

A Brief Overview of Epigenomics

BIOTRAC Lecture

Sean, Soonweng, Cho, PhD

Kennedy Krieger Institute

and

Johns Hopkins University School of Medicine



sean.cho@jhmi.edu



github.com/sean-cho

Outline

What is epigenomics?

Why epigenomics?

How

- DNA methylation
- Histone modifications
- Chromatin features
 - Accessibility
 - HiC
- Integrative epigenomics
 - Imputation
 - Chromatin state
- Integrative analyses
- Study Design
 - Limitations and considerations
 - Balancing and Blocking

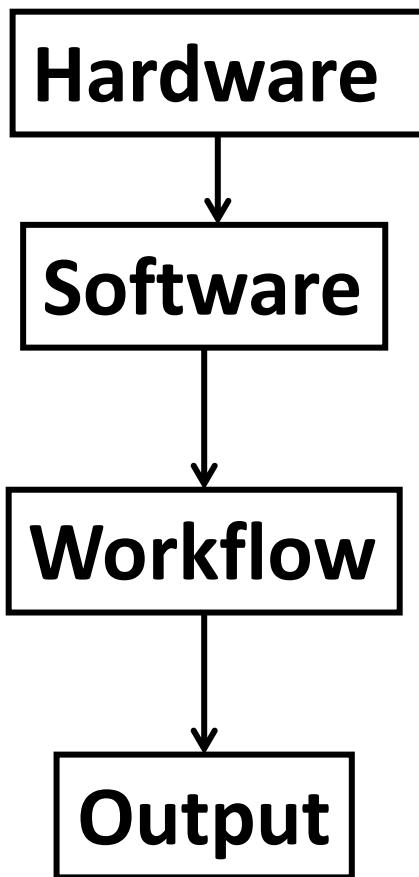
The future of epigenomics

- Nanopore/MinION sequencing
- Single cell epigenomics

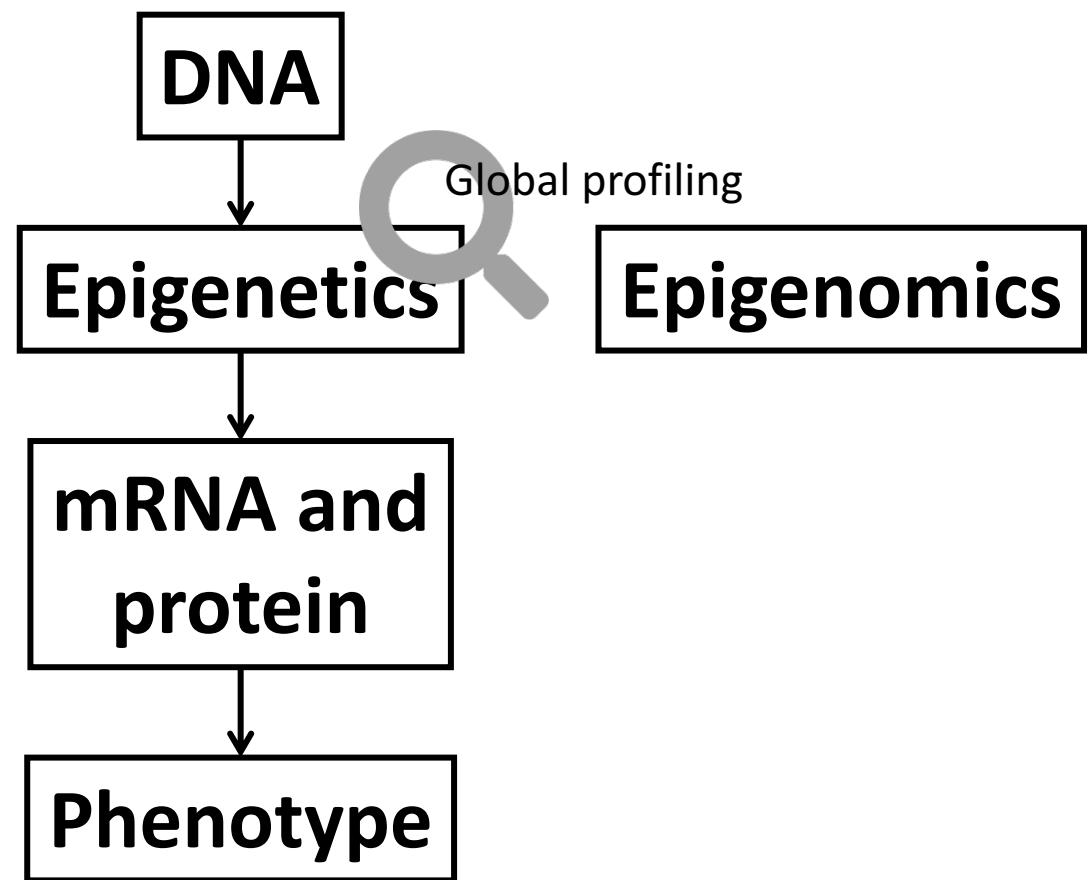
What is epigenetics?

