



DNA methylation data analysis

Yujia Qin

Bioinformatics Core, JABSOM

Department of Quantitative Health Science

University of Hawaii

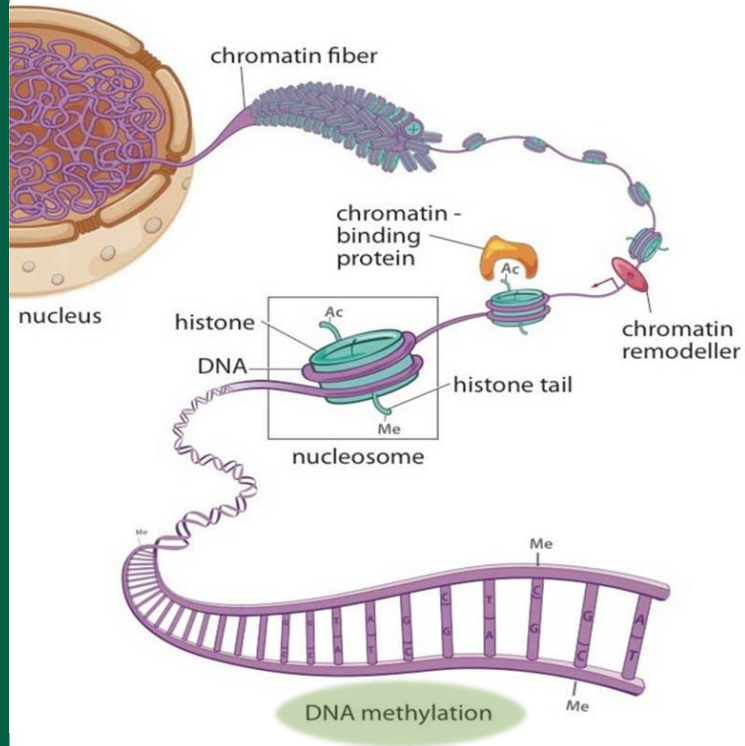
2022



1. DNA methylation

Introduction

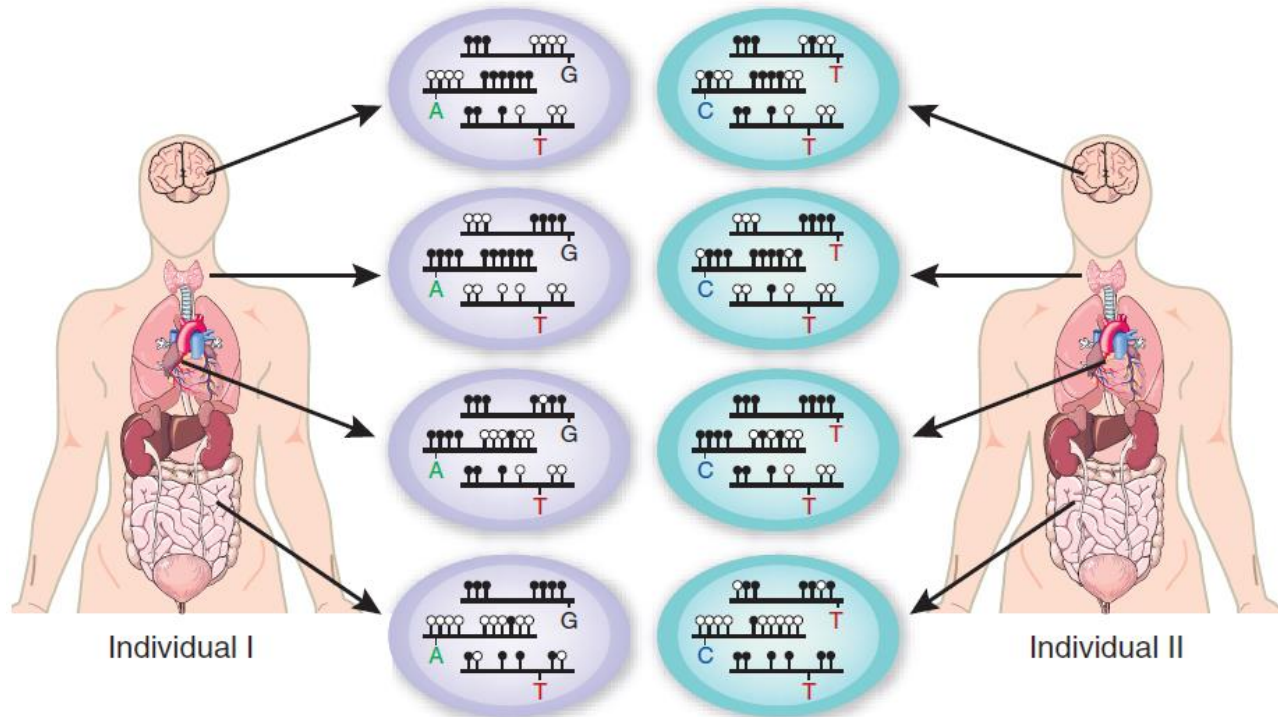
What is DNA methylation?



- **DNA methylation** is a major epigenetic mechanism involving direct chemical modification to the DNA
- **Epigenetics** is the study of the mechanisms that cause changes in gene expression but are not due to the change in the DNA sequences, including DNA methylation, histone modification, non-coding RNA expressions



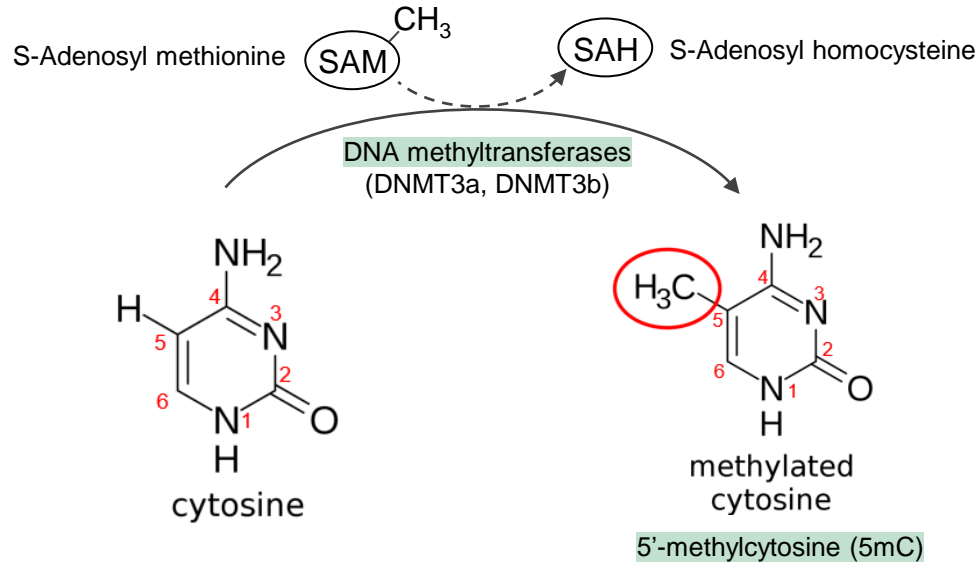
Tissue-specific DNA methylation and epigenetic heterogeneity among individuals



Kim Caesar



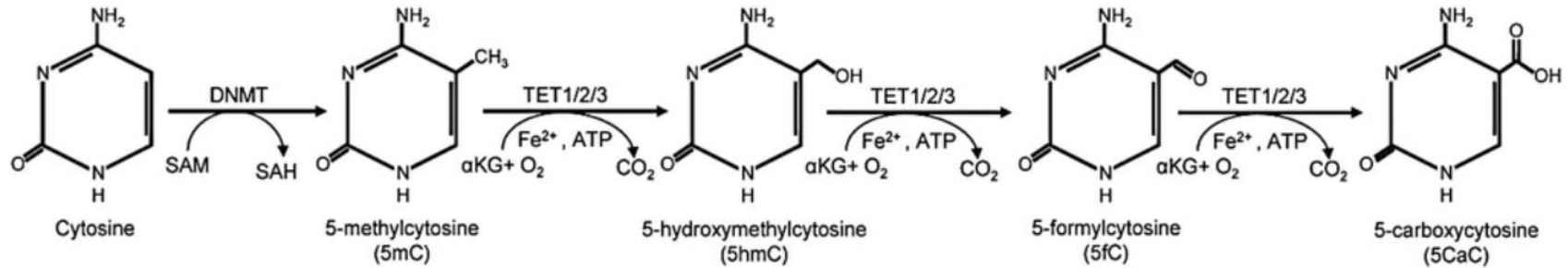
What is DNA methylation?



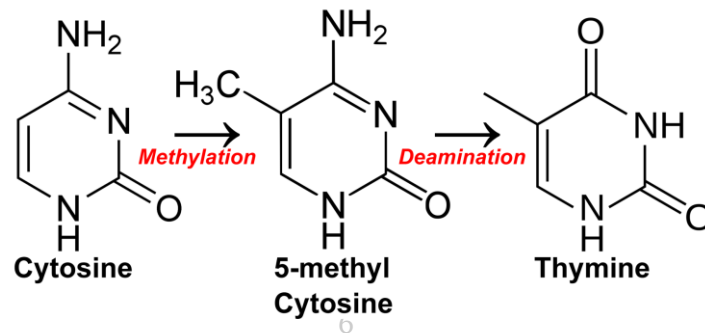
- Most common form of DNA methylation is 5-methylcytosine (5mC)
- Three types of DNA methyltransferases (DNMT's): DNMT1, DNMT3a, DNMT3b
- 5mC is inherently mutagenic: spontaneously undergo deamination, leading to C->T transitions.



