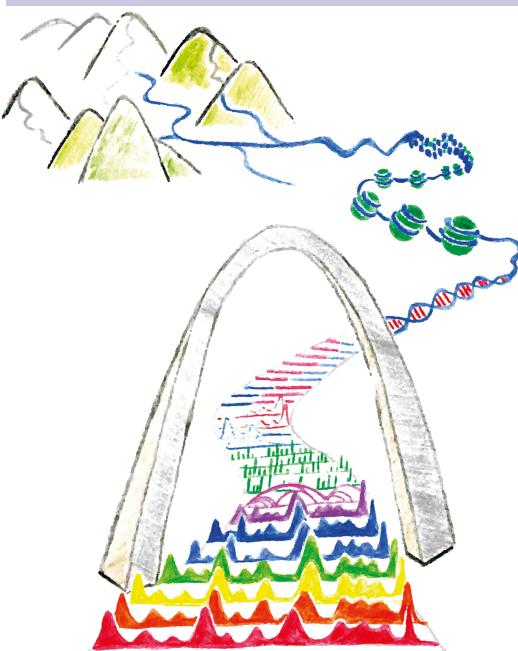


系统生物学与生物信息学  
海外学者短期讲学系列课程

**Current Topics in Epigenomics**

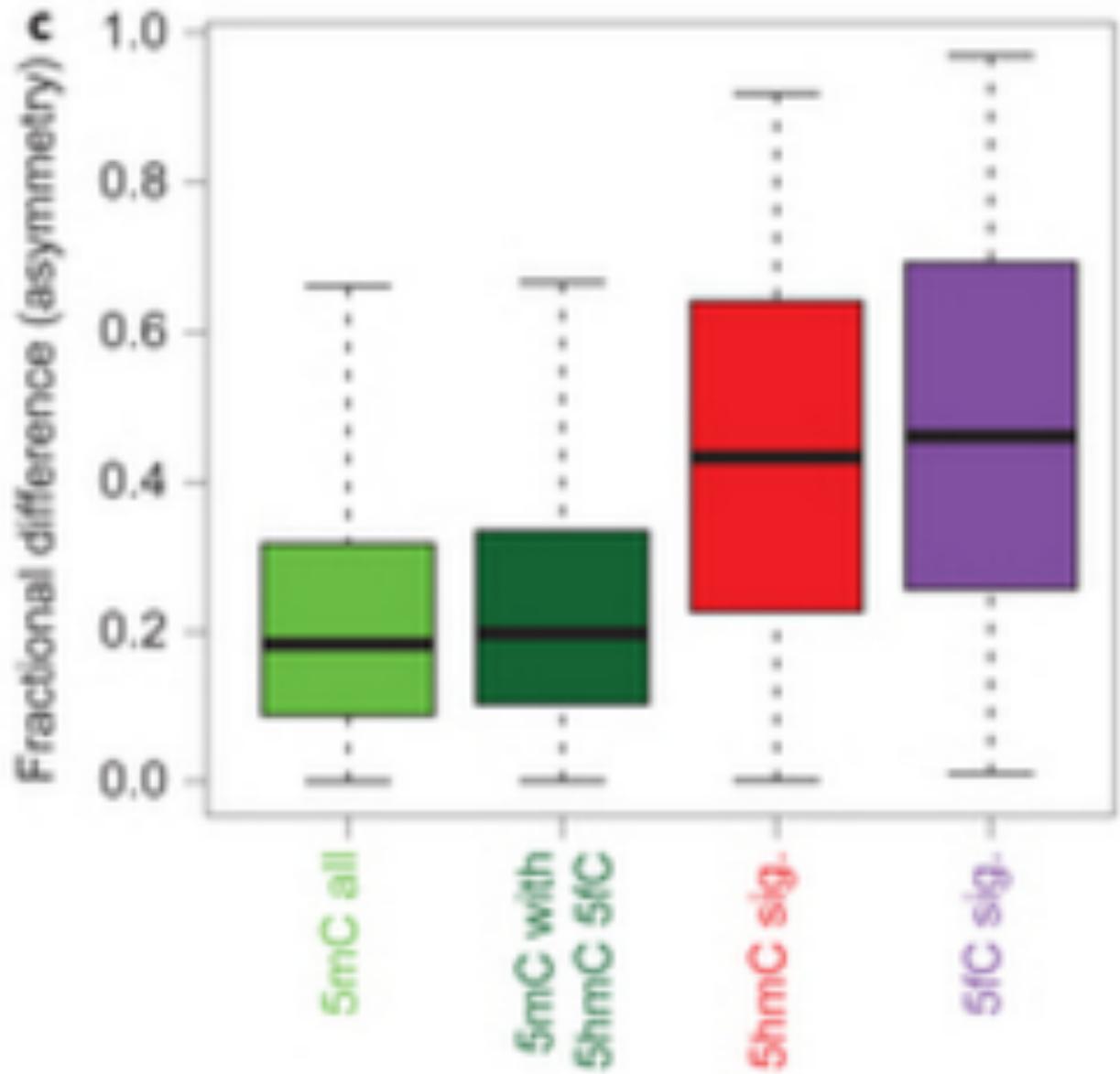
表观基因组学前沿



Ting Wang  
Department of Genetics  
Center for Genome Sciences and Systems Biology  
Washington University School of Medicine

Tsinghua University  
April 15-27

# **Epigenomic technology**



## Review

# Evolution of Epigenetic Regulation in Vertebrate Genomes

Rebecca F. Lowdon,<sup>1,\*</sup> Hyo Sik Jang,<sup>1</sup> and Ting Wang<sup>1,\*</sup>

**Empirical models of sequence evolution have spurred progress in the field of evolutionary genetics for decades. We are now realizing the importance and complexity of the eukaryotic epigenome. While epigenome analysis has been applied to genomes from single-cell eukaryotes to human, comparative analyses are still relatively few and computational algorithms to quantify epigenome evolution remain scarce. Accordingly, a quantitative model of epigenome evolution remains to be established. We review here the comparative epigenomics literature and synthesize its overarching themes. We also suggest one mechanism, transcription factor binding site (TFBS) turnover, which relates sequence evolution to epigenetic conservation or divergence. Lastly, we propose a framework for how the field can move forward to build a coherent quantitative model of epigenome evolution.**

### Trends

Epigenome evolution is characterized by variable conservation and divergence across the genome; within a clade (here vertebrates), rates of conservation or divergence are highly genome feature-specific.

TFBS turnover can mediate epigenome conservation or divergence.

Developmental genes are enriched in loci with divergent chromatin features, suggesting that rapid epigenome evolution may contribute novel regulatory mechanisms for lineage-

## **Questions on Epigenome Evolution:**

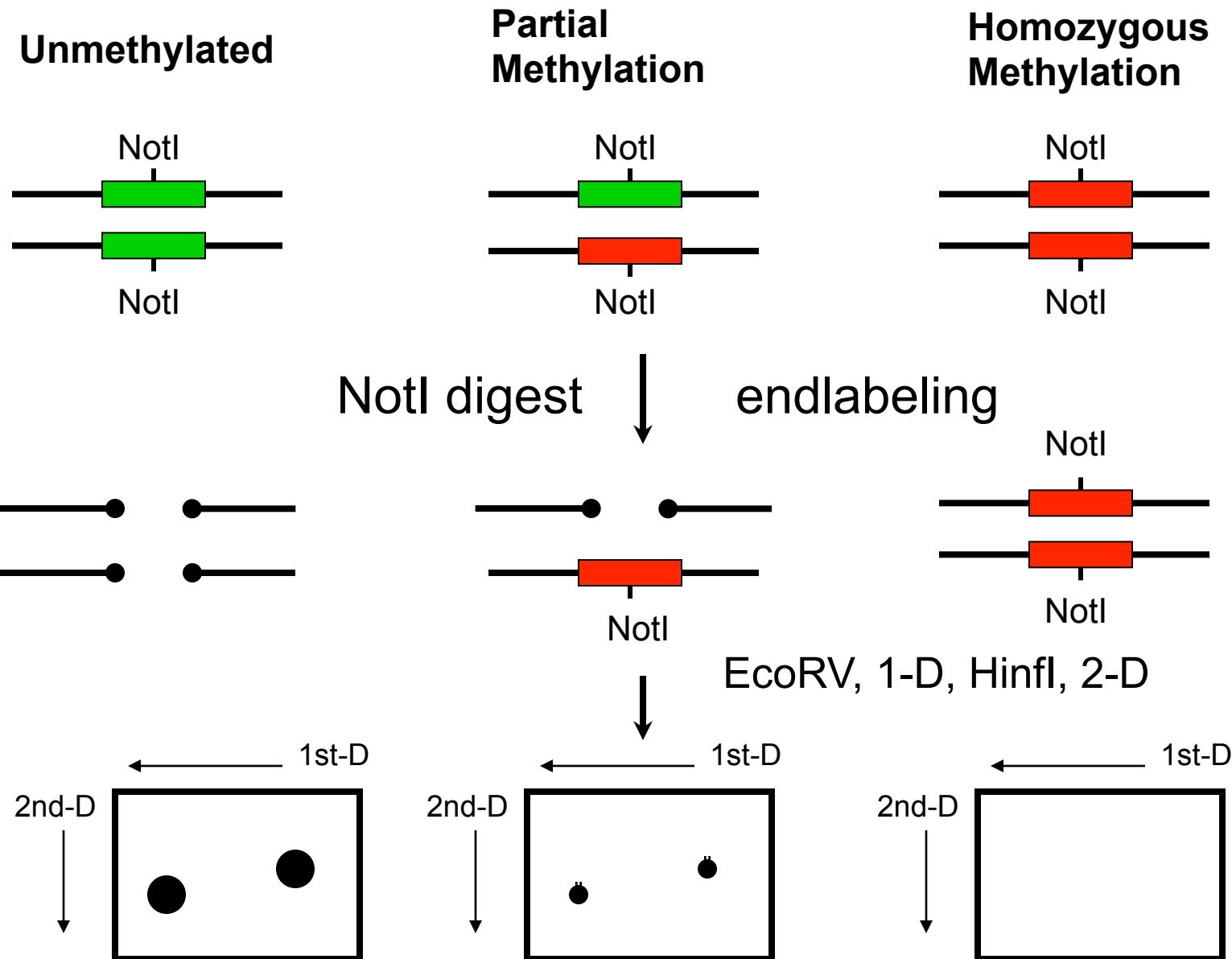
- 1) 细胞治疗是否能启动甲基化。我们实验室做细胞治疗，有时注射细胞后，发现比对照有增重，注射的细胞而非质粒，肯定不能改变基因，那这种增重究竟是不是epigenetic的。如何检测？直接抽血？
- 2) 现已知的表观遗传进化的现象，以及它们在表观遗传进化中所起的作用
- 3) epigenetic revolution的研究，取得了哪些重要或有趣的成果？举个例子？
- 4) epigenetic revolution的研究的难点在哪？比较同一物种内的epigenetic和/或genetic差异，而不是物种之间的差异，是否对我们研究基因调控等方面更直接和更有帮助？

## **Questions on Big Bio Data:**

- a) 基于已知数据库，对于新的数据，我们用怎样的方法将这些数据归类并识别的。比如基因组里有很多的信息是我们无法归类的，那么我们应该对这些不确定的数据怎样处理，归类到编码区，非编码区，启动子或增强子等等？
- b) 对于big data processing的入门学习，我应该从哪些方面入手？
- c) 怎样将多重数据如表达，mutation, 甲基化，Hi-C, chip-seq以及其他数据进行有效整合来阐述表观中重要的assumptions（比如表达与表观遗传修饰之间的定量关系）？
- d) 表观的修饰通常都是genome-wide的形式进行，怎么样通过genome-wide的大数据分析来阐述如reader, writer, eraser间作用具有特异性(specifity)和选择性(selectivity)?
- e) 在大数据时代，数据一致性将会是个问题。不同实验室发布的数据之间一致性到底有多大？single cell得到的数据和mix cell 得到的数据一致性有多大？如果出现不一致，怎么从数据分析的角度评判哪个数据好？
- f) 在同一种细胞群体中和一个细胞的增殖生长周期中，表观遗传学的各检测方法和检测的指标的波动（方差）有多大？能否充分利用现有的大数据信息对这个方差做出估计？

# **How to detect epigenetic marks?**

# Restriction Landmark Genome Scanning (RLGS)

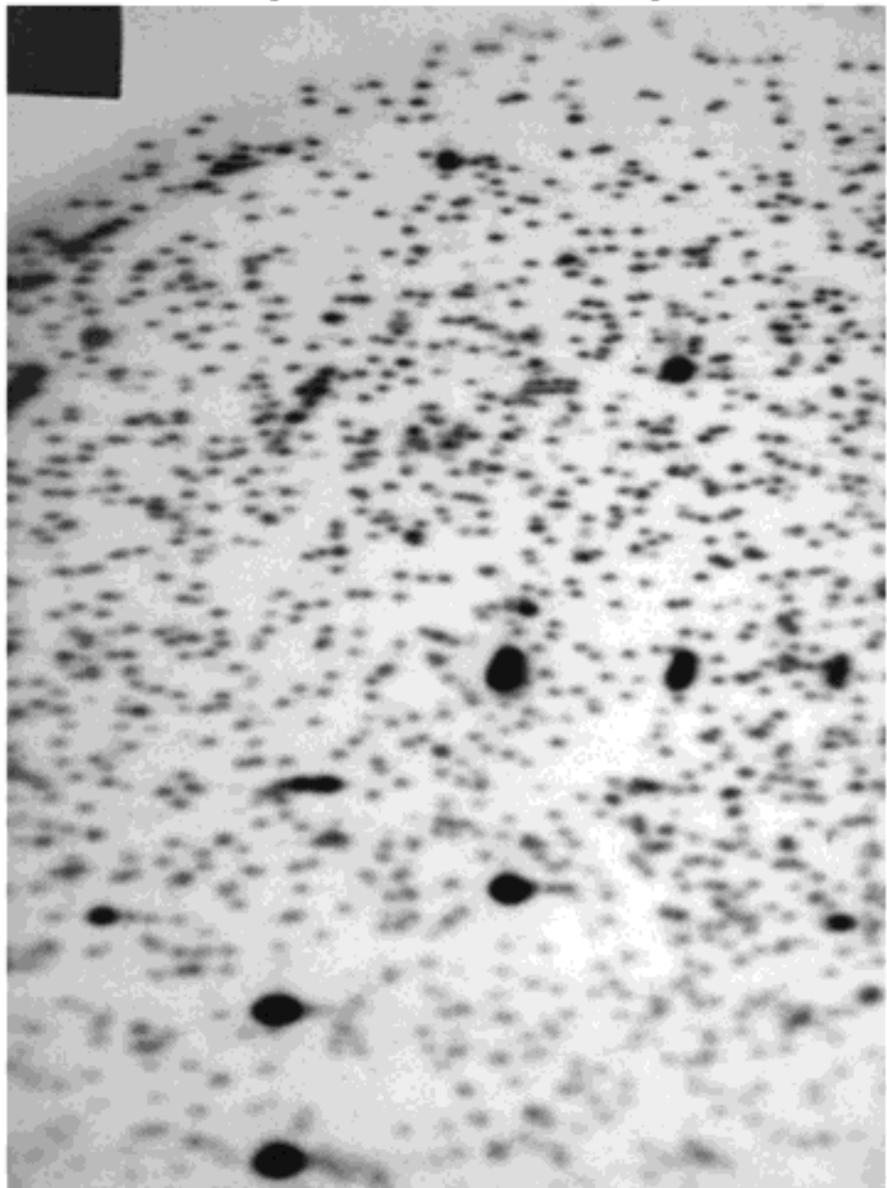


← 1st dimension Not I / EcoR V —

1 kb

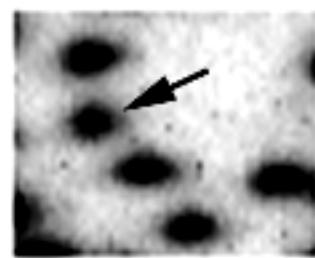
2 kb

— 2nd dimension Hinf I ↓

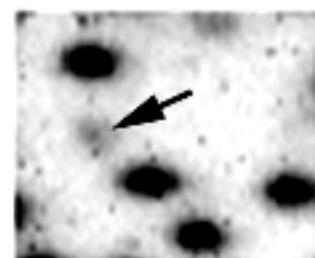


-0.8 kb

Normal



Tumor



-0.3

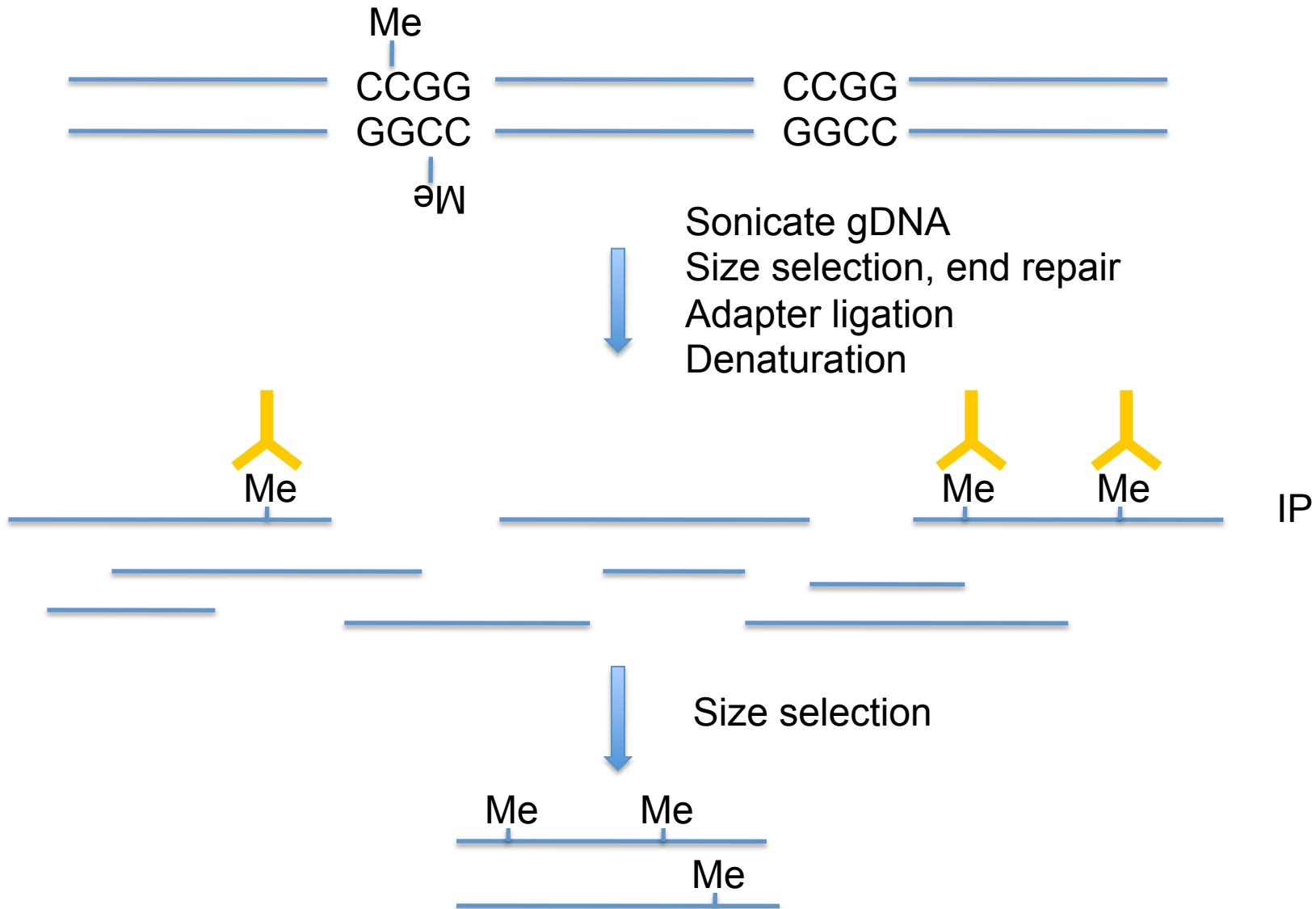
-0.15 kb

# Scaling to high coverage, high resolution

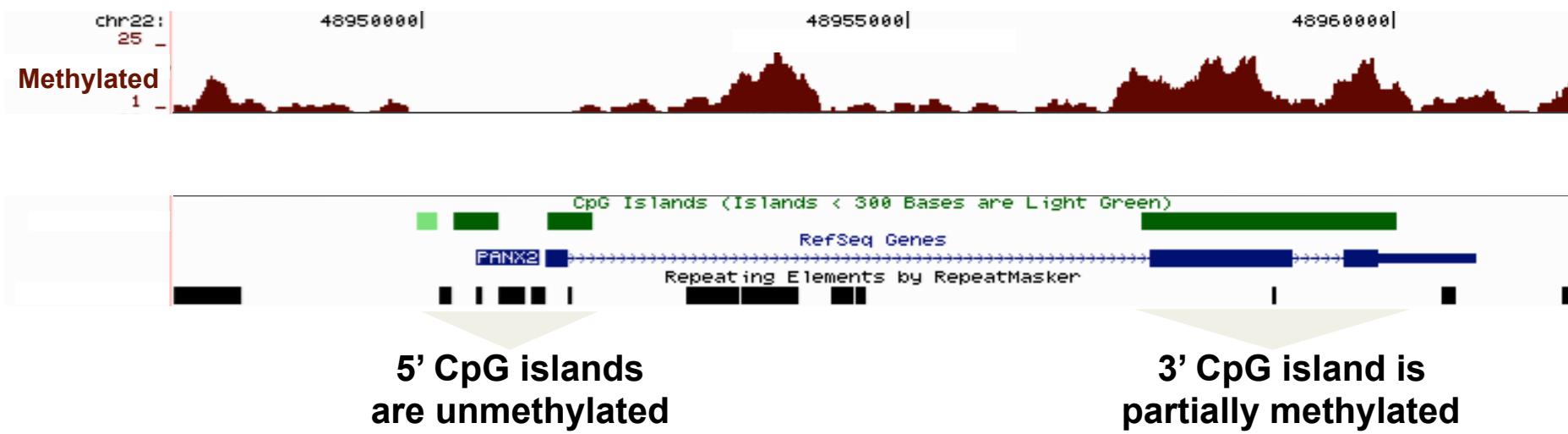
- Enrichment based methods
  - MeDIP-seq
  - MBD-seq/MethylMiner
- Restriction enzyme based methods
  - MRE-seq
  - HELP, Methyl-MAPS, Methyl-seq
- Bisulfite based methods
  - MethylC-seq
  - RRBS, bisulfite padlock
- Direct reading of modified nucleotides
  - SMRT
  - Nanopore sequencing