Epigenetics

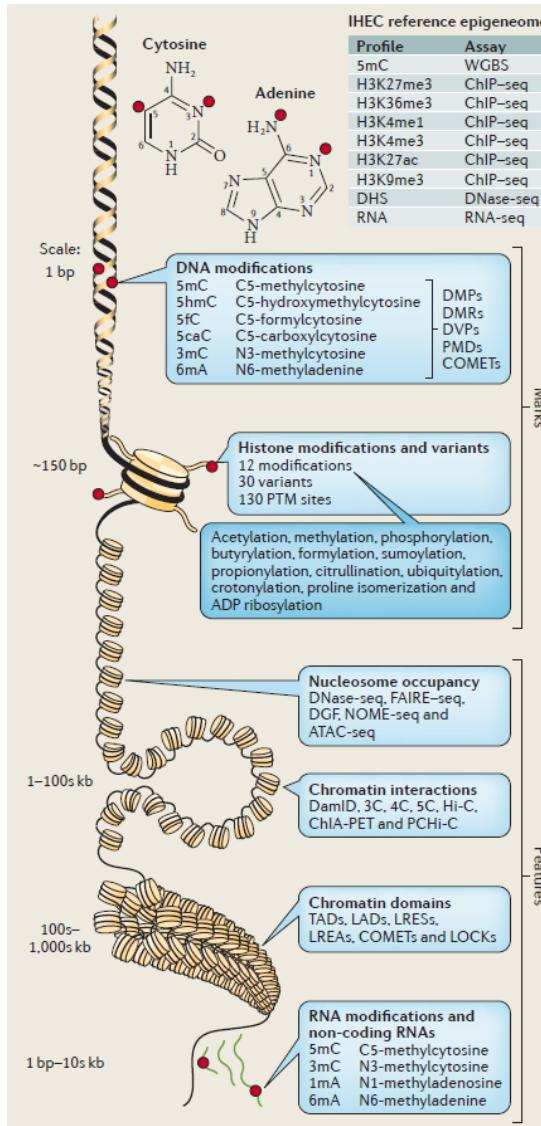
Computer
Science



Biology



Epigenetics → epigenomics



DNA methylation

WGBS

Histone modifications

H3K4me3

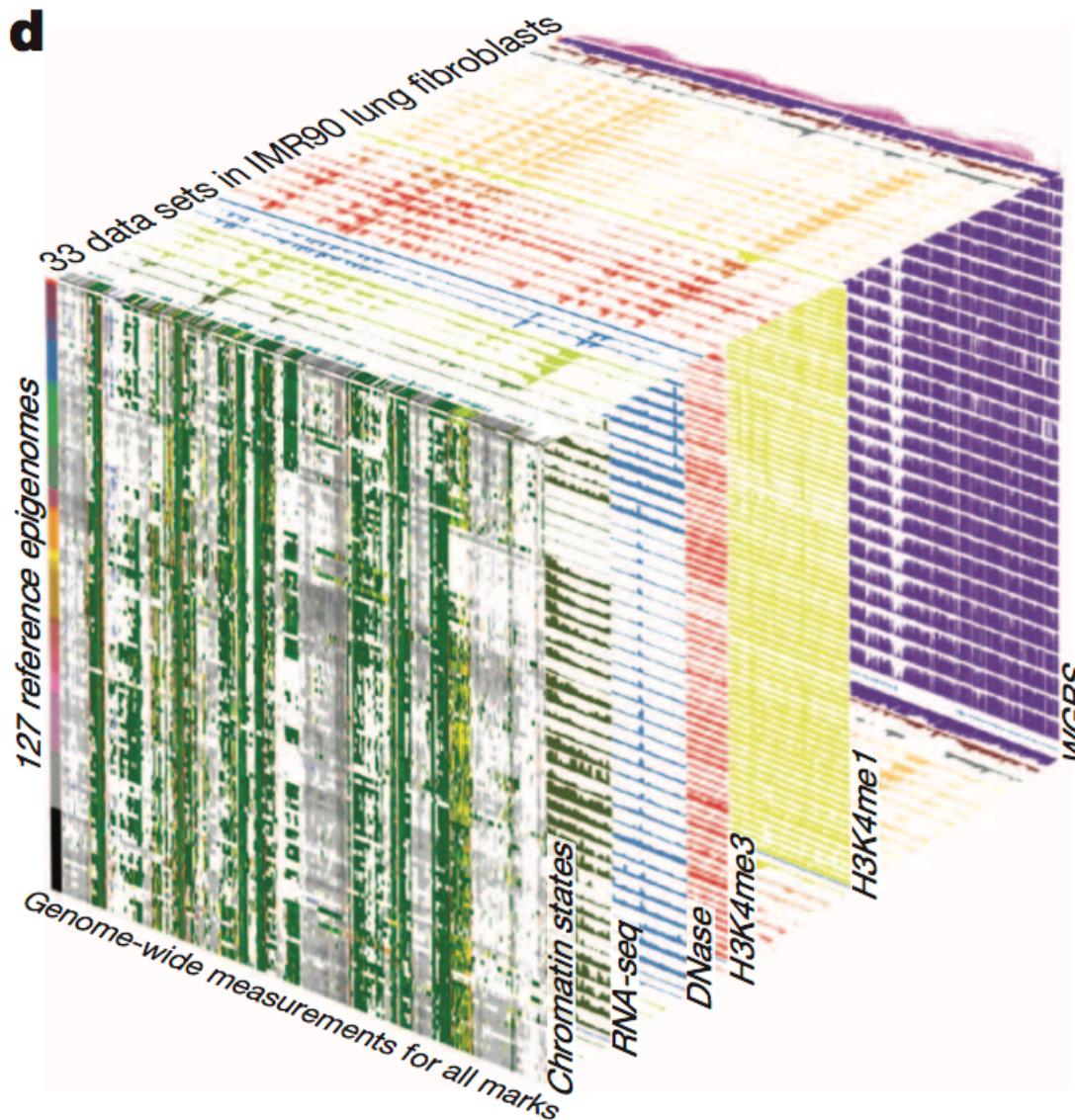
Chromatin features

DNase

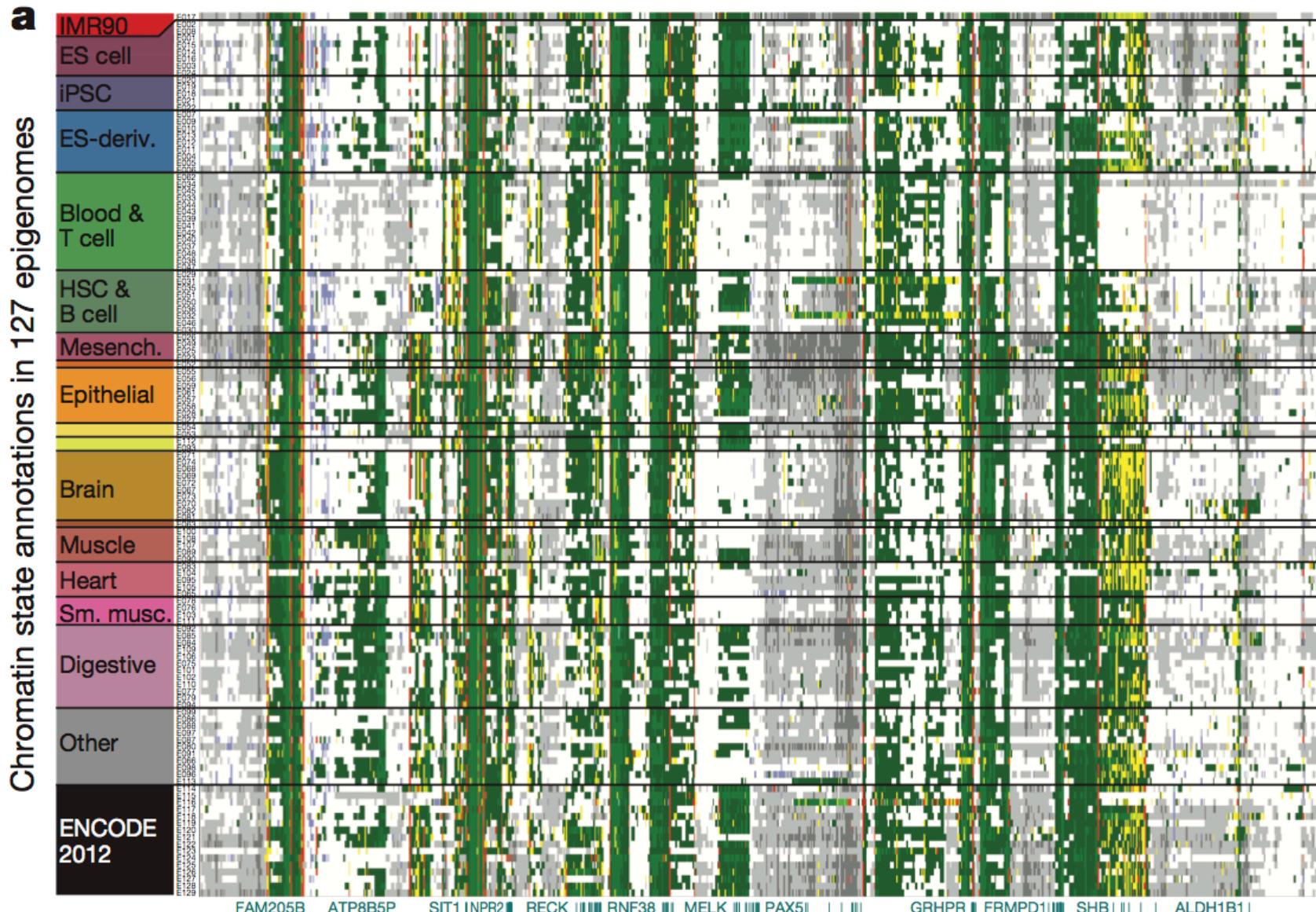
Non-coding RNA

RNA-seq

Integrated into cell specific epigenomes
that govern the function, state, and fate of cells