Oxidized Derivatives of 5-MethylCytosine



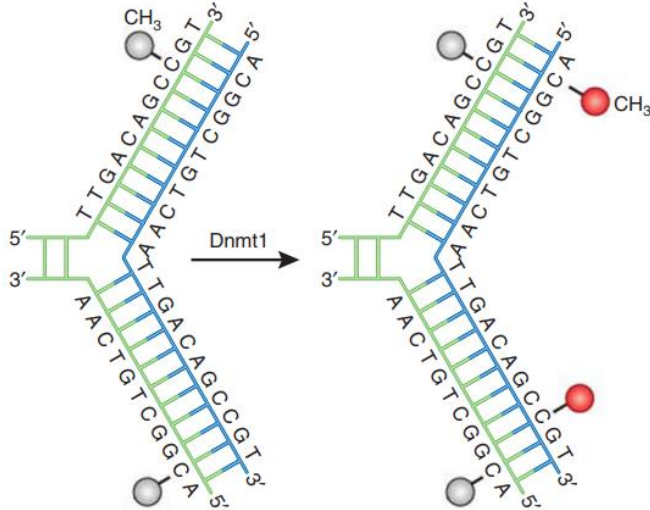
Cytosine Deamination



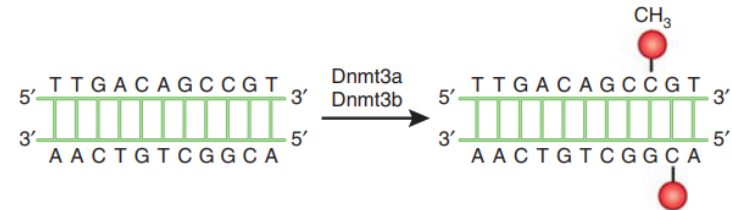


DNA methylation enzymes

DNMT1: maintenance DNA methyltransferase:



DNMT3a, DNMT3b: de novo DNA methyltransferase:





DNA methylation **locations**

- DNA methylation primarily happens at CpG sites.
- **CpG dinucleotides/sites**
 - DNA regions where a cytosine is followed by a guanine in the 5'-3' direction (5'-Cytosine-phosphate-Guanine-3')
 - CG dinucleotides present on single strand of DNA
 - Under-representation
- **CpG island** are regions with a high frequency of CpG sites
 - Length > 200bp
 - GC percentage > 50%
 - Observed-to-expected CpG ratio > 60%
- CpG island '**shores**' are regions of comparatively low CpG density, located approximately 2 kb from CpG islands



CpG island

CpG sites

```
CCCGGGTCGGGGGGAAGAAGCCCTCAAAGGCAAGGCCATCCGCGA
GAAGGCCAGCCCCCGCCGCTGCAAGCCAGGCGCGCGCTCCCGCTG
GGCTGCTCCCTCGGGCCCTGCAAGCCCTCTCTGCTACTTTGGAGCGCTTC
CTCAAGCCCTCTCTCAAGCCCGCGCCAGGCTCCCGCGCGCAAGCTGGGG
ATCTCGGCGCAATAAAGGAAGAAAGGGCGCGCGCTACCGCGCGCAAGTGC
GTGGGCGAGCAGAGCTCAAGCCCTCTCTCAAGCGCGCGCGCGCGCGCG
ACAGCTGCTGGGCTGCAAGGAGCGGTAAGCGCGCGAGCAAGAGAGGCGCG
CTTGACTCGCACTTTTGTGCGTTGCGAAGCTCTCTCTCAAGTGGTGGTGC
AATGCGAGCGCGCTCTTAAATCAAGTGGCGCTCAAGAGCTCAAGTAAAG
GTACAGGCGCTTCGCGCGCGAGTGCCTCGCCCTCAAGCGCGCTCGCCCT
CGGGGATGCCCAAGCCCTGTGGCGCTTCAGCGCTCCCGCGCGAGGCG
CGCTCGGGCTGCGCTGGCTCTTCAGCAAGCGCGCGCGCGCGCGCGCGCG
TGCACTGTGTGAGCAAGCGCGCTTCAGCGCTTCAGCGCGCGCGCGCGCG
GGCGTGGTATTCAGTGCAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGGCGCTAGCGAGTGGAGCGCGCTGGCGCGCGCGCGCGCGCGCGCGCGCG
CGGGAAGCCCTGTCTTTCCCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGTCAACAAGCGCTGTCTTGGGTCAAGGACATCTCGCGCTCTGAAAG
ACCCCGCGCTCTTCGCGCGCGCACTCGCGCTCTGGCGCGCGCGCGCGCG
CGAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGCACTCTAAGGCGCTTCAGCTCTGGCGCGCGCGCGCGCGCGCGCGCG
TCTGCGCTCTGCTTGGGGGAGGGCGCTTTGGGGTCTTCAAGGGCGCGCG
GGAAGGCGCGCGCTGTCTTGGGTCCCGCGGAAGGGTGTGAGATTGAGGCC
CGAAGGCGCGCGCGCTGTGCAAGCGCTCTTCGCGCGCGCGCGCGCGCGCG
```

```
TGCCCCAGCTGAATCCAGCGGCGCGCGCGCTGCTGCTTGTCTTGTCTTCT
CGAGCTGGTATTAAGCGCGCTGCGCGCGCGCGCGCGCGCGCGCGCGCG
TCAGGAAATGCCCAAGCGGAAAGGAGGCGCGCGCGCGCGCGCGCGCGCG
CCAAAGAGGTGGCGCGCGGGAAGCAAGTGTCTTCTGGCGCTTCTGCTCTCT
CTAGGCTGTGACAGCGCACTCTCTGGAGCACTGCGCTGAGGAAGCGCGAG
CTCTGTGTGAGCGCAAGCACTGCGCAAGCGCTCTCTCTCAAGCTCTGCGAG
GAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CCCCCAAGCACTCAGCTCAGGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GAAGCGCTGTGAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
TCTTCAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
TGAGCTCCCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
TGAAGCTGTGACCTCTGGAAGCAAGCGCTCGCGGTAGGTGATGGGTAAAC
ATTCTCTAAATGGTGAAGTCACTGCGCTCTTCAAGCTGCGCGCGCGCG
TAAGCGCTGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGTCACTCTCTGCTCACTTCAAGCGCGCGCGCGCGCGCGCGCGCGCGCG
TTTGGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AAAGCGGGGGAAGCTGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGTAAGGGAAGGGAAGGCGCTGGGTGGCGCGCGCGCGCGCGCGCGCGCG
GCCAGCTTGGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGTGAAGTGGCTGTGTGTGCTGCTGCGCGCGCGCGCGCGCGCGCGCGCG
ACATAGCTGGGCTGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGAGCTGAGGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GGGCGCTGAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CAAGAAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GCTTGGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGGCGCGCTGGTGGGAGTGGGAGCGCGCGCGCGCGCGCGCGCGCGCGCG
CTTACTGTGTGCGAGGCTGCTGCTGGCGCGCGCGCGCGCGCGCGCGCG
CCAGTGGCTTCCCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CTGCAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GTAAGCTGCTGTGAGCTGCTGGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGTGAAGCTGGTGGAGGCTGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CCCTGTAAGCTTCTCTCTCTCTGCGCGCGCGCGCGCGCGCGCGCGCG
CCCTGTAAGCTTCTCTCTCTCTGCGCGCGCGCGCGCGCGCGCGCGCG
AGCGCAAGCTCTCAGCTGAGTCCCGCGCGCGCGCGCGCGCGCGCGCG
GCAAGTGAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GCTTCAAGCTTCTCTCTCTCTCTGCGCGCGCGCGCGCGCGCGCGCGCG
GCTTCAAGCTTCTCTCTCTCTCTGCGCGCGCGCGCGCGCGCGCGCGCG
```

Under-representation:
CG suppression

GpC sites

```
CCCGGGTCGGGGGGAAGAAGCCCTCAAAGGCAAGGCCATCCGCGA
GAAGGCCAGCCCCCGCCGCTGCAAGCCAGGCGCGCGCTCCCGCTG
GGCTGCTCCCTCGGGCCCTGCAAGCCCTCTCTGCTACTTTGGAGCGCTTC
CTCAAGCCCTCTCTCAAGCCCGCGCCAGGCTCCCGCGCGCAAGCTGGGG
ATCTCGGCGCAATAAAGGAAGAAAGGGCGCGCGCTACCGCGCGCAAGTGC
GTGGGCGAGCAAGCGCTCAAGCCCTCTCTCAAGCGCGCGCGCGCGCGCG
ACAGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTG
CTGCACTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTG
AAATGGAAGCGCGCTTAAATCAAGTGGCGCTCAAGAGCTCAAGTAAAG
GTACAGGCGCTTCGCGCGCGAGTGCCTCGCCCTCAAGCGCGCTCGCCCT
CGGGGATGCCCAAGCCCTGTGGCGCTTCAGCGCTCCCGCGCGAGGCG
CGCTCGGGCTGCGCTGGCTCTTCAGCAAGCGCGCGCGCGCGCGCGCGCG
TGCACTGTGTGAGCAAGCGCGCTTCAGCGCTTCAGCGCGCGCGCGCGCG
GGCGTGGTATTCAGTGCAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGGCGCTAGCGAGTGGAGCGCGCTGGCGCGCGCGCGCGCGCGCGCGCGCG
CGGGAAGCCCTGTCTTTCCCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGTCAACAAGCGCTGTCTTGGGTCAAGGACATCTCGCGCTCTGAAAG
ACCCCGCGCTCTTCGCGCGCGCACTCGCGCTCTGGCGCGCGCGCGCGCG
CGAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGCACTCTAAGGCGCTTCAGCTCTGGCGCGCGCGCGCGCGCGCGCGCG
TCTGCGCTCTGCTTGGGGGAGGGCGCTTTGGGGTCTTCAAGGGCGCGCG
GGAAGGCGCGCGCTGTCTTGGGTCCCGCGGAAGGGTGTGAGATTGAGGCC
CGAAGGCGCGCGCGCTGTGCAAGCGCTCTTCGCGCGCGCGCGCGCGCGCG
```

```
TCAGGAAATGCCCAAGCGGAAAGGAGGCGCGCGCGCGCGCGCGCGCG
CCAAAGAGGTGGCGCGCGGGAAGCAAGTGTCTTCTGGCGCTTCTGCTCT
CTAGGCTGTGACAGCGCACTCTCTGGAGCACTGCGCTGAGGAAGCGCGAG
CTCTGTGTGAGCGCAAGCACTGCGCAAGCGCTCTCTCTCAAGCTCTGCGAG
GAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CCCCCAAGCACTCAGCTCAGGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GAAGCGCTGTGAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
TCTTCAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
TGAGCTCCCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
TGAAGCTGTGACCTCTGGAAGCAAGCGCTCGCGGTAGGTGATGGGTAAAC
ATTCTCTAAATGGTGAAGTCACTGCGCTCTTCAAGCTGCGCGCGCGCG
TAAGCGCTGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGTCACTCTCTGCTCACTTCAAGCGCGCGCGCGCGCGCGCGCGCGCGCG
TTTGGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AAAGCGGGGGAAGCTGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGTAAGGGAAGGGAAGGCGCTGGGTGGCGCGCGCGCGCGCGCGCGCGCG
GCCAGCTTGGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGTGAAGTGGCTGTGTGTGCTGCTGCGCGCGCGCGCGCGCGCGCGCGCG
ACATAGCTGGGCTGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGAGCTGAGGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GGGCGCTGAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CAAGAAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GCTTGGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
AGGCGCGCTGGTGGGAGTGGGAGCGCGCGCGCGCGCGCGCGCGCGCGCG
CTTACTGTGTGCGAGGCTGCTGGCGCGCGCGCGCGCGCGCGCGCGCG
CCAGTGGCTTCCCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CTGCAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GTAAGCTGCTGTGAGCTGCTGGCGCGCGCGCGCGCGCGCGCGCGCGCG
CGTGAAGCTGGTGGAGGCTGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
CCCTGTAAGCTTCTCTCTCTCTCTGCGCGCGCGCGCGCGCGCGCGCGCG
CCCTGTAAGCTTCTCTCTCTCTCTGCGCGCGCGCGCGCGCGCGCGCGCG
AGCGCAAGCTCTCAGCTGAGTCCCGCGCGCGCGCGCGCGCGCGCGCG
GCAAGTGAAGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCGCG
GCTTCAAGCTTCTCTCTCTCTCTGCGCGCGCGCGCGCGCGCGCGCGCG
```

