
| RefSeq genes | 21,231 (99%) |
| CpG islands | 26,658 (96%) |
| CpG island shores (0–2 kb from CGI) | 26,249 (92%) |
| CpG island shelves (2–4 kb from CGI) | 24,018 (86%) |
| HMM islands[a] | 62,600 |
| FANTOM 4 promoters (High CpG content)[a] | 9426 |
| FANTOM 4 promoters (Low CpG content)[a] | 2328 |
| Differentially methylated regions (DMRs)[a] | 16,232 |
| Informatically-predicted enhancers[a] | 80,538 |
| DNAse hypersensitive sites | 59,916 |
| Ensemble regulatory features[a] | 47,257 |
| Loci in MHC region | 12,334 |
| HumanMethylation27 loci | 25,978 |
| Non-CpG loci | 3091 |

# 450 Array

- 50bp single stranded DNA oligos ("probes") attached to silica beads

- 2 detection channels: red and green

- Hybrid of 2 different probe designs
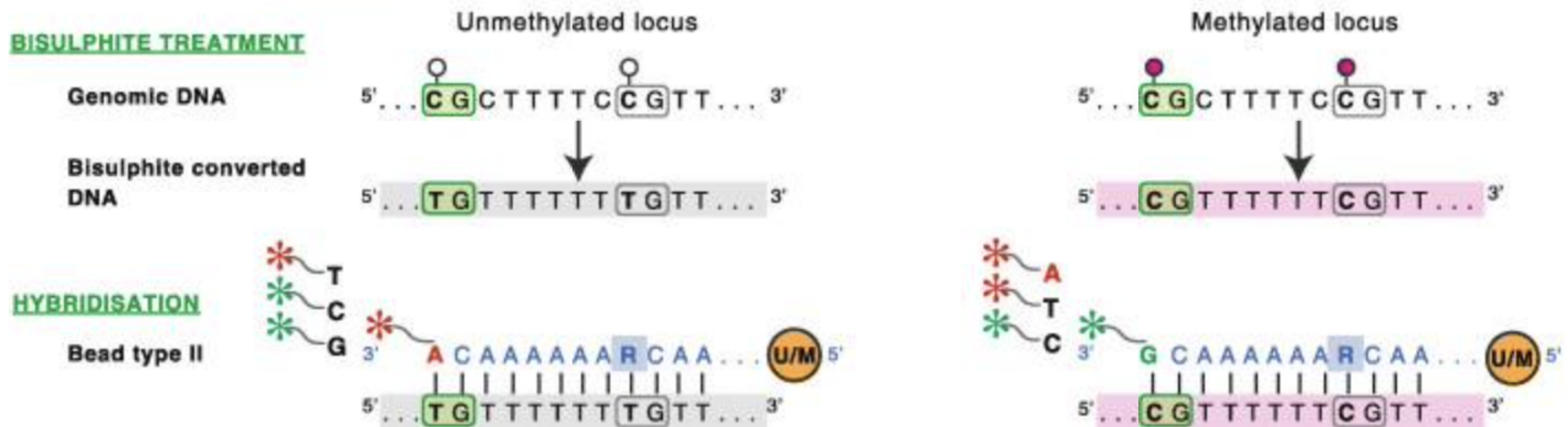
Epigenomics Data analysis 2023: Methylation

# Infinium: Type I vs II design

- Type I: single color detection, two beads

# Infinium: Type I vs II design

- Type II: two color detection, single bead

# Infinium: Type I vs II design

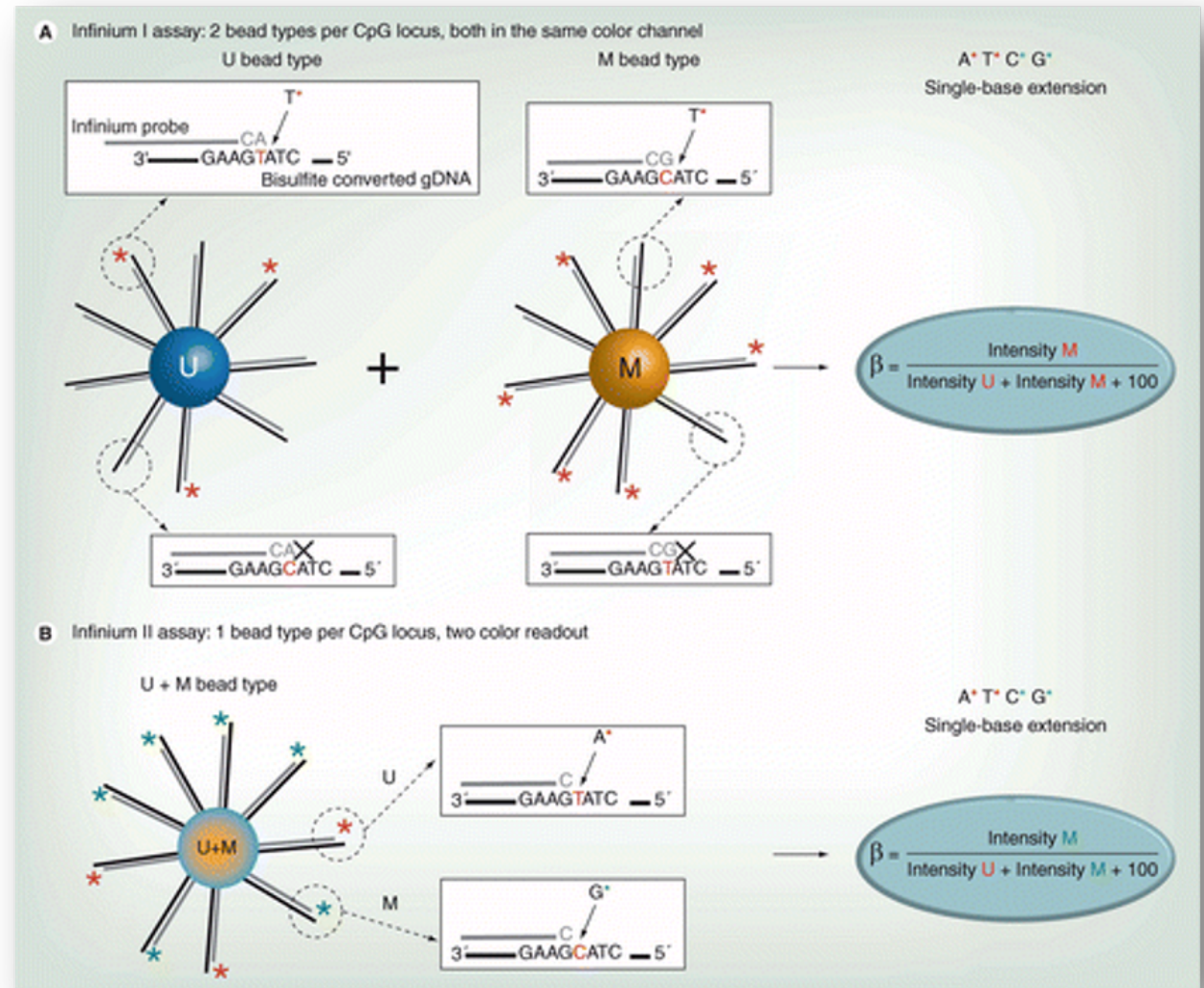| Type I | Type II |
| --- | --- |
| Same chemistry as 27K | New from 450K on |
| 2 beads/CpG | 1 bead/CpG (fits more) |
| Better for CpG dense regions | better for less dense regions |
| More stable/reproducible | lower dynamic range |

# From red/green to methylation level

- Intensities are used to estimate *Beta* values; for both probe designs

$$beta = M/(M + U + 100)$$

- *Beta* value between 0 and 1
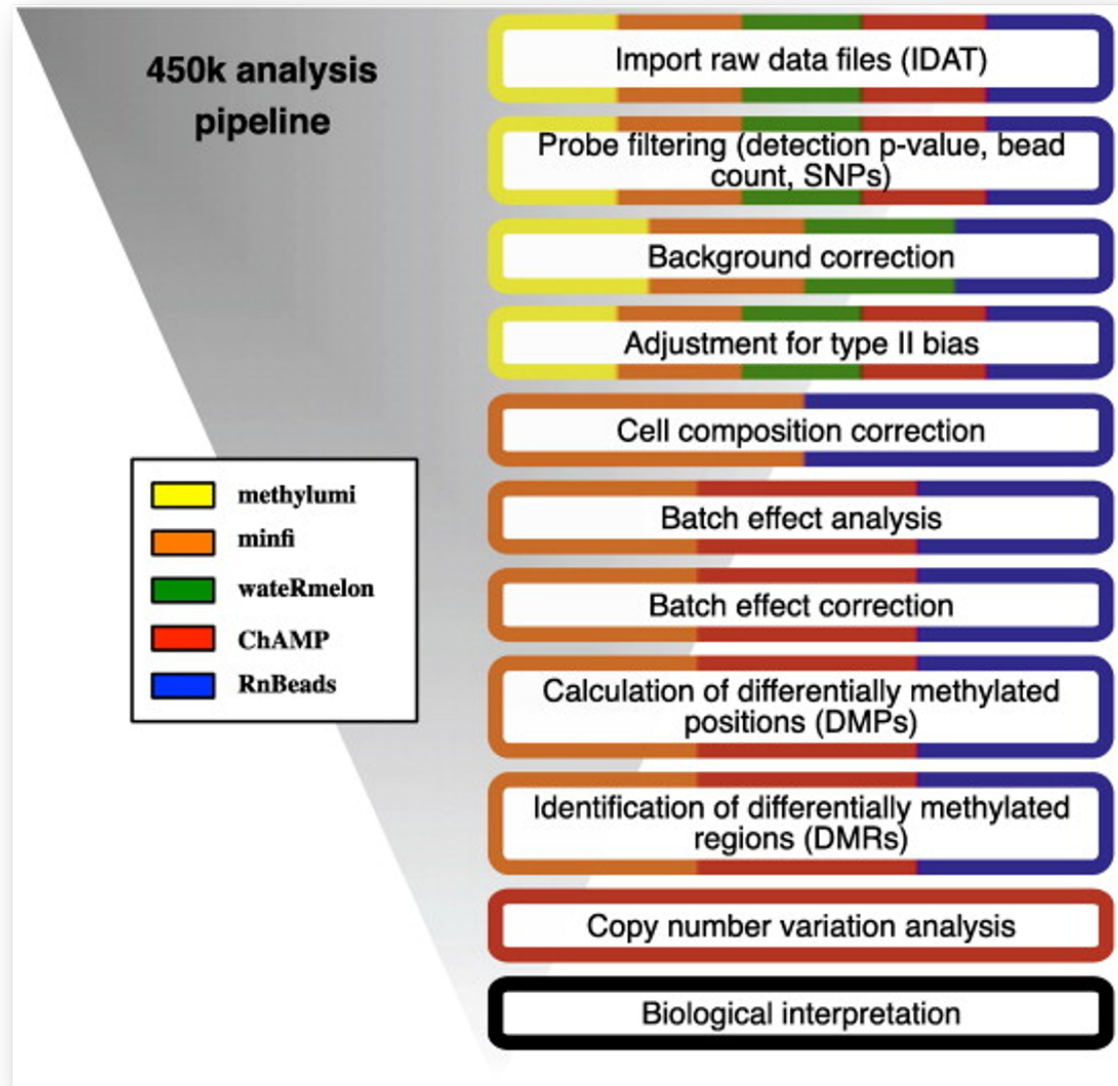- Easily interpretable, but related M-value has better statistical properties

