# Bioinformatics Workshop

Session #16

Bisulfite Sequencing and Analysis
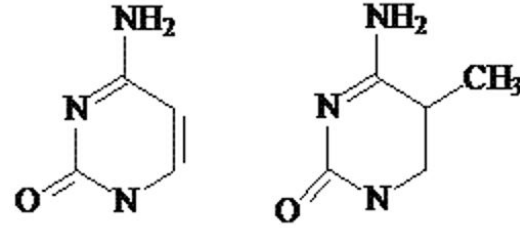
Chris Miller

# Epigenetics



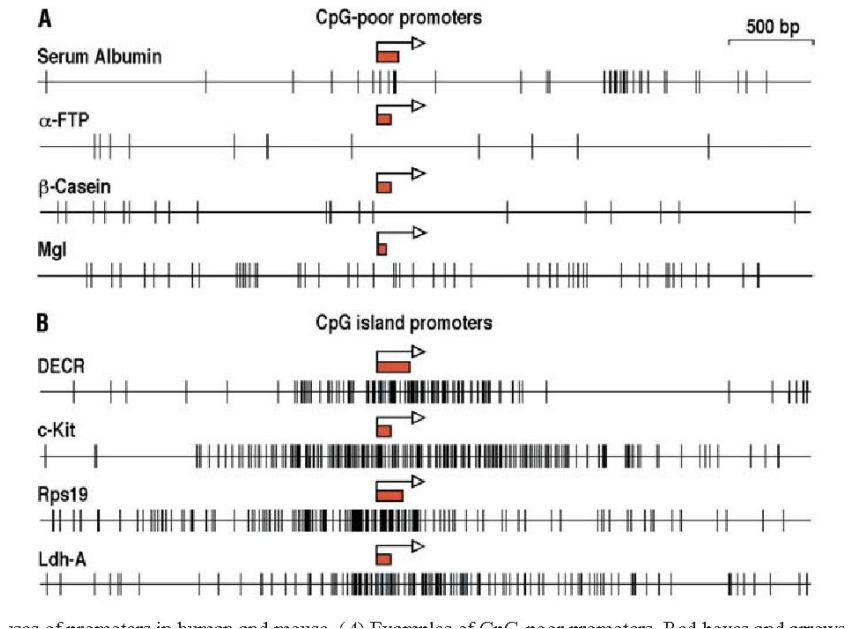The two main components of the epigenetic code

**DNA methylation**
Methyl marks added to certain DNA bases repress gene activity.

**Histone modification**
A combination of different molecules can attach to the 'tails' of proteins called histones. These alter the activity of the DNA wrapped around them.

Histone tails

Histones

Chromosome

# DNA Methylation

- Mostly happens at CpGs

- About 25 million CpGs in human genome

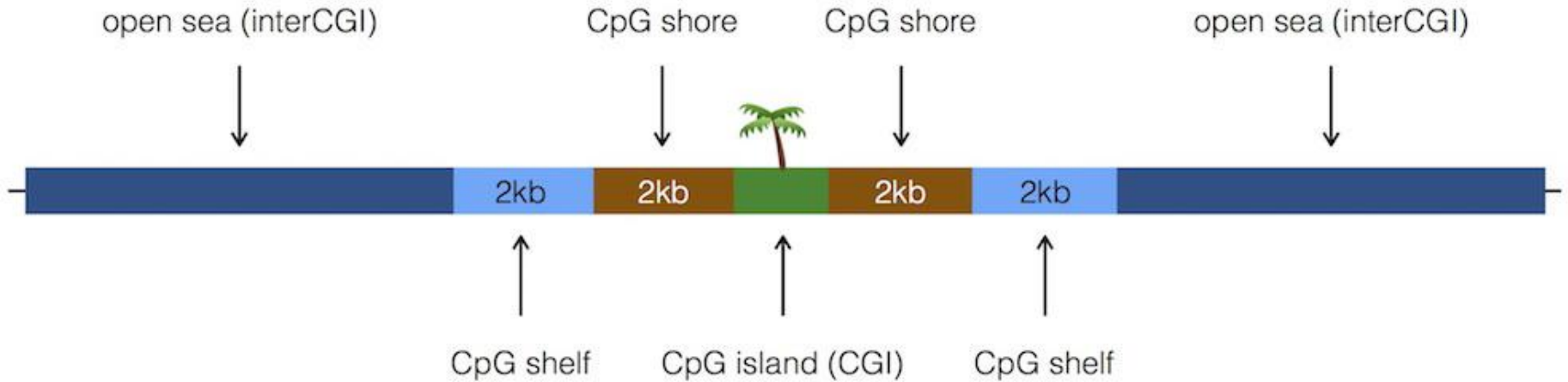● https://en.wikipedia.org/wiki/CpG_site#/media/File:CpG_vs_C-G_bp.svg

# DNA Methylation

- CpG Islands

- Length >= 200 bp
  GC% > 50%
  o/e CpG ratio > 60%

- Selective pressure/
  Evolutionary constraint

A — CpG-poor promoters. Serum Albumin, α-FTP, β-Casein, Mgl. 500 bp.
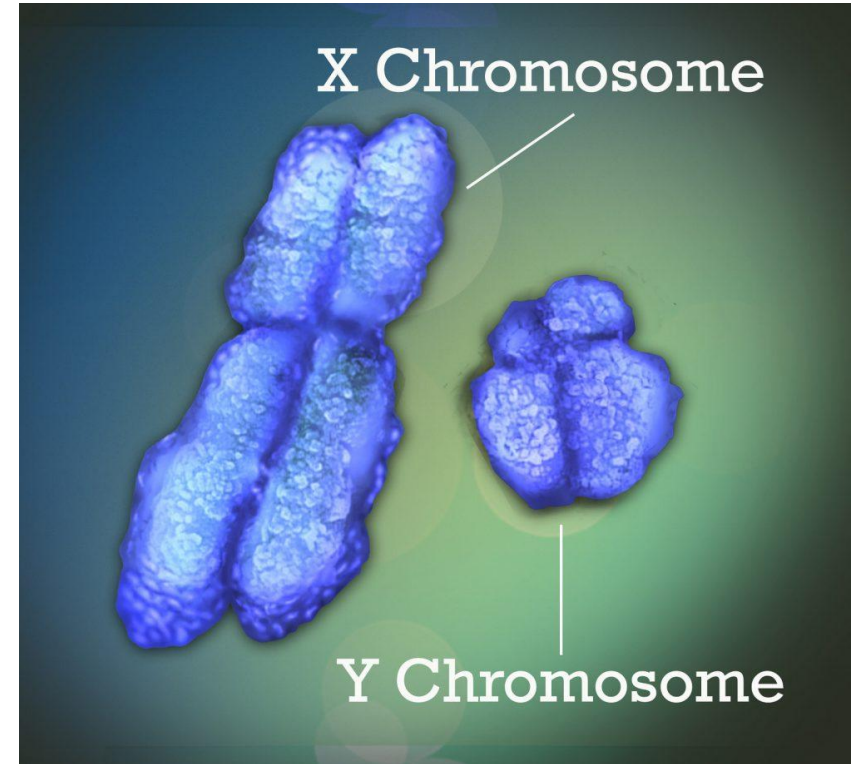B — CpG island promoters. DECR, c-Kit, Rps19, Ldh-A.

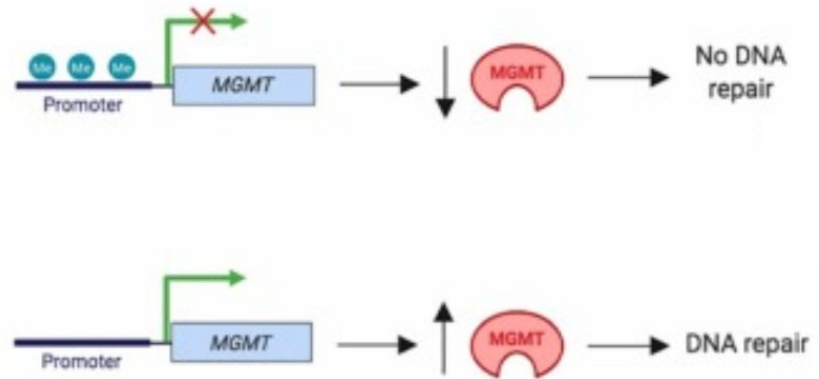# Islands, shores, and shelves

# What does DNA methylation do?

- The short answer: It depends!