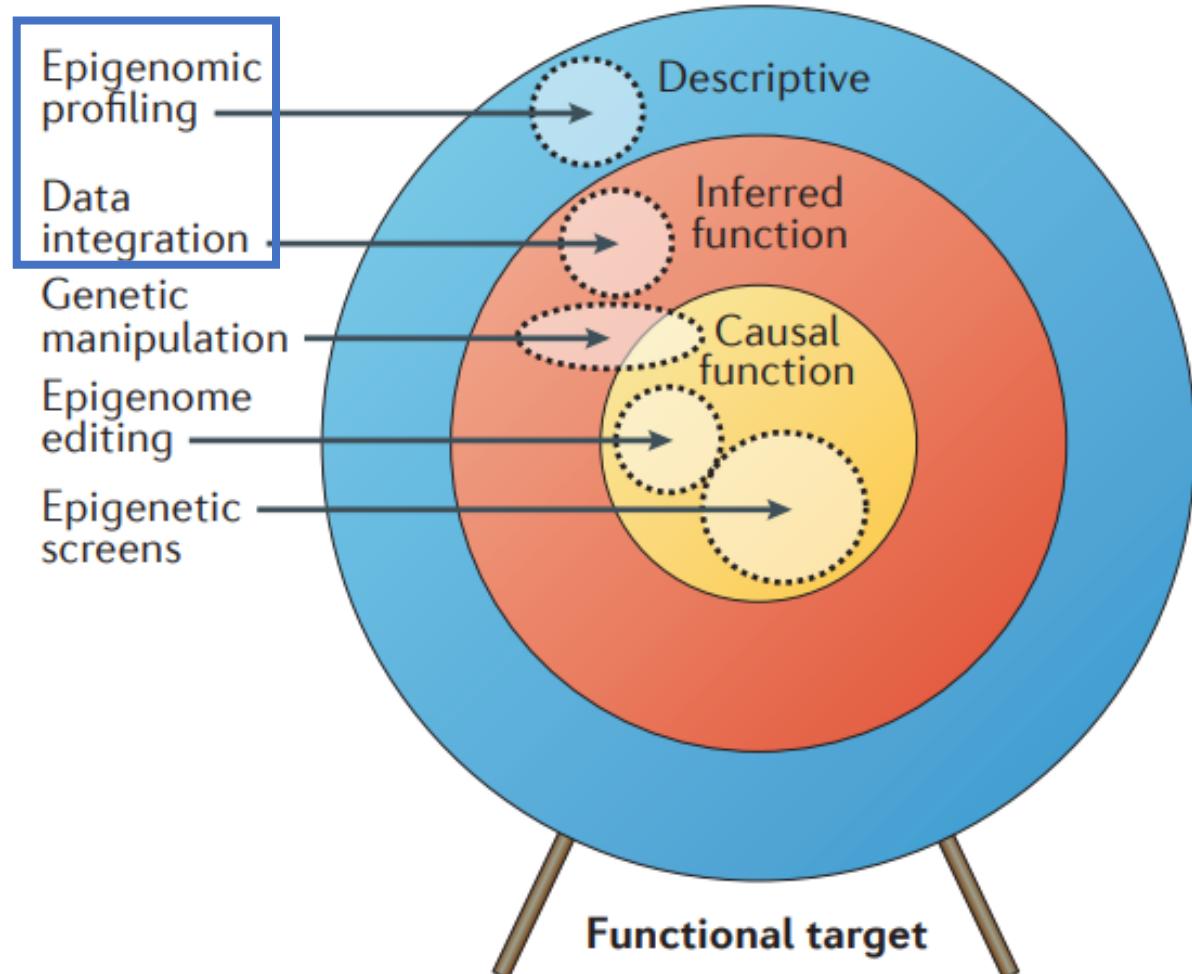
Roadmap Epigenomics Consortium



Why epigenomics?

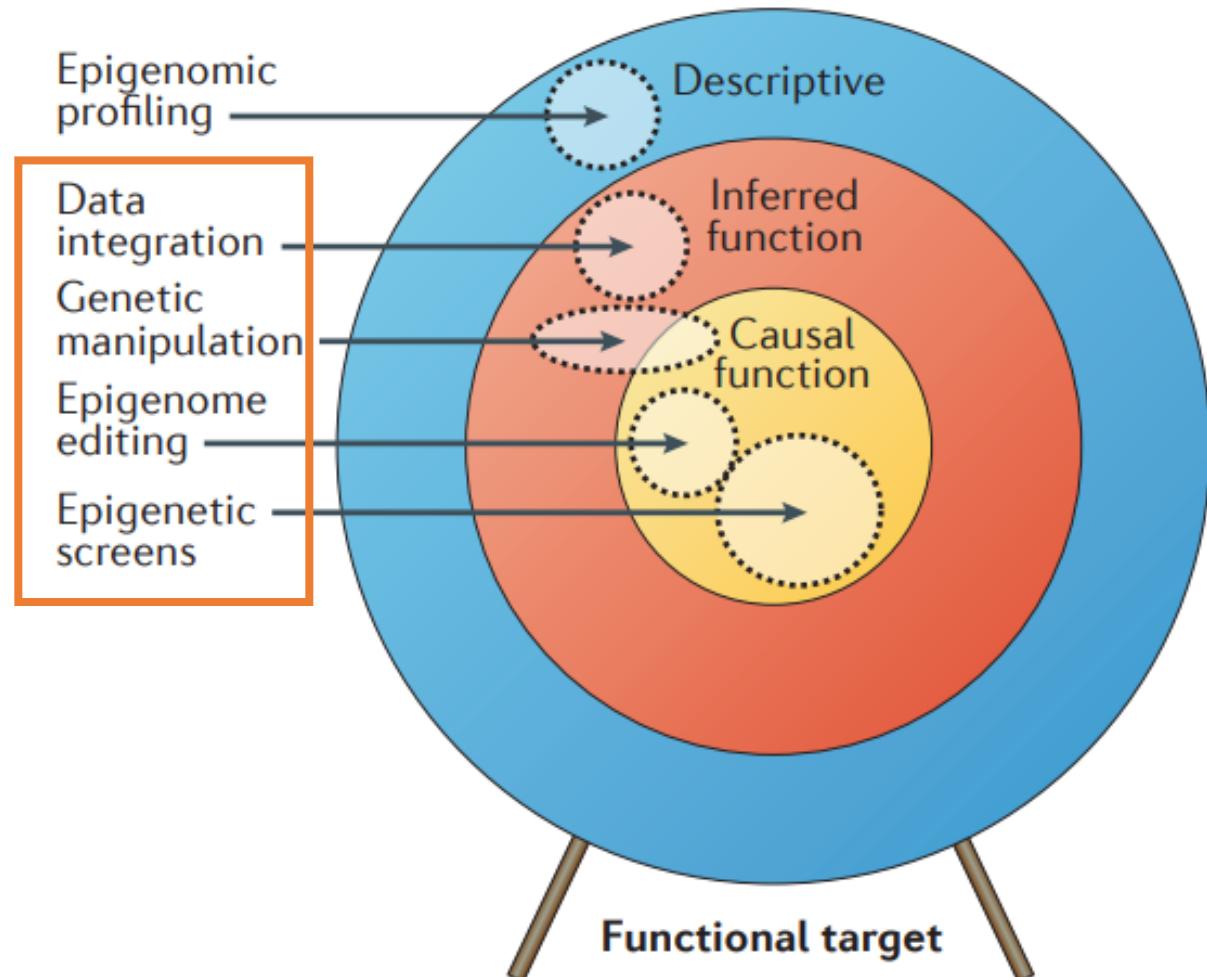
Why epigenomics?