$$Mvalue = log2(M/U)$$

# Analysis Workflow

- Typical analysis consists of different steps…

- Many tools for analyzing Illumina arrays

- R package minfi

```
1  library(minfi)
```
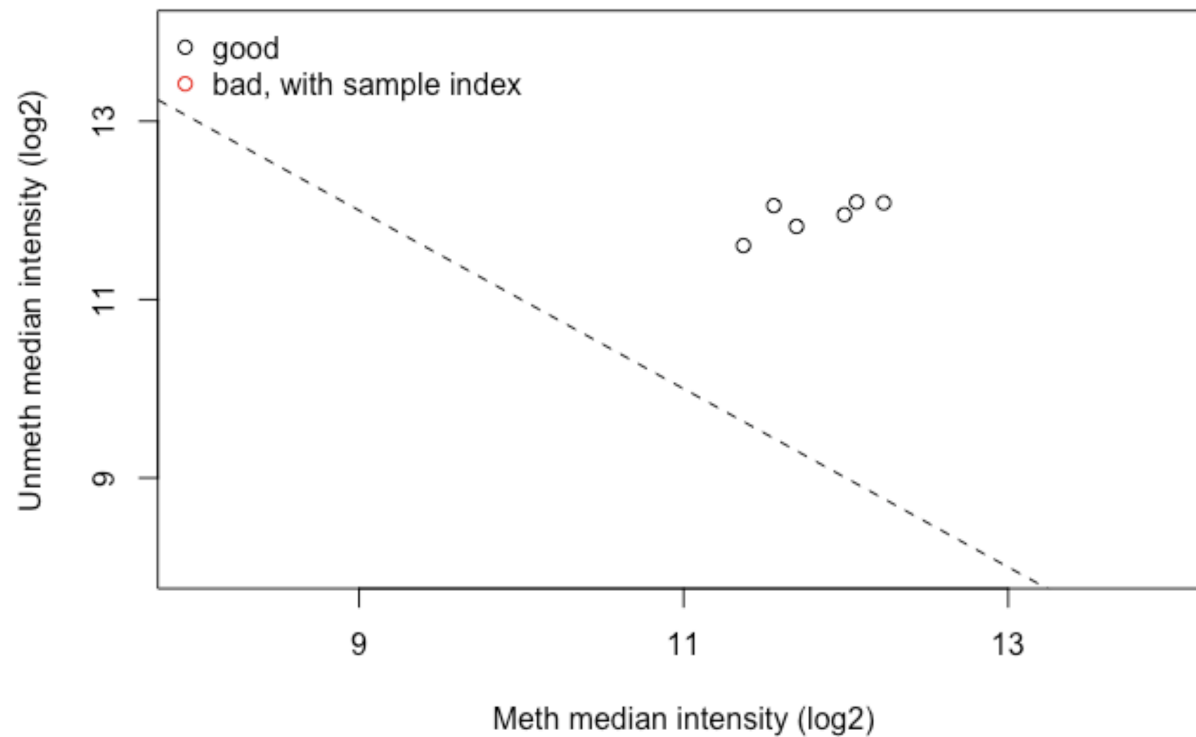
# Import data

- IDAT files; slide scanner output

  - 5859594006_R01C01_Grn.idat

```r
1  dataDirectory <- "/sw/courses/epigenomics/DNAmethylation/array_data/"
2  # read in the sample sheet for the experiment
3  targets <- read.metharray.sheet(dataDirectory, pattern="SampleSheet.csv")
4  # read in the raw data from the IDAT files
5  rgSet <- read.metharray.exp(targets=targets)
6  # Go from intensity data to methylation levels
7  MSet <- preprocessRaw(rgSet)
```

# Initial Quality Control
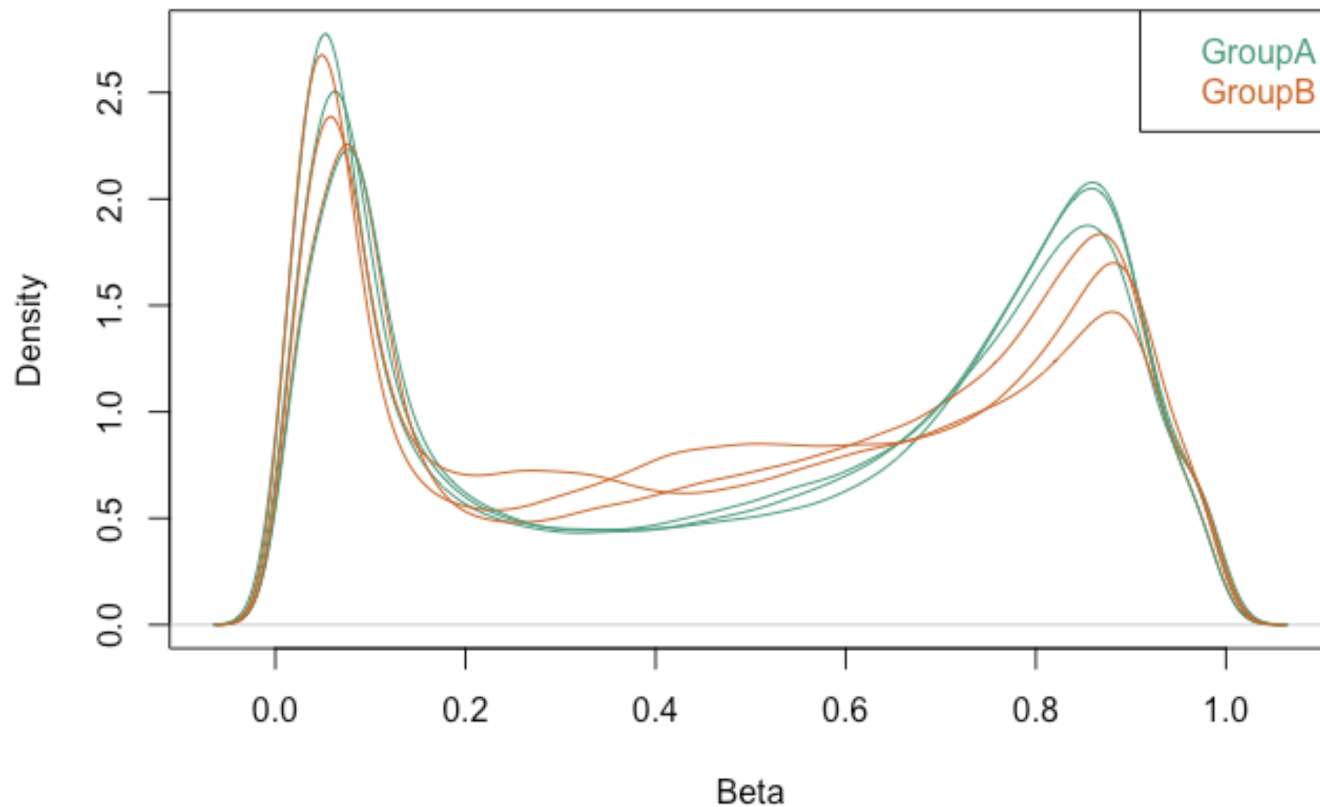
- Plot median intensity in M vs U

```
1  qc <- getQC(MSet)
2  plotQC(qc)
```

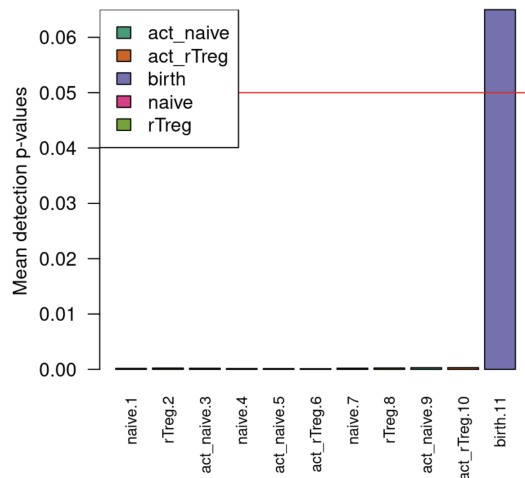# Initial Quality Control

- *Beta* value density distribution

```
1  densityPlot(MSet, sampGroups = phenoData$Sample_Group)
```