- X-chromosome inactivation

- Silencing of transposable elements

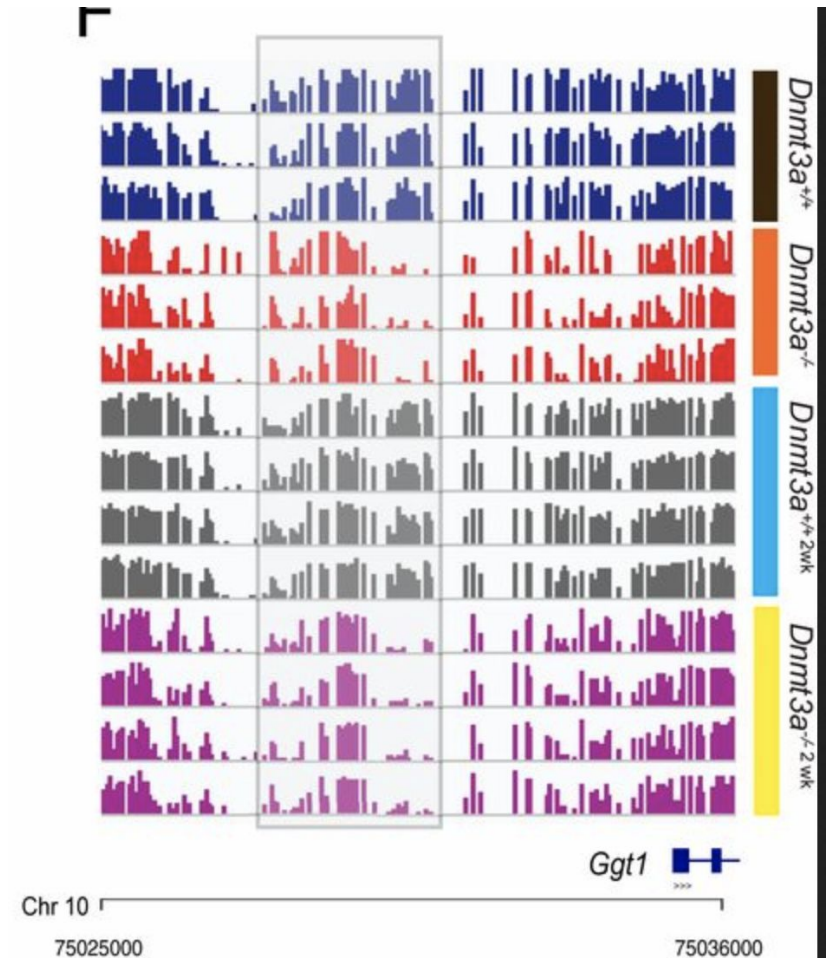- Cellular differentiation

- Cancer - hypo/hypermethylation


X Chromosome

Y Chromosome

# MGMT and Temozolomide

- TMZ is an alkylating agent - damages DNA, causes cell death

- MGMT "cleans up" the damage

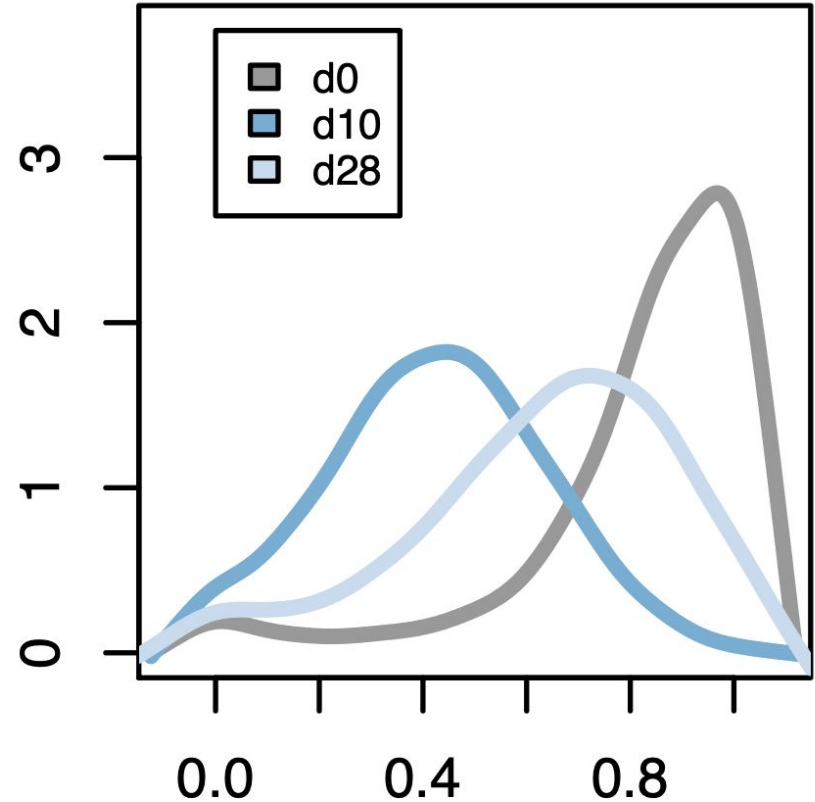- Methylation of the MGMT promoter is linked to better outcomes!

# Methylation Patterns
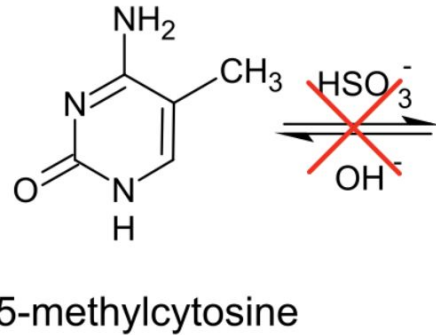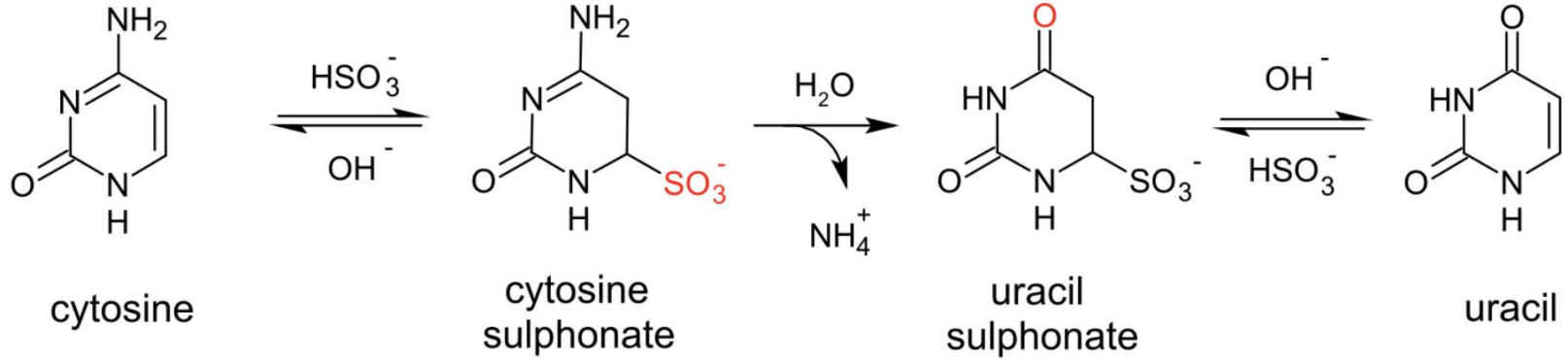
- Methyltransferases that act locally
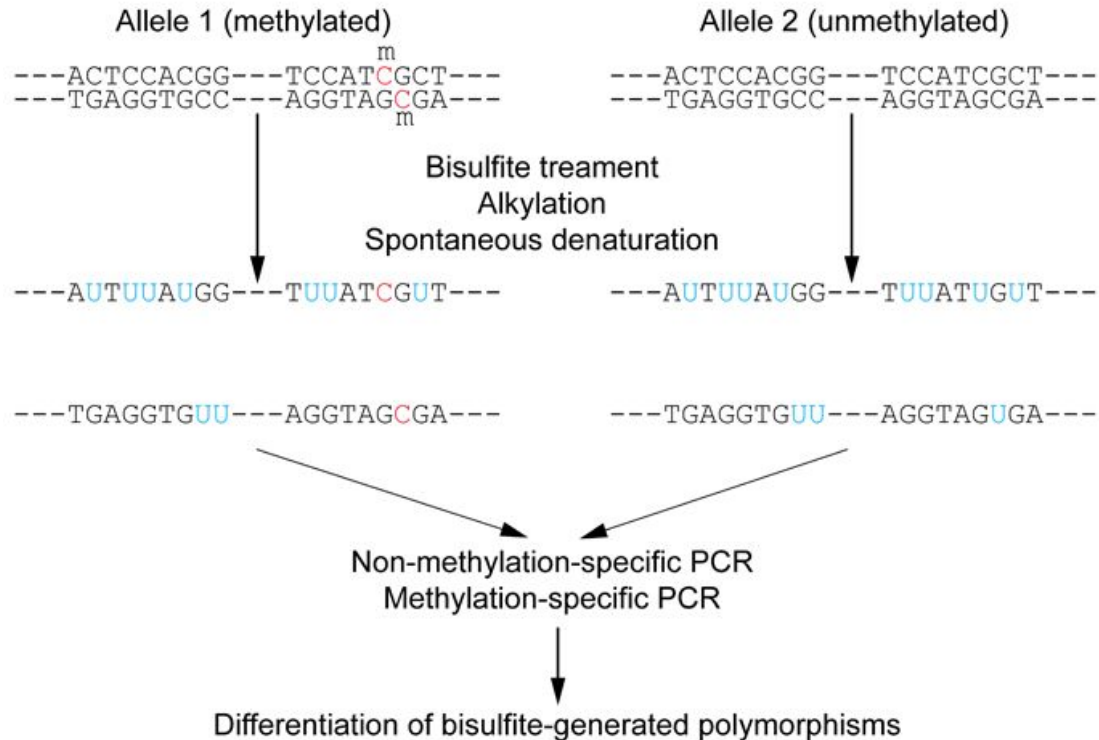
# Methylation Patterns