Hypothesis generating



Why epigenomics?

Hypothesis testing



Case 01: Identification of differentially methylated loci and biologically relevant clusters in breast cancer

EMBO
Molecular Medicine

Research Article
Epigenetic portraits of human breast cancers

DNA methylation profiling reveals
a predominant immune component
in breast cancers

Sarah Dedeurwaerder^{1†}, Christine Desmedt^{2†}, Emilie Calonne¹, Sandeep K. Singhal²,
Benjamin Haibe-Kains^{2,3}, Matthieu Defrance¹, Stefan Michiels², Michael Volkmar¹, Rachel Deplus¹,
Judith Luciani¹, Françoise Lallemand², Denis Larsimont⁴, Jérôme Toussaint², Sandy Haussy²,
Françoise Rothé², Ghizlane Rouas², Otto Metzger², Samira Majjaj², Kamal Saini², Pascale Putmans¹,
Gérald Hames⁵, Nicolas van Baren⁶, Pierre G. Coulie⁵, Martine Piccart⁷,
Christos Sotiriou^{2***,†}, François Fuks^{1*,†}

Overview

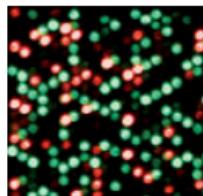
MAIN SET OF PATIENTS:
123 breast tissues



Discovery cohort

VALIDATION SET OF PATIENTS:
125 breast tissues

**Independent
validation cohort**



**14, 475 genes
(27,578 CpGs)
investigated per sample**

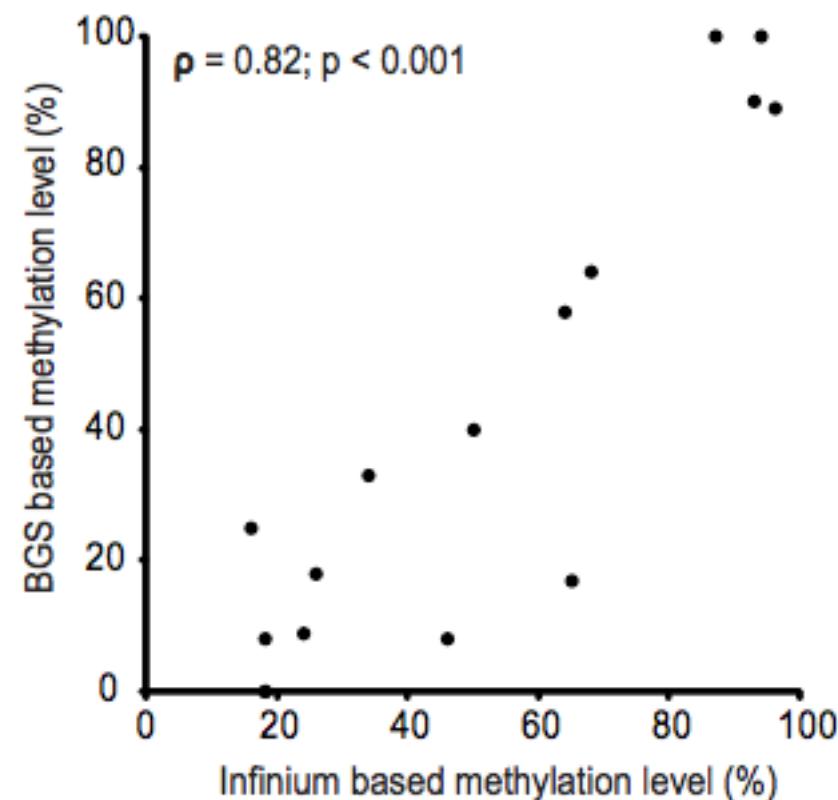
Important, especially for classification and biomarkers studies to validate that discovery is applicable to other general datasets.

Identification of differentially methylated loci

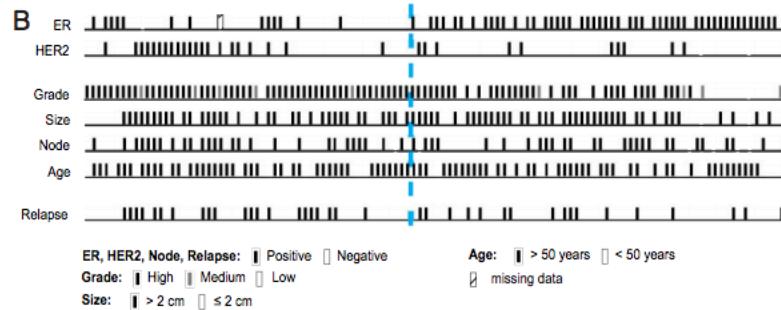
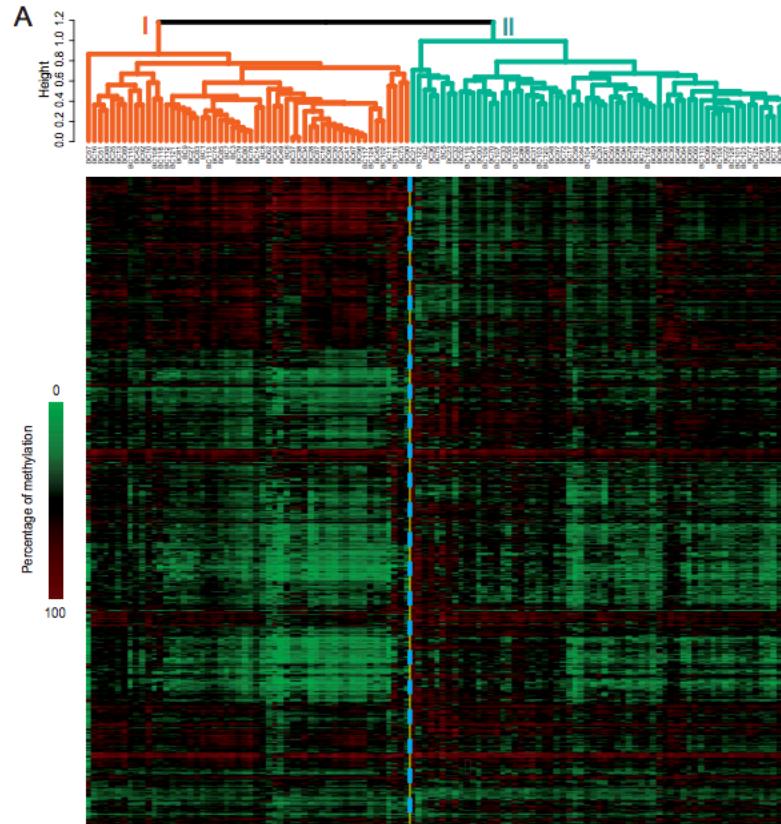
Gene	Illumina ID	Infinium methylation (frequency, %)	Reported methylation data (frequency, %; technique)*	Correlation Infinium vs. reported methylation data
<i>RASSF1A</i>	cg00777121	71	70; MSP / 56; MSP / 56; MPS	++
	cg08047457	72	65; MSP	++
	cg21554552	70	65; MSP	++
<i>CCND2</i>	cg25425078	9	46; MSP / 28; MSP / 55; MSP	+
<i>APC</i>	cg16970232	39	45; MSP / 28; MSP / 39; MSP / 49; MSP	++
	cg20311501	35	45; MSP / 28; MSP / 39; MSP / 49; MSP	++
<i>RARβ2</i>	cg27486427	12	17; PS / 0; PS	++
	cg26124016	4	23; MSP	+
<i>CDH13</i>	cg08747377	17	33; MSP	++
<i>SDHB</i>	cg24305835	0	0; MS-HRM	++
	cg03861428	0	0; MS-HRM	++
<i>FH</i>	cg06806184	0	0; MS-HRM	++