In mammalian genome,
70-80% CpGs are
methylated.

60% of mammalian
promoters have CpG
islands, but they are
mostly unmethylated



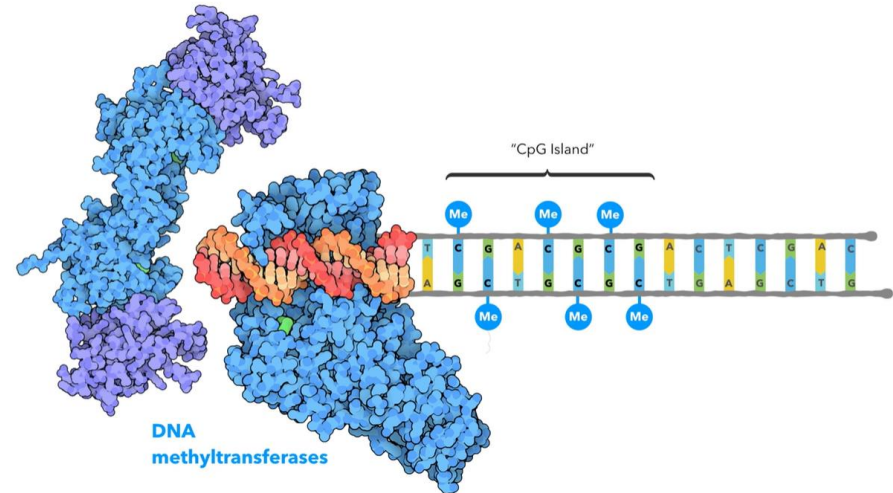
DNA methylation **functions** vs. **locations**

- DNA methylation effects on transcriptional regulation differ depending on the location of the CpG site (intragenic vs promoter region vs enhancer).
 - **Methylated** CpG island **promoters** are associated with **gene repression**.
 - CpG islands, occur and commonly span **promoters** of house-keeping genes. These promoter CpG islands typically remain unmethylated, resulting in active gene expression
 - **Gene bodies** tend to have intermediate CpG densities. Unlike CpG island promoters, extensive exonic or genic methylation is typically associated with **active gene expression**.
- Beyond these regions, the genome has a lower-than-expected frequency of CpG sites which are typically methylated



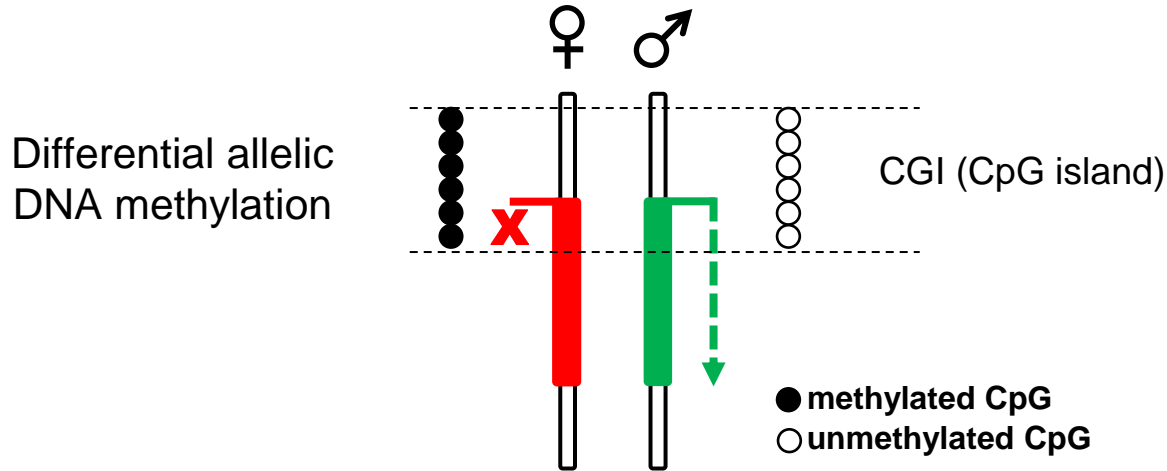
DNA methylation **functions**

- DNA methylation has been shown to be important for:
 - Genomic imprinting
 - Transposable element silencing
 - Stem cell differentiation
 - Embryonic development
 - Inflammation
 - Cancer





Imprinted Genes: mono-allelic expression



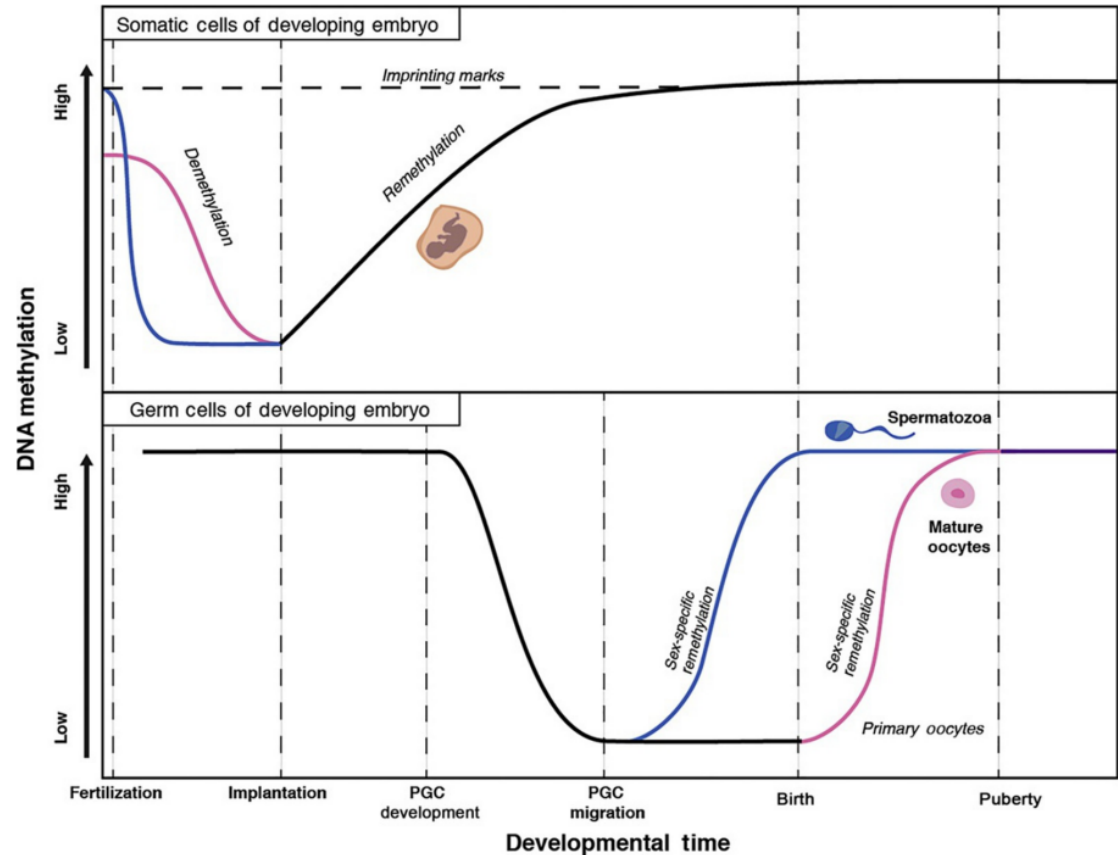
Imprinted Genes: Mono-allelic expression with parent-of-origin specificity.
Have key roles in energy metabolism, placenta functions.