# Initial Quality Control

- Detection p-value: Are the intensities significantly above background?

```
1  # Calculate the detection p-values
2  detP <- detectionP(rgSet)
3  # examine mean detection p-values across all samples to identify any failed
4  barplot(colMeans(detP), las=2, cex.names=0.8, ylab="Mean detection p-values
5  abline(h=0.05,col="red")
```
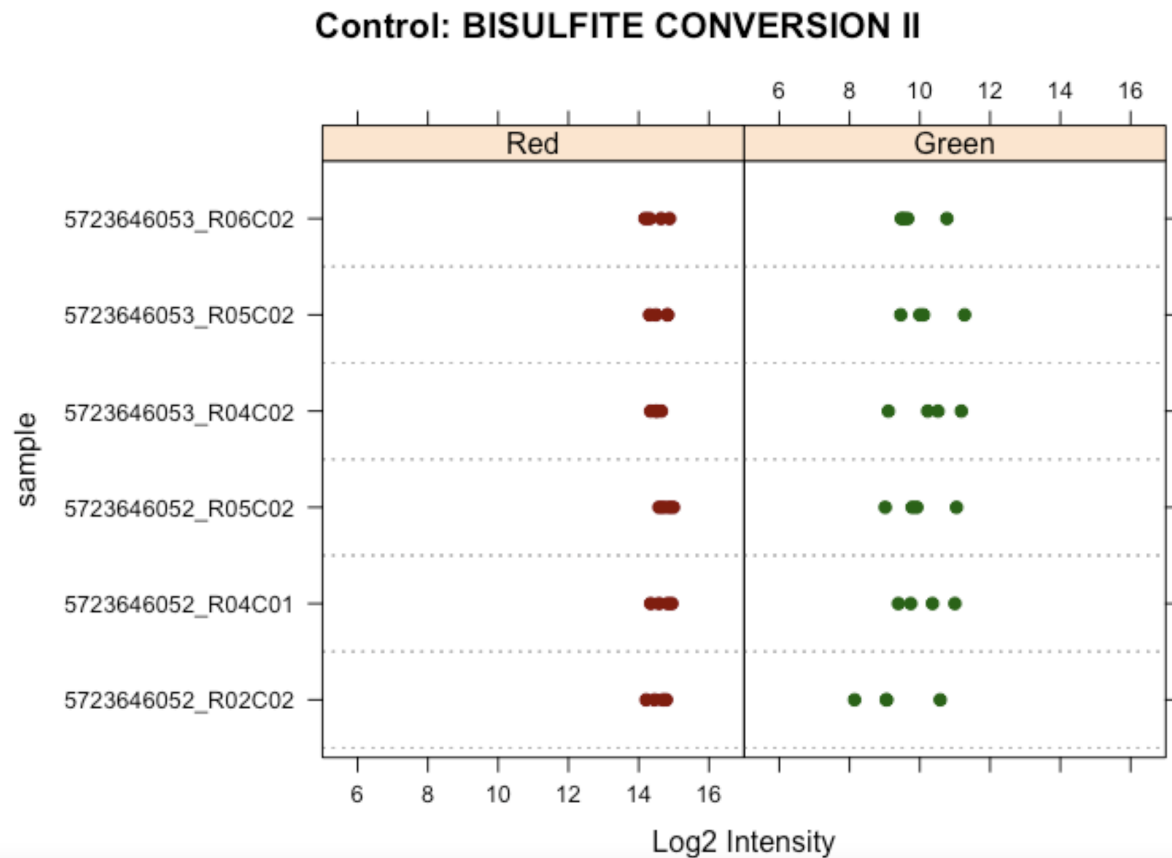


Potentially remove bad samples and/or probes.

# Initial Quality Control

- Several internal control probes for different sample preparation steps (bisulfite conversion, hybridization, …)

```
1  controlStripPlot(RGSet, controls="BISULFITE CONVERSION II")
```



**Control: BISULFITE CONVERSION II**

- Staining control
- Bisulfite conversion
- extension controls
- specificity controls
- hybridization controls
- target removal controls
- negative controls

Description in Illumina manual

# Other considerations…

- Remove X/Y Chromosome CpGs?

- Remove CpG overlapping with known SNP and/or cross reactive probes

- Check sample structure with PCA

Many of the previous plots can be looked at interactively with shinyMethyl.

paper: A comprehensive overview of Infinium HumanMethylation450 data processing.

# Normalization

- Within and across array normalization

**A systematic study of normalization methods for Infinium 450K methylation data using whole-genome...**

Ting Wang

**Between-array normalization for...**

**A systematic assessment of normalization approaches for the Infinium 450K methylation platform**

Michael C Wu, Bonnie R Joubert, Pei-fen Kuan, Siri E Håberg, Wenche Nystad, Shyamal D Peddada & Stephanie J London
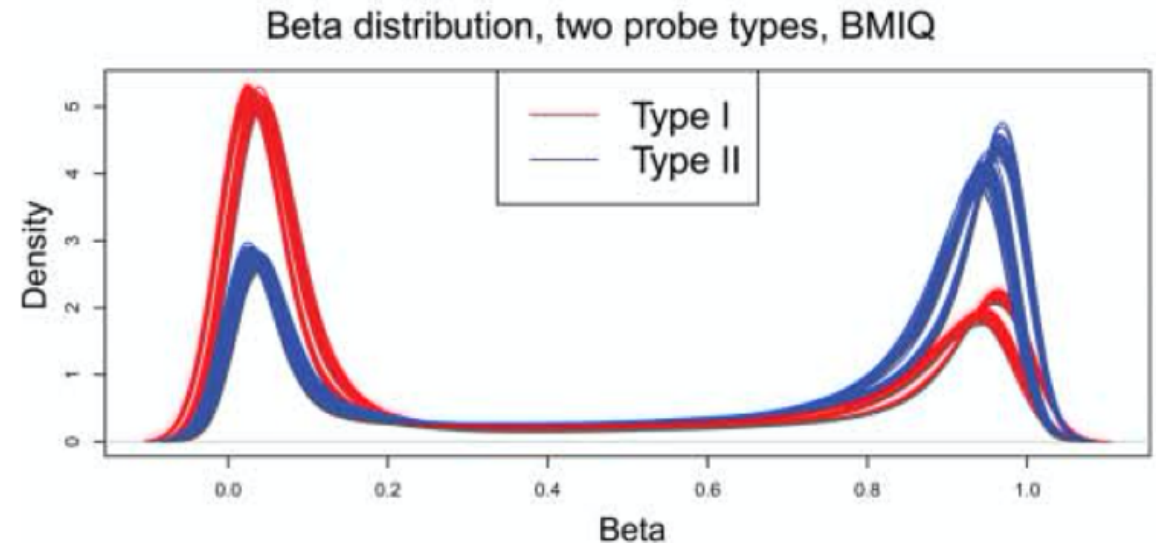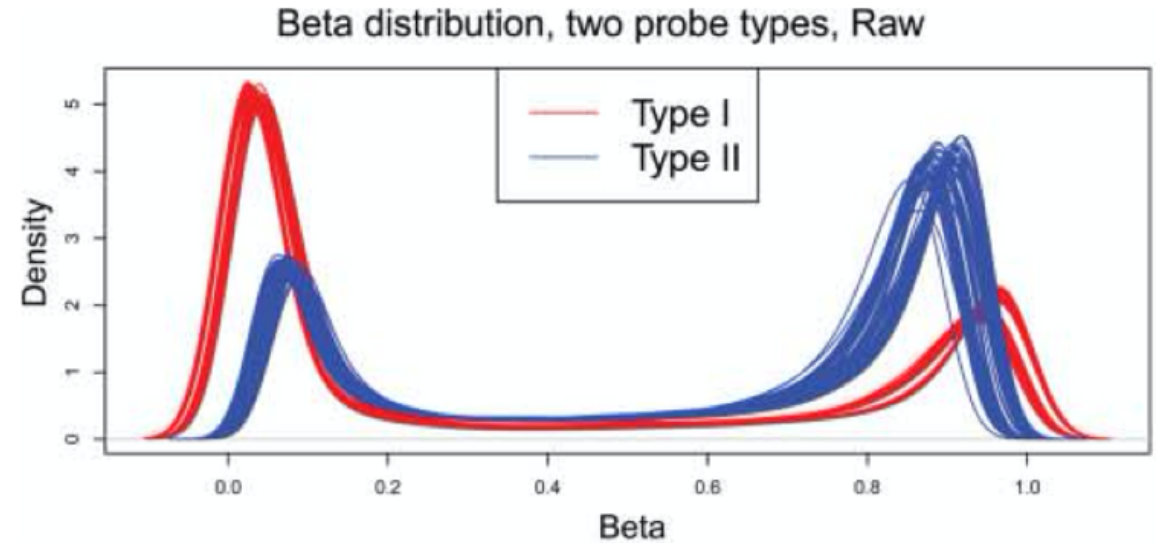
Func
array
stud

Jean-Philippe Fortin[1], Aurélie Labbe[2,3,4], Mathieu Lemire[5], Brent W Zanke[6], Thomas J Hudson[5,7], Elana J Fertig[8], Celia MT Greenwood[2,9,10] and Kasper D Hansen[1,11]*

# Normalization

- Within and across array normalization

- Within array:

    - background correction

    - dye bias adjustment

    - Type I/II bias correction

- Between array:

    - starting material

    - labeling efficiency

- Good overview + described in lab

- An evaluation of processing methods for HumanMethylation450 BeadChip data

# Assess normalization case by case

- Within and across array normalization not always necessary

- Depends on biological signal



Beta distribution, two probe types, Raw
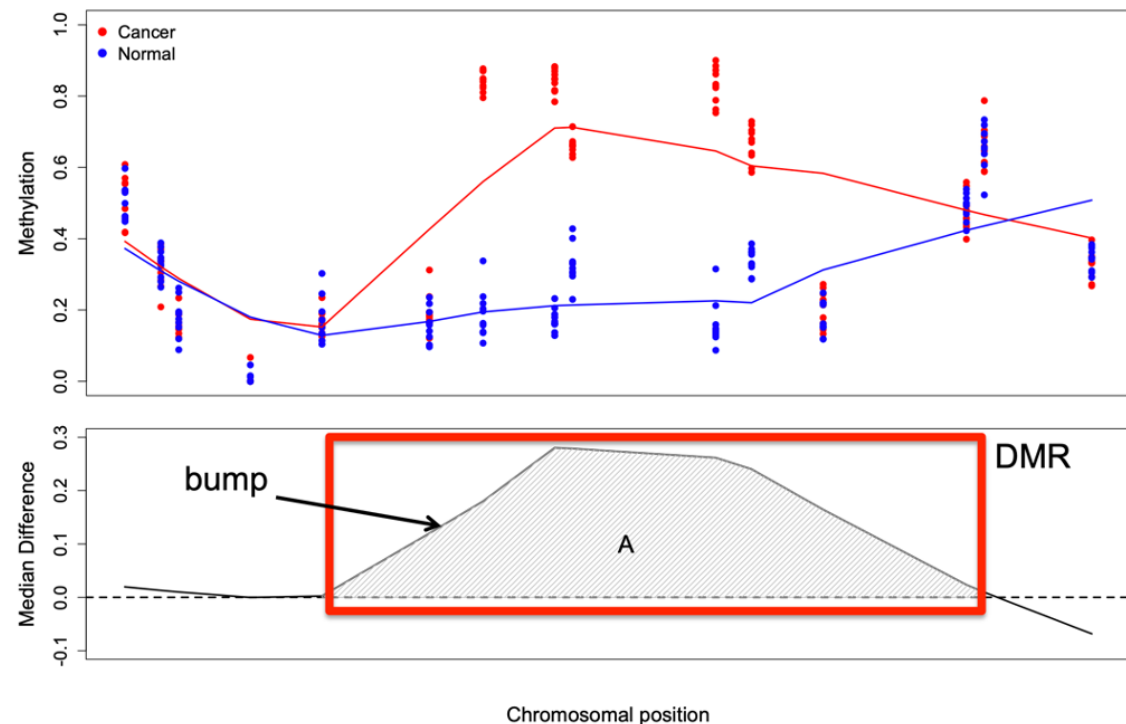
Beta distribution, two probe types, BMIQ

# Differential Methylation

- Identification of systematic differences in methylation between groups of samples (case vs control, smokers vs non-smokers, …)

- Usually starts on a per CpG basis

- Many ways to approach this
    - Questions being asked of data, available information on potential confounders, nature/structure of the data (repeated measures, …)

- Some possible approaches
    - T-test and ANOVA models
    - Wilcoxon rank-sum and Kruskall Wallis test
    - Linear, logistic and Cox regression or mixed effect models

- Use M-values: $M = log2(M/U)$ and *Beta* minimal difference cutoff

# Differential Methylation

- Single CpG often less informative than region (DMR)

- How to define region?

  - Sliding window

  - Heuristic cutoff

  - Functional units

- We will try last two in the lab

# Gene Set Enrichment

- Long list of DMP or DMR…. What does it mean?

- Gene expression -> GO analysis

- Not so straightforward for methylation data!

  - CpG link to genes unclear

  - Directionality
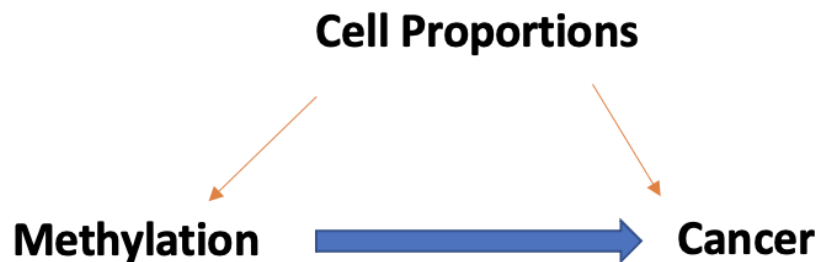
  - Bias! Number of CpG per gene differs

## Gene-set analysis is severely biased when applied to genome-wide methylation data

Paul Geeleher[1,2], Lori Hartnett[3], Laurance J. Egan[3], Aaron Golden[4], Raja Affendi Raja Ali[3] and Cathal Seoighe[2,*]

- missMethyl, methylGSA, BioMethyl

# Cell Type Deconvolution

- Estimates the relative proportion of pure cell types within a sample
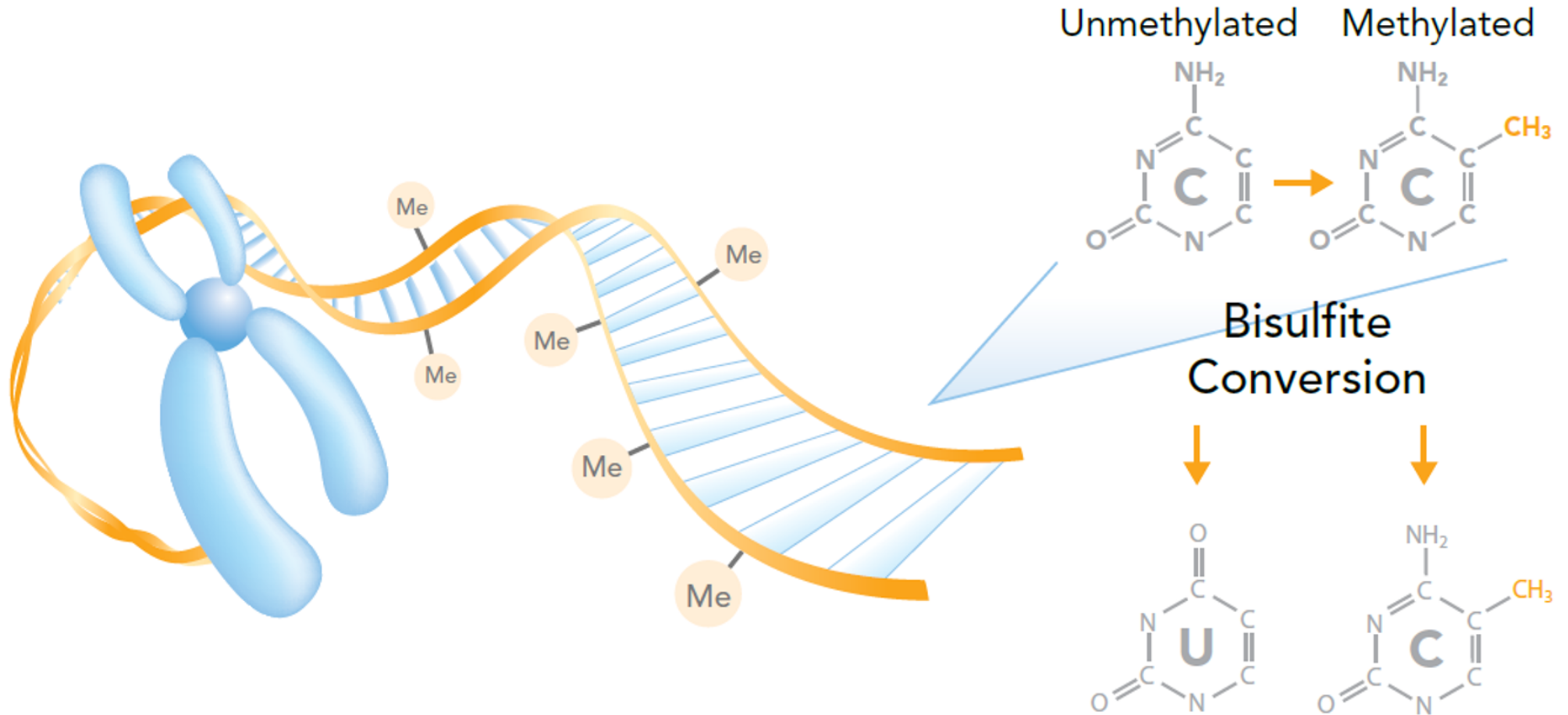
- Most cohort studies use data from blood samples



- Minfi: RGChannelSet returns relative proportions of CD4+ and CD8+ T-cells, NK cells, monocytes, granulocytes and B-cells in each sample

# Datasets

- Small toy data

- IDAT files

- 10 samples: 4 different T-cell types from 3 individuals

  - Naive

  - Treg

  - act_naive

  - act_Treg

- An additional sample has been added from another study GSE51180, to illustrate approaches for identifying poor quality samples.

# Bisulfite Sequencing

# Bisulfite Sequencing



Unmethylated    Methylated

Bisulfite Conversion

# Easy readout… in theory



Epigenomics Data analysis 2023: Methylation

# … but not in reality



- 2 different PCR product and 4 possible different sequence strands from one genomic locus

- Each of these 4 can exist in any possible conversion state

# 3-letter alignment



sequence of interest    TTGGCATGTTTAAACGTT

bisulfite convert read (treat sequence as both forward and reverse strand)

5'...**TT**GG**T**ATGTTTAAA**TGT**T...3'    5'...TT**AA**CAT**A**TTTAAAC**A**TT...3'

(1)    (2)

align to bisulfite converted genomes

(3)    (4)

...**TT**GG**T**ATGTTTAAA**TGT**T...    ...CC**AA**CAT**A**TTTAAAC**A**CT...
...**AA**CC**A**TACAAATTT**A**CAA...    ...GG**TT**GTA**T**AAATTTG**T**GA...
forward strand C -> T converted genome    forward strand G -> A converted genome
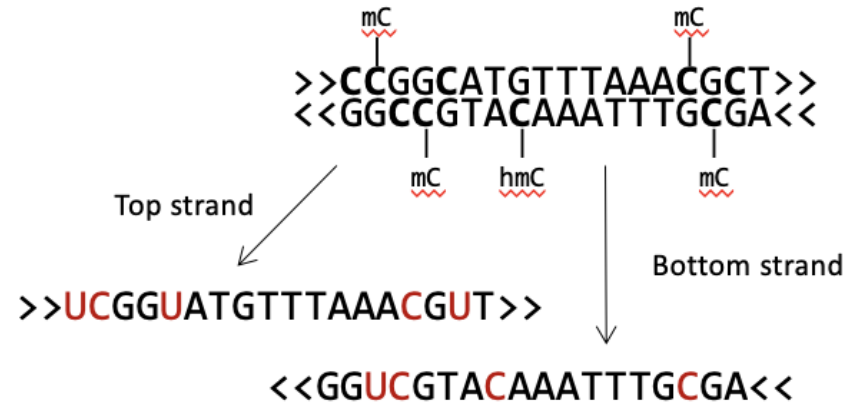(equals reverse strand C -> T conversion)