Identification of differentially methylated loci

Gene	Illumina ID	Infinium methylation (frequency, %)	Reported methylation data (frequency, %; technique)*	Correlation Infinium vs. reported methylation data
<i>RASSF1A</i>	cg00777121	71	70; MS	
	cg08047457	72	65; MS	
	cg21554552	70	65; MS	
<i>CCND2</i>	cg25425078	9	46; MS	
<i>APC</i>	cg16970232	39	45; MS	
	cg20311501	35	45; MS	
<i>RARβ2</i>	cg27486427	12	17; PS	
	cg26124016	4	23; MS	
<i>CDH13</i>	cg08747377	17	33; MS	
<i>SDHB</i>	cg24305835	0	0; MS-	
	cg03861428	0	0; MS-	
<i>FH</i>	cg06806184	0	0; MS-	



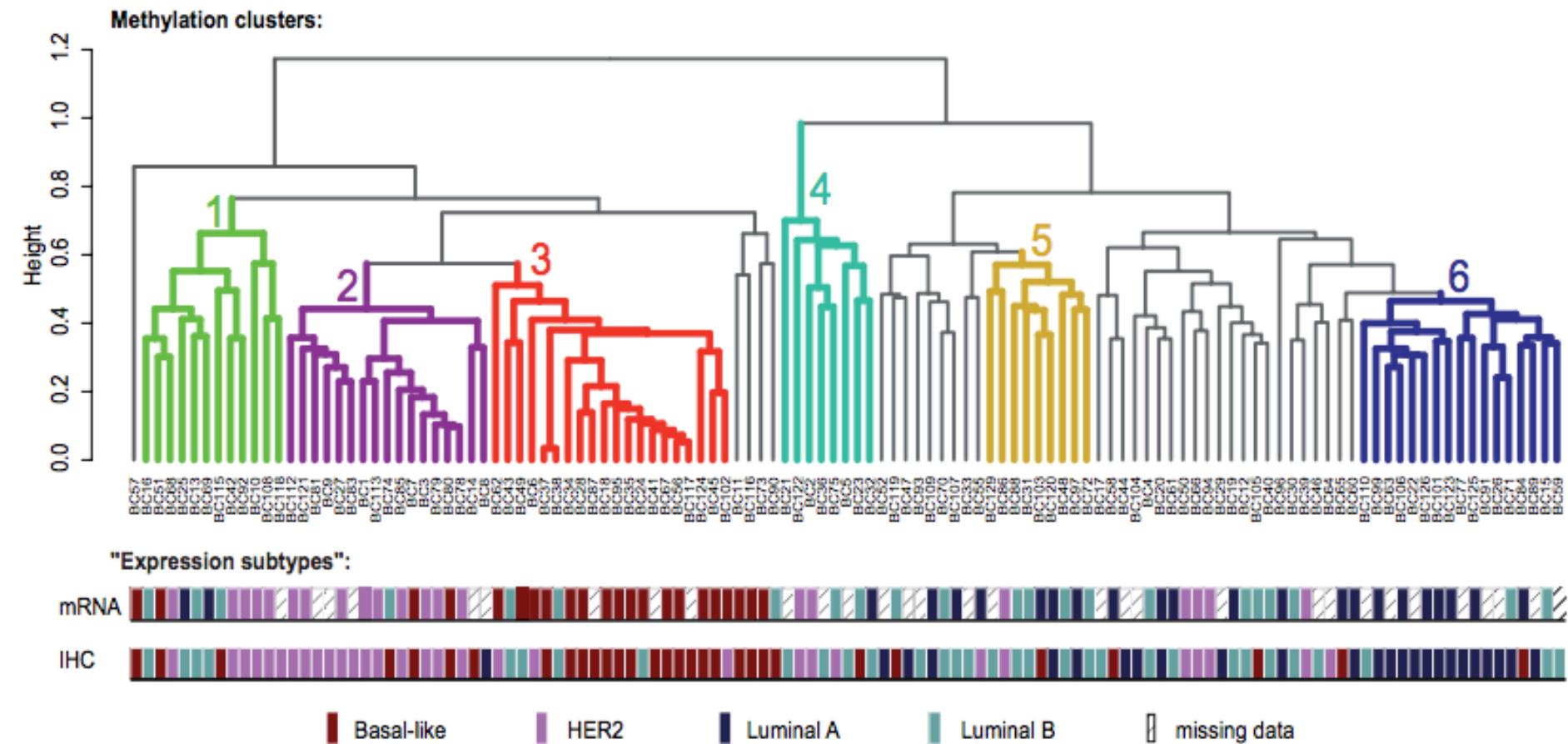
Unsupervised clustering identifies primary clusters by ER status