# **Enriching for methylated DNA targets**

# MeDIP-seq and MBD-seq

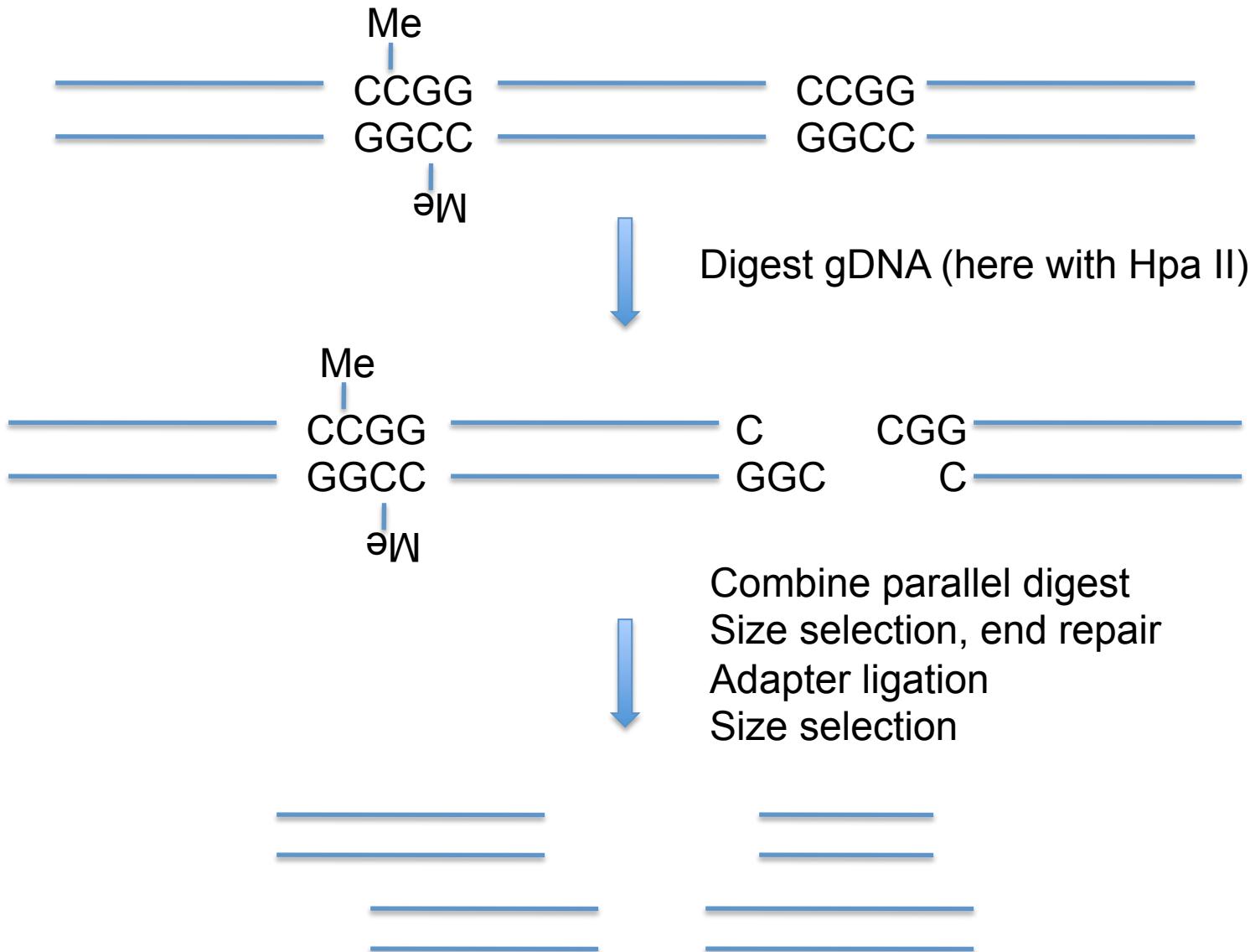


# Typical MeDIP data on a genome browser

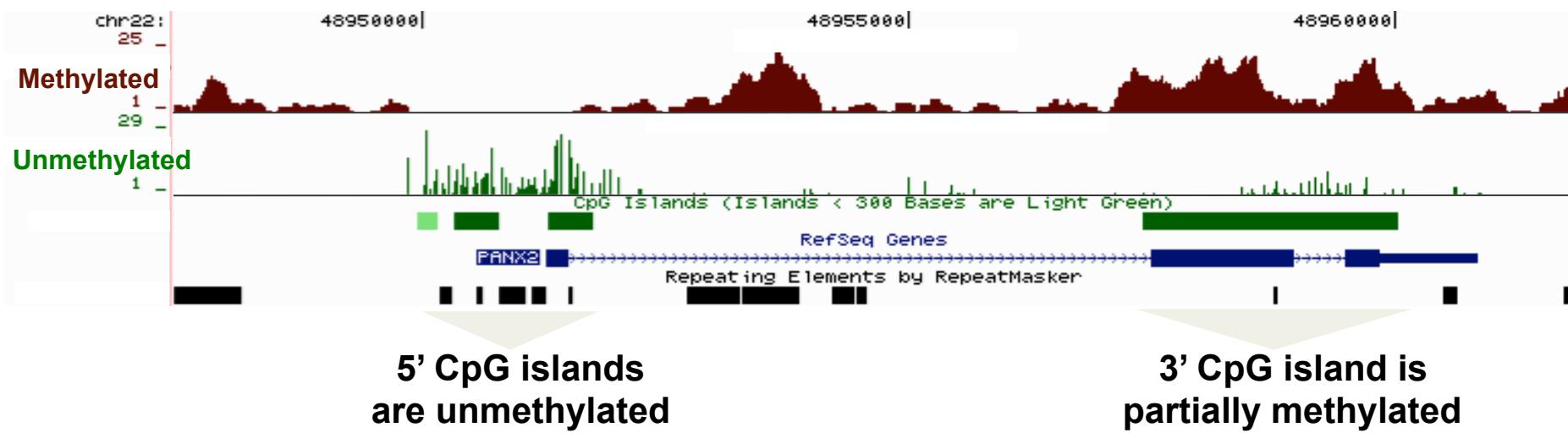


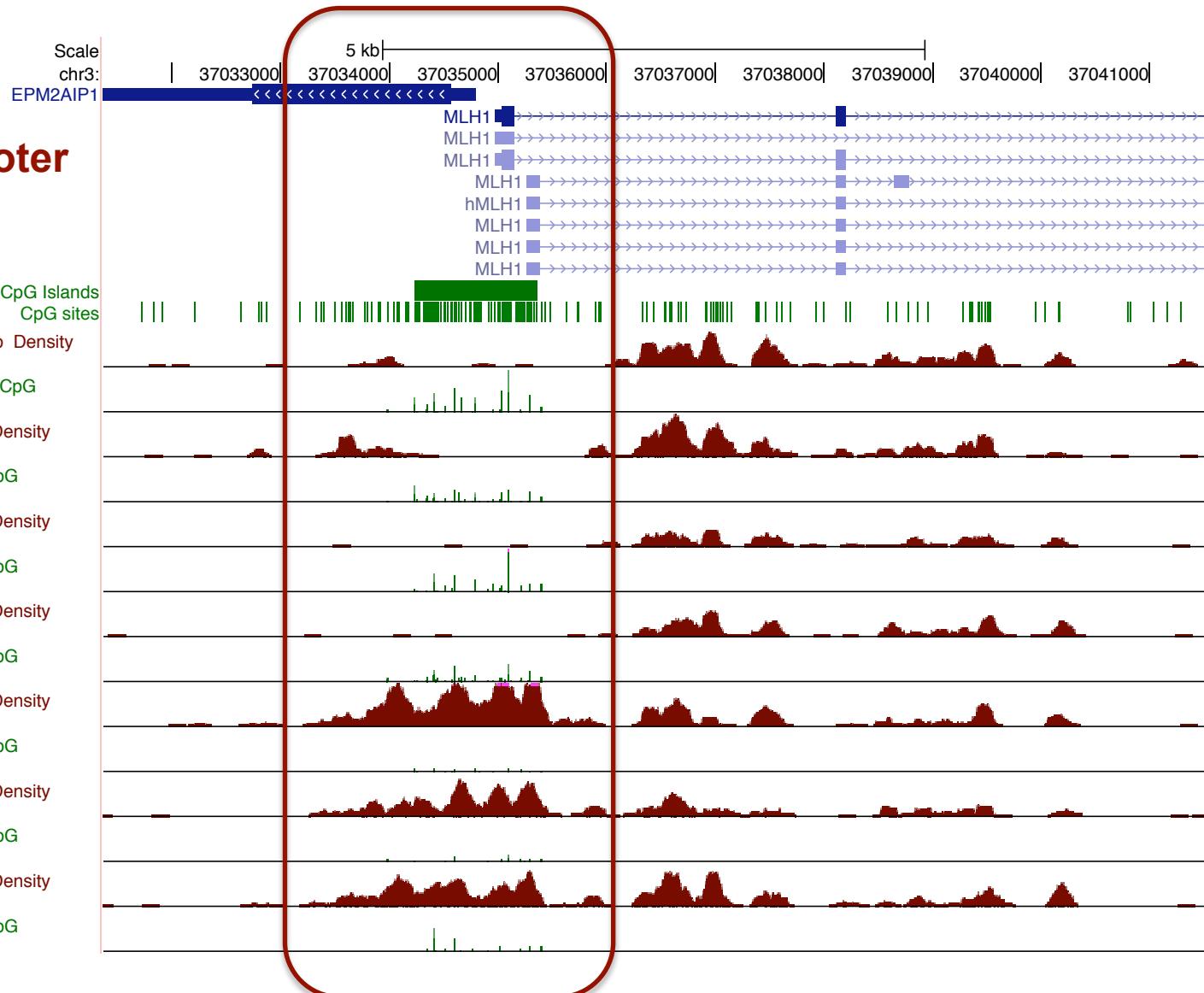
# Taking advantage of methylation dependent restriction enzymes

# MRE-seq

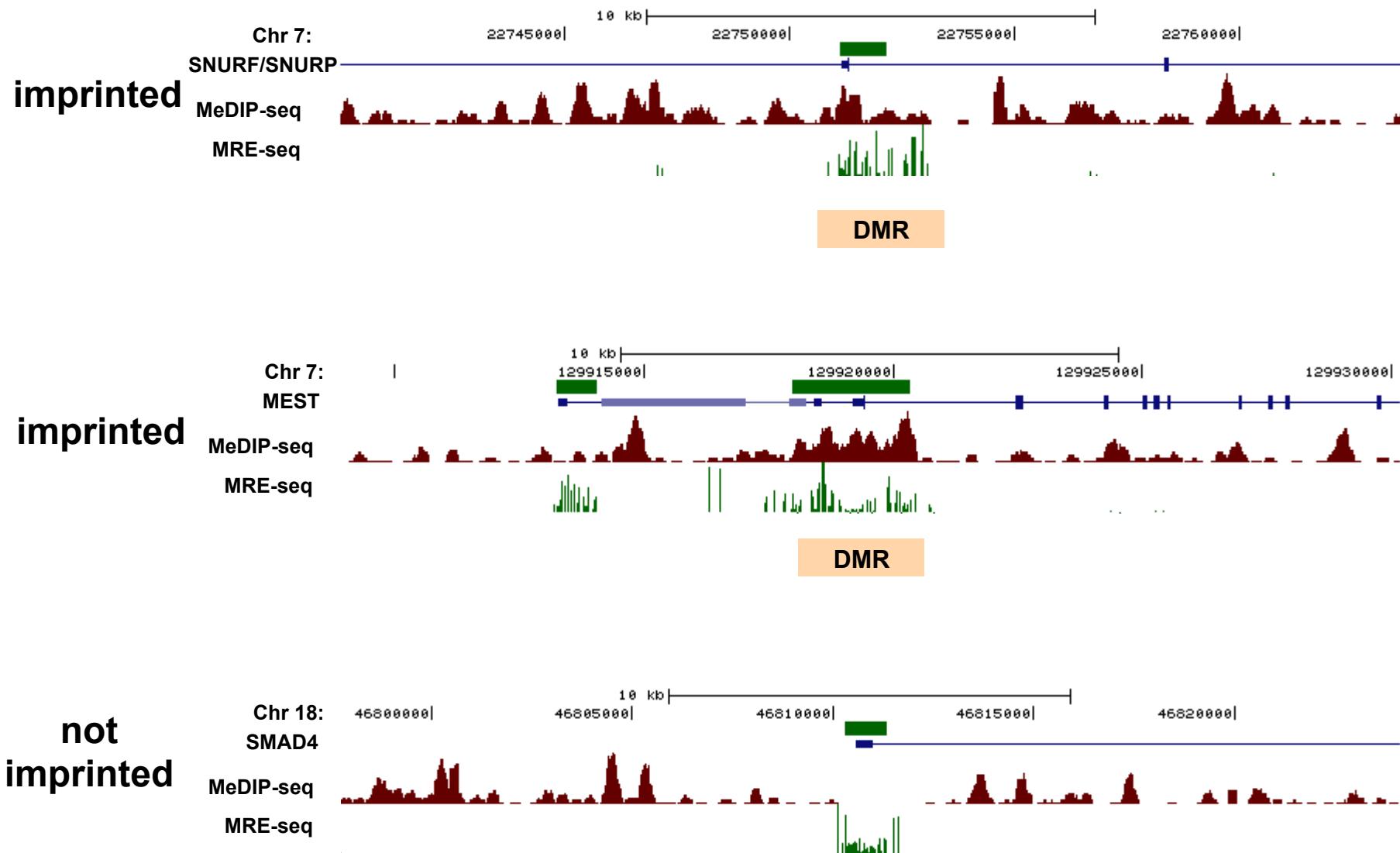


# Typical MeDIP/MRE data on a genome browser

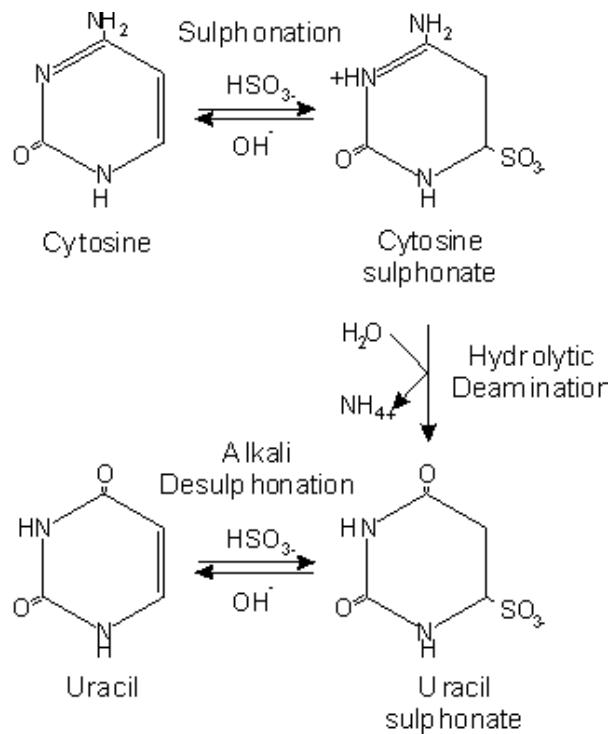




# Allele-specific methylation at imprinted genes



# **Gold standard: bisulfite sequencing**



# Bisulfite sequencing

## Sonication

Watson >>**AC<sup>m</sup>GTTCGCTTGAG**>>  
 Crick <<**TGC<sup>m</sup>AAGCGAACTC**<<

$\text{C}^m$  methylated  
 $\text{C}$  Un-methylated

## Denaturation

Watson >>**AC<sup>m</sup>GTT**C**GCTTGAG**>>  
 Crick <<**TGC<sup>m</sup>AAG**C**GAACTC**<<

## Bisulfite Treatment

BSW >>**AC<sup>m</sup>GTT**U**GTGTTGAG**>>  
 BSC <<**TGC<sup>m</sup>AAG**U**GAATU**<<

## PCR Amplification

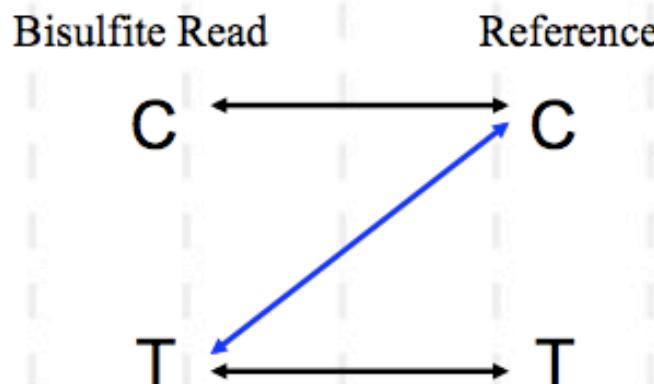
BSW >>**AC<sup>m</sup>GTT**T**GTTGAG**>>  
 BSC <<**TGC<sup>m</sup>AAG**T**GAATT**<< C-poor stands  
 BSWR <<**TG CAAACAAACTC**<<  
 BSCR >>**ACG TTC**A**CTTAAA**>> G-poor stands