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

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

5'...CCGGCATGTTTAAACGCT...3'

read sequence      TTGGCATGTTTAAACGTTA
genomic sequence   CCGGCATGTTTAAACGCTA

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

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

Epigenomics Data analysis 2023: Methylation

# Common library preparations



## 1) Directional libraries
(vast majority of kits, also EpiGnome/Truseq)

## 2) PBAT libraries

## 3) Non-directional libraries
(e.g. single-cell BS-Seq, Zymo Pico Methyl-Seq)
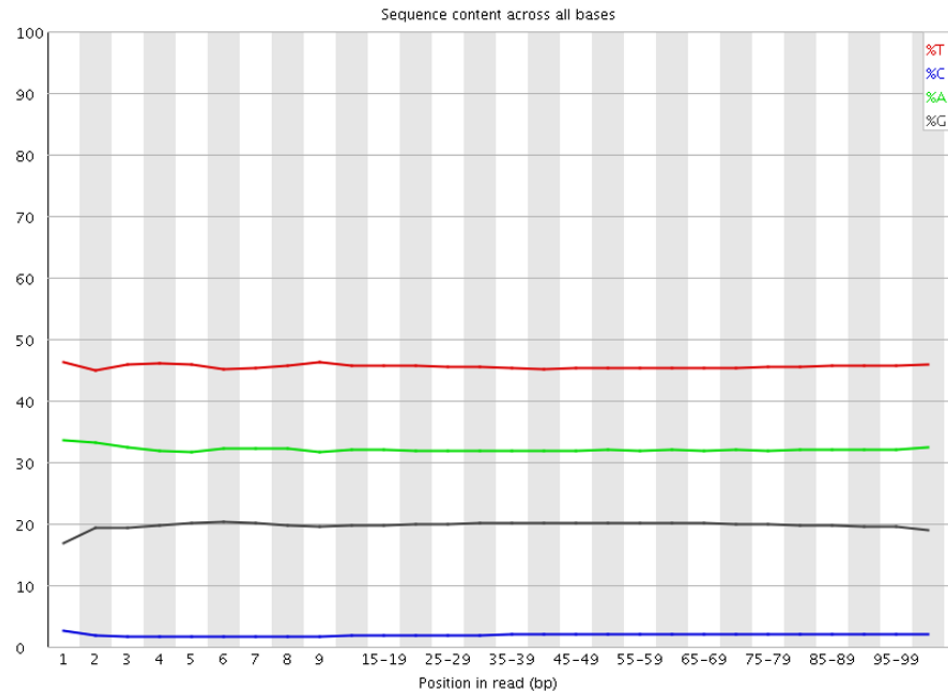
# Quality Control is essential

- Accurate C >T detection

- Pre-alignment

  - Base quality/composition

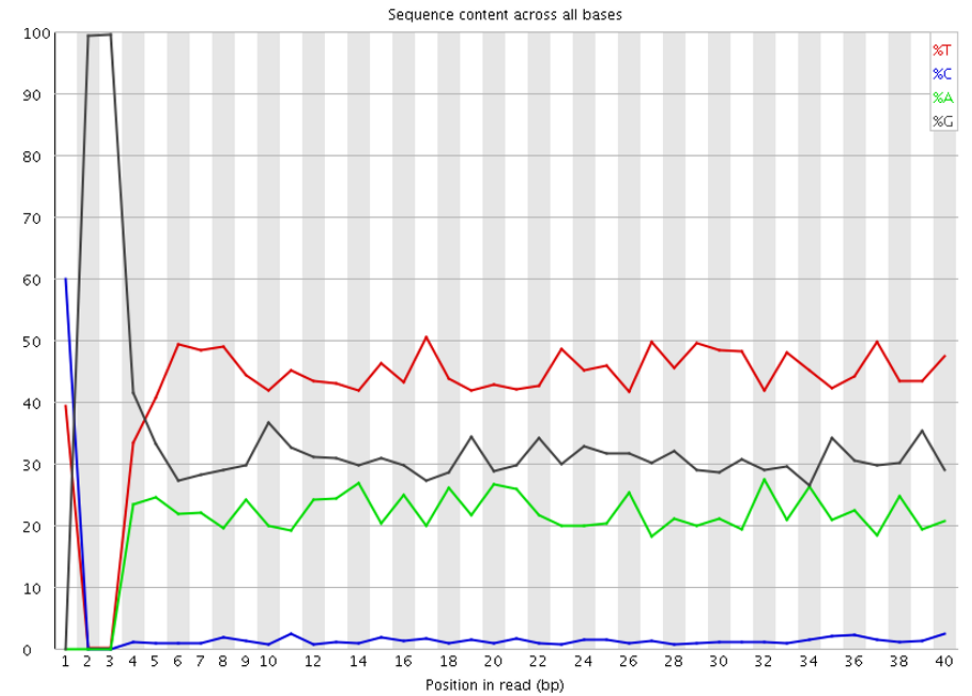  - Duplication levels

  - Adapter removal

# Average Base Quality



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

Error rate

0.1%

1%

10%

Position in read (bp)

Epigenomics Data analysis 2023: Methylation
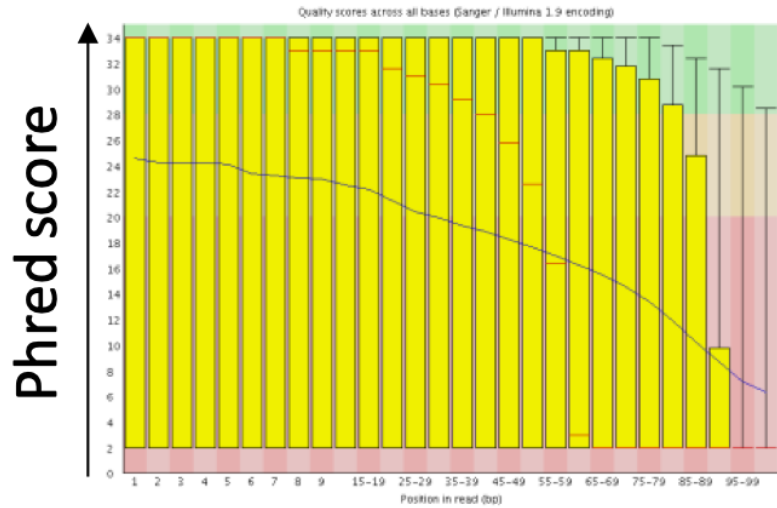
# Base Composition

WGBS

RRBS

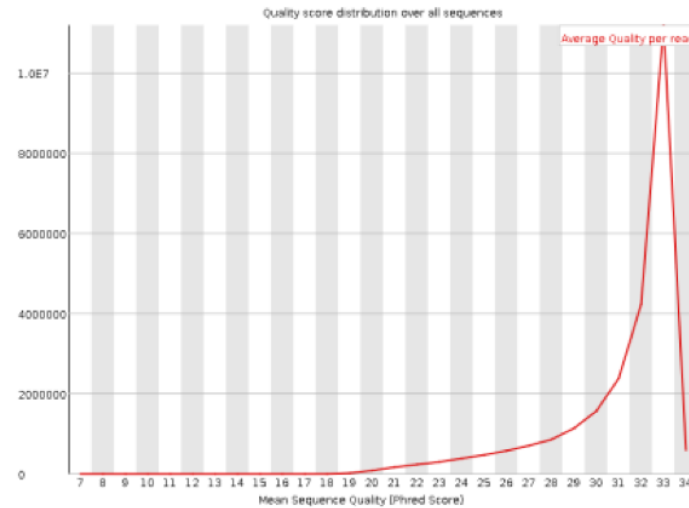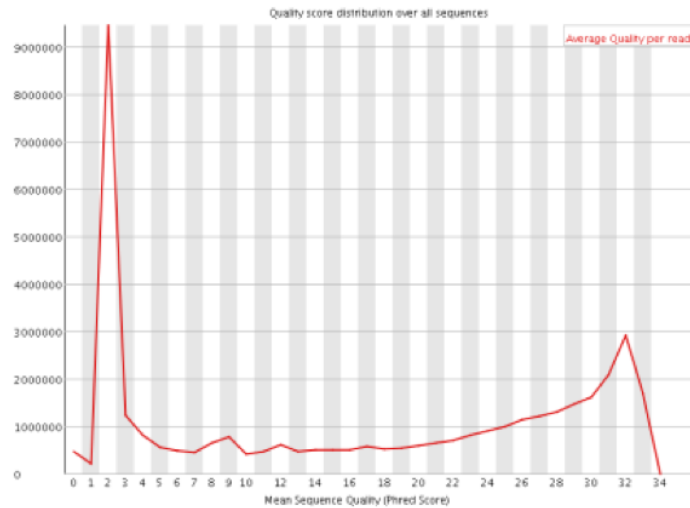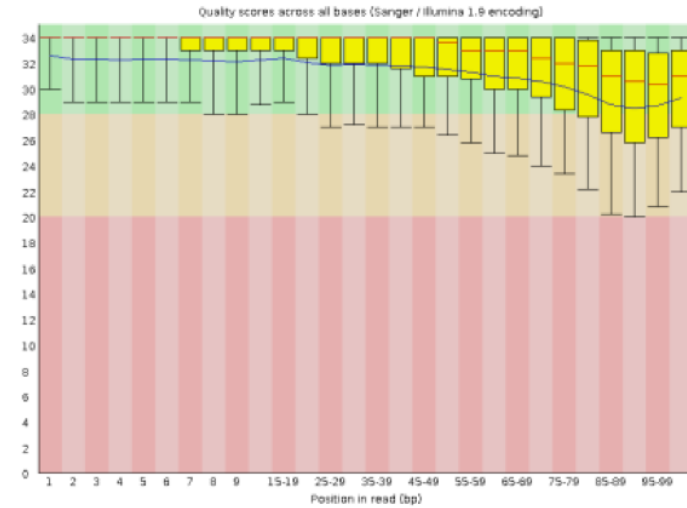# Common bisulfite sequencing QC issues