DNA methylation is reset during reprogramming

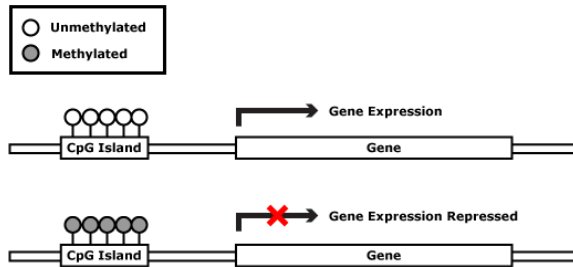
The genome undergoes two waves of global demethylation and re-methylation for the purpose of producing the next generation:

1. After fertilization
2. Germ cell development





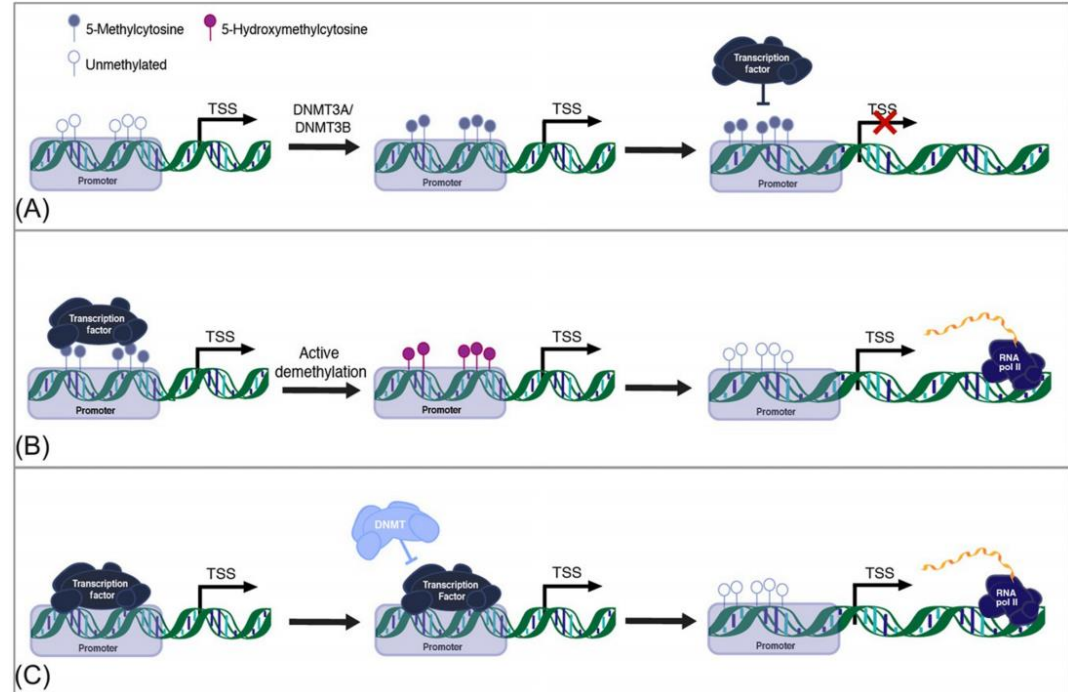
DNA methylation in gene promoter



Silencing of gene expression

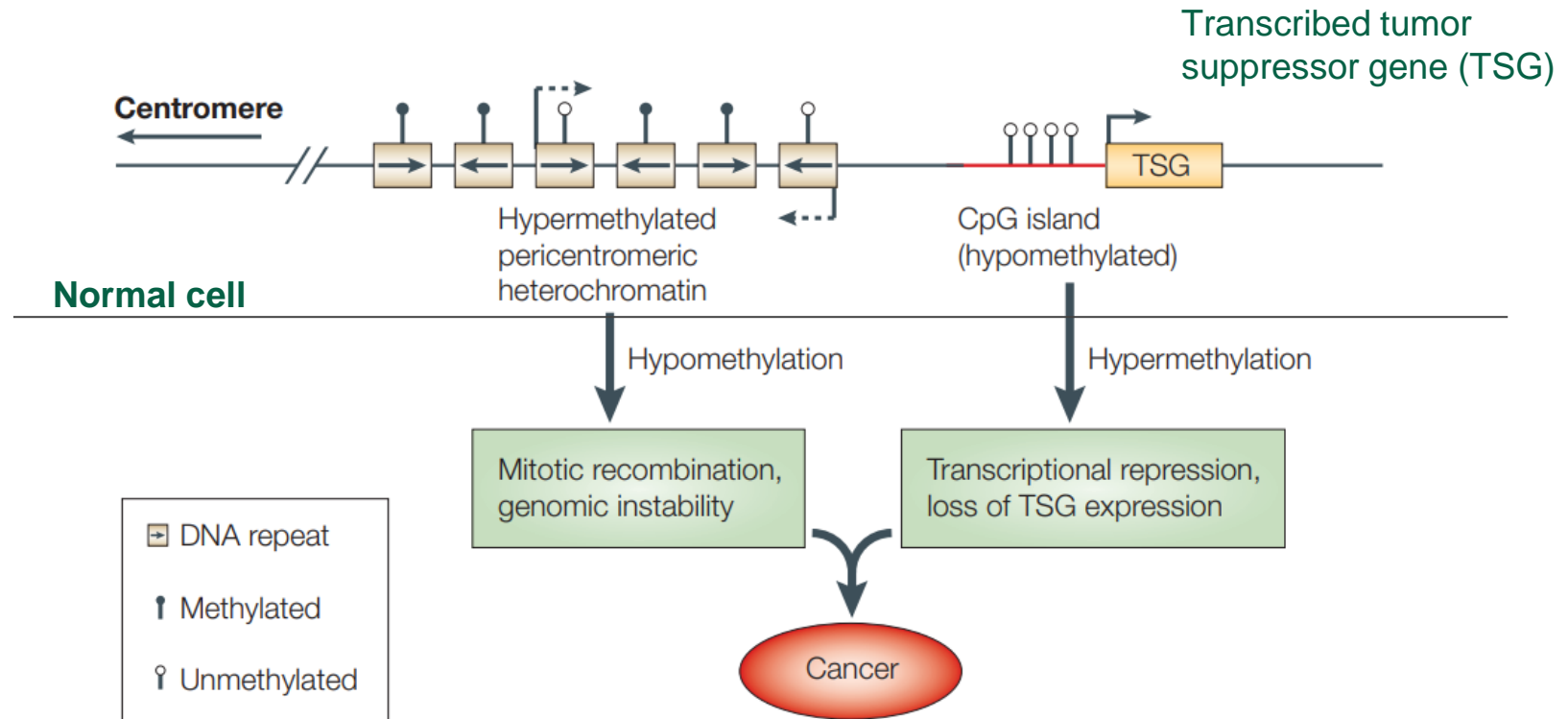
Two possible ways:

- Prevent TF from binding
- Wrap DNA up make it inaccessible.





DNA methylation and cancer



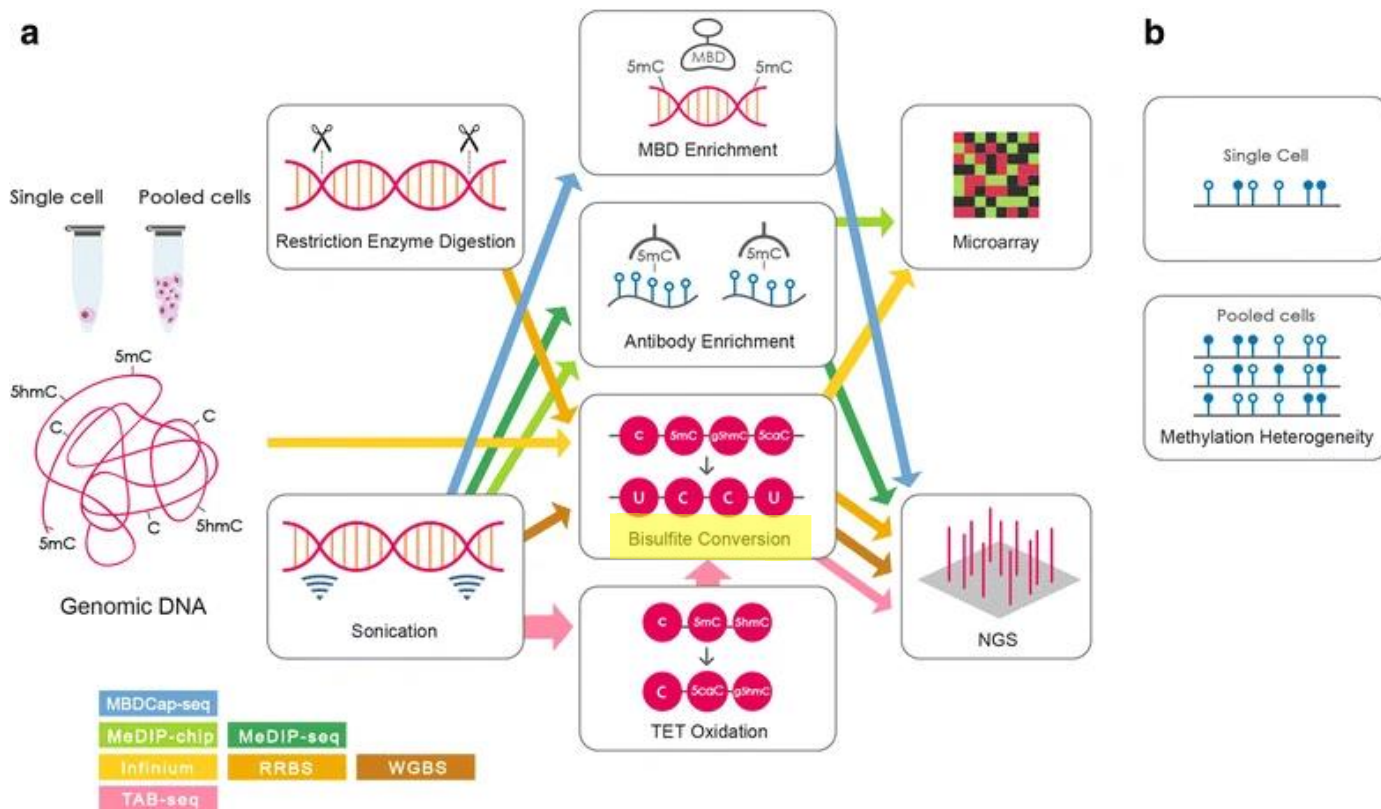


2. How to measure DNA methylation?

Bisulfite conversion



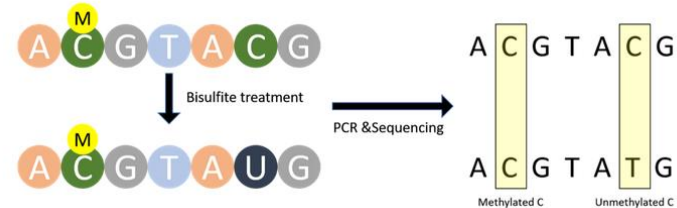
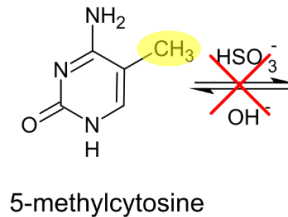
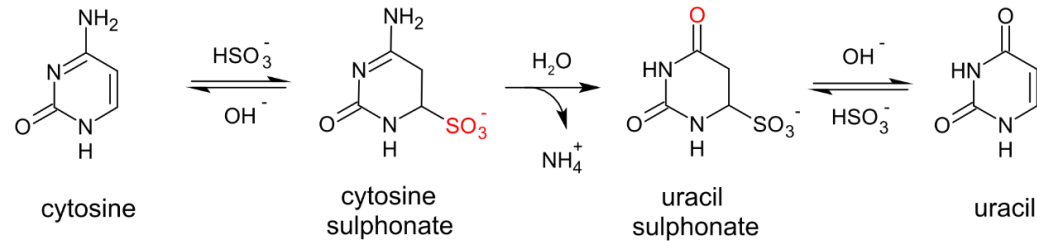
Assays for measuring DNA methylation