## Sequencing and Mapping

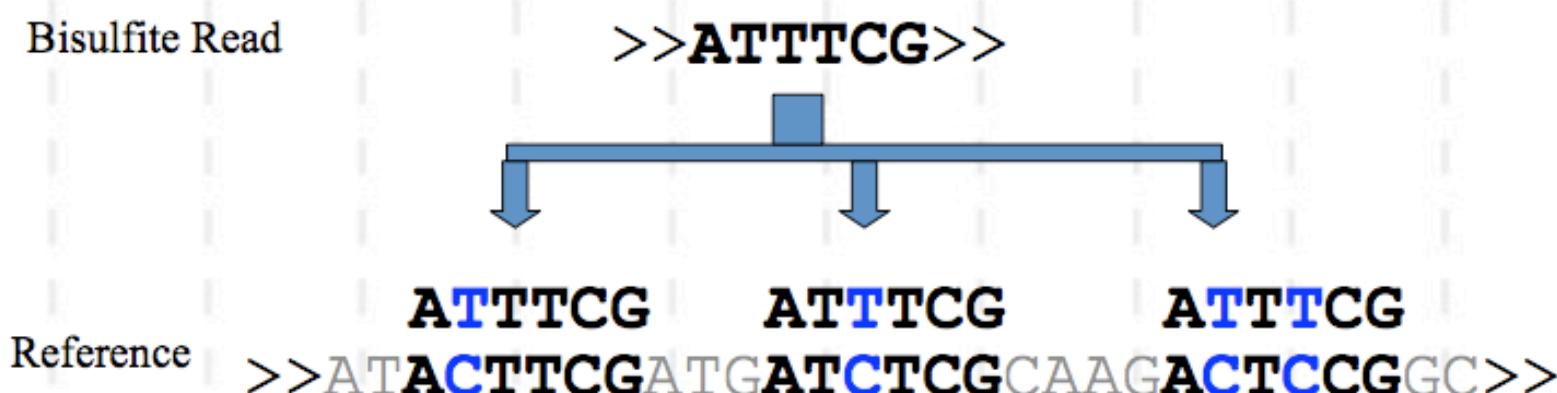
# Methylation call

# Bioinformatics challenges in mapping high throughput bisulfite reads

- Mapping Asymmetry

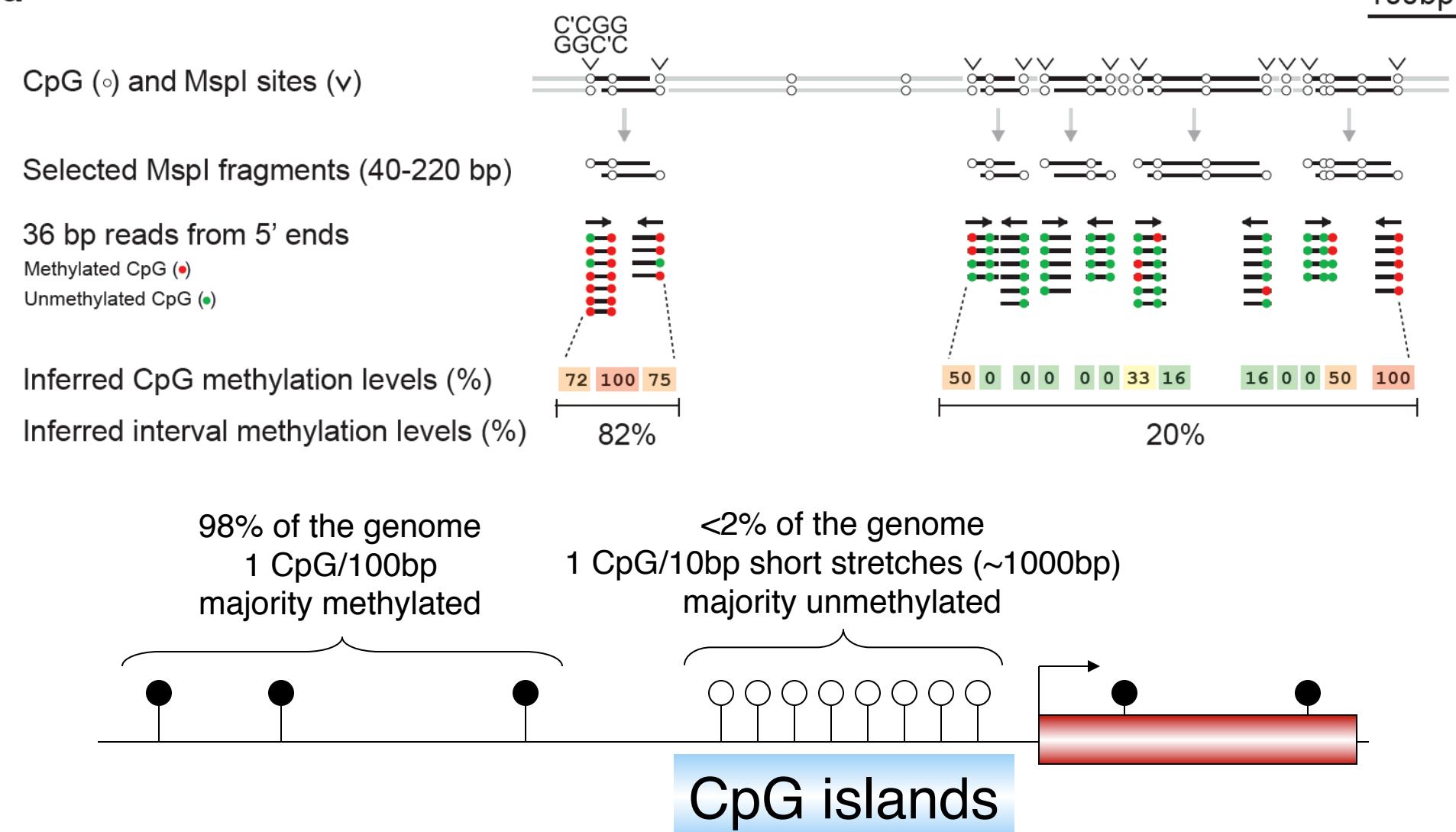


- Increased Mapping Space

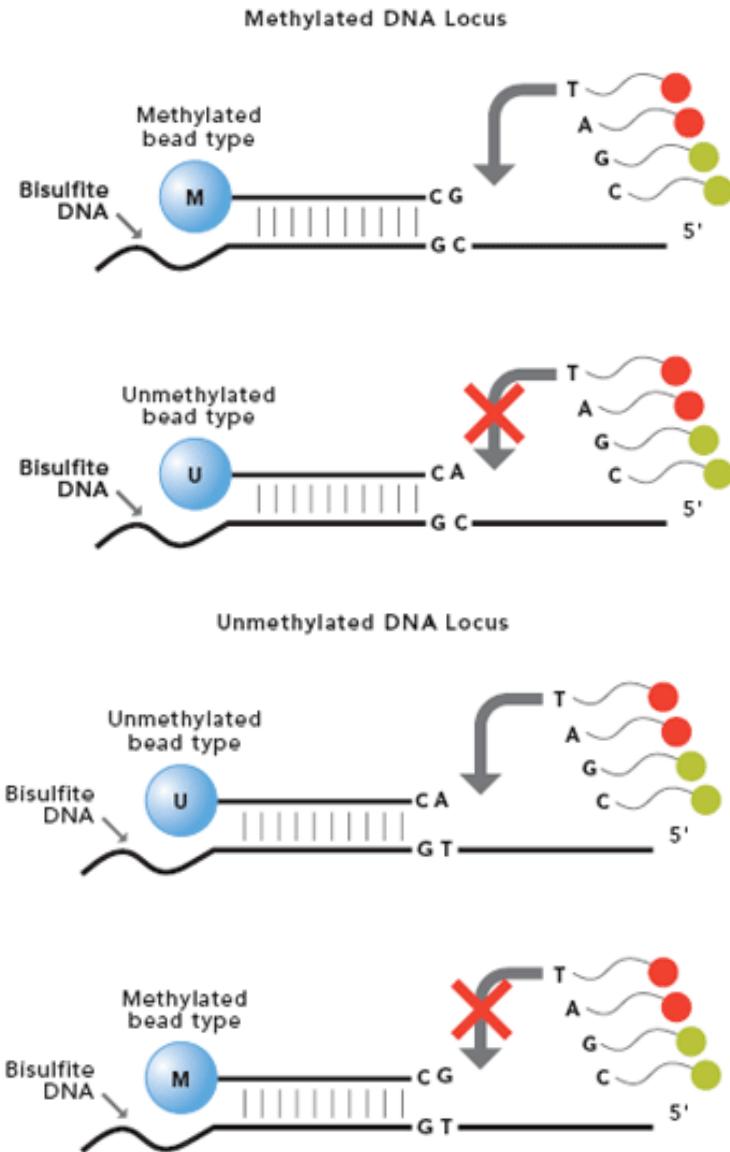


# RRBS: Reduced Representations Allow Enrichment of CpG Dinucleotides

a



# 5-mC detection with the Infinium BeadChip



- **Pros**

- Relatively inexpensive
- Single-base resolution
- Internal quality controls
- Highly reproducible ( $r > 0.98$ )
- PCR-free protocol
- Works with FFPE samples

- **Cons**

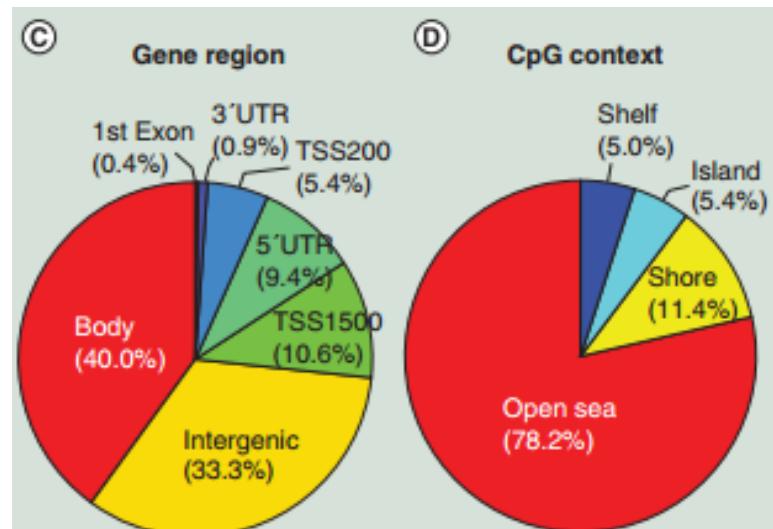
- Only covers a subset of the methylome
- Bisulfite treatment can damage DNA
- Dependent on bisulfite conversion

# 5-mC detection with the Infinium BeadChip

Array	Year released	# of Sites	Targeted sites
HumanMethylation27k	2008	> 27k	<ul style="list-style-type: none"><li>&gt;14K RefSeq genes</li></ul>
HumanMethylation450k	2011	> 450k	<ul style="list-style-type: none"><li>99% of RefSeq genes</li></ul>
Infinium MethylationEPIC	2015	> 850k	<ul style="list-style-type: none"><li>&gt; 90% of sites in the 450k</li><li>FANTOM5 and ENCODE enhancers (350k sites)</li><li>CpG sites outside of CpG islands</li><li>ENCODE open chromatin</li></ul>



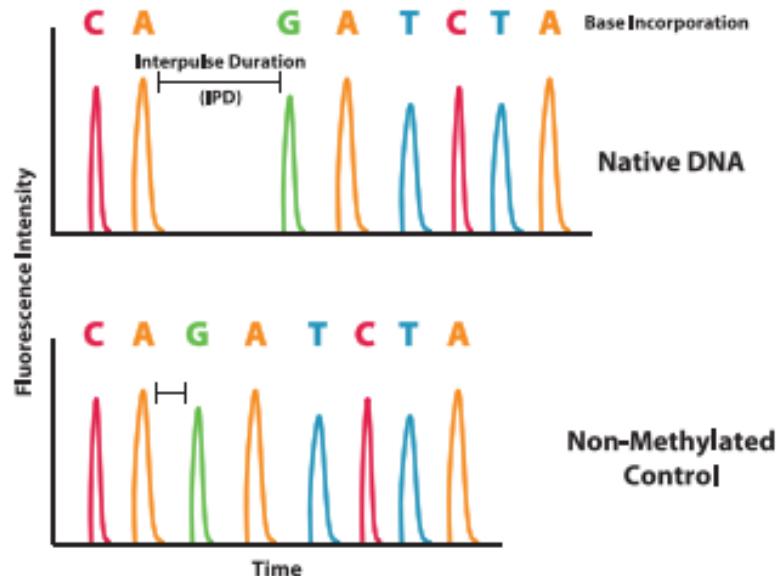
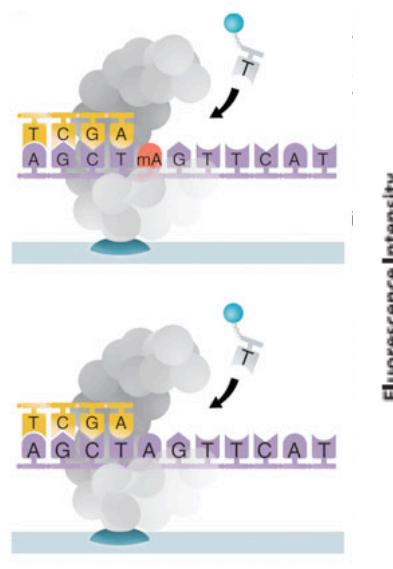
Context of the new EPIC probes



# **Direct detection of modified nucleotides**

# Single-molecule real-time (SMRT) sequencing

SMRT sequencing discriminates between different bases by analyzing variations in polymerase kinetics



- Pros

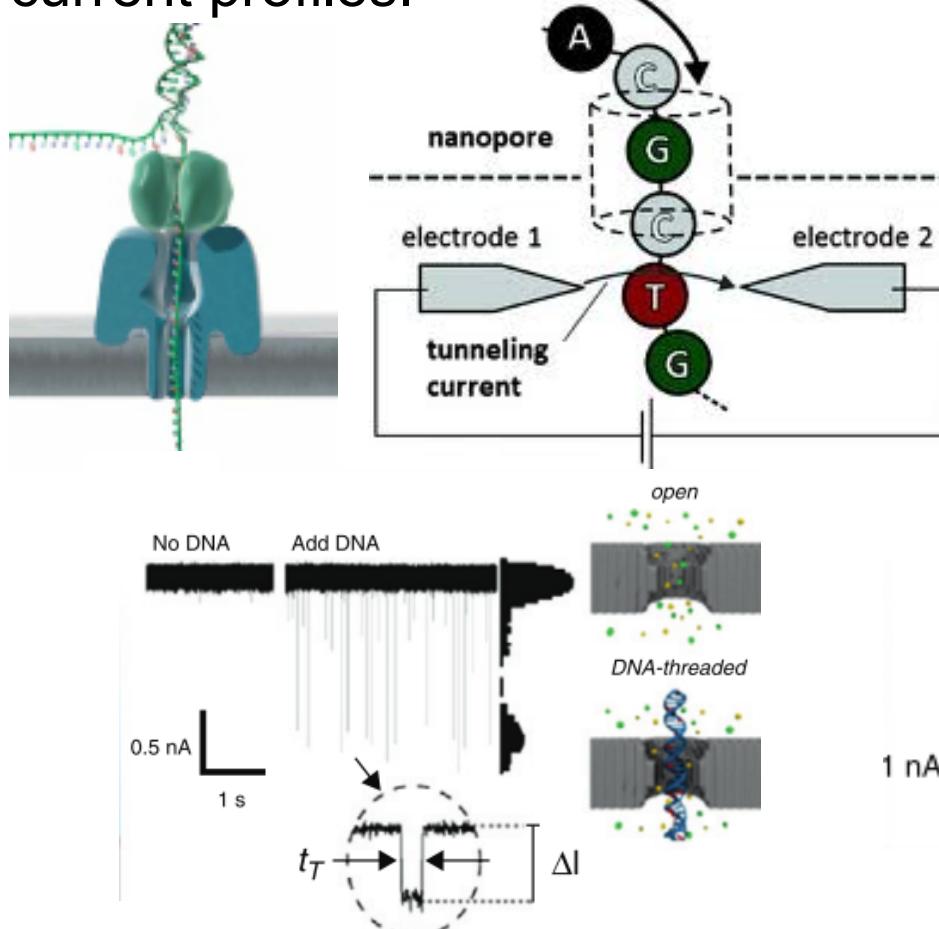
- Single-base resolution
- Measures absolute levels of many modified nucleotides
- “Raw” DNA is used
- Long reads

- Cons

- Suboptimal accuracy
- Low throughput

# Nanopore sequencing

Nanopore amperometry methods can discriminate between C, 5-mC, and 5-hmC due to differences in current profiles.



- **Pros**

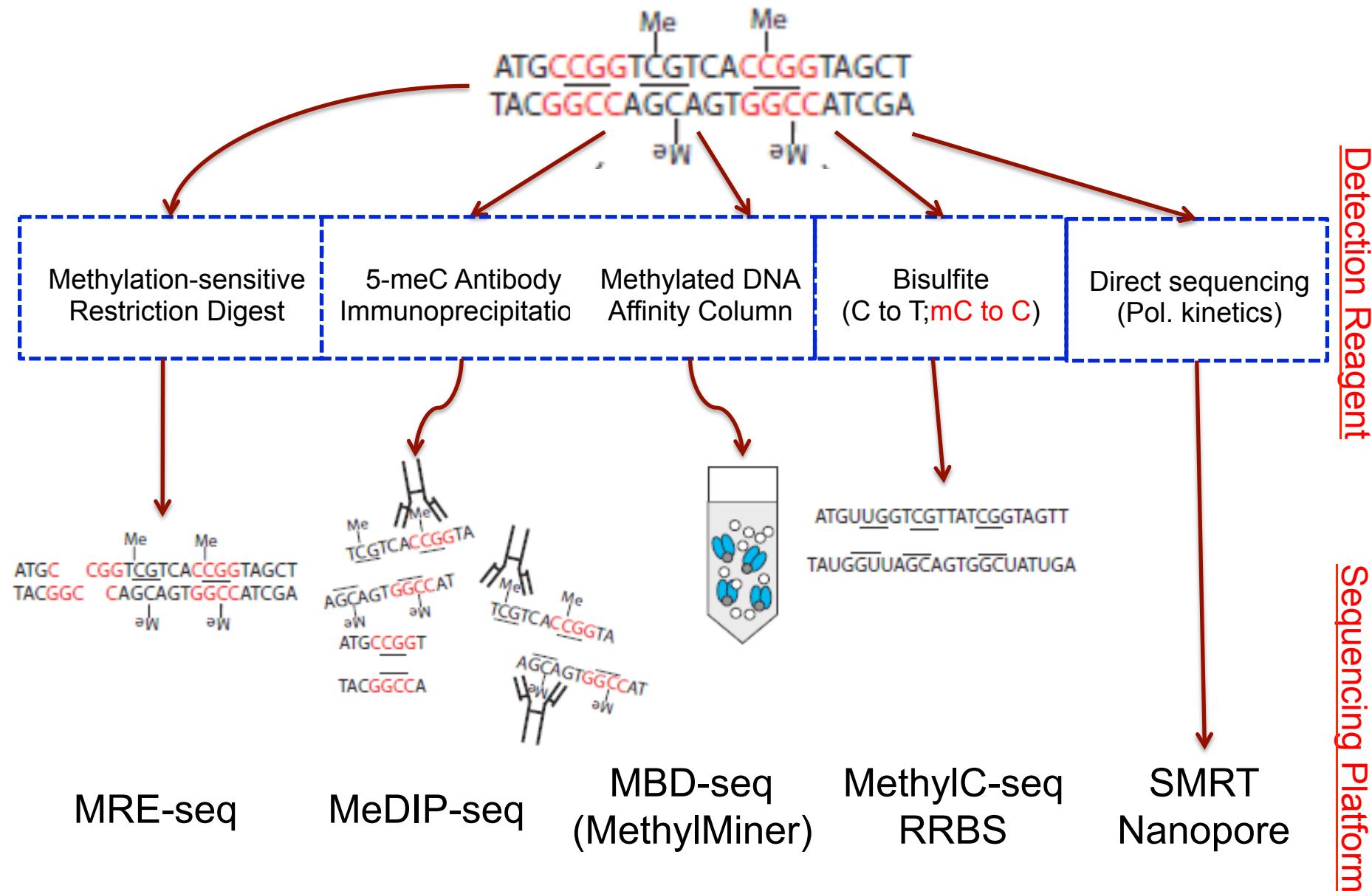
- Single-base resolution
- Measures absolute levels of many modified nucleotides
- “Raw” DNA is used
- Long reads