Six stable, biologically relevant clusters

A

UNSUPERVISED CLUSTERING ON THE MAIN SET

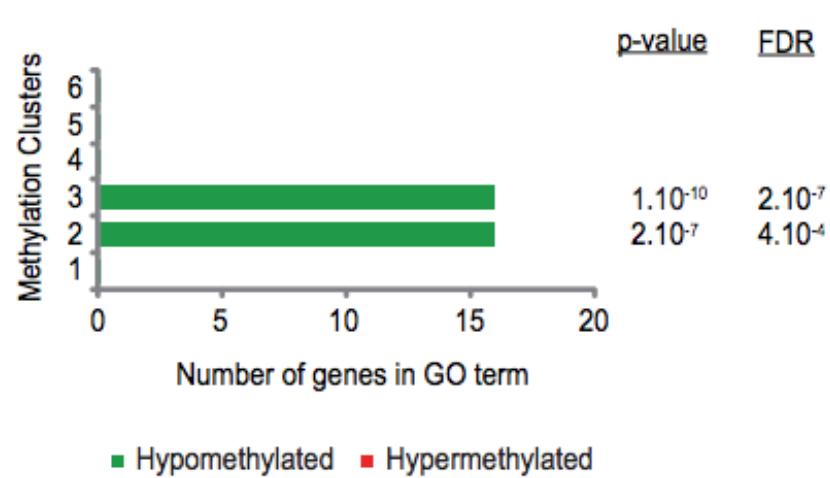


Validated the presence of these clusters in an independent data set

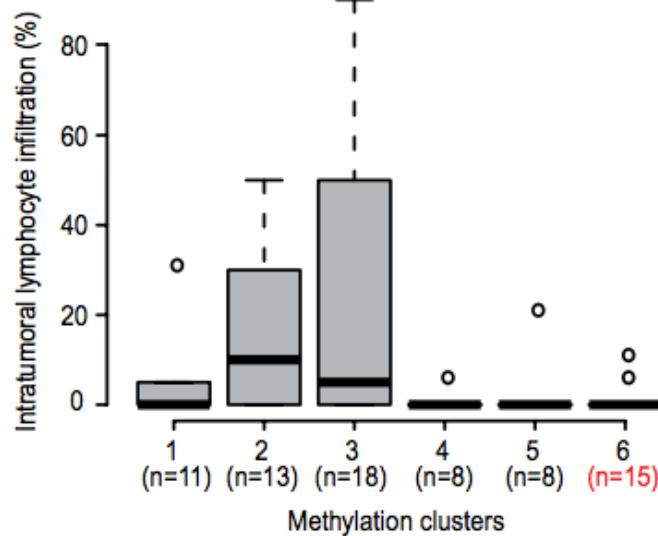
Identification of immune-related clusters

B

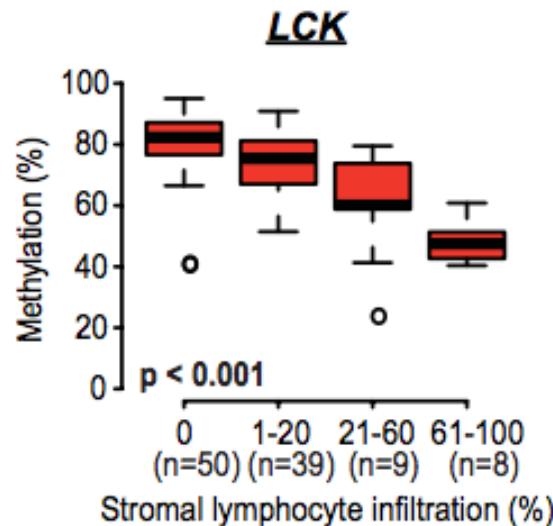
Immune system process



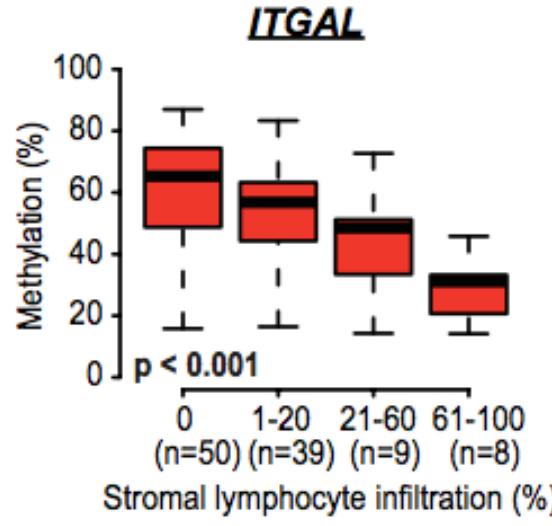
Intratumoral lymphocyte infiltration



LCK



ITGAL

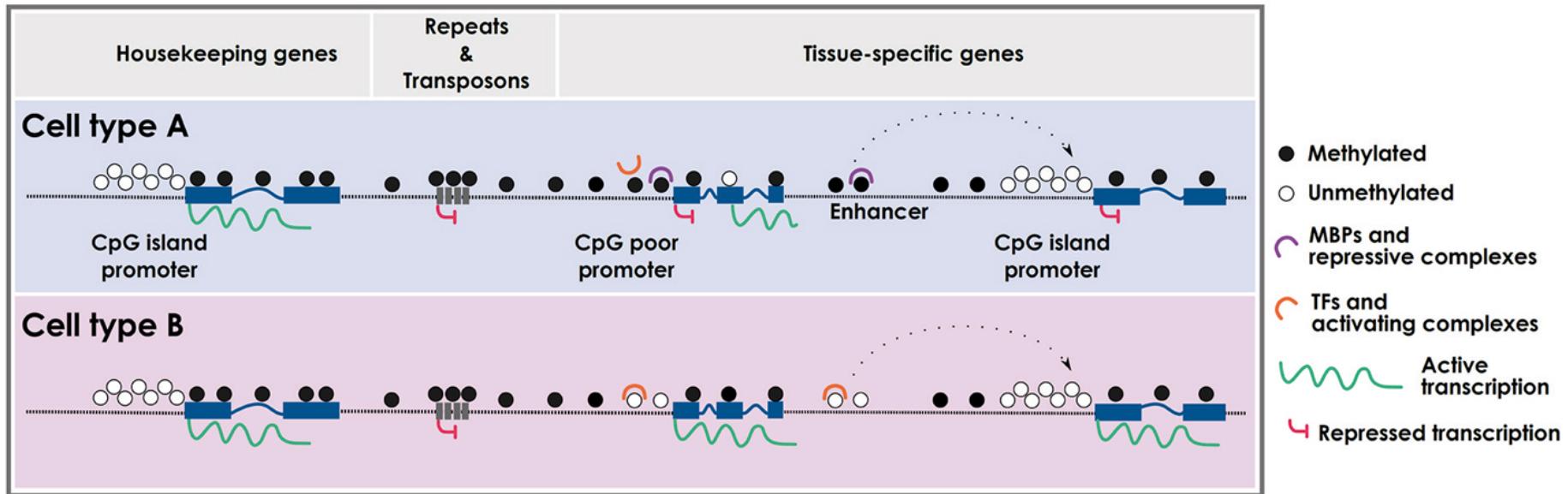


How?