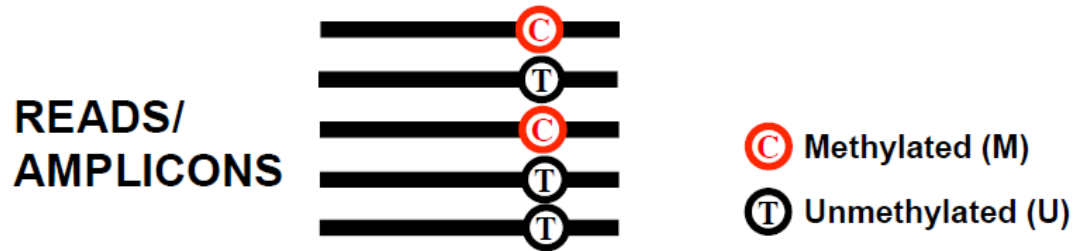
Bisulfite treatment

- Treatment of DNA with **bisulfite** converts cytosine residues to uracil but leaves 5-methylcytosine residues unaffected.





Measuring DNA methylation



$$\% \text{Methylation} = \frac{M}{M + U} \times 100$$



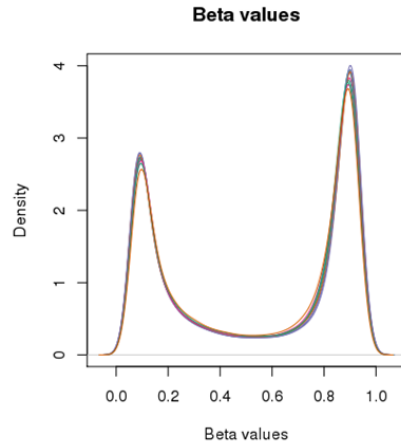
Measuring DNA methylation

- For **microarrays**, there are other measurements:

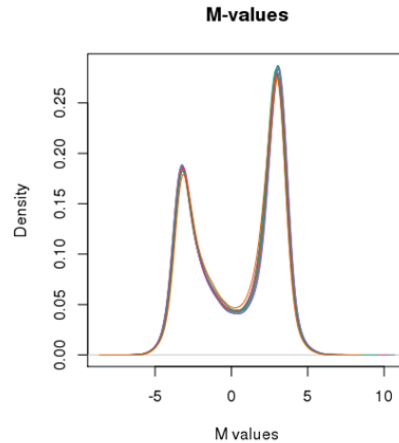
$$\beta = \frac{Meth}{Meth + Unmeth}$$

$$M = \log_2 \frac{Meth}{Unmeth}$$

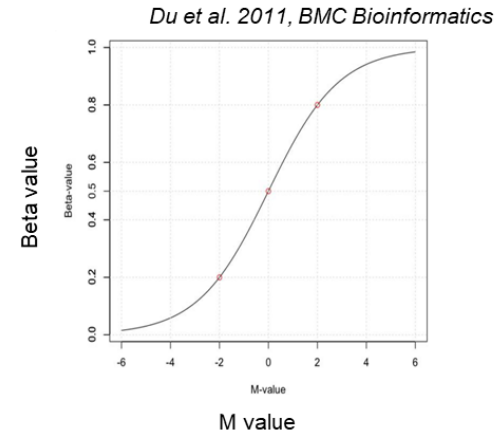
$$M = \log_2 \frac{\beta}{1-\beta}$$



Intuitive, easy to interpret, great for visualisation



Better statistical properties, recommended for statistical testing



Can convert between them via a logit transformation



Illumina Infinium **Human** Methylation BeadChips

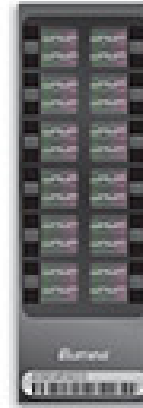
27k array (2009)



1 chip = 12 samples

>27,000 unique CpG sites
measured in each sample

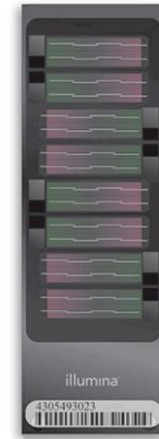
450k array (2011)



1 chip = 12 samples

>450,000 unique CpG sites
measured in each sample

EPIC array (Today)



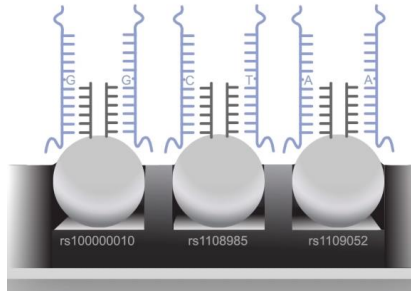
1 chip = 8 samples

>850,000 unique CpG sites
measured in each sample

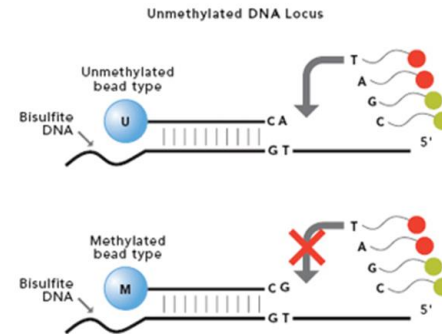
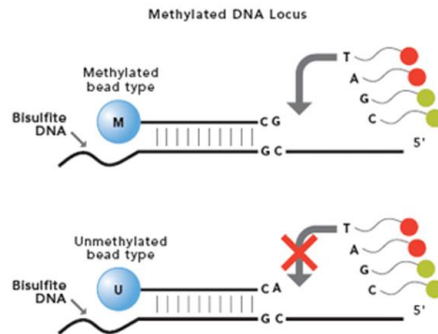
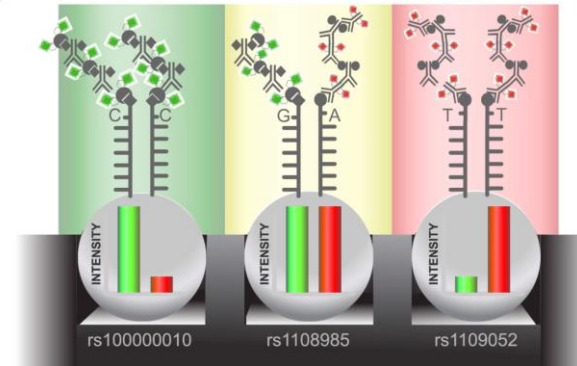


Assays for measuring DNA methylation – bisulfite microarrays

Each probe binds to a complementary sequence.



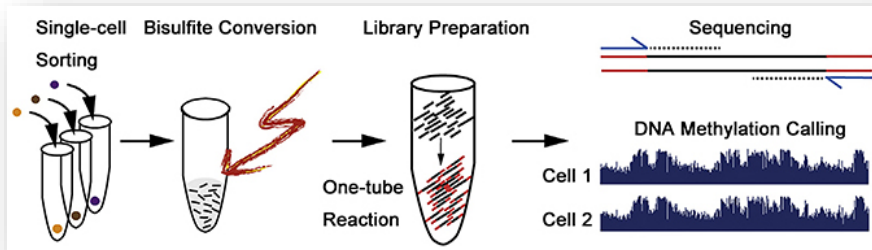
Illumina Infinium Methylation Assay





Whole-genome bisulfite sequencing (WGBS)

- WGBS experiment steps:
 - Sonication
 - End-repair
 - A-tailing
 - Adapter ligation
 - Bisulfite conversion
 - Amplification
 - Sequencing
- Advantages:
 - measure single-cytosine methylation levels genome-wide
 - directly estimate the ratio of molecules methylated rather than enrichment levels





3. WGBS data analysis workflow

The real deal



WGBS data analysis steps



1. Quality control

Quality trimming, adapter trimming

2. Bisulfite sequence alignment

Two strategies: wild-card alignment, three letter alignment

3. Quantification of DNA methylation (methylation calling)

Sequence deduplication, **.bed** file format

4. Visualization

5. Differentially methylated regions (DMRs) detection



Why do we need quality control?

NGS generates highly accurate data, but it can have certain types of errors:

- contamination with adapters
- technical duplication in the library
- failure at specific parts of the flowcell
- PCR duplicates

Reads without proper trimming reads may result in

- low mapping efficiency
- mis-alignments
- errors in methylation calls since adapters are methylated
- base-call errors tend toward 50%

Tools to use: FastQC, Trim Galore!, Trimmomatics ...