- **Cons**

- Suboptimal accuracy
- Low throughput



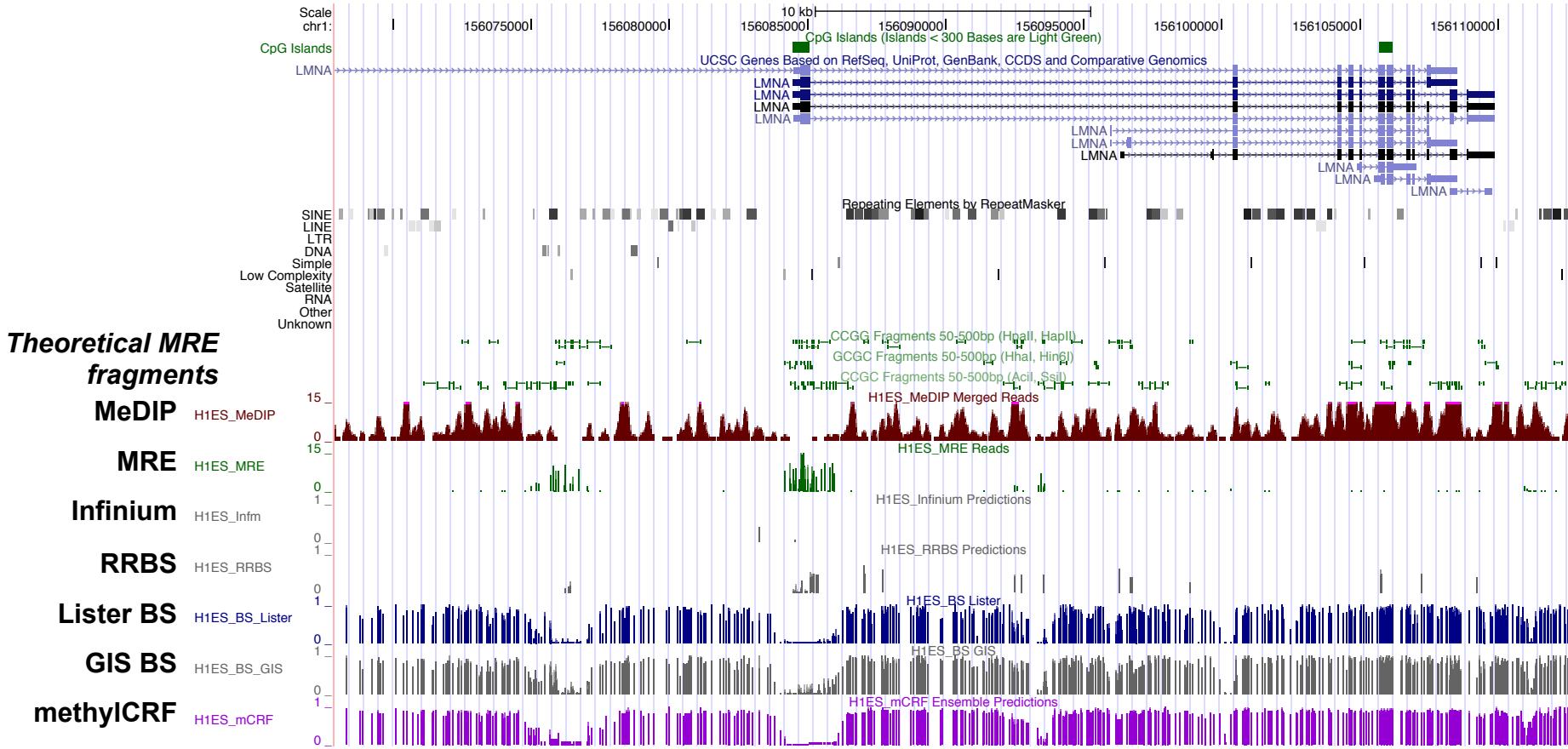
# Modern DNA Methylocomics



# Comparison of Methylome coverage

	CpG coverage	CpG island coverage	Resolution (bp)	Illumina Lanes
<b>Genome total</b>	28 M	28 K	NA	NA
BS shot gun	26 M	27 K	1	207 (2009) 3 (now)
RRBS	0.2-1M	15 K	1	1
MRE-seq and MeDIP-seq	25 M	27 K	1 and 200	8 (2009) half (now)
Golden-Gate	1,500	800	1	NA
Infinium	27,500 480,000 (2012) 850,000 (2016)	12 K 27 K 27 K	1	NA

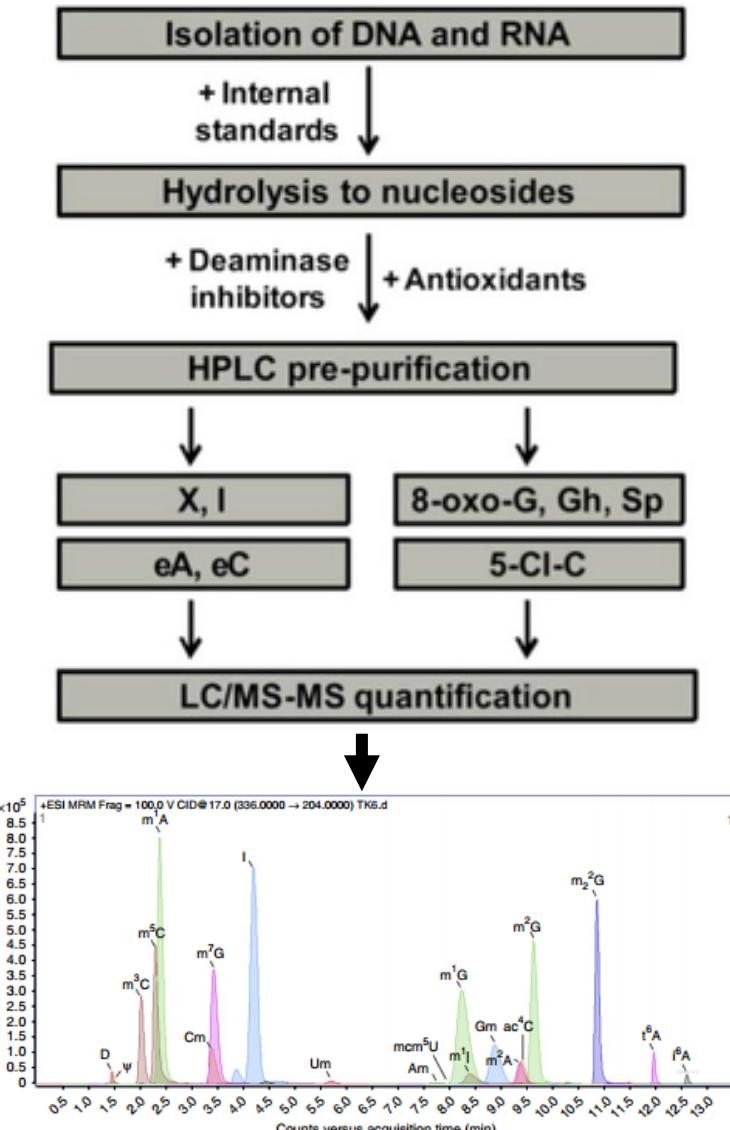
# Predicting single CpG methylation level with Conditional Random Field (methylCRF)



# Technologies for quantifying the methylome & hydroxymethylome

DNA modification	Measurement	Non-targeted enrichment	Targeted enrichment	Whole genome	Arrays
5mC	Absolute (single base)	RRBS, mRRBS, LCM-RRBS or scRRBS	Microdroplet PCR Bisulphite, Patch PCR, mTACL, BSPP, LHC-BS (pre- and post-conversion) or RSMA	WGBS, T-WGBS or PBAT	Infinium BeadChip
	Relative (peak)	MRE-seq, MeDIP-seq, MBD-seq or MethylCap-seq			CHARM or MeKL-ChIP
5mC oxidation derivatives	Absolute (single base)	RRHP  Reduce representation sequencing with TAB-seq, oxBs-seq, CAB-seq, fCAB-seq or redBS-seq	Locus-specific sequencing with TAB-seq, oxBs-seq, CAB-seq, fCAB-seq or redBS-seq	TAB-seq, oxBs-seq, CAB-seq, fCAB-seq, redBS-seq	
	Relative (peak)	DIP-seq, anti-CMS, hMe-Seal, fC-Seal, GLIB, JBP1, EpiMark or Aba-seq			

# Liquid chromatography and mass spectrometry (LC-MS)



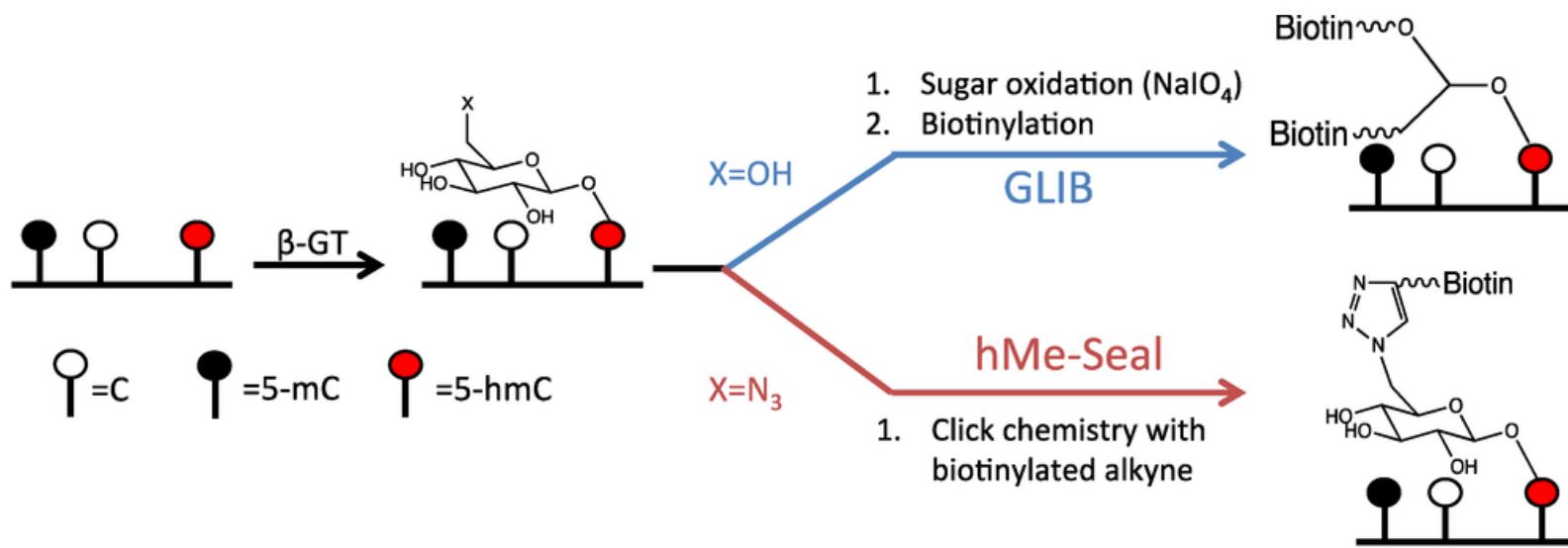
- Pros

- Well established method
- A gold standard
- High sensitivity

- Cons

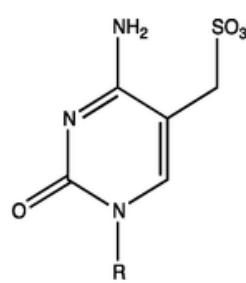
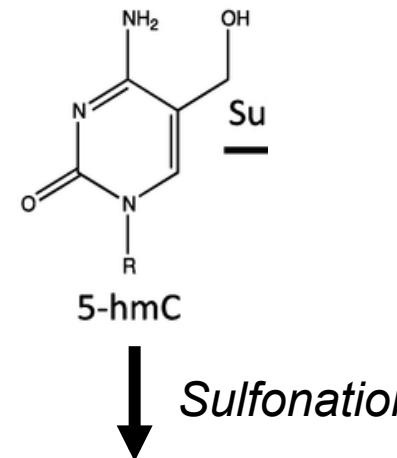
- Requires extensive experience
- Only measures global levels

# 5-hmC detection with GLIB- & hME-Seal-seq



- **Pros**
  - Covalent chemical labeling and biotin-based purification is sensitive and specific
- **Cons**
  - Biotinylated adduct can affect PCR so amplification-free library and sequencing preferred

# 5-hmC detection with CMS-seq

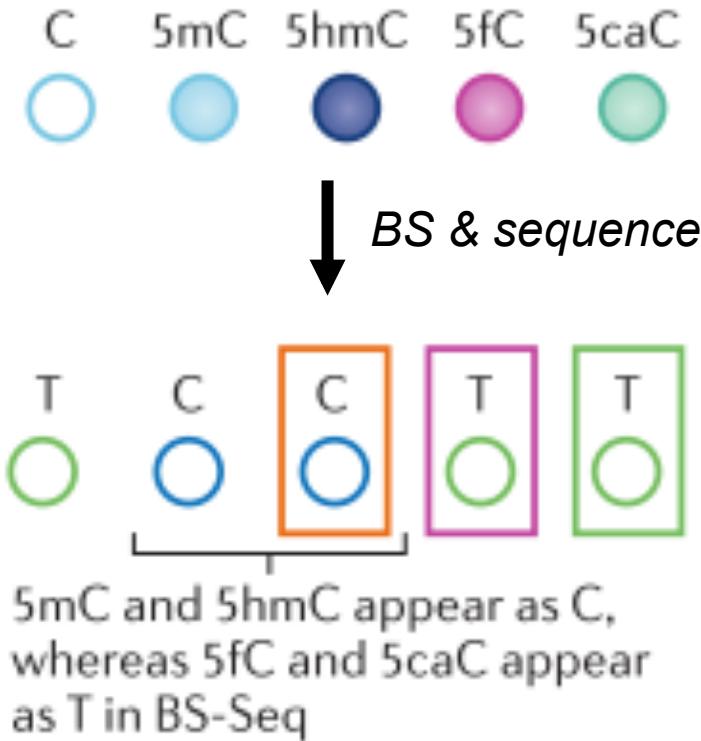


**5-methylsulphonate  
(CMS)**

An arrow labeled "Pull-down with anti-CMS antibody & sequence" points downwards from the CMS structure, indicating the final step in the detection process.

- **Pros**
  - Anti-CMS antibodies are more sensitive and less density dependent than anti-5-hmC antibodies
- **Cons**
  - Low resolution
  - Limited to specific sequences

# Whole genome bisulfite sequencing (WGBS)



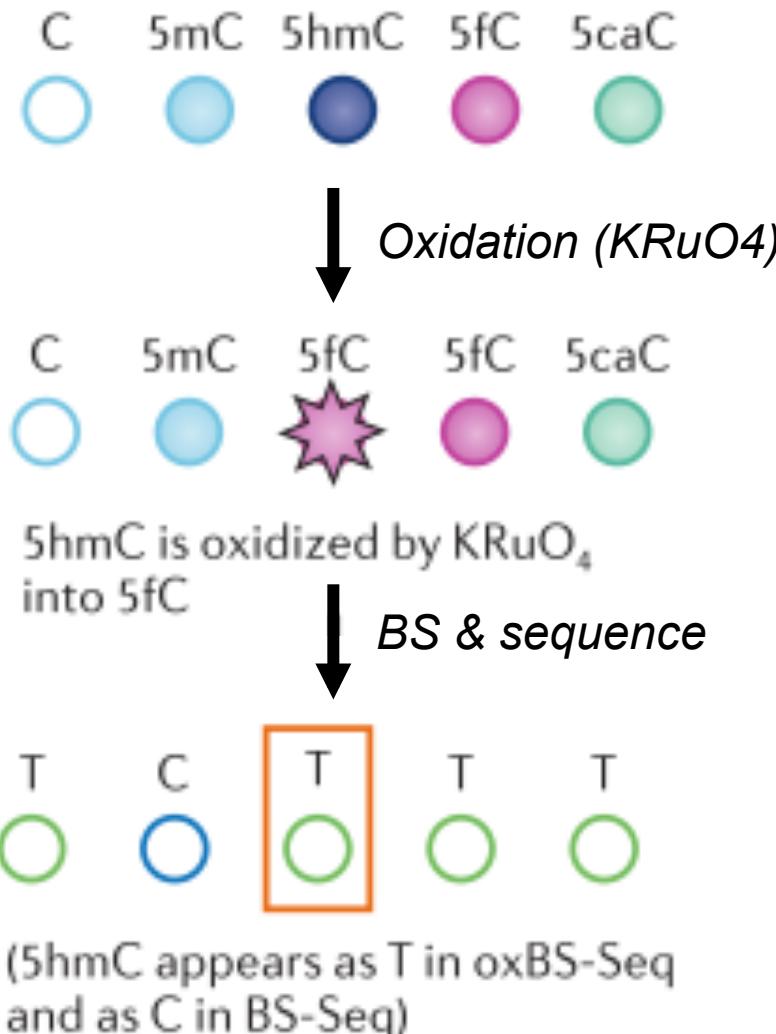
- Pros

- “Gold standard”
- Single-base resolution
- Profiles a large % of the “CpG-ome”
- Measures absolute levels

- Cons

- \$\$\$
- Does not distinguish between 5-mC and 5-hmC
- Does not distinguish between C, 5-caC, & 5-fC
- Bisulfite treatment can damage DNA
- Dependent on bisulfite conversion
- Reduced complexity of sequencing reads creates alignment challenges
- Requires a large amount of DNA

# 5-mC & 5-hmC detection with oxBs-Seq



- Pros

- A gold standard
- Single-base resolution
- Measures absolute levels of 5-hmC and 5-mC

- Cons

- \$\$\$
- Requires a large amount of DNA
- Bisulfite treatment can damage DNA
- Dependent on bisulfite & oxidative efficiencies
- Requires a large amount of DNA
- Reduced complexity of sequencing reads creates alignment challenges