- Methyltransferases that act locally

- Other alterations (or treatments) that act globally

# Bisulfite sequencing



cytosine  →(HSO₃⁻ / OH⁻)→  cytosine sulphonate  →(H₂O, −NH₄⁺)→  uracil sulphonate  →(OH⁻ / HSO₃⁻)→  uracil

5-methylcytosine  (HSO₃⁻ / OH⁻ reaction blocked)

https://en.wikipedia.org/wiki/File:Bisulfite_conversion.svg

# Bisulfite sequencing

# Bisulfite sequencing

# Whole-genome Bisulfite Sequencing (WGBS)

- Need a special aligner - has to expect many C > T mismatches!

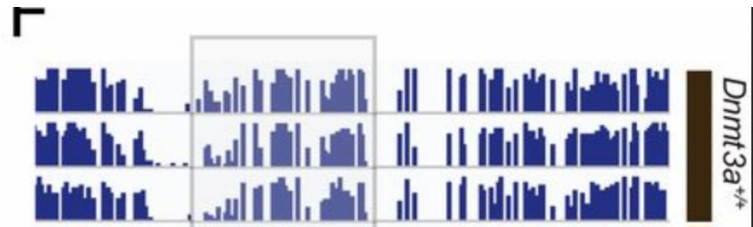- BSMAP
- bismark
- BWA-meth
- biscuit

# Methylation calling

- Determine methylation fraction at each site in the genome

    - Count the Cs and Ts, taking strandedness into account

    - Some tools account for SNPs while doing this
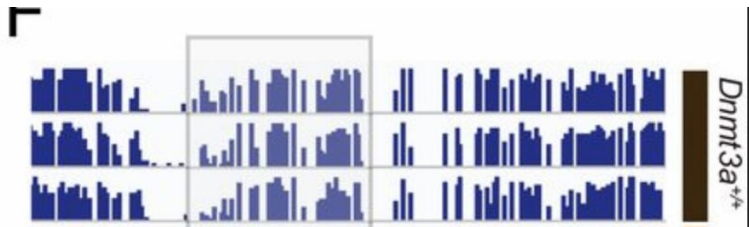
-

# Methylation calling

- Determine methylation fraction at each site in the genome

    - Count the Cs and Ts, taking strandedness into account

    - Some tools account for SNPs while doing this

- Why isn't every position 0%, 50% or 100%?
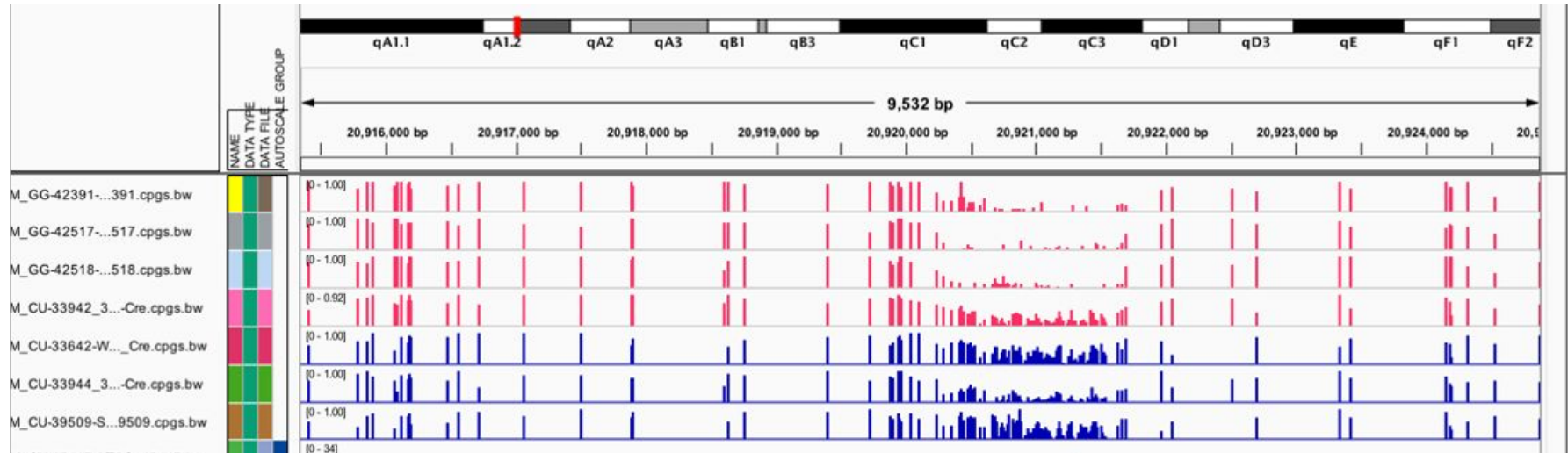


*Dnmt3a+/+*

# Methylation calling

- Determine methylation fraction at each site in the genome

    - Count the Cs and Ts, taking strandedness into account

    - Some tools account for SNPs while doing this

- Why isn't every position 0%, 50% or 100%?
    - we're sequencing a population of cells!