1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



4. Visualization



5. DMRs detection



Sequencing data files - FASTQ

Pair-end reads

sample_1_R1_001.fastq

sample_1_R2_001.fastq

sample_2_R1_001.fastq

sample_2_R2_001.fastq

FASTQ Format

Identifier	———	@HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Sequence	———	TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTTNNNNNNNNNNNTAGTTTCTTGAGA
+ sign & identifier	———	+HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Quality scores	———	efcfffffcfeefffcfffffdddf`feed]`_]_Ba^__[YBBBBBBBBBBRTT\]] [] dddd`

Base T
phred Quality] = 29

1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



4. Visualization



5. DMRs detection

Q-score: “Phred quality scores”

These scores represent the likelihood of the base being called wrong.

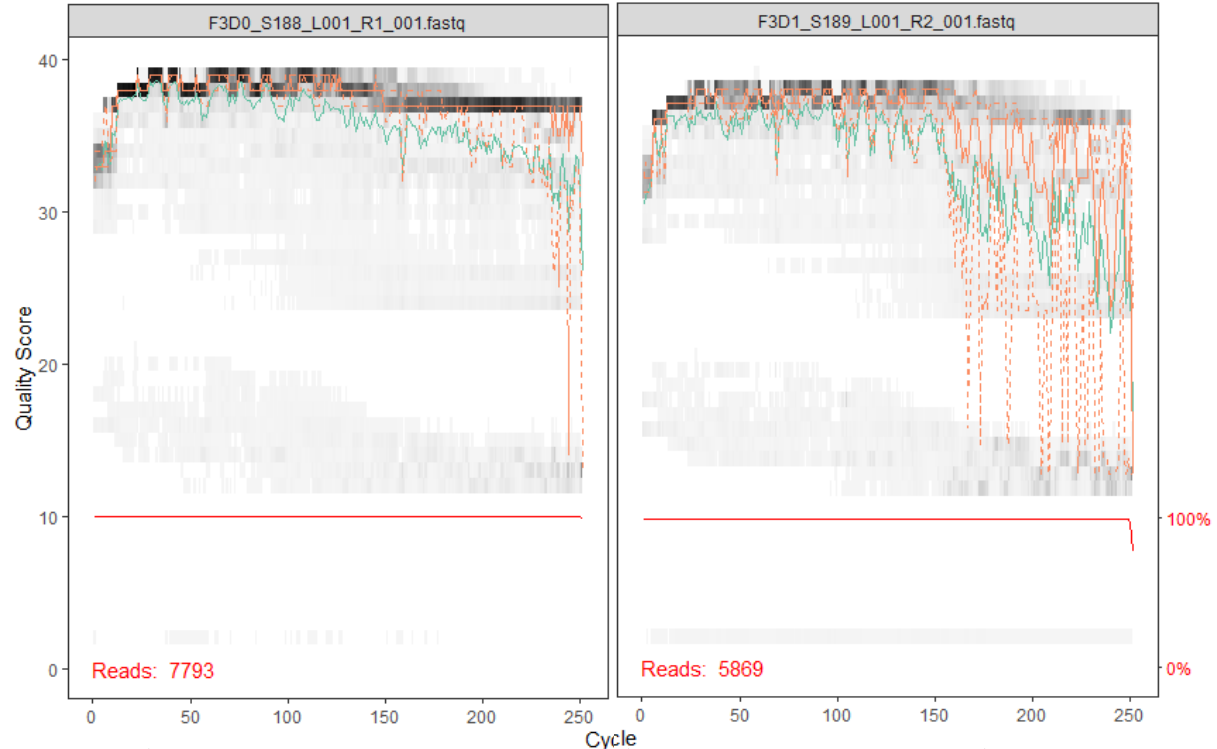
$$Q_{\text{phred}} = -10 \log_{10} e,$$

e is the probability that the base is called wrong.

When $Q = 30$, $e = 0.1\%$

For Illumina, $Q > 30$, the base call quality value is good enough

Quality control – FASTQC



1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



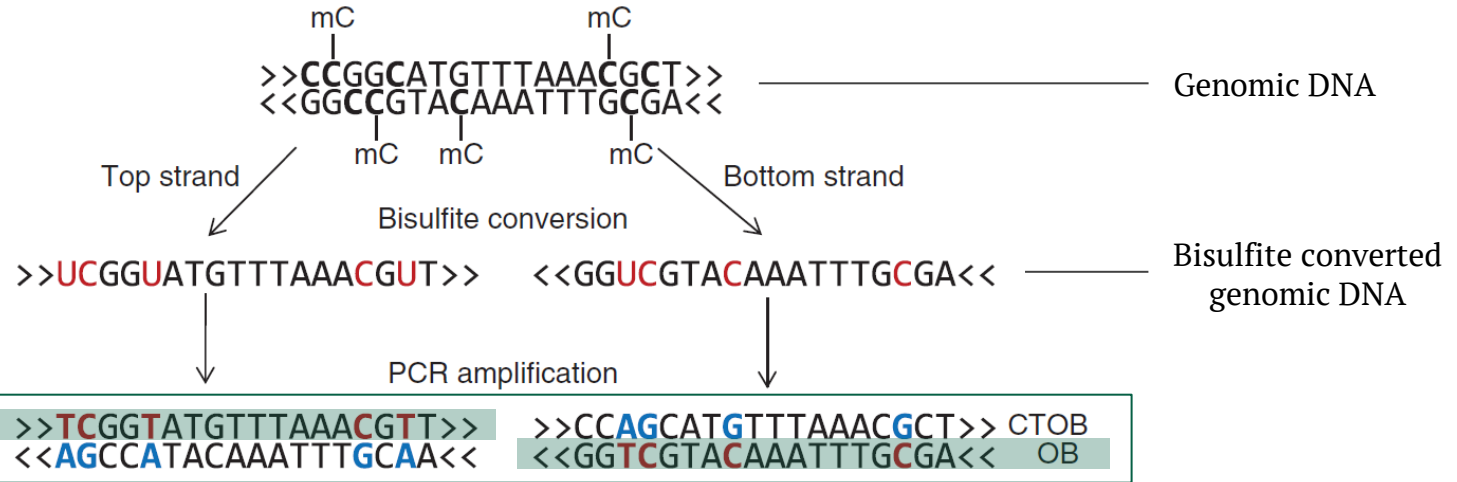
4. Visualization



5. DMRs detection



Effect of bisulfite treatment of DNA



1. Quality control

2. Bisulfite alignment

3. Methylation Quantification

4. Visualization

5. DMRs detection



Bisulfite sequence alignment – two strategies

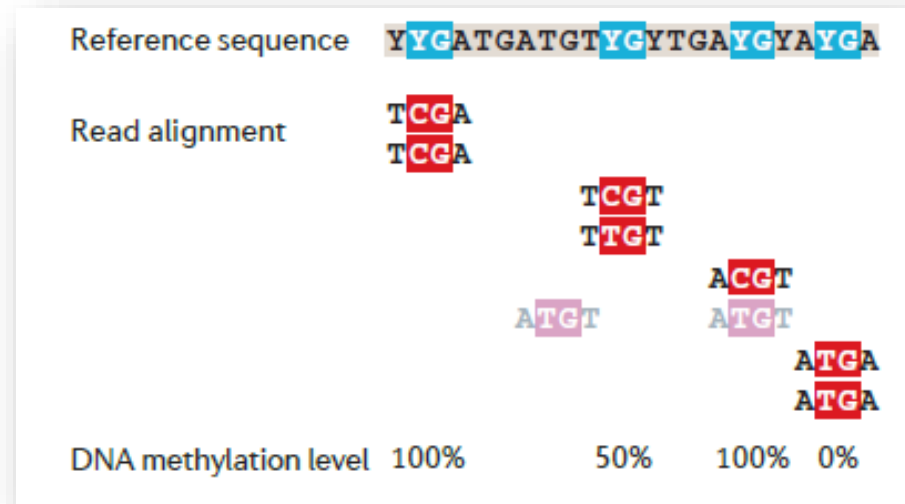
Genomic DNA sequence **C****CG**ATGATGT**CG**CTGA**CG**CA**CG**A
 DNA methylation level 100% 50% 50% 0%

Bisulphite-sequencing reads **A****CGT**, **A****TGA**, **A****TGA**, **A****TGT**,
T**CGA**, **T****CGA**, **T****CGT**, **T****TGT**

1. Wild-card aligner

Replace Cs in the genomic DNA sequence (reference) by the wild-card letter Y, which matches both Cs and Ts in the read sequence

The **DNA methylation level**: the percentage of aligning Cs among all uniquely mapped reads



1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



4. Visualization



5. DMRs detection



Bisulfite sequence alignment – two strategies

Genomic DNA sequence **C****CG**ATGATGT**CG**CTGA**CG**CA**CG**A
 DNA methylation level 100% 50% 50% 0%

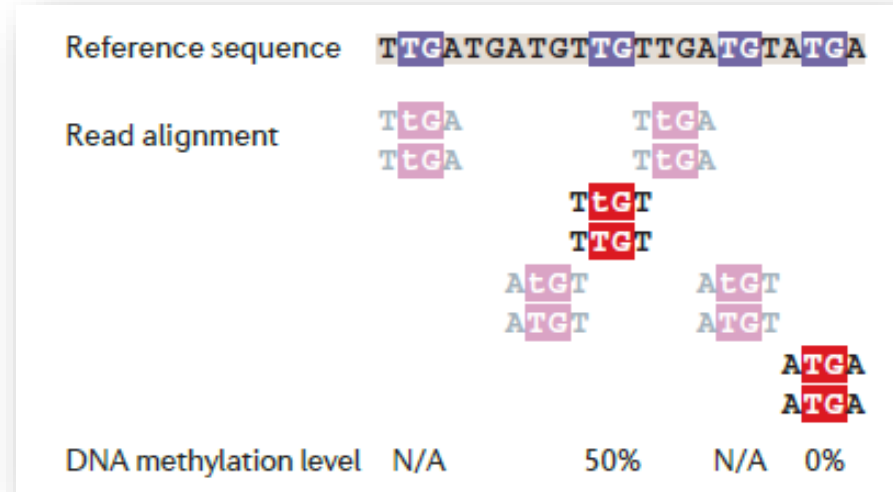
Bisulphite-sequencing reads **A****CGT**, **A****TGA**, **A****TGA**, **A****TGT**,
T**CGA**, **T****CGA**, **T****CGT**, **T****TGT**

2. Three-letter aligner

Converting all Cs into Ts in the reads (in lower case t) and genomic DNA sequences

Carry out the alignment on a three-letter alphabet (A, G and T) using standard aligner, such as Bowtie.