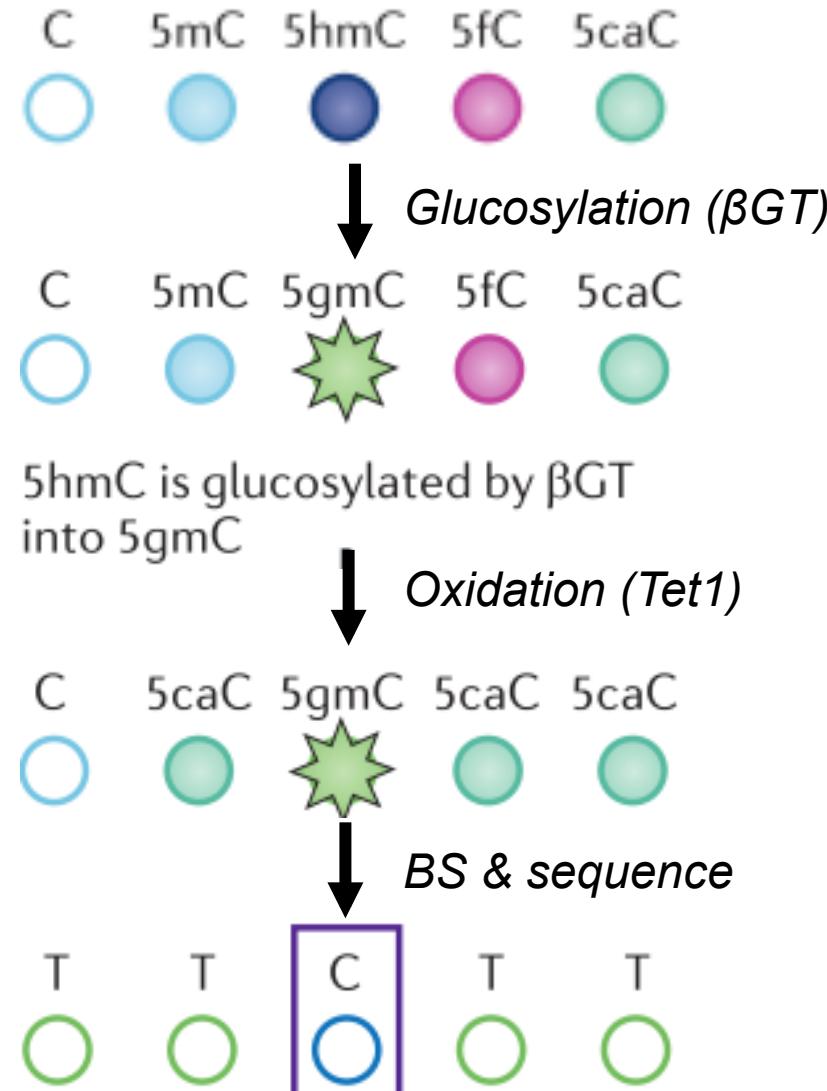
oxBS-seq =Oxidative bisulfite sequencing

KRuO<sub>4</sub> = potassium perruthenate

BS = Bisulfite treatment

Figures adapted from Plongthongkum et al., *Nature Reviews Genetics* (2014) | Booth et al., *Science* (2012)

# 5-mC & 5-hmC detection with TAB-Seq



- **Pros**

- Single-base resolution
- Enzymatic treatments are more mild than oxBs-seq
- Measures absolute levels of 5-hmC and 5-mC

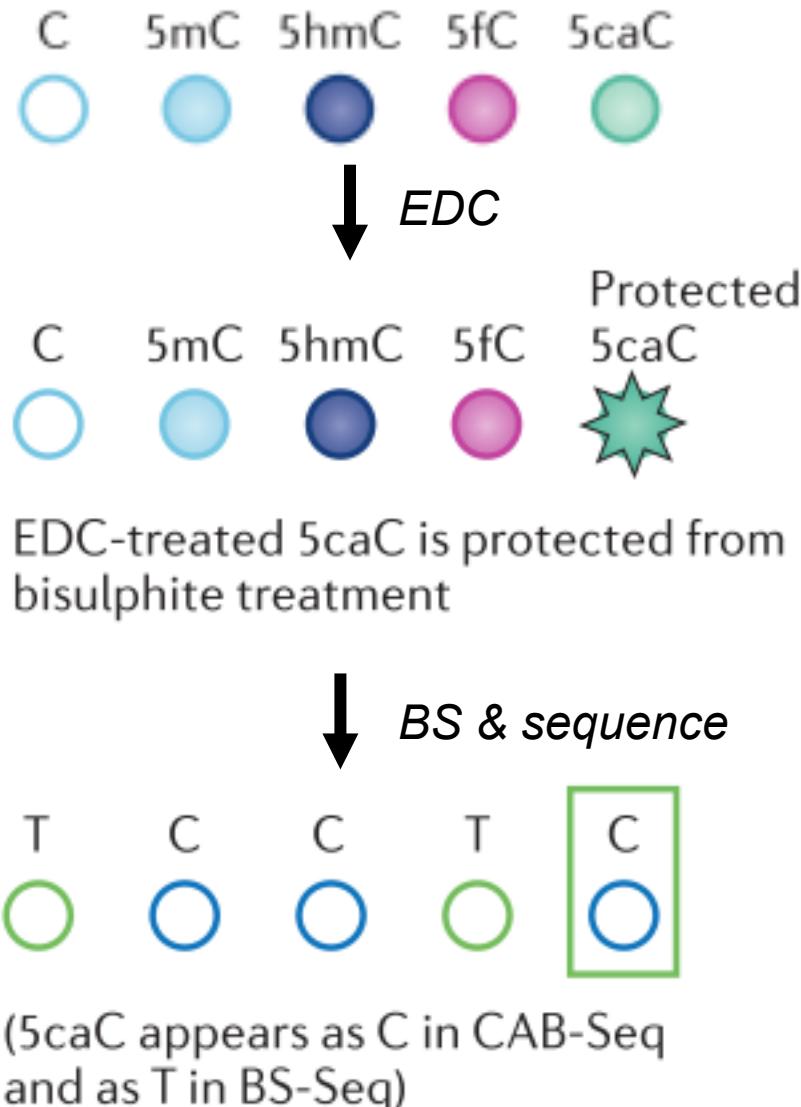
- **Cons**

- \$\$\$
- Bisulfite treatment can damage DNA
- Dependent on bisulfite and glycosylation efficiencies
- Requires a large amount of DNA
- Reduced complexity of sequencing reads creates alignment challenges

TAB-seq = Tet-assisted bisulfite sequencing

BS = Bisulfite treatment

# 5-caC detection with CAB-Seq



- Pros

- Single-base resolution

- Cons

- \$\$\$

- Bisulfite treatment can damage DNA

- Dependent on bisulfite conversion

- Requires a large amount of DNA

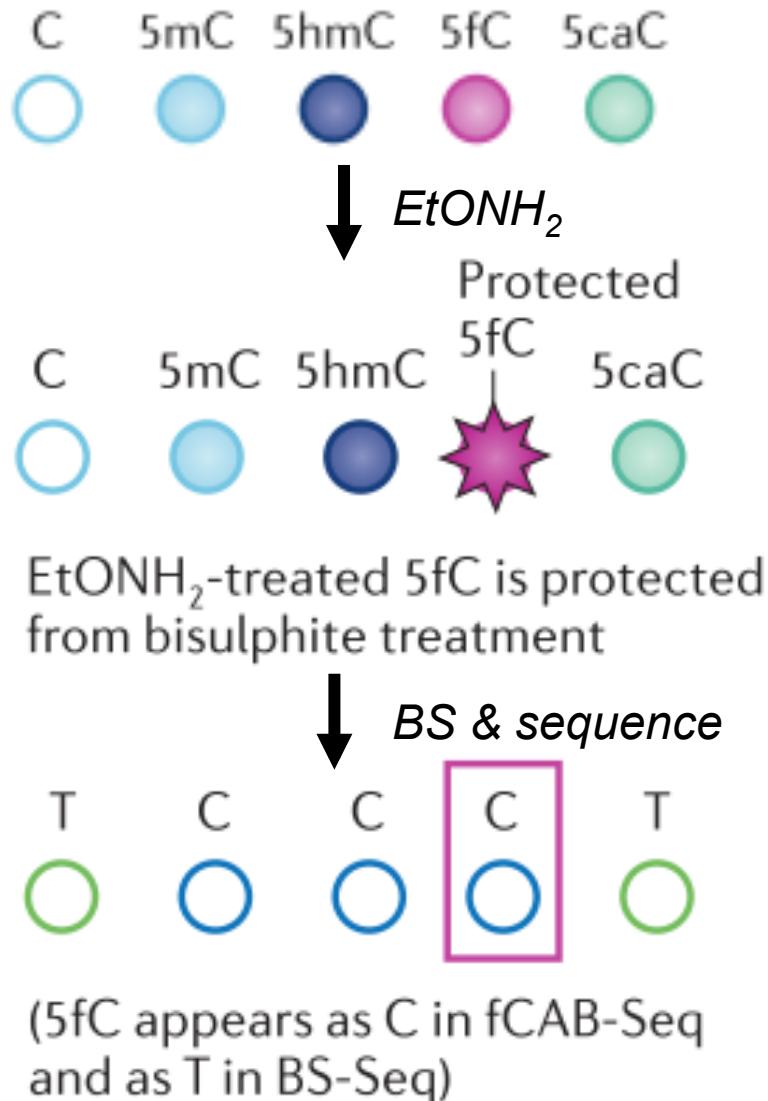
- Reduced complexity of sequencing reads creates alignment challenges

CAB-seq = Chemical modification-assisted bisulfite sequencing

BS = Bisulfite treatment

Figures adapted from Plongthonkum et al., *Nature Reviews Genetics* (2014) | Lu et al., *Journal American Chemical Society* (2013)

# 5-fC detection with fCAB-Seq



- **Pros**

- Single-base resolution

- **Cons**

- \$\$\$
- Bisulfite treatment can damage DNA
- Dependent on bisulfite conversion
- Requires a large amount of DNA
- Reduced complexity of sequencing reads creates alignment challenges

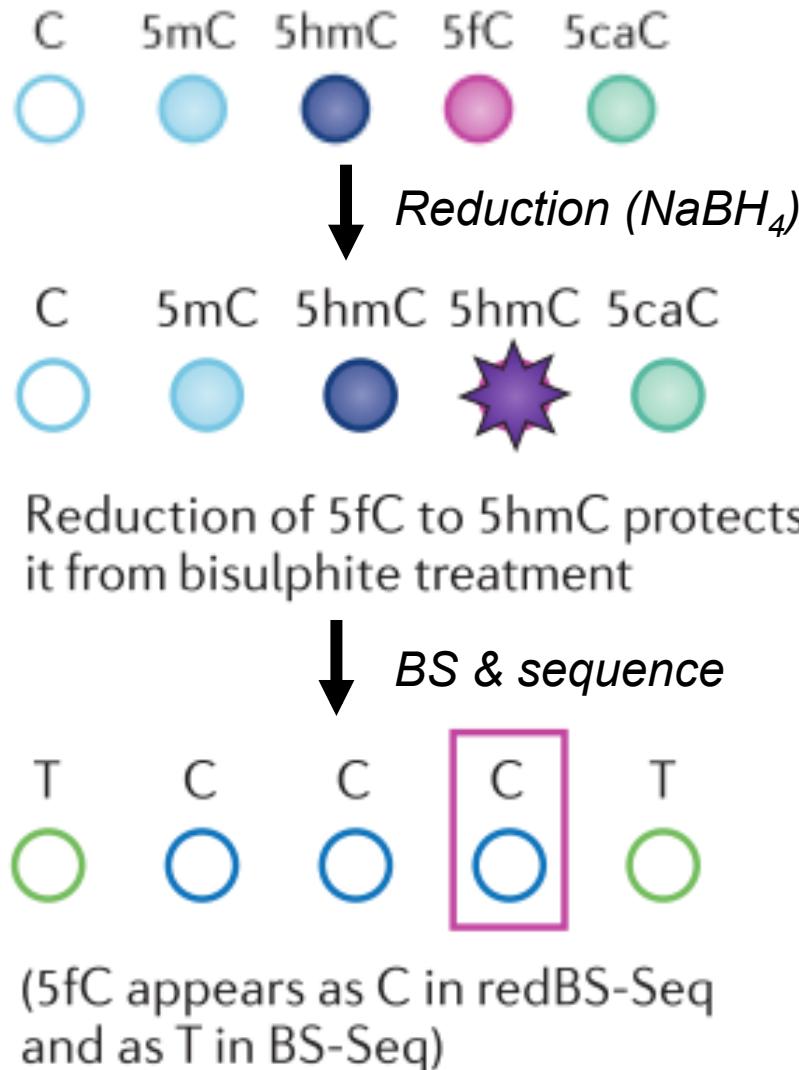
fCAB-seq = 5fC Chemical modification-assisted bisulfite sequencing

$EtONH_2$  = O-ethylhydroxylamine

BS = Bisulfite treatment

Figures adapted from Plongthonkum et al., *Nature Reviews Genetics* (2014) | Song et al., *Cell* (2013)

# 5-fC detection with redBS-Seq



- **Pros**

- Single-base resolution

- **Cons**

- \$\$\$
- Bisulfite treatment can damage DNA
- Dependent on bisulfite conversion
- Requires a large amount of DNA
- Reduced complexity of sequencing reads creates alignment challenges

redBS-seq = Reduced bisulfite sequencing

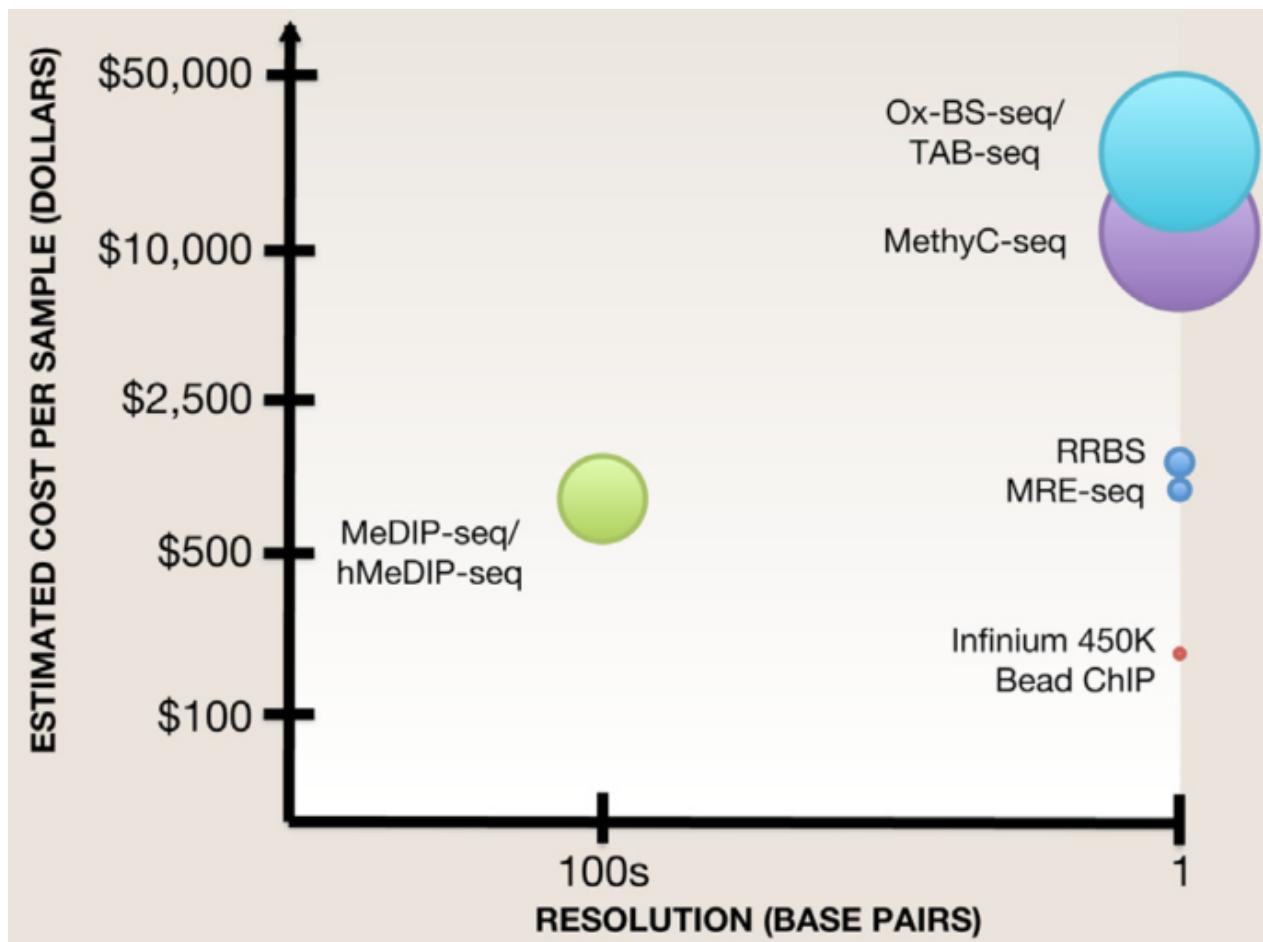
$\text{NaBH}_4$  = sodium borohydride

BS = Bisulfite treatment

Figures adapted from Plongthongkum et al., *Nature Reviews Genetics* (2014) | Booth et al., *Nature Chemistry* (2014)

# Comparison of DNA 5-mC and 5-hmC mapping technologies

Major considerations when selecting a profiling method are: 1) cost, 2) resolution, & 3) genome coverage



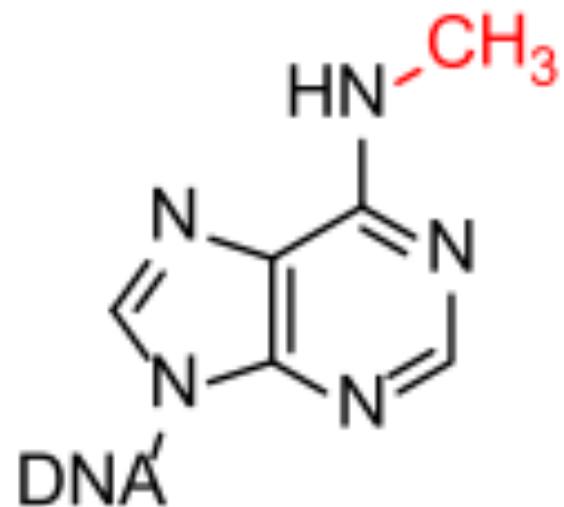
# Sources of bias in 5-mC and 5-hmC mapping technologies

	Method	Amplification	Effect of copy number variation on quantification	Incomplete treatment (chemical or enzyme digestion)	Background signals	Batch to batch variation	Cross hybridization	Sensitivity to sequence context
5mC assays	BS-seq	Low*	None	Low	None	None	None	Medium <sup>†</sup> (RRBS only)
	Low input BS-seq	High <sup>§</sup>	None	Low	None	None	None	Medium (RRBS only)
	BS-based arrays	Low	None	Low	None	Medium	Medium	None
	Non-BS-seq	Low	Medium	Low	High	None	None	Medium (MeDIP/MBD)
	Low input non-BS-seq	High	Medium	None	High	None	None	Medium (MeDIP/MBD)
	Non-BS-arrays	Low	Medium	None	High	Medium	Medium	None
5mC oxidation derivatives assays	BS-based seq	Low	None	Low	None	None	None	None
	Non-BS-based seq	Low	Medium	Low	High	None	None	Medium
	RE-based 5hmC seq	Low	None	Low	None	None	None	Medium
All	Third generation sequencing	Low	None	None	None	None	None	None

\*Low, can be ignored in most cases. †Medium, should be considered when performing analyses. § High, must be considered when performing analyses.