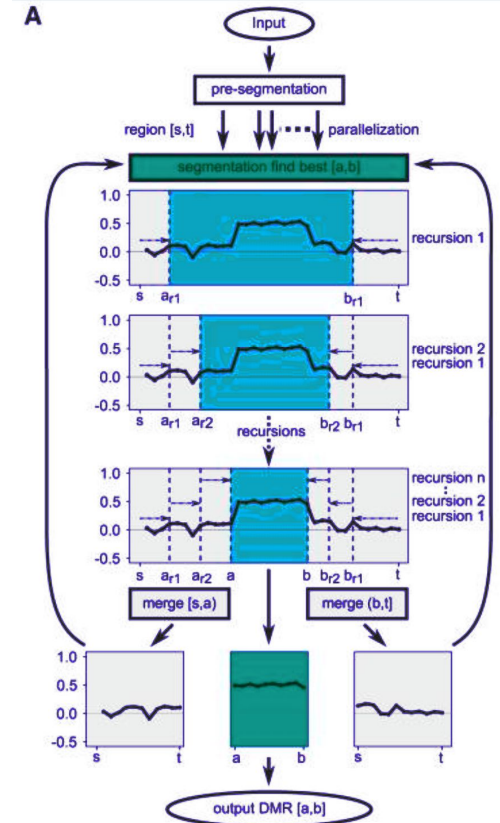# Workflow/File formats

- Aligning:  FASTQ > BAM/CRAM

- Pileup:  BAM/CRAM > VCF
    - (entries for every site, allele frequencies)

- VCF > bedgraph
    - chr, start, stop, beta_value  (methylation fraction)

- bedgraph > bigwig  (for visualization in IGV)

- We have a workflow for this!
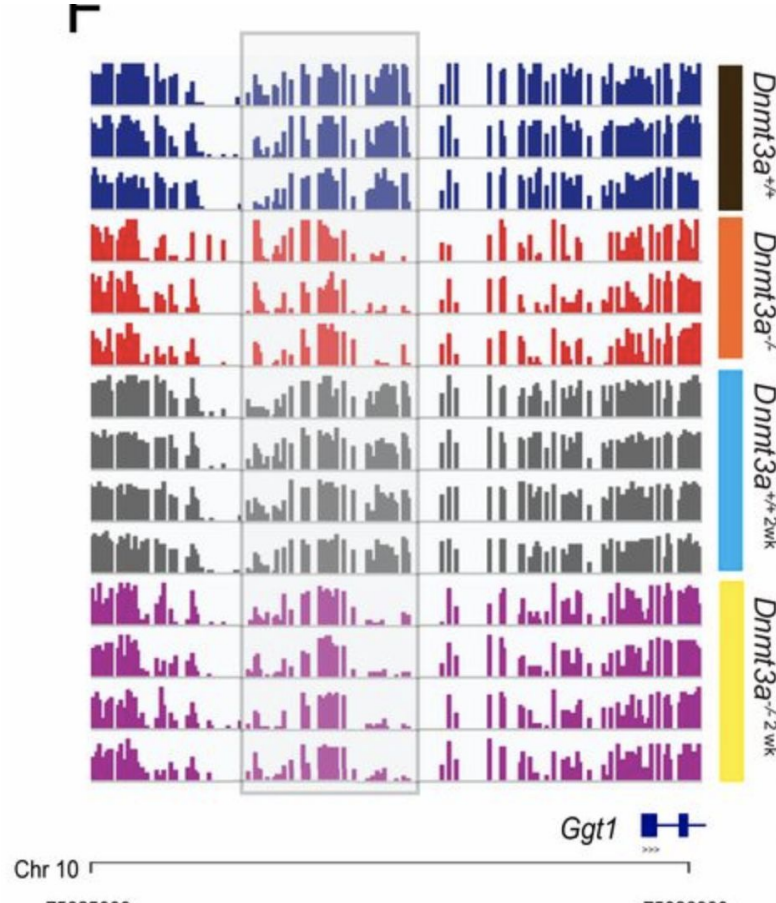
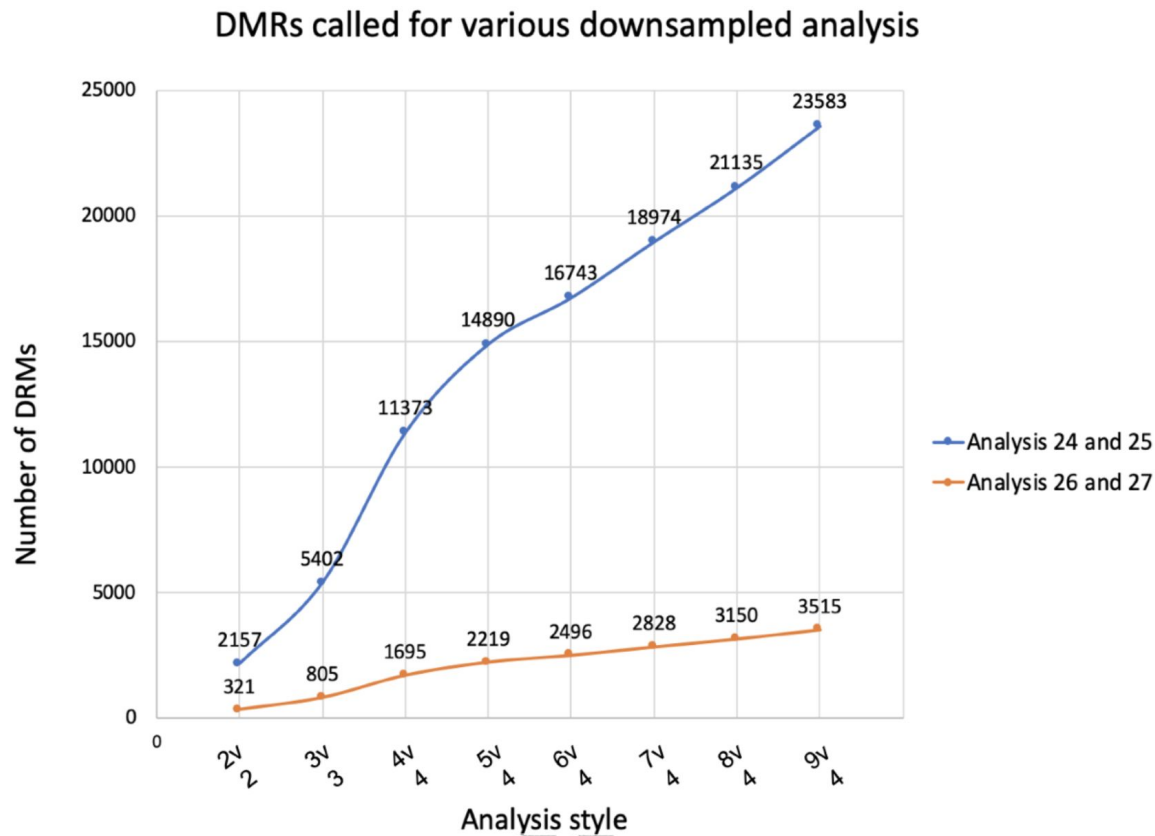# IGV visualization