Not observed in ChIP or RNA-Seq

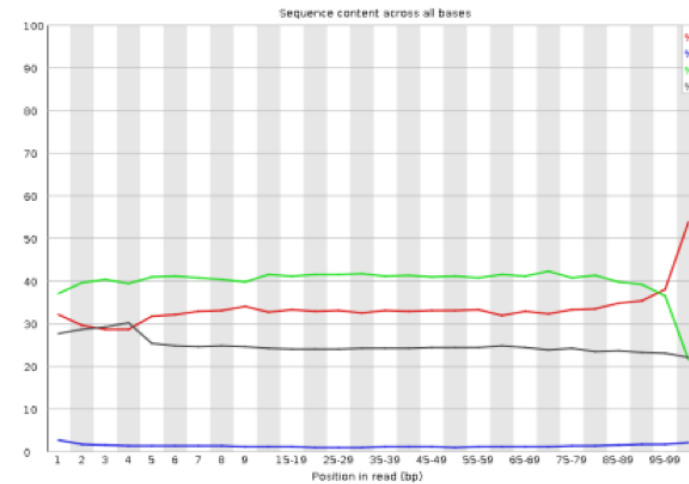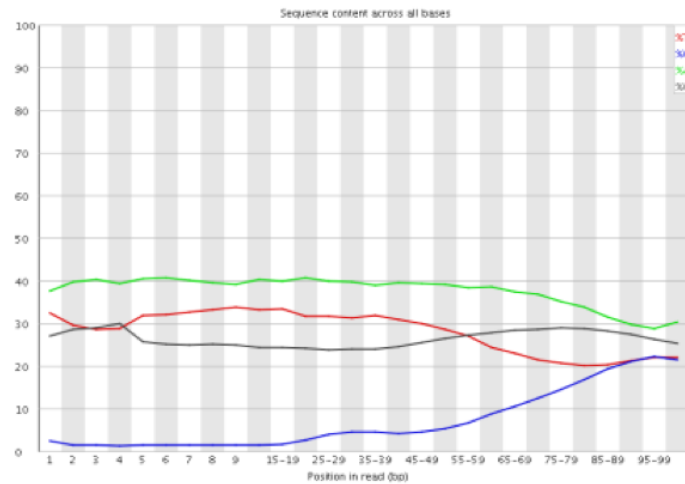# Remove poor quality basecalls

**before trimming**          **after trimming**



Epigenomics Data analysis 2023: Methylation

# Remove adapter contamination



Epigenomics Data analysis 2023: Methylation

# Summary Adapter/Quality trimming

Important to trim, if not:

- Low mapping efficiency
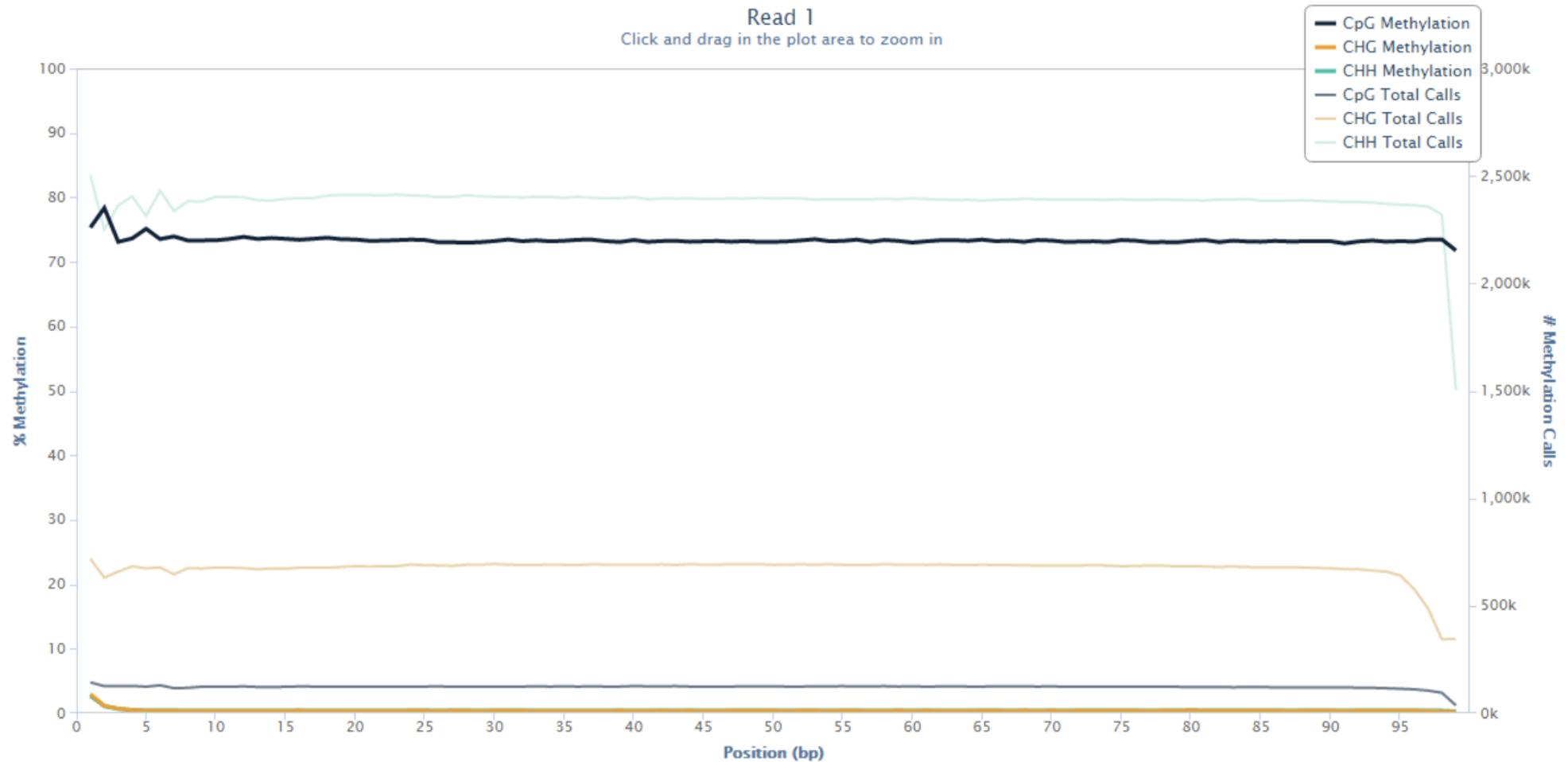
- misalignments

- errors in methylation calls (adapters are methylated)

- basecall errors

# Quality Control is essential

- Post-alignment

  - Incomplete conversion? non-CpG should be near 100%

  - Degradation? Check alignment rates and insert length

  - Average methylation levels

  - PCR bias? Deduplicate?

# M-bias

Average methylation levels across the entire length of the read

# M-bias

## Average methylation levels across the entire length of the read



```
1  bismark_methylation_extractor --ignore_r2 2 --gzip sample1_bismark_bt2_pe.bam
```

Epigenomics Data analysis 2023: Methylation
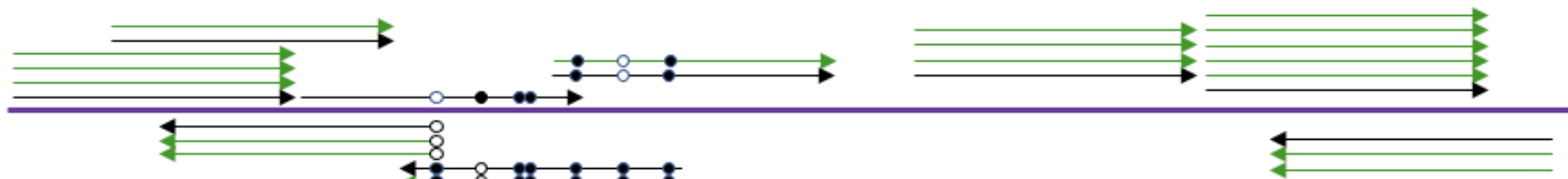
# Sequence duplication

**Complex/diverse library:**

**Duplicated library:**

percent methylation      55   17   100/100   100   71   100

deduplication

percent methylation      33   50   100/100   100   50   100

Epigenomics Data analysis 2023: Methylation

# Deduplication?

Advisable for large genomes and moderate coverage

- Unlikely to sequence several genuine copies

- Should have sufficient coverage, even after dedup

NOT advisable for RRBS or other target enrichment methods - high coverage expected

# Workflow



- nf-core pipeline: methylseq (see Thursday)
- Preprocessing + QC
  - 2 aligners: Bismark or bwa/meth/MethylDackel
  - QC: qualimap, preseq an multiqc
- Output ready for downstrea analysis