Plongthongkum et al., *Nature Reviews Genetics* (2014)

# N<sup>6</sup>-methyladenine



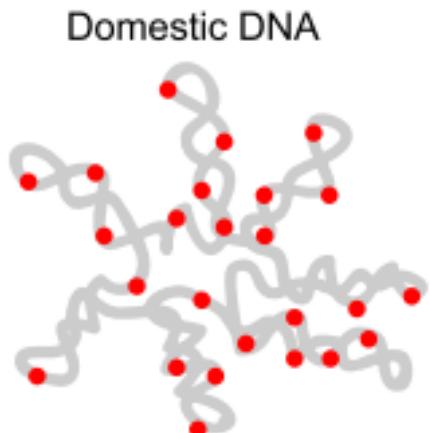
# N<sup>6</sup>-methyladenine

## Bacteria

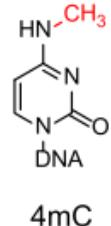
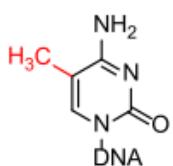
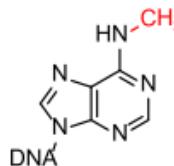
Foreign DNA



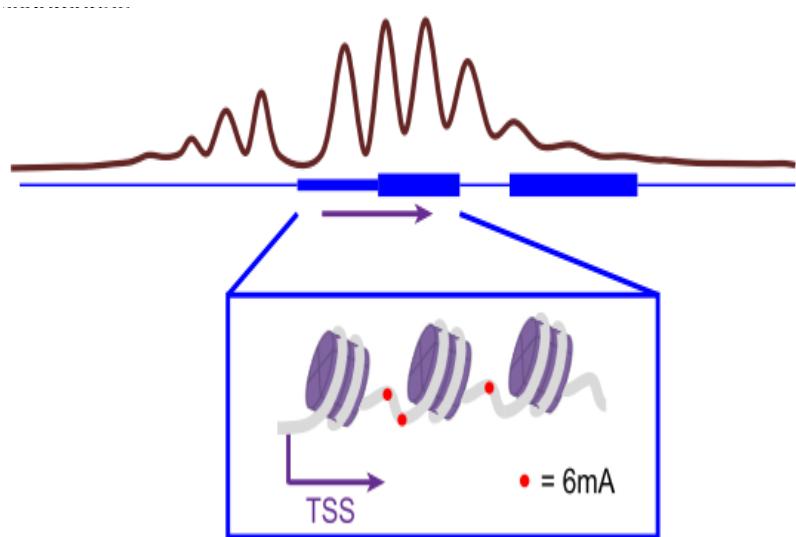
Restriction and  
modification system;  
Package of  
phage DNA



• = 6mA or 5mC or 4mC



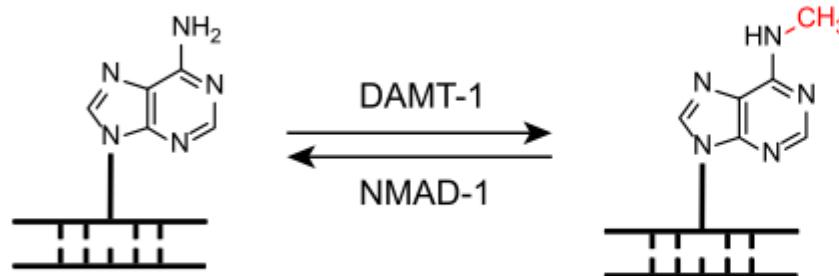
## C. reinhardtii



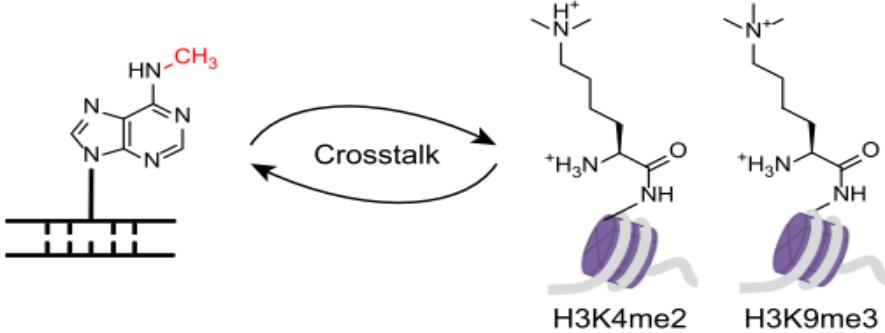
- 6-mA depleted at TSS
- 6-mA enriched at linker regions between nucleosomes
- 6-mA marks actively transcribed genes

# N<sup>6</sup>-methyladenine (cont.)

## C. elegans

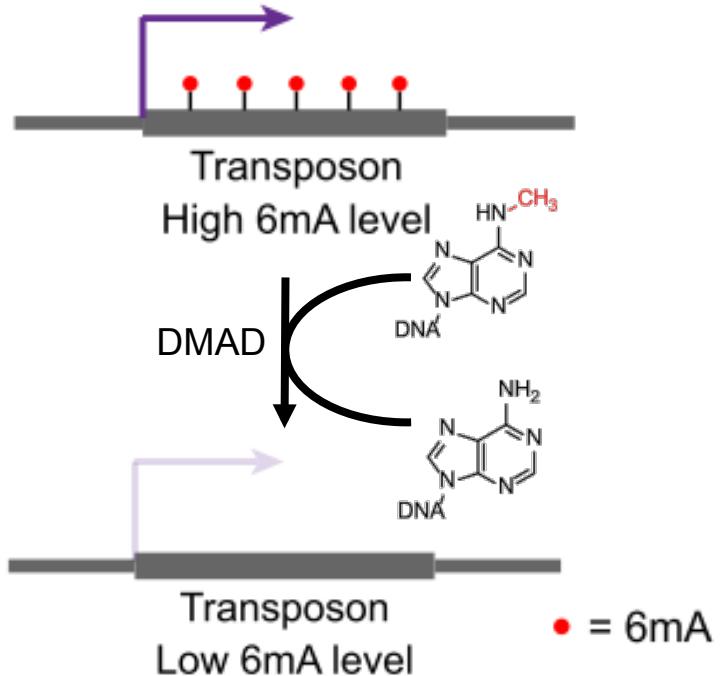


- DAMT-1 is a 6-mA writer
- NMAD-1 is a 6-mA eraser



- 6-mA may play a role regulating gene expression
- *C. elegans* lacks 5-mC

## D. melanogaster

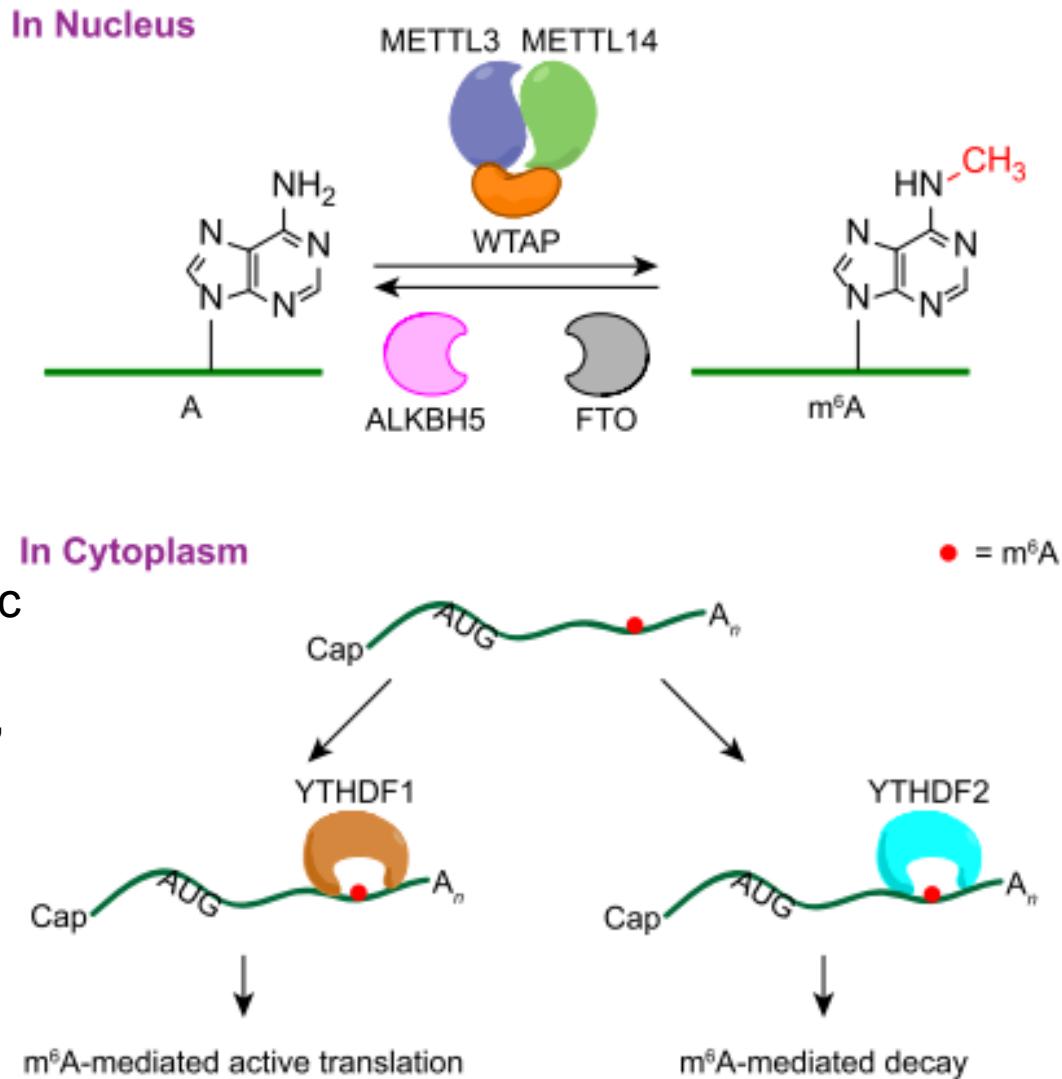


- DMAD, a Tet homolog, may be an eraser of 6mA
- Dynamic 6-mA demethylation associated with transposon expression
- 6-mA plays a critical roles in development

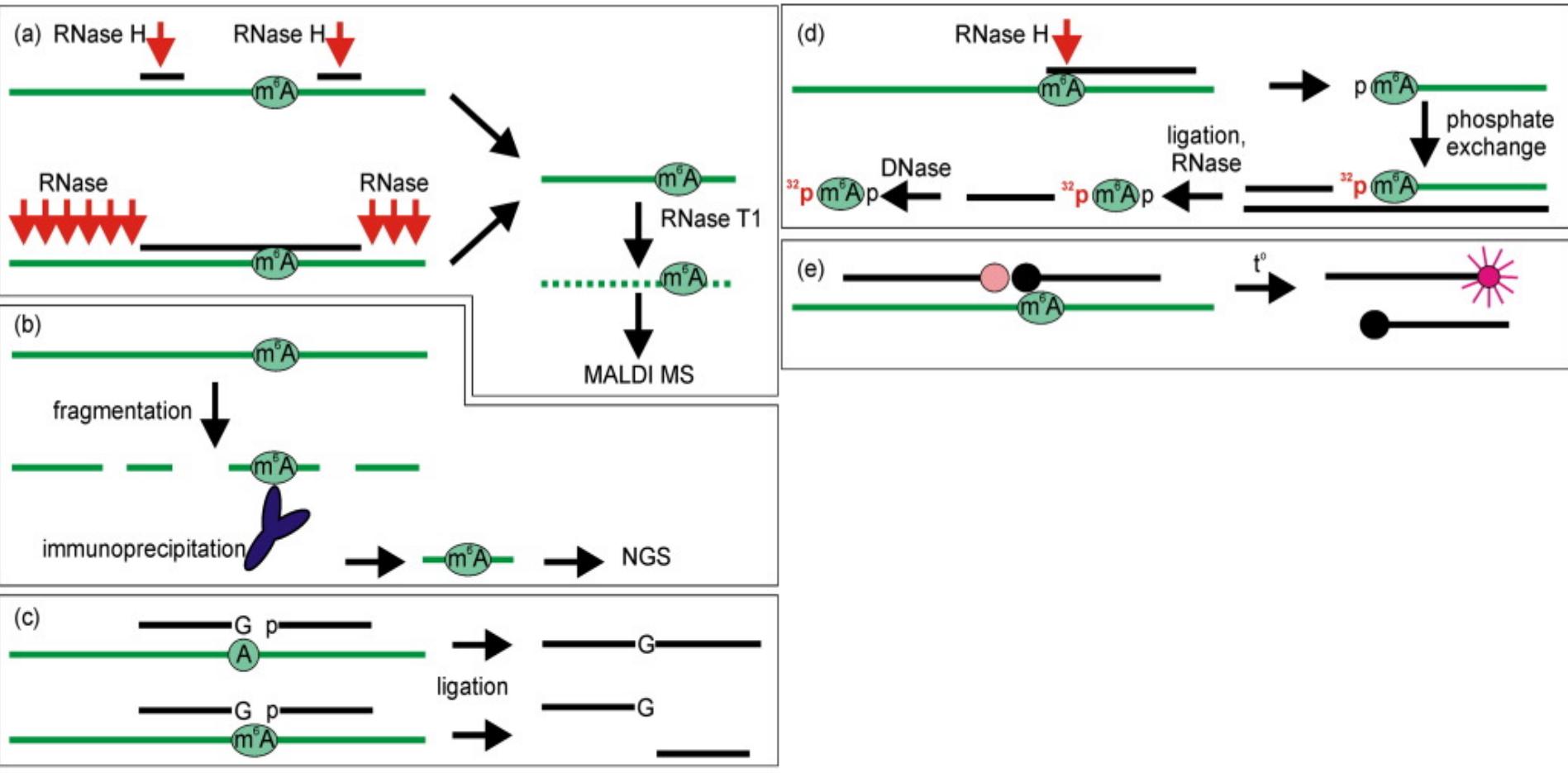
# **RNA modifications**

# N<sup>6</sup>-methyladenosine (m<sup>6</sup>A) in mRNA

- Discovered in eukaryotic and viral mRNAs in 1974
- Most prevalent modification in mRNAs and lncRNAs in higher eukaryotes
- m<sup>6</sup>A plays a critical role in development & disease
  - Deficiencies in methyltransferases is embryonic lethal in mouse and leads to developmental defects in yeast, flies, zebrafish & plants
  - m<sup>6</sup>A dysregulation associated with obesity, cancer, and other diseases in humans