The **DNA methylation level**: the percentage of aligning Cs among all uniquely mapped reads





Bismark primary alignment output (BAM file)

Read 1

chromosome

position

```
HISEQ2000-06:366:C3G4NACXX:3:1101:1316:2067_1:N:0: 99 16 71322125 255 100M =
71322232 207
NTTATTTAGTTTTTTAGGGTTTGTGTGTAGGAGTGTGGGAATTATGTTTTTATGGTTGATATTTATTTAAAGTGAGTATAAATTATATATATTTTTTT
#1=DDDDDAFFHIIA:<FGHCCFEGHD?CFFBBBGEHHGHII<FEHIIIII==DE??EHHFEEEEEEEC>;>66;@CDEEDCEEEEEEDDCBB
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
```

Read 2

methylation call

5'-TTGGCATGTTTAAACGTT-3'

5'...ccggcatgtttaaacgct...3'

z

Z

x

X

h

H

bisulfite read

genomic sequence

methylation call

z unmethylated C in CpG context

Z methylated C in CpG context

x unmethylated C in CHG context

X methylated C in CHG context

h unmethylated C in CHH context

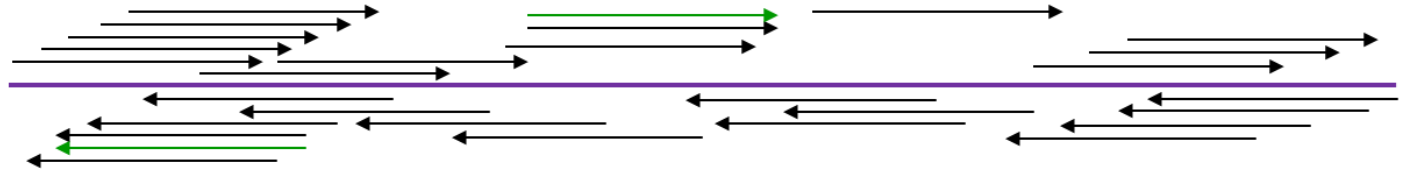
H methylated C in CHH context



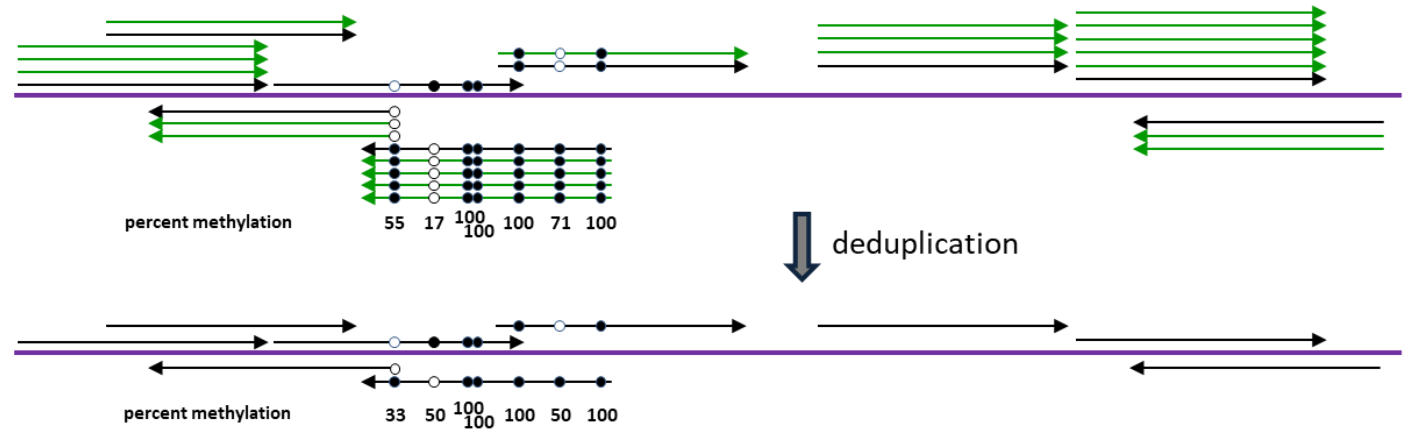


Sequence duplication

Complex/diverse library:



Duplicated library:



1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



4. Visualization



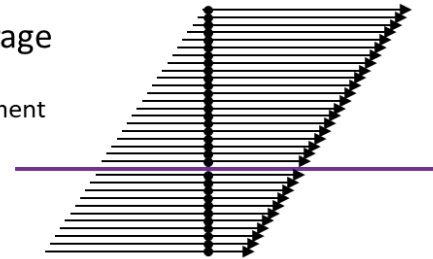
5. DMRs detection



Sequence Deduplication - considerations

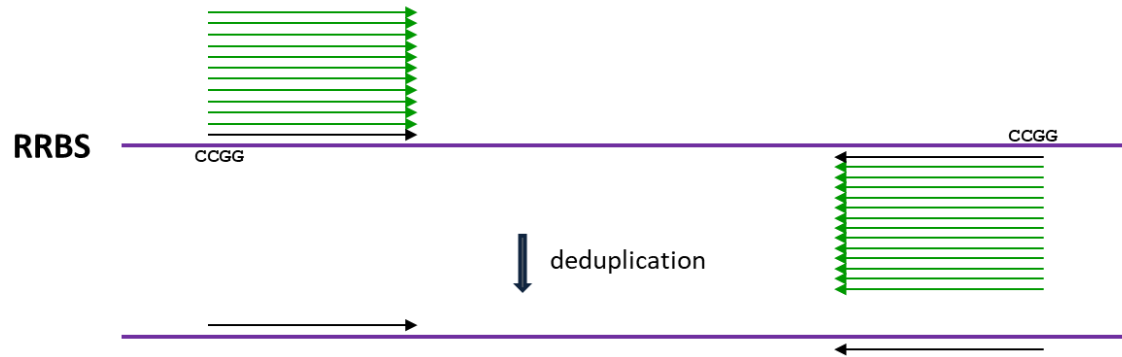
Advisable for large genomes and moderate coverage

- unlikely to sequence several genuine copies of the same fragment amongst >5bn possible fragments with different start sites
- maximum coverage with duplication may still be (read length)-fold (even more with paired-end reads)



NOT advisable for RRBS or other target enrichment methods

where higher coverage is either desired or expected



1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



4. Visualization



5. DMRs detection



Methylation extraction – output `.bed/bedGraph` file

`.bed/bedgraph` file format (UCSC)

Following the track definition line are the track data in four column BED format:

```
chromA  chromStartA  chromEndA  dataValueA
chromB  chromStartB  chromEndB  dataValueB
```