# MethylKit: R package



Can read Bismark coverage files as input

# Descriptive statistics

Coverage file Bismark

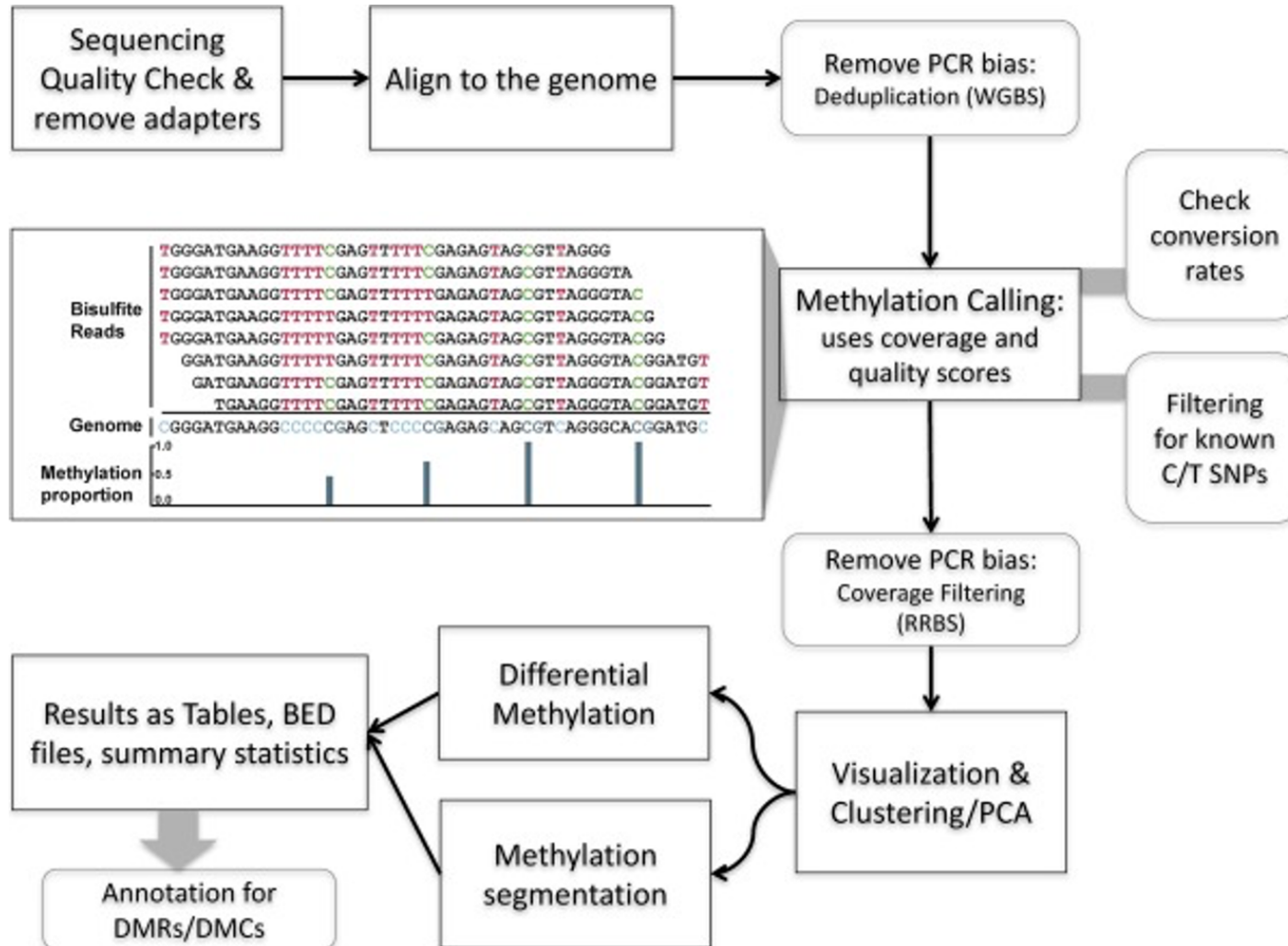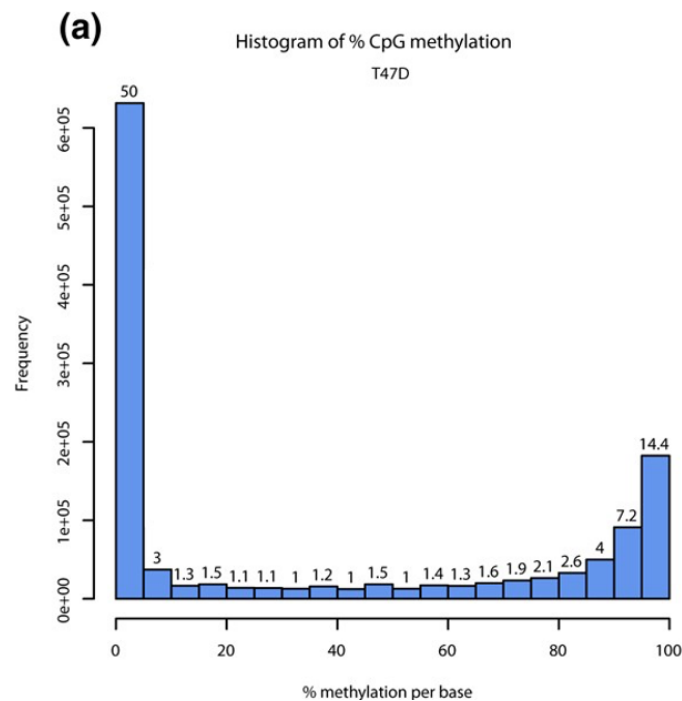| Chr | Start | End | Methylation Prop. | # mC | # C |
|---|---|---|---|---|---|
| chr8 | 3052997 | 3052997 | 0.00000 | 0 | 1 |
| chr8 | 3052998 | 3052998 | 53.26087 | 49 | 43 |
| chr8 | 3068732 | 3068732 | 57.14286 | 8 | 6 |
| chr8 | 3068733 | 3068733 | 100.00000 | 11 | 0 |
| chr8 | 3089948 | 3089948 | 100.00000 | 5 | 0 |
| chr8 | 3089984 | 3089984 | 100.00000 | 5 | 0 |



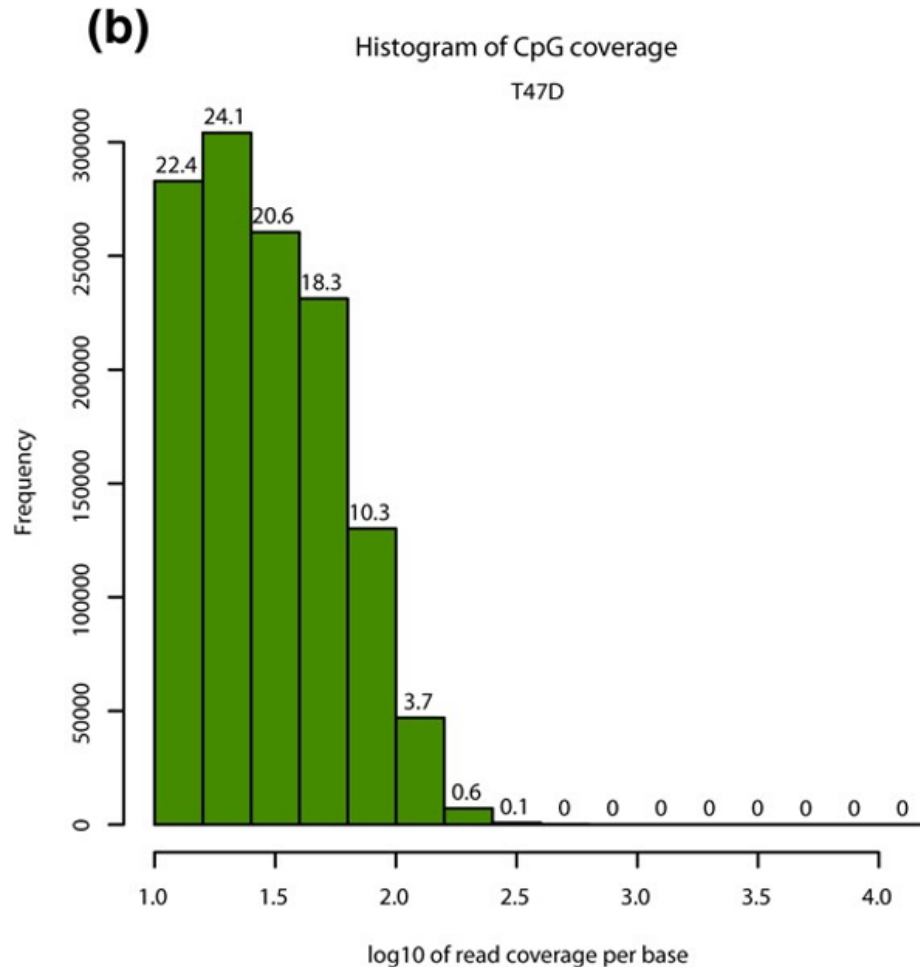(a) Histogram of % CpG methylation T47D

```
1   # Define the list containing the bismark coverage files.
2   file.list <- list(
3       "/sw/courses/epigenomics/DNAmethylation/biseq_data/P6_1.bismark.cov.gz"
4       "/sw/courses/epigenomics/DNAmethylation/biseq_data/P6_4.bismark.cov.gz"
5       "/sw/courses/epigenomics/DNAmethylation/biseq_data/P8_3.bismark.cov.gz"
6       "/sw/courses/epigenomics/DNAmethylation/biseq_data/P8_6.bismark.cov.gz"
7
8   # read the listed files into a methylRawList object making sure the other
9   # parameters are filled in correctly.
10  myobj <- methRead(file.list,
11          sample.id=list("Luminal_1","Luminal_2","Basal_1","Basal_2"),
12          pipeline = "bismarkCoverage",
13          assembly="mm10",
14          treatment=c(1,1,0,0),
15          mincov = 10
16          )
17
18  # Get a histogram of the methylation percentage per sample
19  # Here for sample 1
20  getMethylationStats(myobj[[1]], plot=TRUE, both.strands=FALSE)
```

Epigenomics Data analysis 2023: Methylation

# Descriptive statistics

## Coverage Distribution

```
1  # Get a histogram of the read coverage per sample
2  getCoverageStats(myobj[[1]], plot=TRUE, both.strands=FALSE)
```



(b)

**Histogram of CpG coverage**
T47D

- Secondary peak towards the right -> PCR duplication?