Epigenomics methods

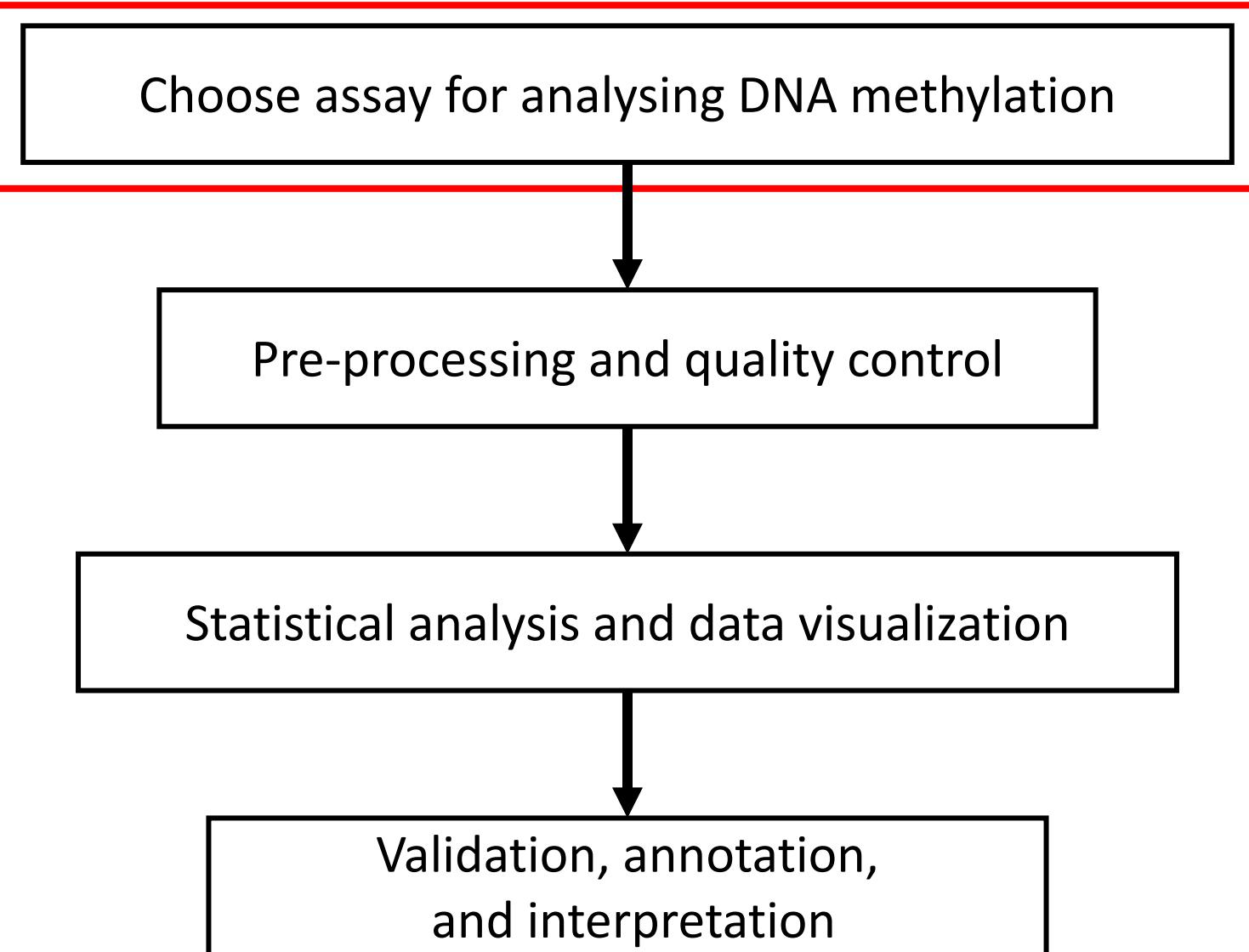
DNA Methylation

“Typical” 5mC DNA methylation in mammalian genome

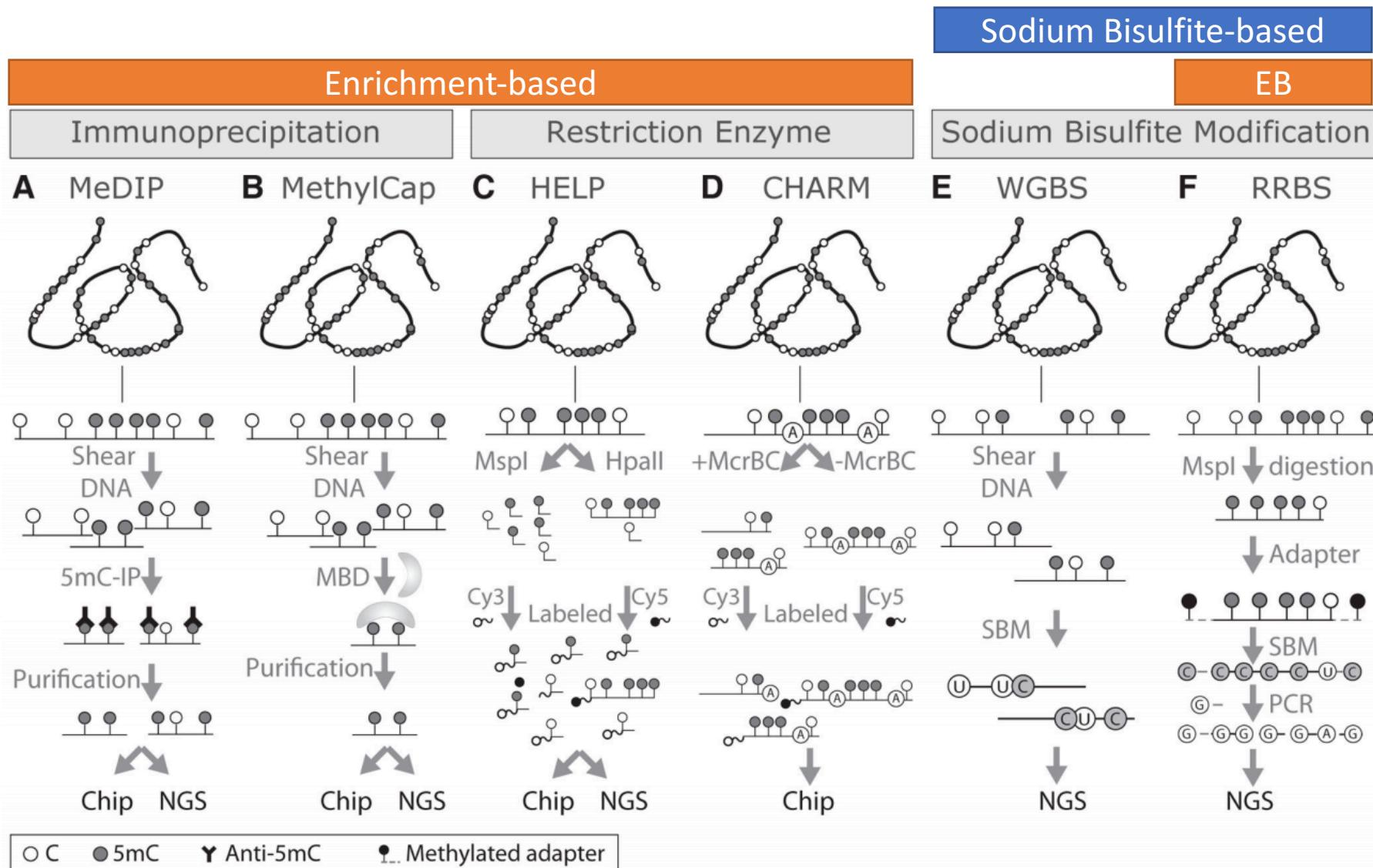


- CpG island promoter and enhancer methylation is inversely proportional to gene expression
- Gene body methylation is correlated with gene expression of highly transcribed genes
- CpG methylation is typical in repeats and transposons
- Changes in these patterns, both globally and loci-specific, have been observed in disease progression

Anatomy of a DNA methylome project



Epigenomics methods for studying 5-mC DNA methylation



Choose assay for analysing DNA methylation



Pre-processing and quality control

Analysis ready data!