Example of methylation extraction output (bedGraph/coverage 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
```

Chromosome start end Methylation
percentage meth unmeth

1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



4. Visualization



5. DMRs detection



Decide early on which data to use

Methylation contexts

- CpG: Only generally relevant context for mammals
- CHG: Only known to be relevant in plants
- CHH: Generally unmethylated

Methylation strands

- CpG methylation is generally symmetric
- Normally makes sense to merge OT / OB strands

1. Quality control



2. Bisulfite alignment



3. Methylation Quantification




4. Visualization



5. DMRs detection



Always start by looking at your data ...



Integrative Genomics Viewer

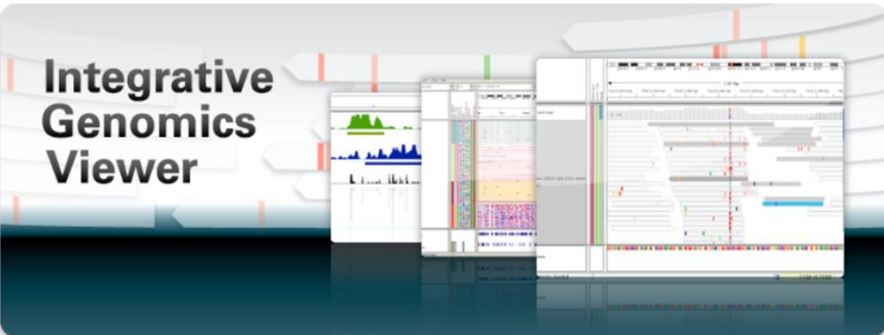
- Home
- Downloads
- Documents
 - IGV User Guide
 - Tutorial Videos
 - File Formats
 - Hosted Genomes
 - FAQ
 - Release Notes
 - Credits
- Contact

Search website

search

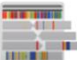
© 2013-2018
Broad Institute
and the Regents of the
University of California

Home



Integrative Genomics Viewer

Overview



The **Integrative Genomics Viewer (IGV)** is a high-performance, easy-to-use, interactive tool for the visual exploration of genomic data. It supports flexible integration of all the common types of genomic data and metadata, investigator-generated or publicly available, loaded from local or cloud sources.

IGV is available in multiple forms, including:

- the original **IGV** - a Java desktop application,
- IGV-Web** - a web application,
- igv.js** - a JavaScript component that can be embedded in web pages (*for developers*)

Citing IGV

To cite your use of IGV in your publication, please reference one or more of:

James T. Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, Jill P. Mesirov. [Integrative Genomics Viewer](#). *Nature Biotechnology* 29, 24–26 (2011). (Free PMC article [here](#)).

Helga Thorvaldsdóttir, James T. Robinson, Jill P. Mesirov. [Integrative Genomics Viewer \(IGV\): high-performance genomics data visualization and exploration](#). *Briefings in Bioinformatics* 14, 178-192 (2013).

1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



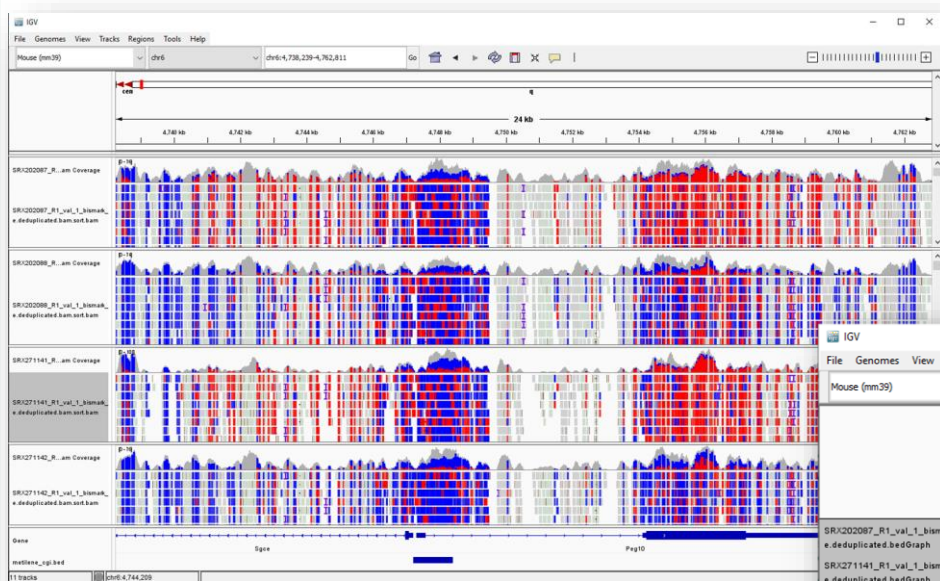
4. Visualization



5. DMRs detection

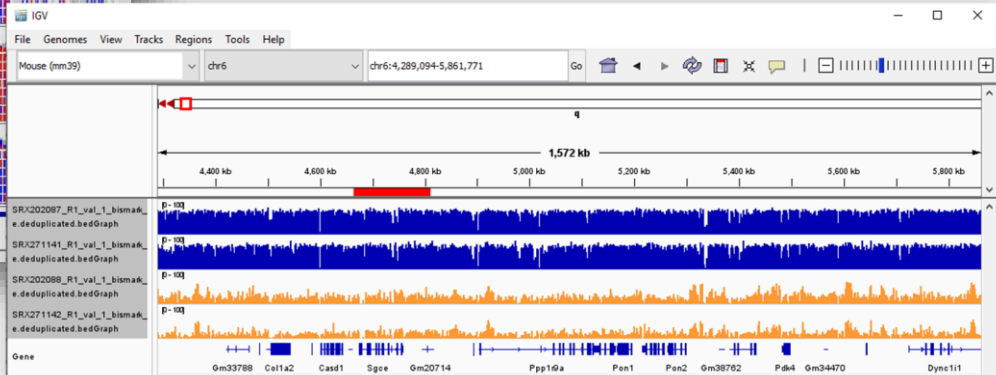


Visualization of DNA Methylation: IGV



View the alignment (bisulfite colored mode)

View the methylation profiles



1. Quality control

2. Bisulfite alignment

3. Methylation Quantification

4. Visualization

5. DMRs detection



Differentially Methylated Regions (DMRs)

- **DMRs:** regions that exhibit consistently different DNA methylation levels between sample groups (e.g., cases vs. control).
- DMRs can be a single C (differentially methylated cytosine, DMC) or as large as an entire gene locus.
- Size: a few hundred to a few thousand bases

1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



4. Visualization



5. DMRs detection



Differentially Methylated Regions (DMRs)

DMR detection:

- Basic: t-test, Wilcoxon rank-sum test
- Advanced: mixture models, Shannon entropy, feature selection, logistic M values

Test at large number of genomic loci

- Correction for multiple hypothesis testing: false discovery rate (FDR: q-value)
- Only strongest single-CpG difference tend to remain significant

To improve the statistical power

- Larger genomic regions
- Pre-selected set of candidate genomic regions

1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



4. Visualization



5. DMRs detection



Resolution levels

Single CpG

a

Genomic DNA sequence		CG	...	CG	CG	...	CG	CG	...	CG	...	CG	...	CG	...	CG	CG
Cases	Sample 1	3%		6%				80%		57%				1%		0%		1%		1%		42%			78%
	Sample 2	2%		0%				50%		74%				0%		1%		0%		0%		38%			85%
	Sample 3	0%		1%				95%		86%				2%		0%		0%		0%		41%			67%
Controls	Sample 4	0%		2%				8%		1%				12%		3%		15%		8%		36%			72%
	Sample 5	1%		4%				5%		2%				15%		5%		33%		11%		39%			94%
	Sample 6	0%		2%				13%		1%				19%		2%		24%		22%		33%			92%

Tilling region

Larger

b

Single-CpG analysis	CG1	CG2		CG3	CG4		CG5	CG6	CG7	CG8	CG9		CG10
Higher in cases (q value)	0.333	0.993		0.085	0.068		0.993	0.993	0.993	0.993	0.196		0.993
Higher in controls (q value)	0.993	0.732		0.993	0.993		0.070	0.104	0.104	0.110	0.993		0.351

c

Genome-wide tiling analysis	Tiling region 1		Tiling region 3		Tiling region 5		Tiling region 7		Tiling region 8
		Tiling region 2		Tiling region 4		Tiling region 6			
Higher in cases (q value)	0.549		0.048*		0.988	0.988	0.549		0.988
Higher in controls (q value)	0.768		0.993		0.067	0.067	0.299		0.299

Annotated region

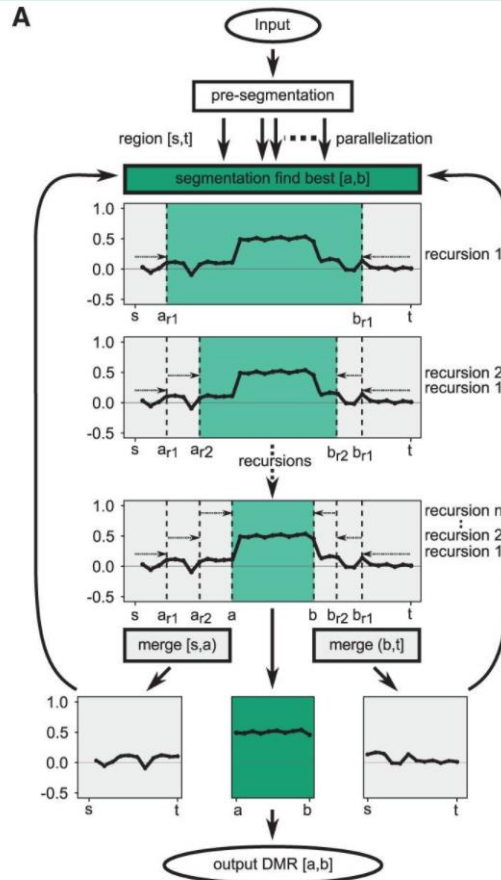
Largest

d

Annotated genome analysis	Enhancer		Promoter region		First exon
Higher in cases (q value)	0.024*		0.986		0.353
Higher in controls (q value)	0.993		0.045*		0.299



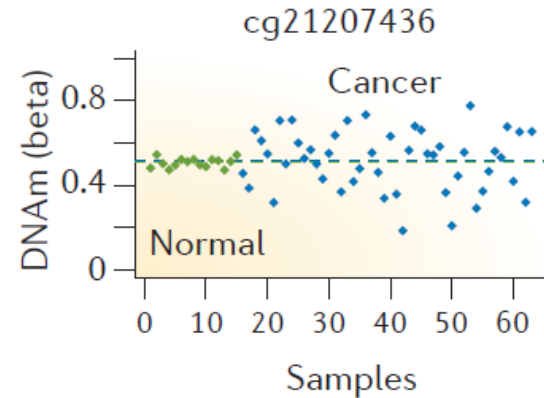
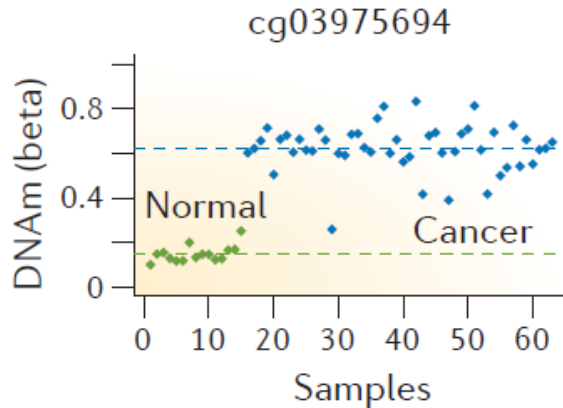
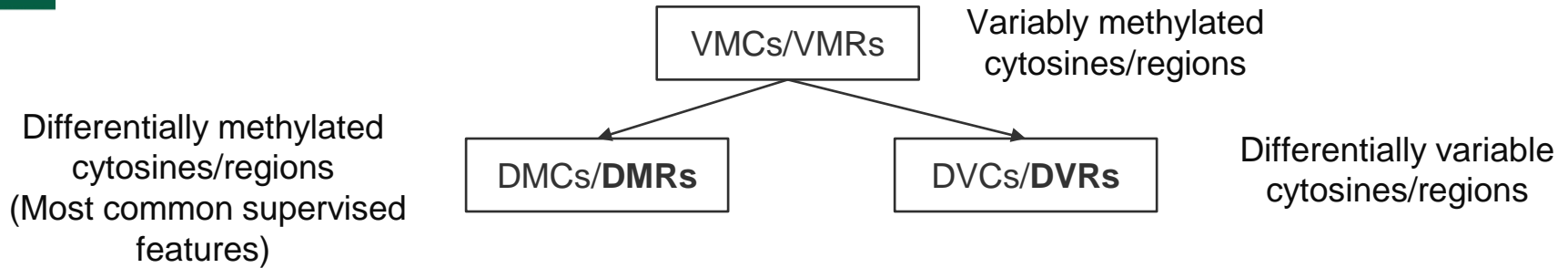
DMR detection - metilene



- The objective problem of finding DMRs has two dimensions:
 - find a genomic region
 - the individuals of two groups are significantly distinct in their methylation levels.
- De-novo DMR detection



Differentially variability



1. Quality control



2. Bisulfite alignment



3. Methylation Quantification



4. Visualization



5. DMRs detection



The end.