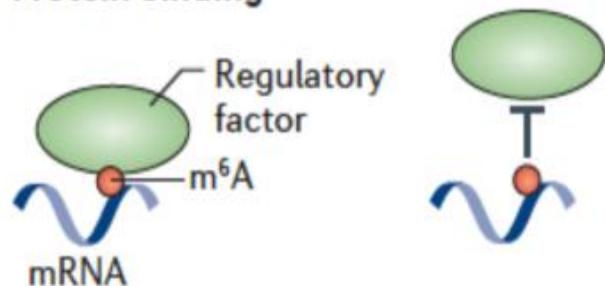


# Methods for detecting $\text{m}^6\text{A}$ in RNA

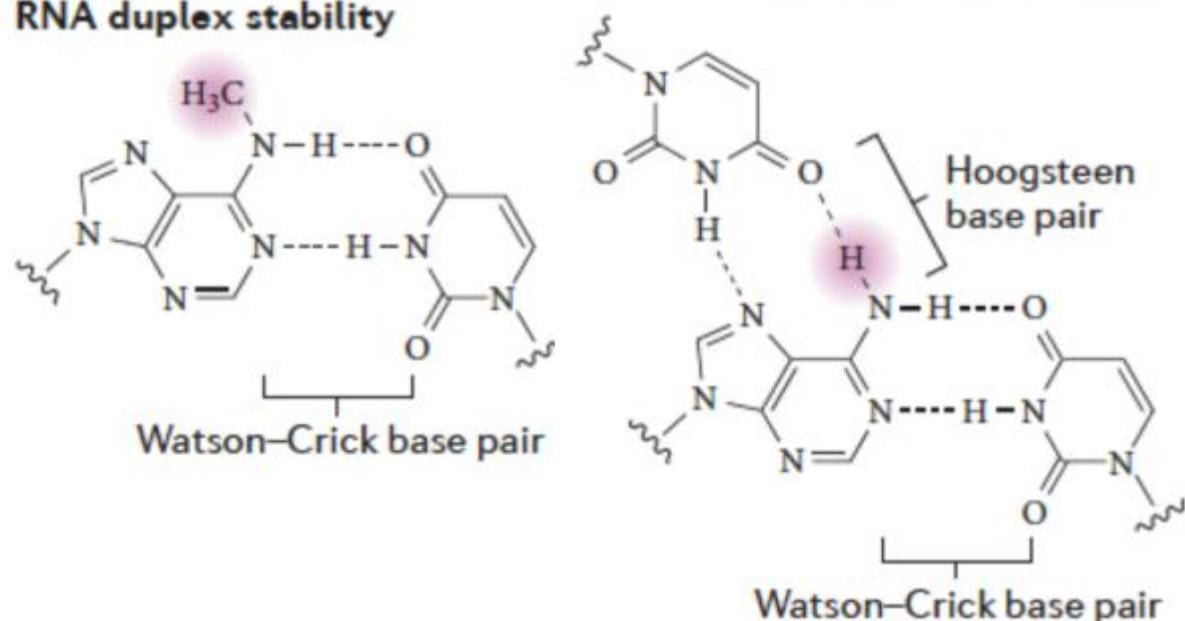


# Mechanisms and functions of m<sup>6</sup>A

## a Protein binding



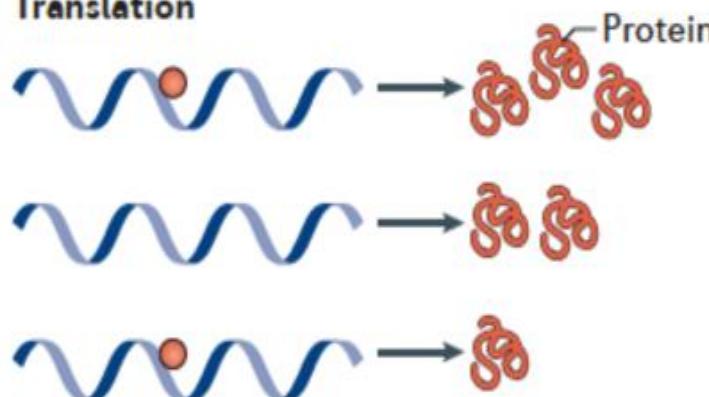
## b RNA duplex stability



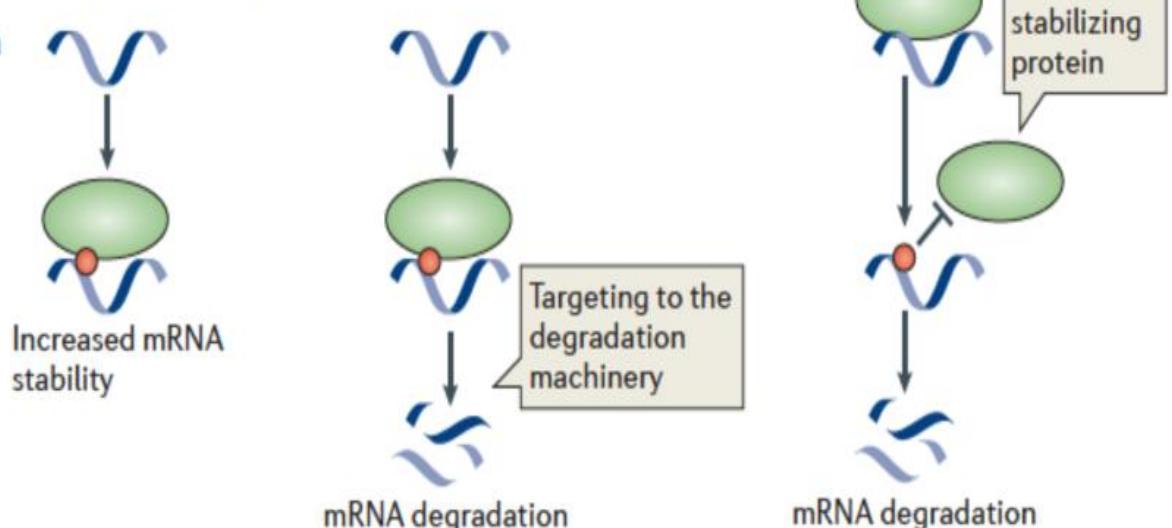
## c Splicing



## d Translation

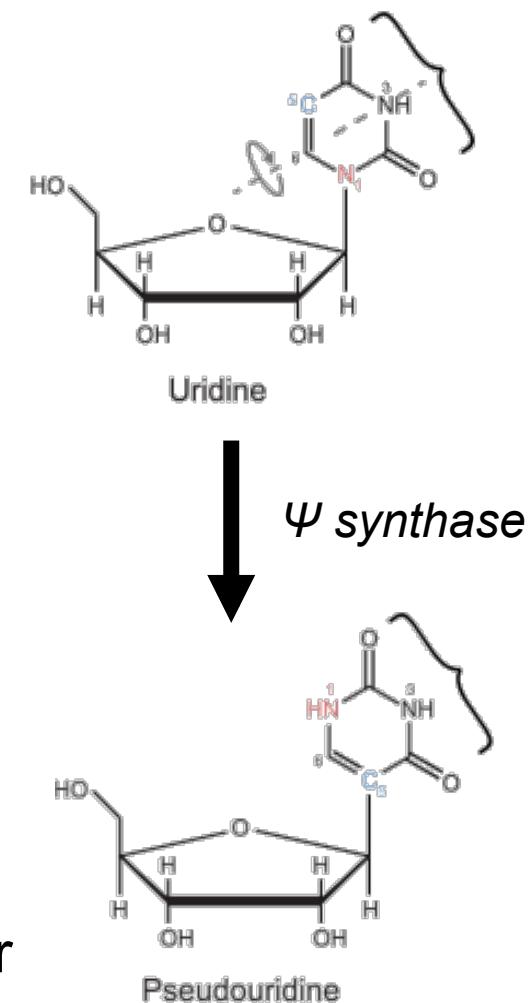


## e Stability and degradation



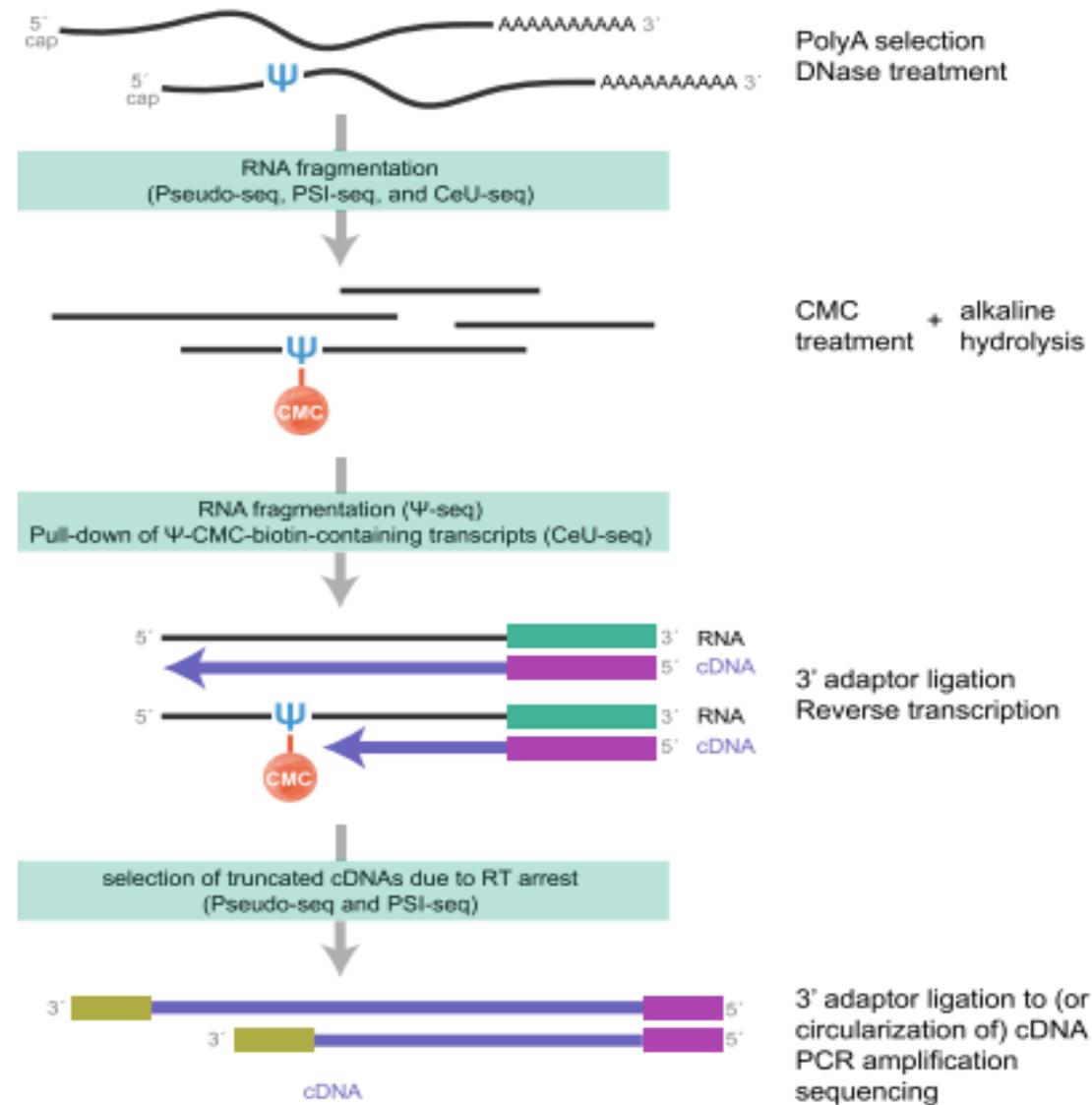
# Other RNA modifications: Pseudouridine ( $\Psi$ )

- Most abundant of over 100 post-transcription modifications
  - Several hundred to thousands of  $\Psi$  sites genome-wide
- Found on rRNAs, snRNAs, mRNAs, & tRNAs
- Functions
  - rRNA: required for ribosome biogenesis and translation fidelity
  - snRNA: alters the structural stability of snRNAs and plays critical roles in snRNP assembly, spliceosome formation, and splicing
  - mRNA: may play a role in response to cellular stress and differentiation



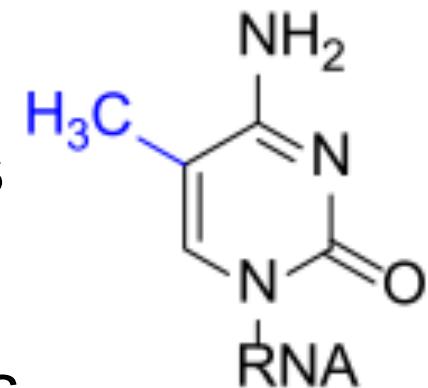
# Methods for detecting $\Psi$

- Methods to map  $\Psi$  modifications recently invented
  - Pseudo-seq (Carlie et al., 2014)
  - $\Psi$ -seq (Schwartz et al., 2014)
  - PSI-seq (Lovejoy et al., 2013)
  - CeU-seq (Li et al., 2015)



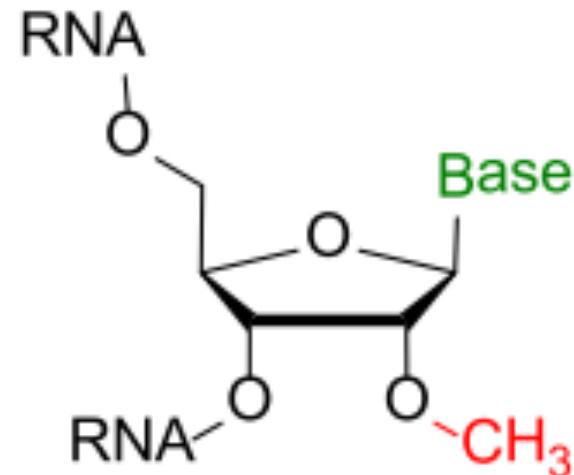
# Other RNA modifications: 5-methylcytidine

- Found on mRNAs, tRNAs, & rRNAs
- Functions
  - tRNAs: increases the stability of tRNAs
- Plays a critical role in development & disease
  - Deficiencies RNA cytosine methyltransferases (RCMT) associated with embryonic lethality in mouse
  - Dysregulation of RCMTs associated with cancer infertility, & mental retardation in humans



# Other RNA modifications: 2'-O-methylation

- Found on miRNAs, snRNAs, tRNAs, & rRNAs
- Functions
  - miRNAs: increases the stability of miRNAs
  - snRNAs: alters the structural stability of snRNAs & plays critical roles in snRNP assembly, spliceosome formation, and splicing



# tRNA modifications

- Modifications in tRNAs affect tRNA stability, protein translation, and play key roles in stress response and immune recognition
- In *S. cerevisiae*, there are ~25 different tRNA modifications
  - Dysregulation leads to lethality, poor growth, & temperature sensitivity
- Methods to map tRNA modifications recently invented
  - ARM-seq (Cozen et al., 2015)
  - DM-tRNA-seq (Zheng et al., 2015)

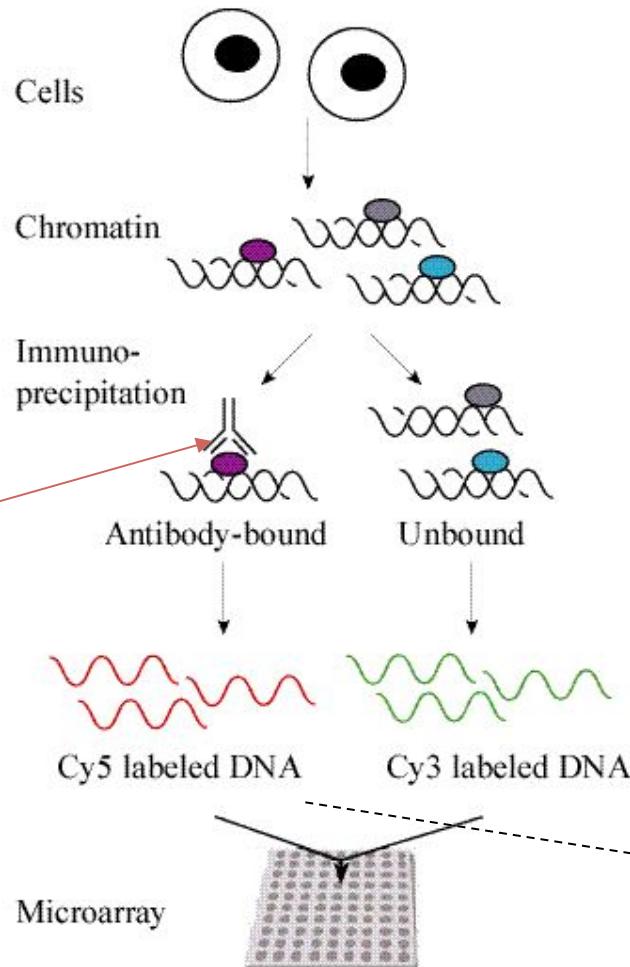
# Summary

- Nucleic acid modifications is an exciting and evolving area of research
- Nucleic acid modifications expand the complexity of gene regulation
- DNA 5-mC is found in almost all kingdoms of life
- Dysregulation of DNA 5-mC is associated with cancer and other human diseases
- DNA 6-mA is a recently discovered epigenetic mark and may play a complementary role to 5-mC
- RNA modifications are diverse and have many functions
- Many technologies have been developed for mapping DNA and RNA modifications

# Technologies for Interrogating Chromatin States

ChIP-chip

Antibody specific  
to one type of  
histone  
modification

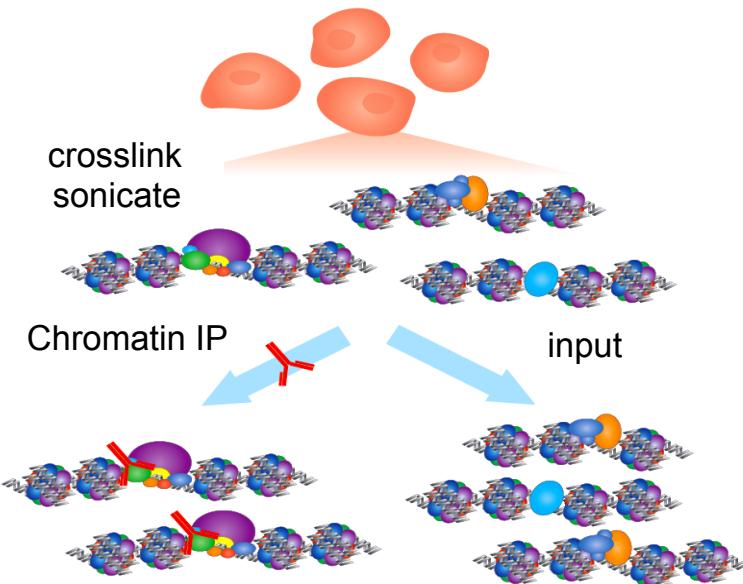


Histone Modifications

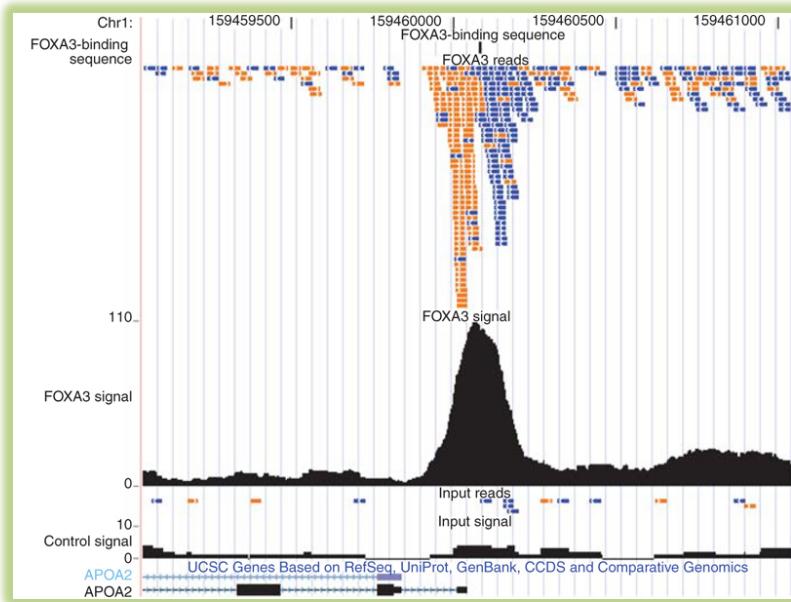
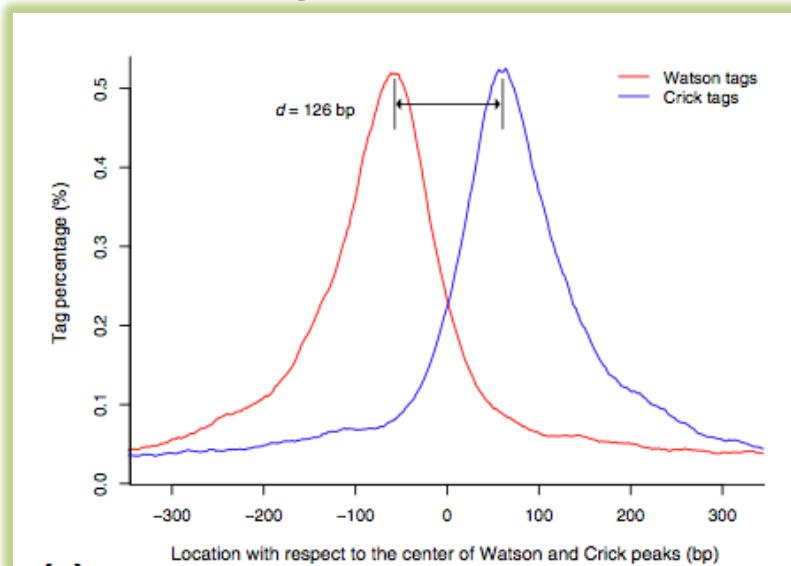
ChIP-seq

Deep sequencing

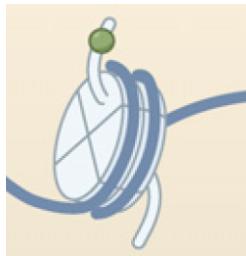
# ChIP-seq



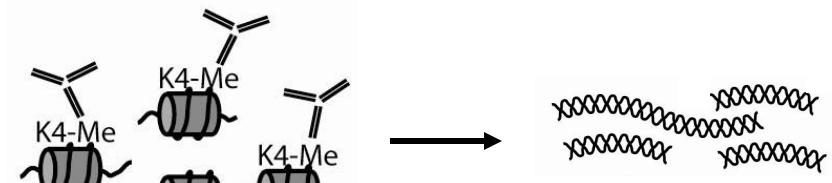
Zhang et al. Genome Biol 2008



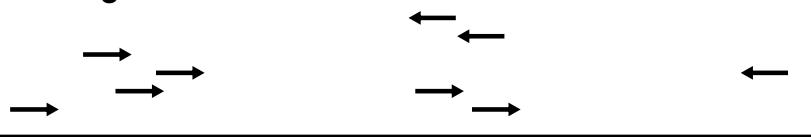
# Chromatin-IP Sequencing



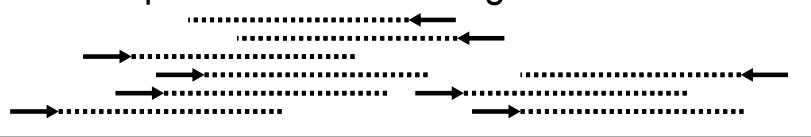
K4me1  
K4me2  
K4me3 } "active"  
K27me3 } "repressive"



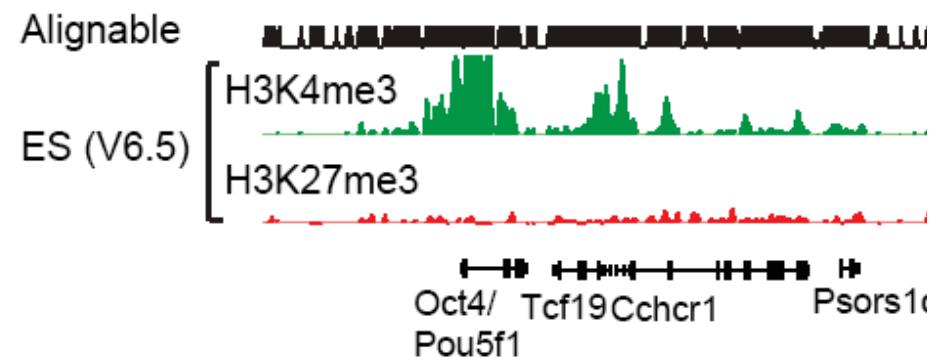
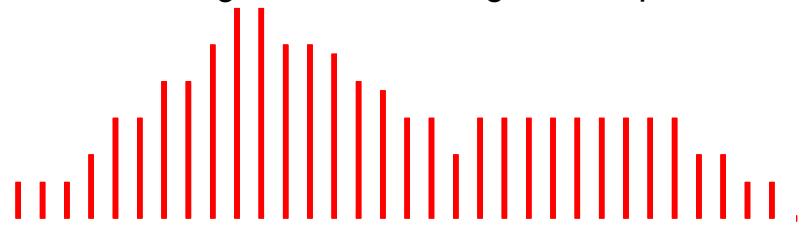
A. Align reads