# Differentially methylated regions



- Comparing two groups to find changes

- Finding DMRs is a segmentation problem
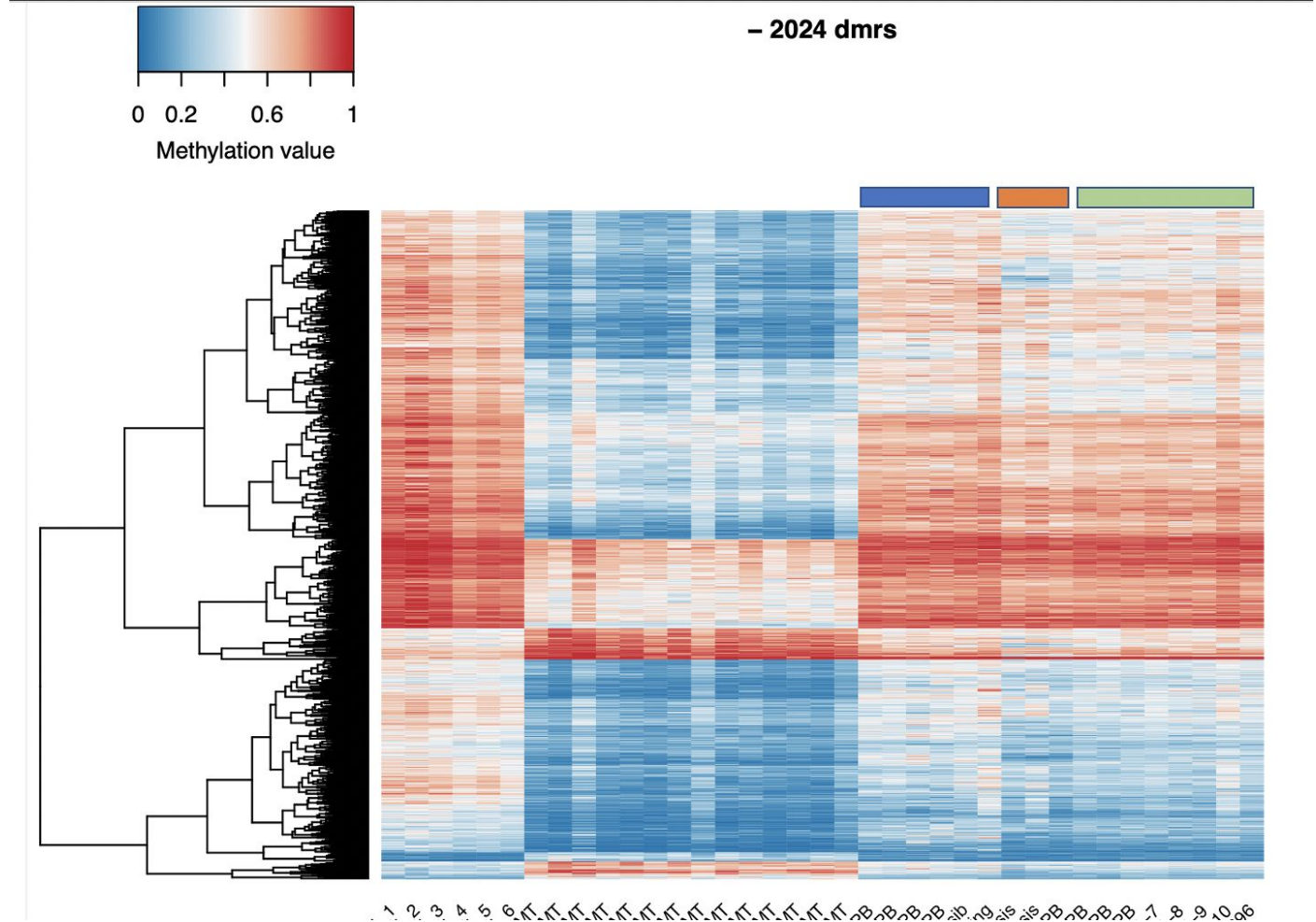
- We use a tool called metilene
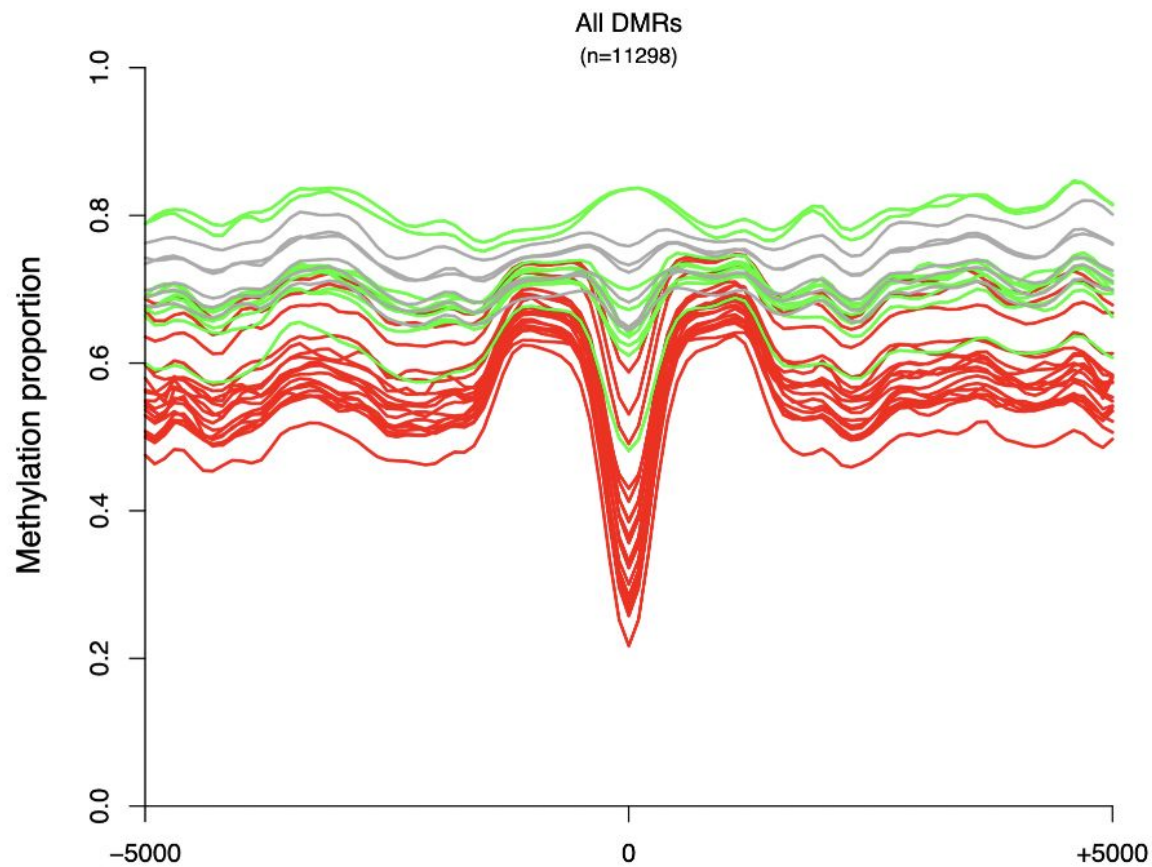
# Differentially methylated regions
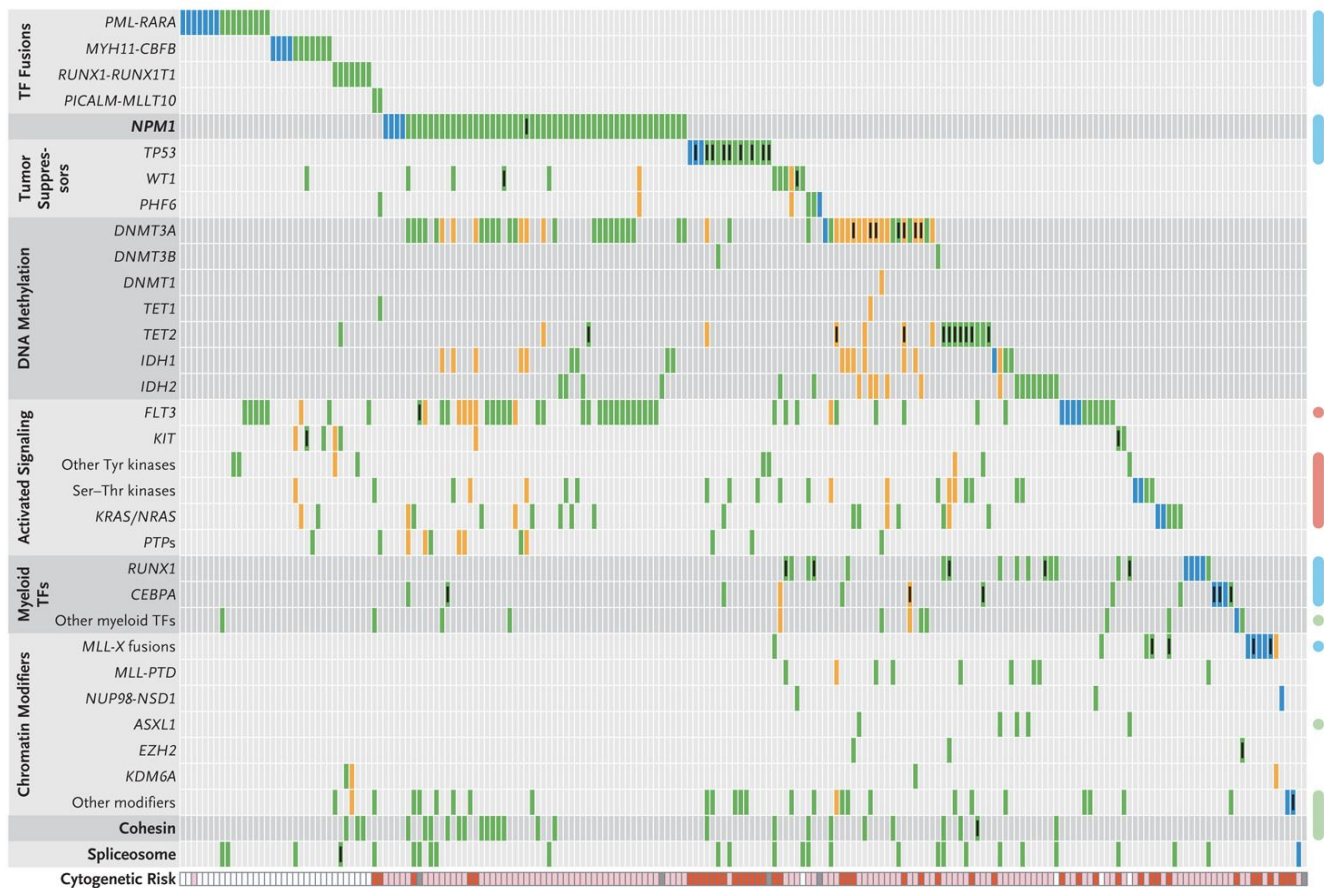
# Number of samples matters!