- Filter cutoff?

```
1  myobj.filt <- filterByCoverage(myobj,
2                                  lo.count=10,
3                                  lo.perc=NULL,
4                                  hi.count=NULL,
5                                  hi.perc=99.9)
```
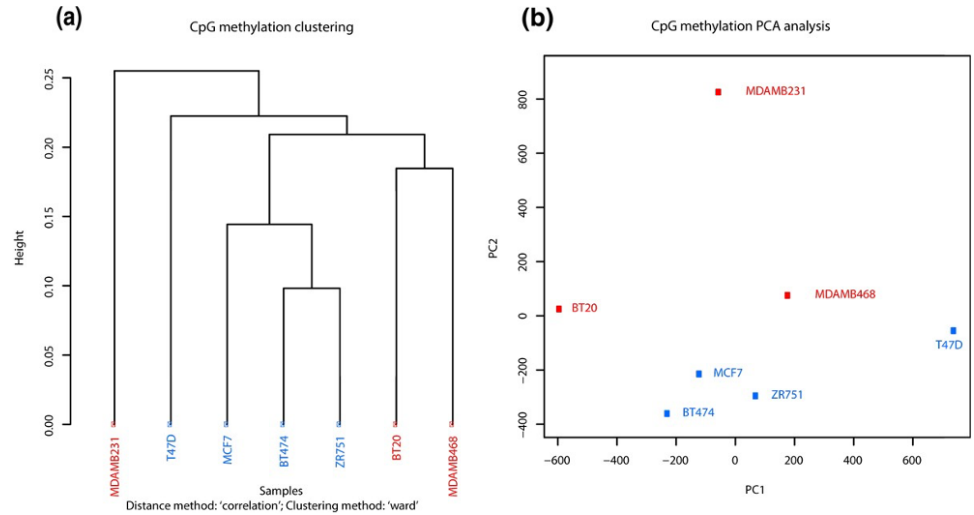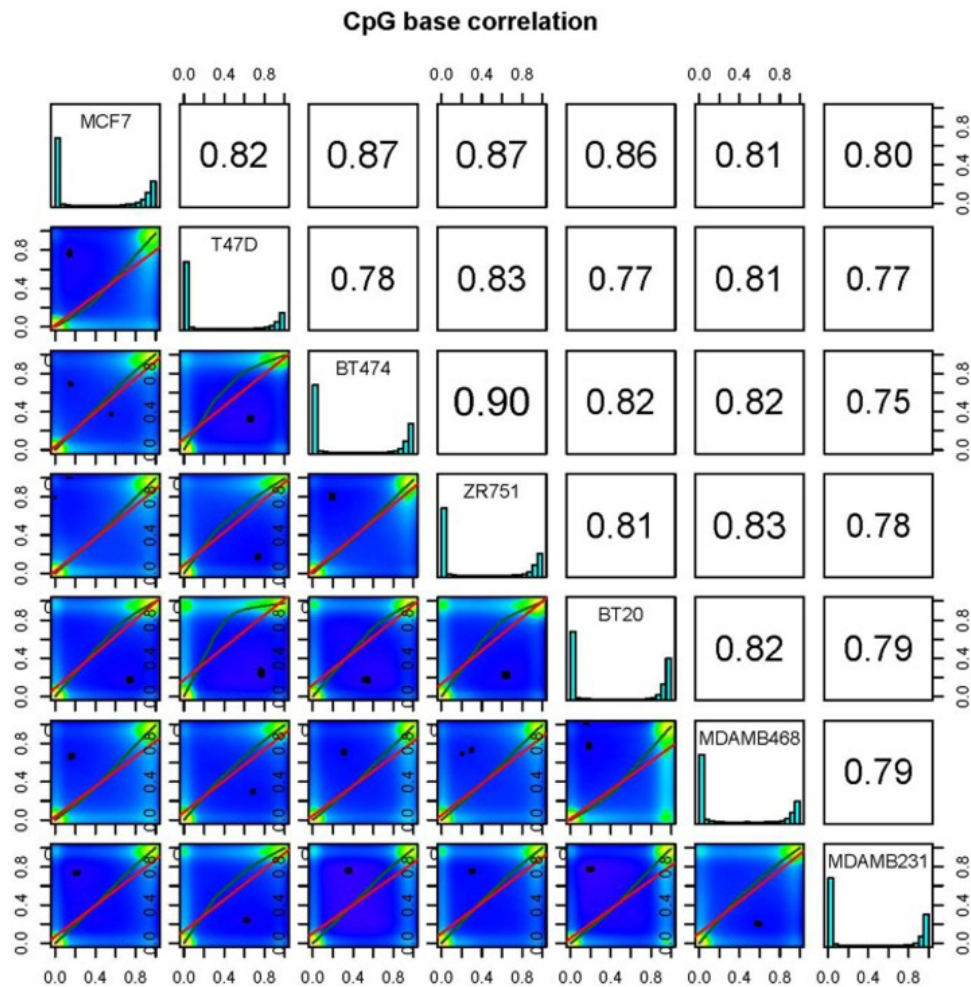
# Filtering

## Remove CpG that have no variation

```r
1   # get percent methylation matrix
2   pm=percMethylation(meth)
3
4   # calculate standard deviation of CpGs
5   sds=matrixStats::rowSds(pm)
6
7   # Visualize the distribution of the per-CpG standard deviation
8   # to determine a suitable cutoff
9   hist(sds, breaks = 100)
10
11  # keep only CpG with standard deviations larger than 2%
12  meth <- meth[sds > 2]
```

## Remove SNP overlap

```r
1   # give the locations of 2 example SNPs
2   mut <- GRanges(seqnames=c("chr21","chr21"),
3           ranges=IRanges(start=c(9853296, 9853326),
4                           end=c( 9853296,9853326)))
5
6   # select CpGs that do not overlap with mutations
7   meth <- meth[!as(meth,"GRanges") %over% mut, ]
```

# Sample Structure



**CpG base correlation**



(a) CpG methylation clustering
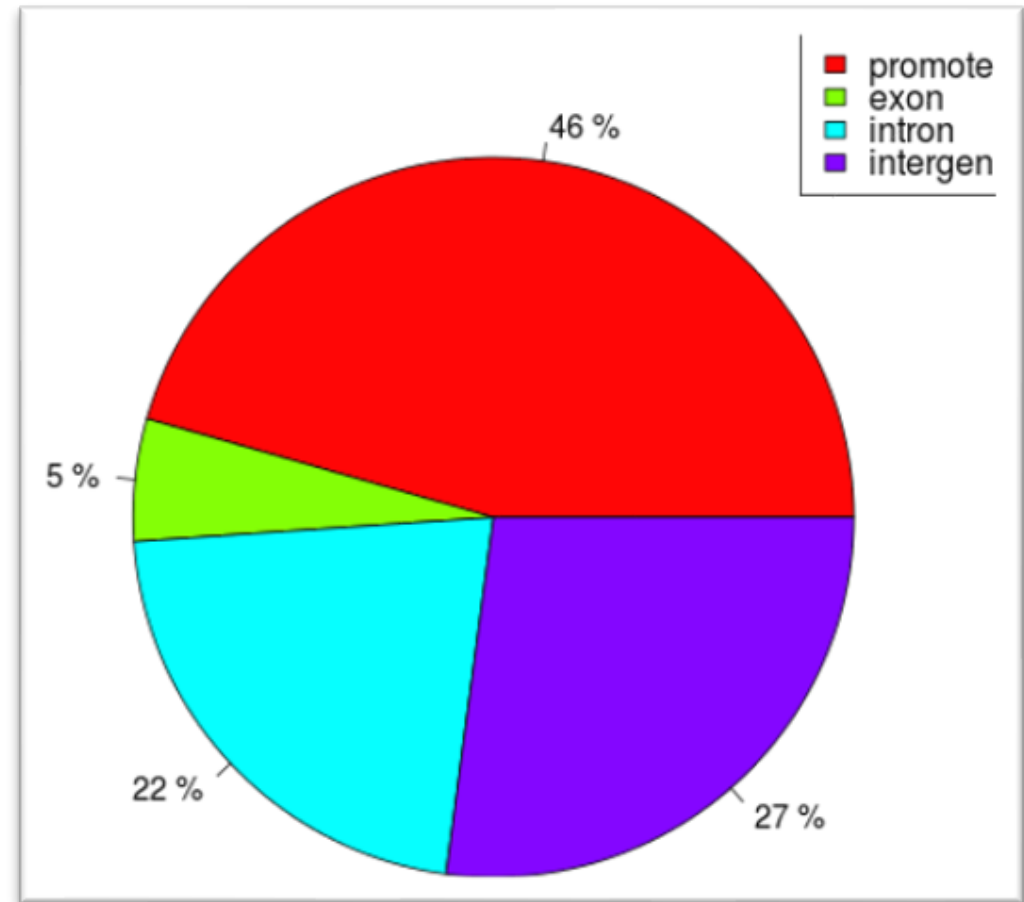
(b) CpG methylation PCA analysis

```
1  getCorrelation(meth,plot=TRUE)
2  clusterSamples(meth, dist="correlation",
3                 method="ward", plot=TRUE)
4  PCASamples(meth)
```

# Differential Methylation

- Many choices; often calculated by comparing proportion in methylated Cs in a test relative to control

- No replicates: Fisher's exact test

- Replicates:

  - linear regression

  - logistic regression (works with [0-1] data)

  - Beta-binomial (count data)

- Regression models can add covariantes/confounders

- Aggregate in regions (see lab)

# Annotate results

- How to interpret the DMR/DMPs?

- Integrate with genome annotation datasets

  - Where in relation to gene/regulatory region?

- Genomation R packge: toolkit for annotation

- Lab: basic annotation transcripts and CpG islands

- Requires some knowledge of R (GenomicRanges package)

# Remarks

- Normalization somewhat less important for bisulfite sequencing (Fisher's exact is sensitive to sequencing depth)

- Gene enrichments is as difficult as for arrays, not many implemented methods (rGREAT, Goseq)

# Lab

- Small dataset of mammary gland cells in mouse

- 4 samples: 2 luminal, 2 basal

- Bismark coverage files

| Chr | Start | End | Methylation Prop. | # mC | # C |
|-----|-------|-----|-------------------|------|-----|
| chr8 | 3052997 | 3052997 | 0.00000 | 0 | 1 |
| chr8 | 3052998 | 3052998 | 53.26087 | 49 | 43 |
| chr8 | 3068732 | 3068732 | 57.14286 | 8 | 6 |
| chr8 | 3068733 | 3068733 | 100.00000 | 11 | 0 |
| chr8 | 3089948 | 3089948 | 100.00000 | 5 | 0 |
| chr8 | 3089984 | 3089984 | 100.00000 | 5 | 0 |