B. Infer positions of ChIP fragments

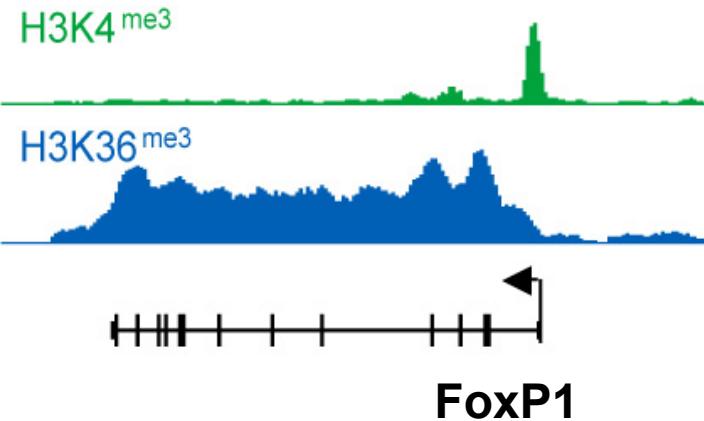


C. Count fragments at each genomic position

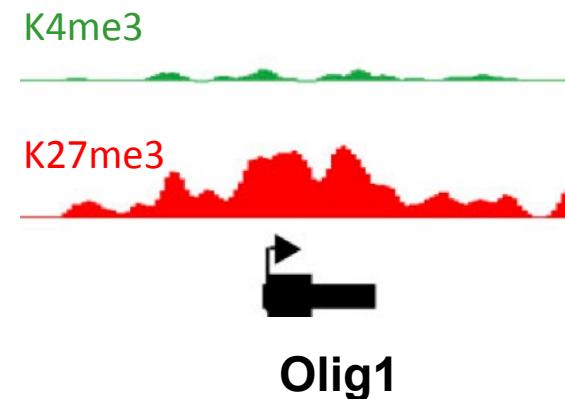


# Histone methylation and transcriptional state

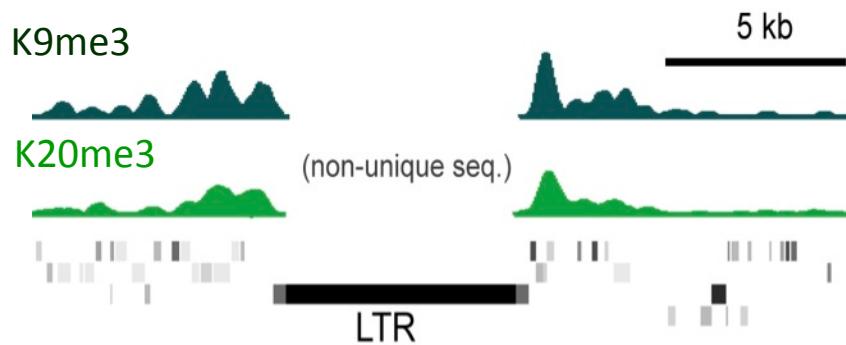
Transcribed gene



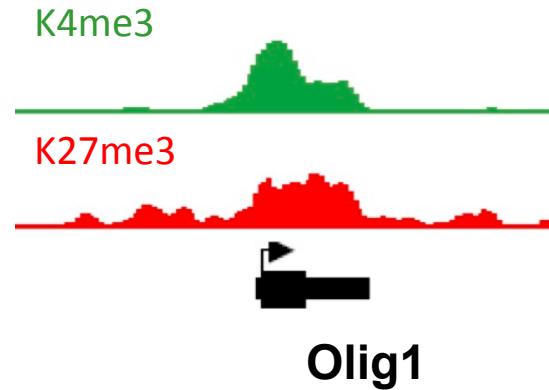
Silent developmental gene



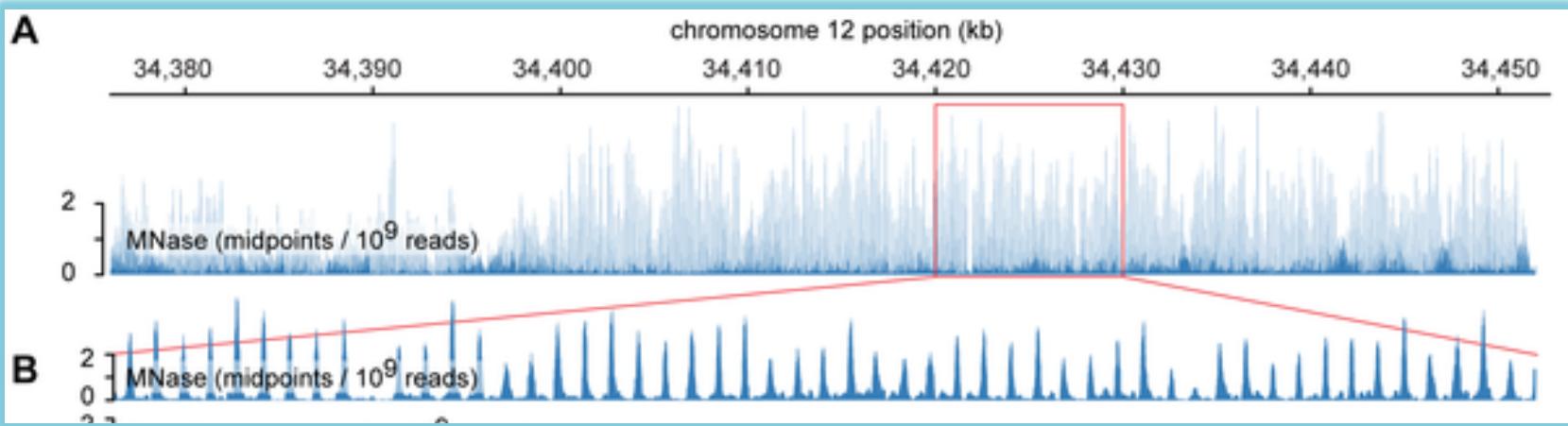
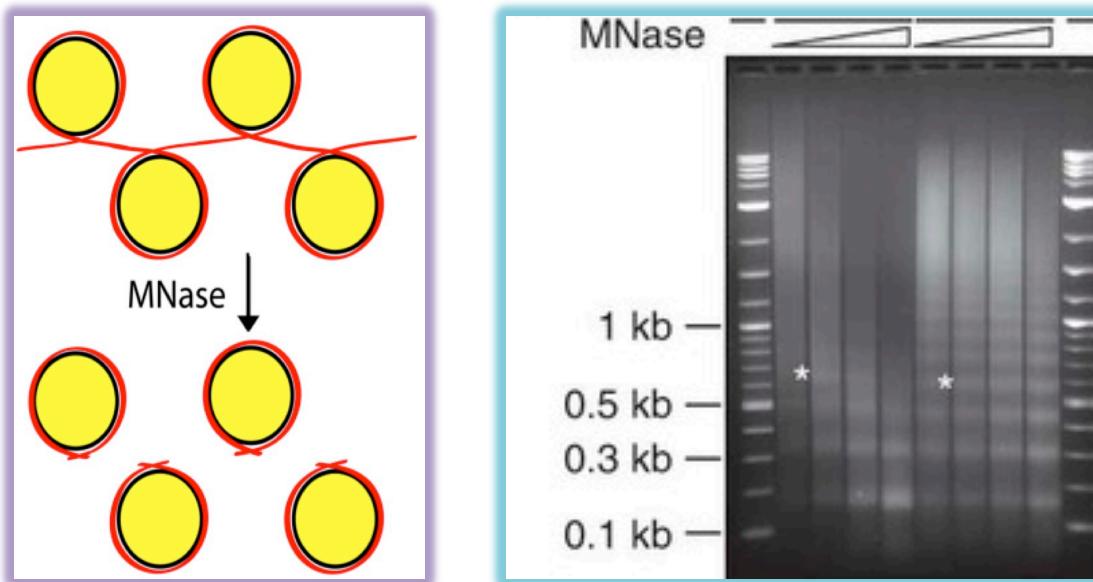
Constitutive heterochromatin



'Poised' developmental gene



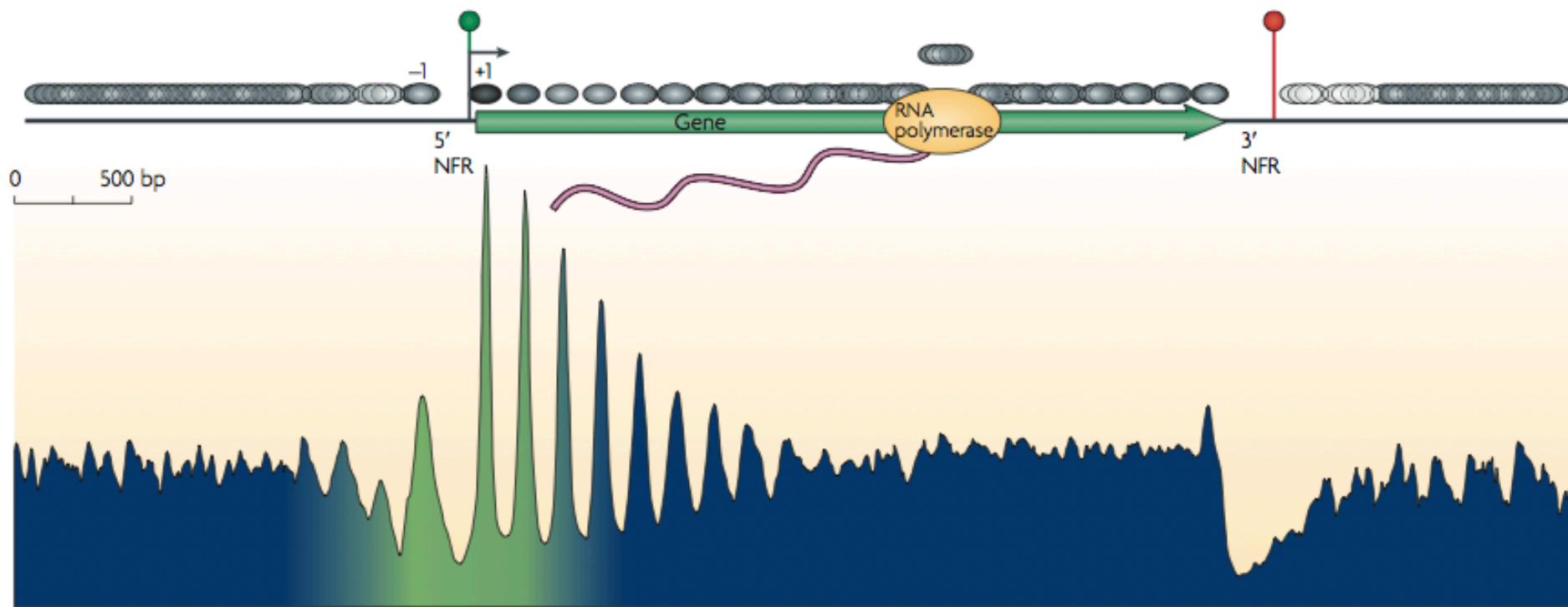
# Mapping nucleosome positions



Gaffney DJ, McVicker G, Pai AA, Fondufe-Mittendorf YN, et al. (2012) Controls of Nucleosome Positioning in the Human Genome. PLoS Genet 8(11): e1003036. doi:10.1371/journal.pgen.1003036

<http://www.plosgenetics.org/article/info:doi/10.1371/journal.pgen.1003036>

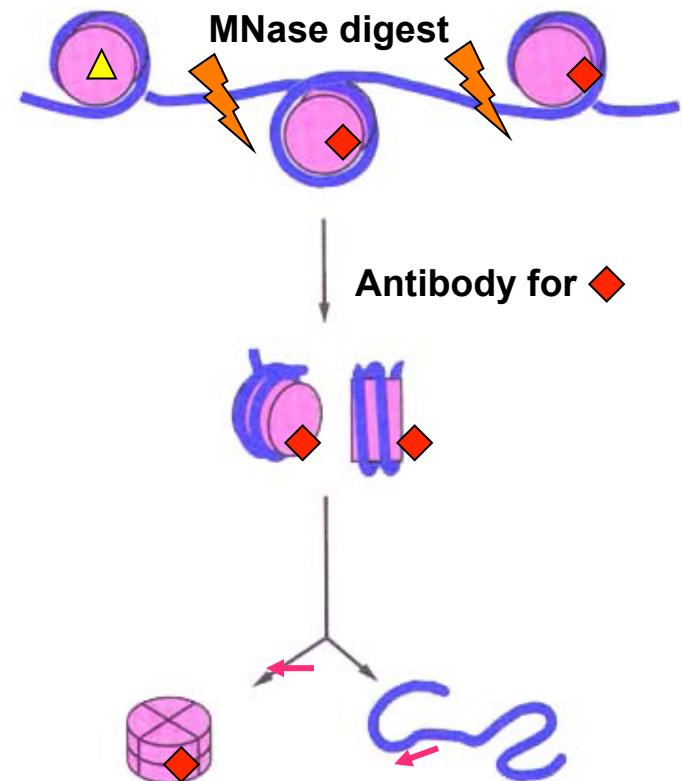
# Genomic distribution of nucleosomes



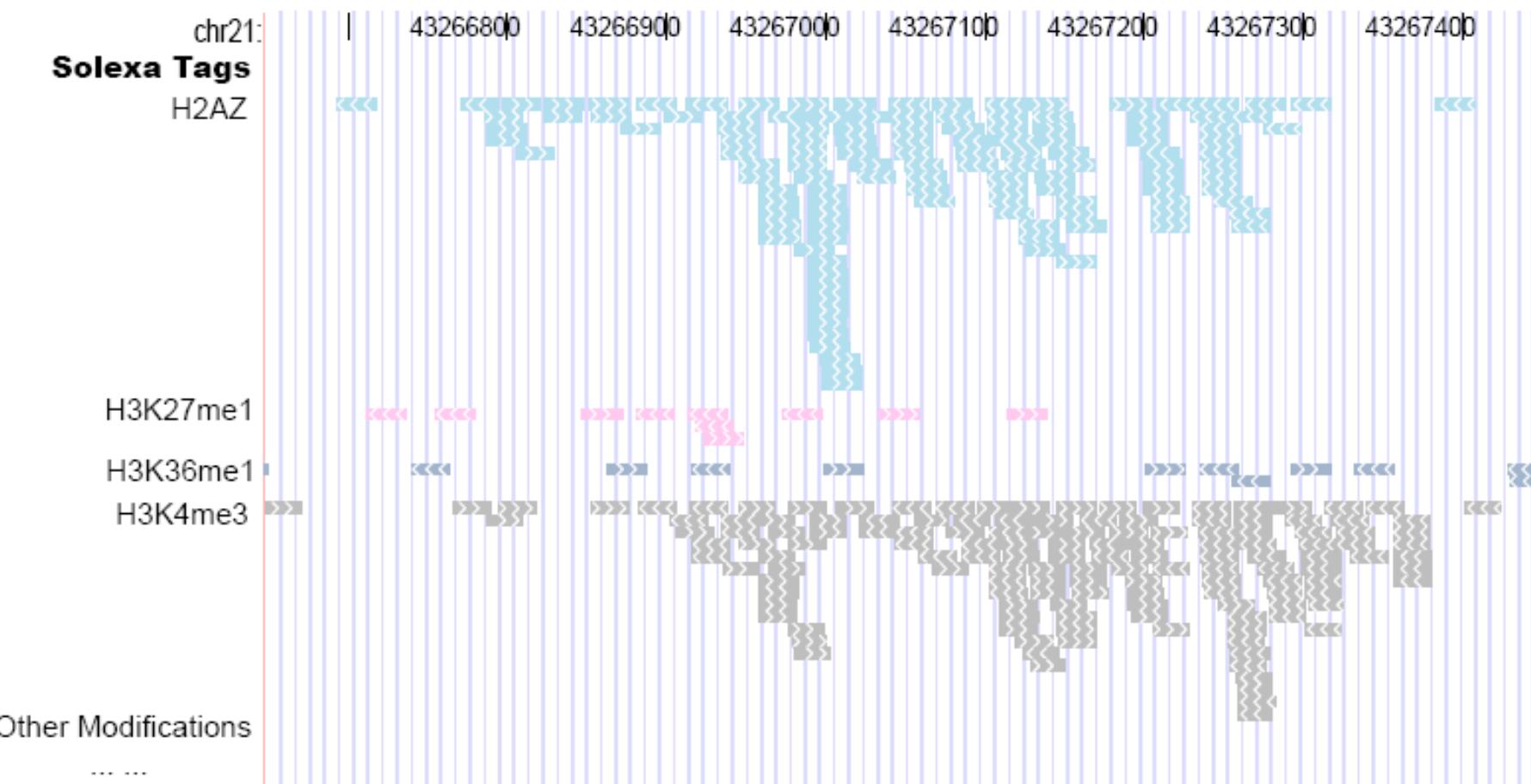
- The presence of NFRs demonstrated that open promoter states are stable and common, even at genes that are transcribed so infrequently

# Nucleosome Positioning from Histone ChIP-seq

- Barski et al, Cell 2007
  - Nucleosome resolution ChIP-seq of 21 histone marks in CD4<sup>+</sup> T-cells
  - Total 185.7 M 25 nt tags sequenced
  - Analysis not at nucleosome resolution to map nucleosomes at specific regions

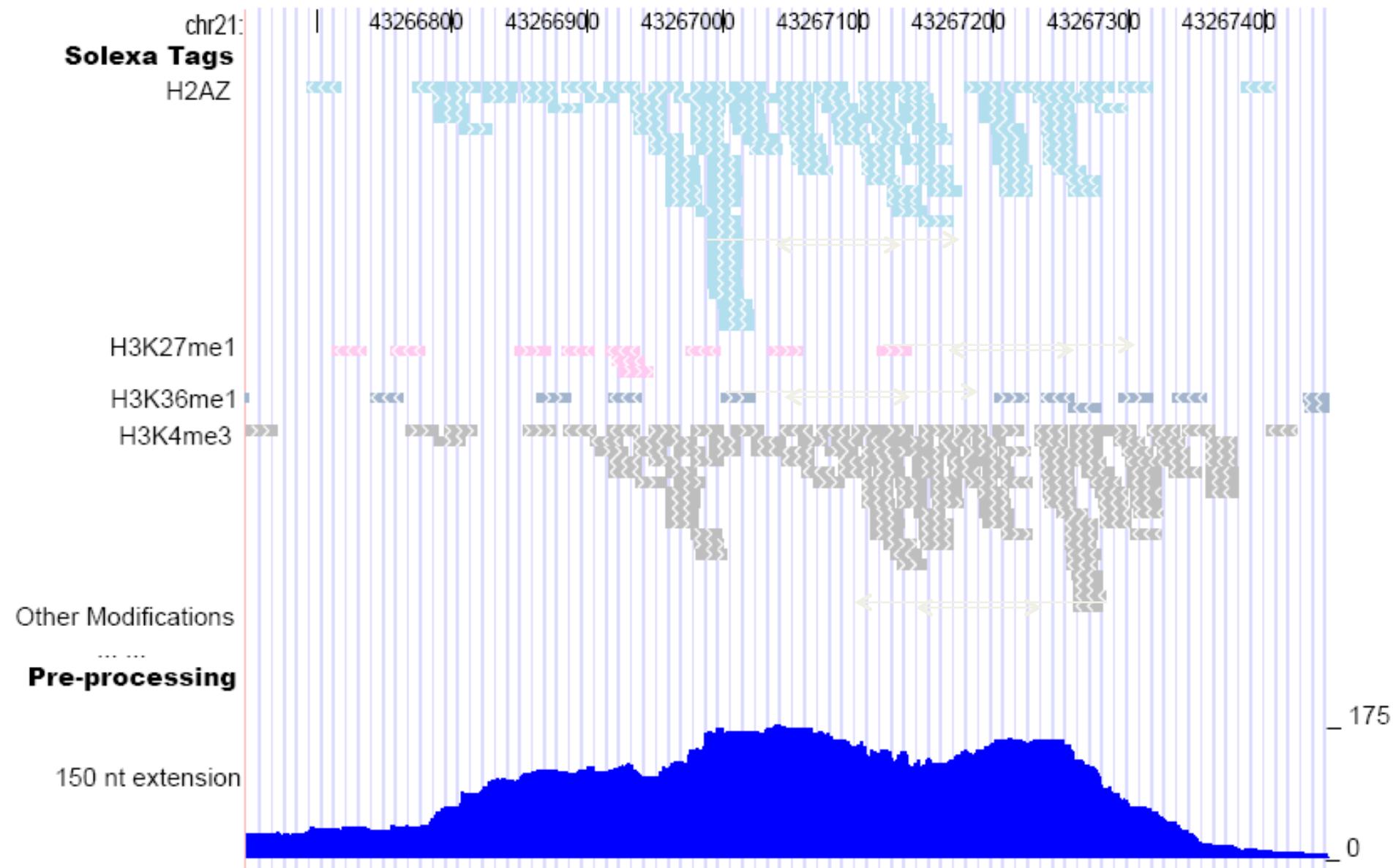


# Combine Tags From All ChIP-Seq



# Extend Tags 3' to 150 nt

## Check Tag Count Across Genome



# Take the middle 75 nt