DMRs called for various downsampled analysis

Number of DRMs

- Analysis 24 and 25
- Analysis 26 and 27

Analysis style

Blue line (Analysis 24 and 25): 2157, 5402, 11373, 14890, 16743, 18974, 21135, 23583

Orange line (Analysis 26 and 27): 321, 805, 1695, 2219, 2496, 2828, 3150, 3515

X-axis: 2v2, 3v3, 4v4, 5v4, 6v4, 7v4, 8v4, 9v4

# Heatmaps

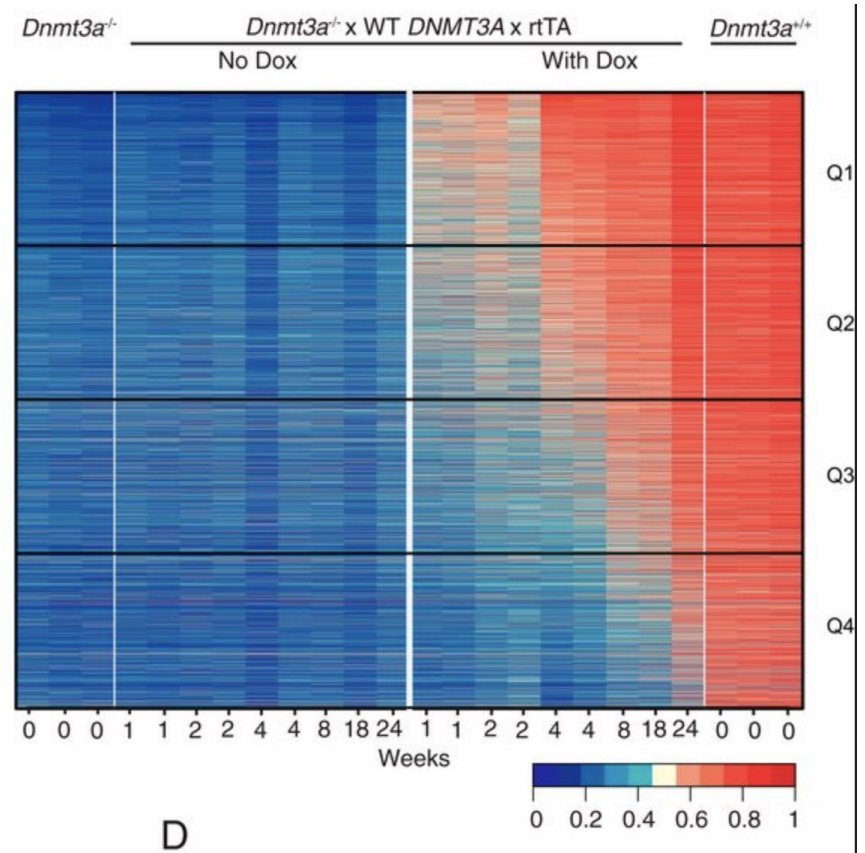# Canyon Plots



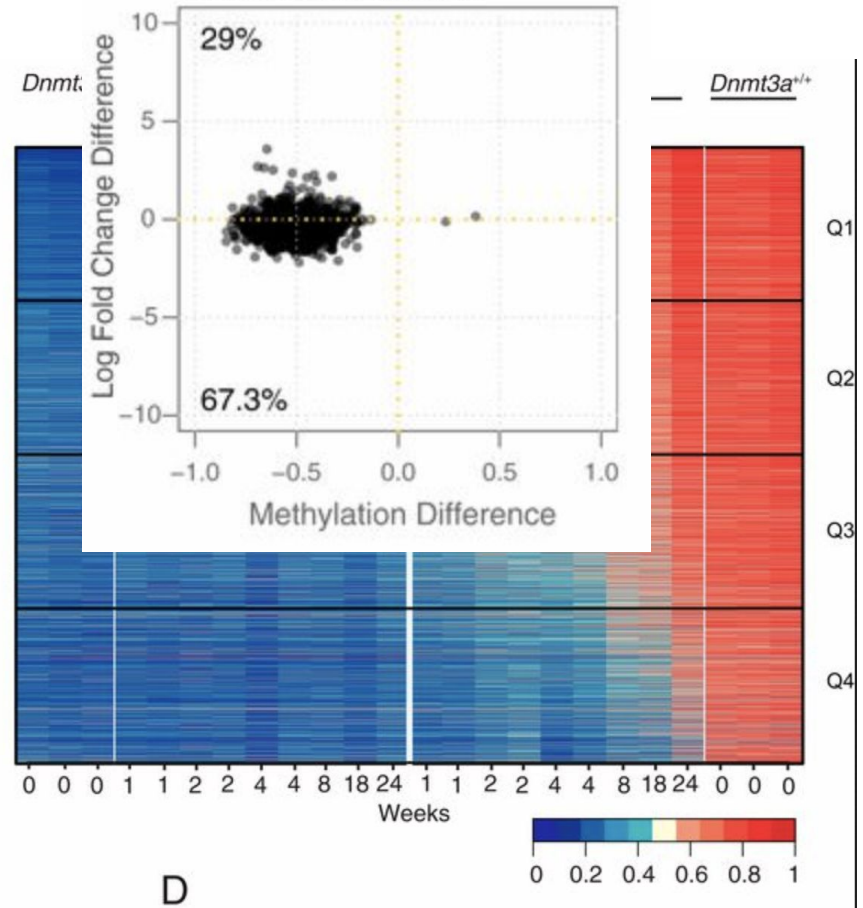All DMRs
(n=11298)

# DNMT3A deficiency

# DNMT3A deficiency

- Mouse models (and human data)

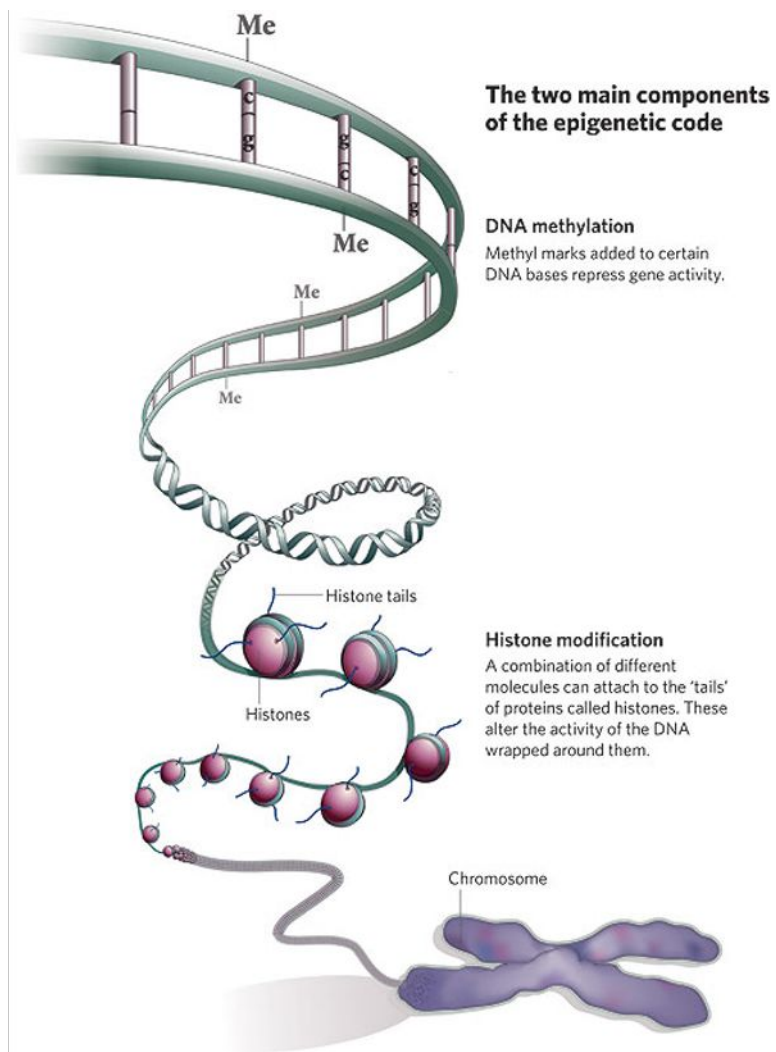- Looking at context, effects, and reversibility

# DNMT3A deficiency

- Mouse models (and human data)

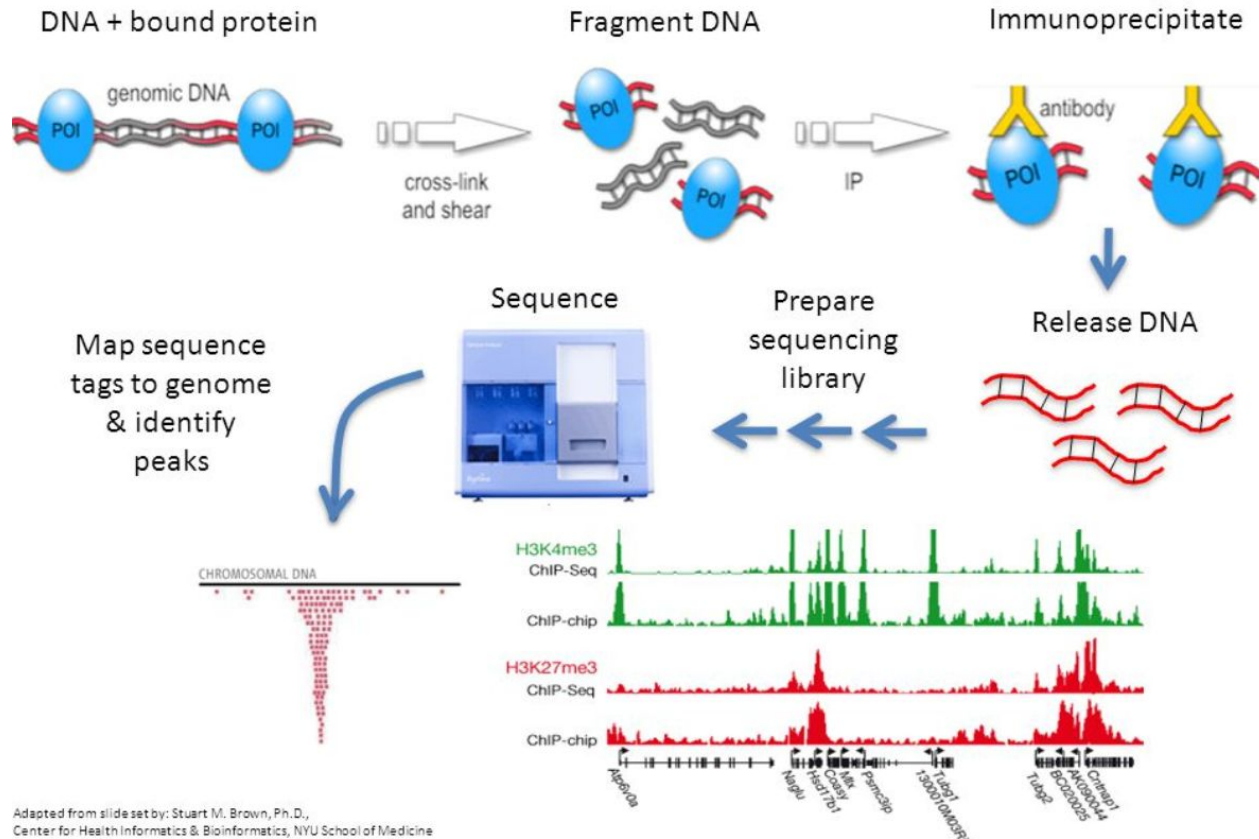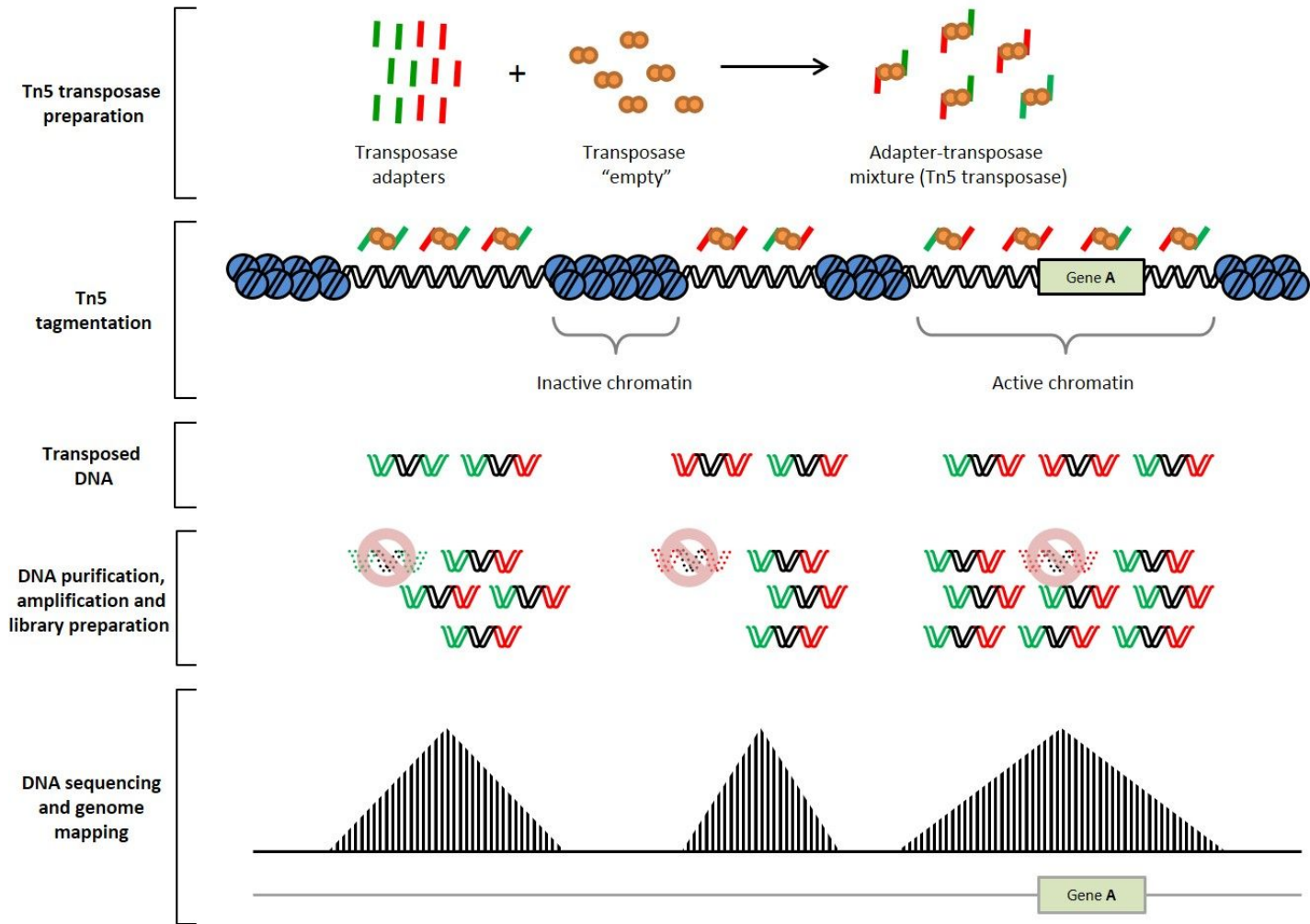- Looking at context, effects, and reversibility

# ChIP-seq/ATAC-seq

- Alterations of DNA state or accessibility

- Wrapped around histones

- Bound by transcription factors

- etc



**The two main components of the epigenetic code**

**DNA methylation**
Methyl marks added to certain DNA bases repress gene activity.

**Histone modification**
A combination of different molecules can attach to the 'tails' of proteins called histones. These alter the activity of the DNA wrapped around them.

# ChIP-seq



Adapted from slide set by: Stuart M. Brown, Ph.D.,
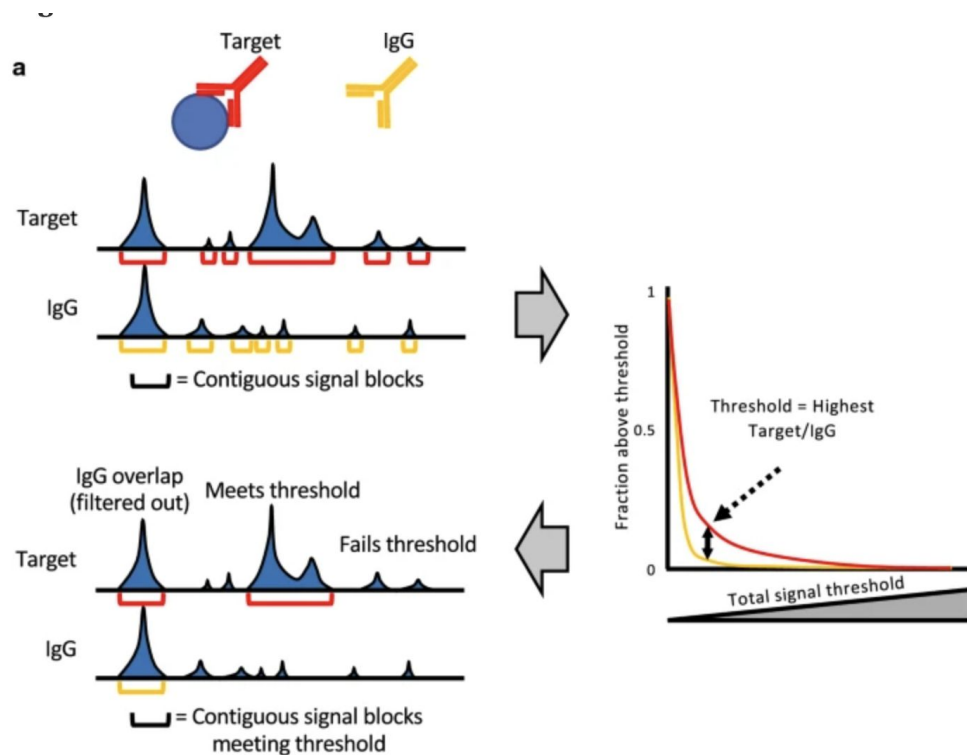Center for Health Informatics & Bioinformatics, NYU School of Medicine
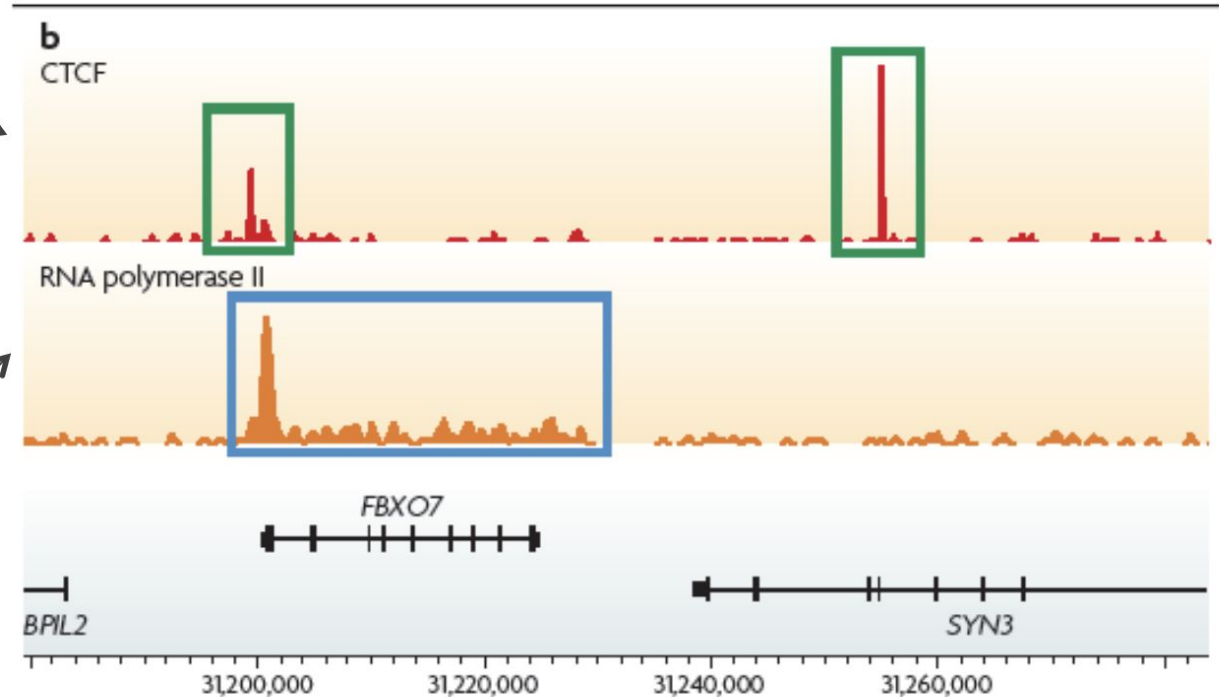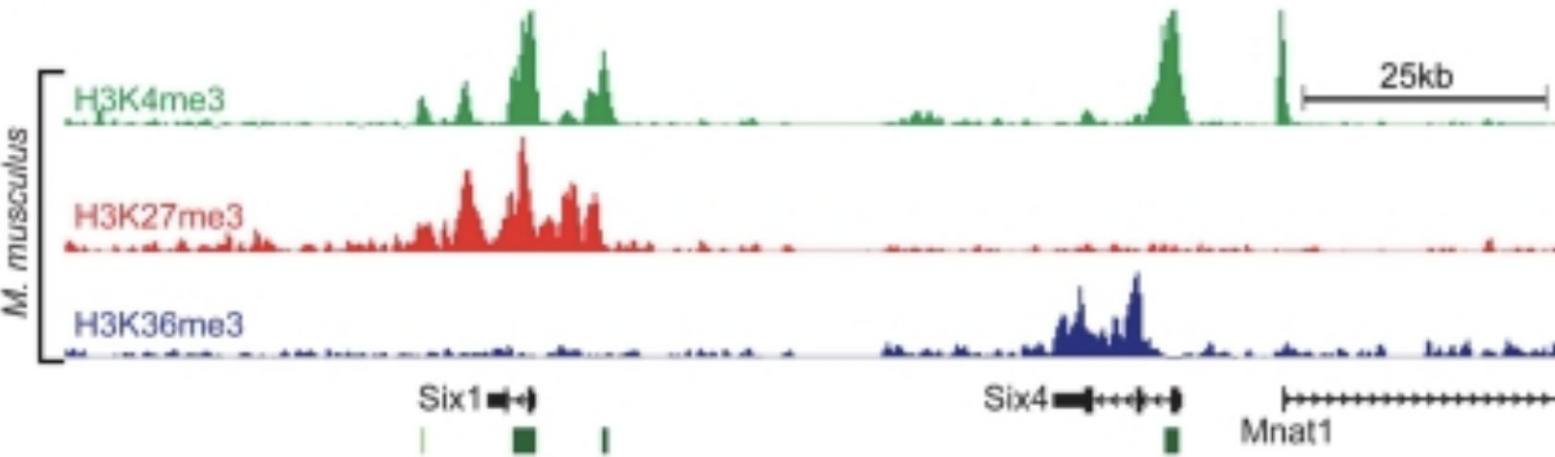
ATAC-seq

# Peak-calling

MACS2, HOMER, SEACR, etc

# Proteins bind in different ways

Transcription factor – tight, highly-peaked binding region

RNA PolII – enriched at TSS but bound throughout gene body

# Interpretation