Statistical analysis and data visualization

Validation, annotation,
and interpretation

Goal: Identify differentially methylated probe/site

Linear model

Observed signal for probe p

Coefficients for probe p

$$E(Y_p) = X\beta_p$$

Design matrix

$$\sim \beta_D + \beta_1 + \cdots + \beta_n + \varepsilon$$

Identifier	Disease	Gender	...	Age
Case01	1	1	...	21
Case02	1	0	...	25
Ctrl01	0	1	...	22
Ctrl02	0	0	...	27

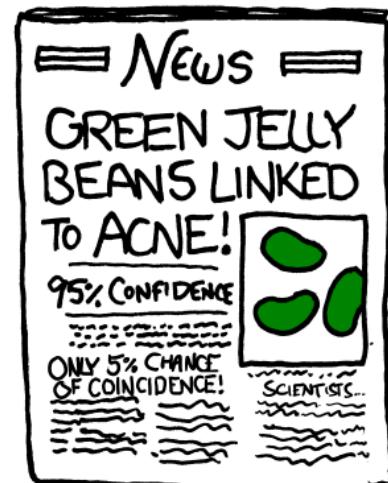
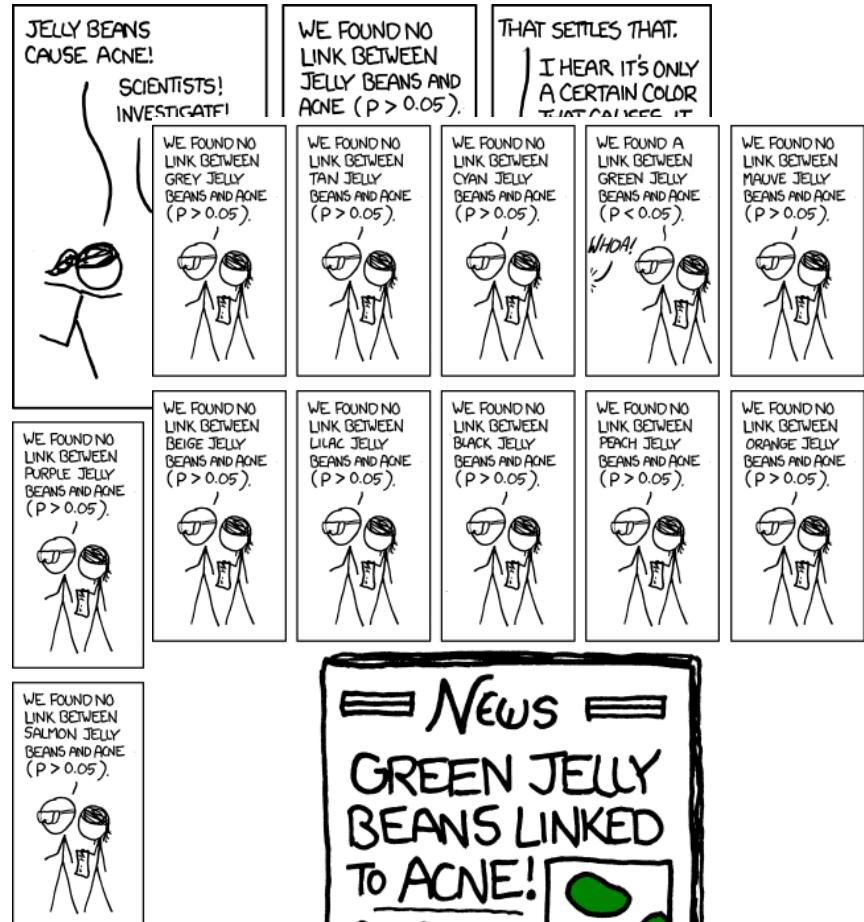
Multiple hypothesis correction

$$E(Y_1) = X\beta_1$$

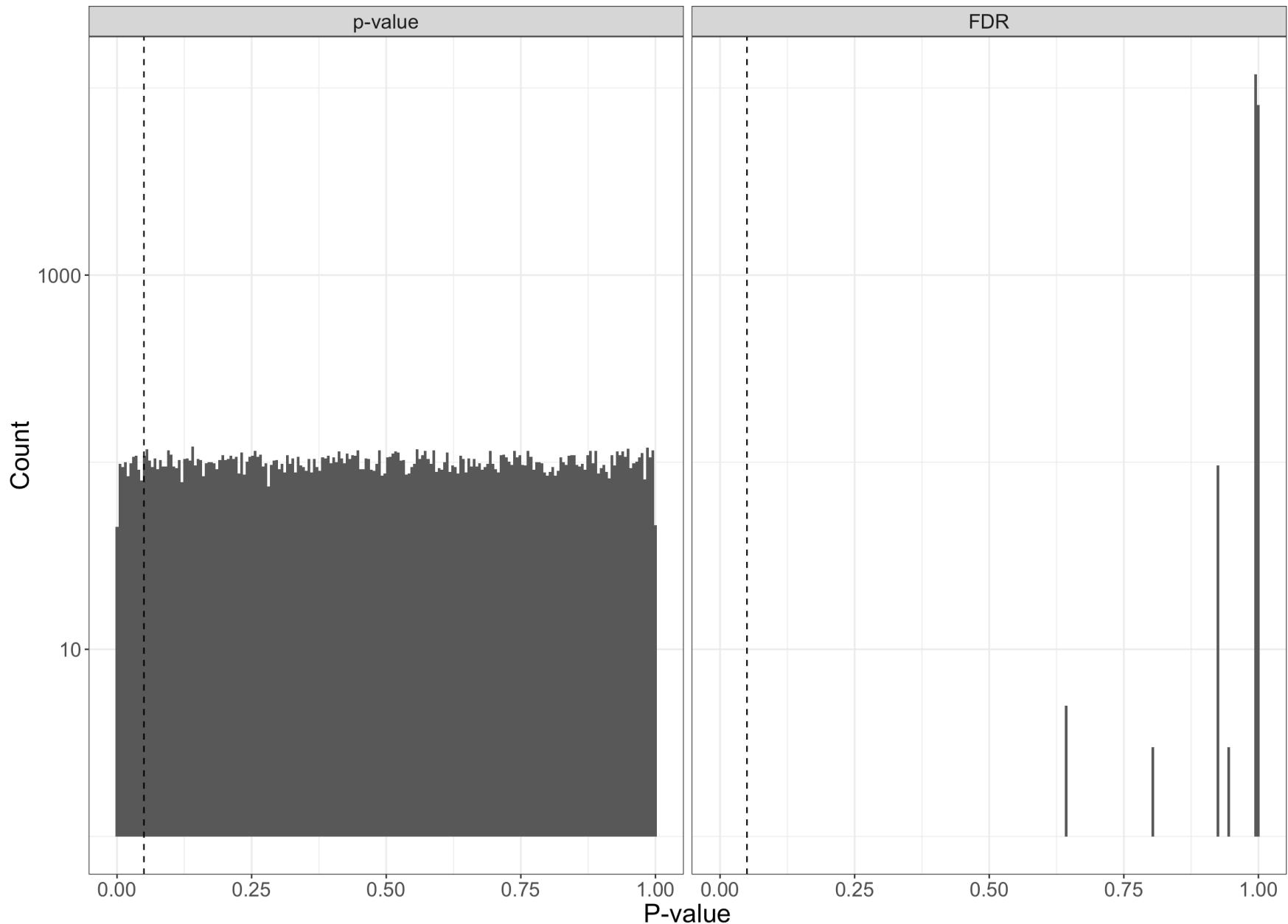
...

$$E(Y_n) = X\beta_n$$

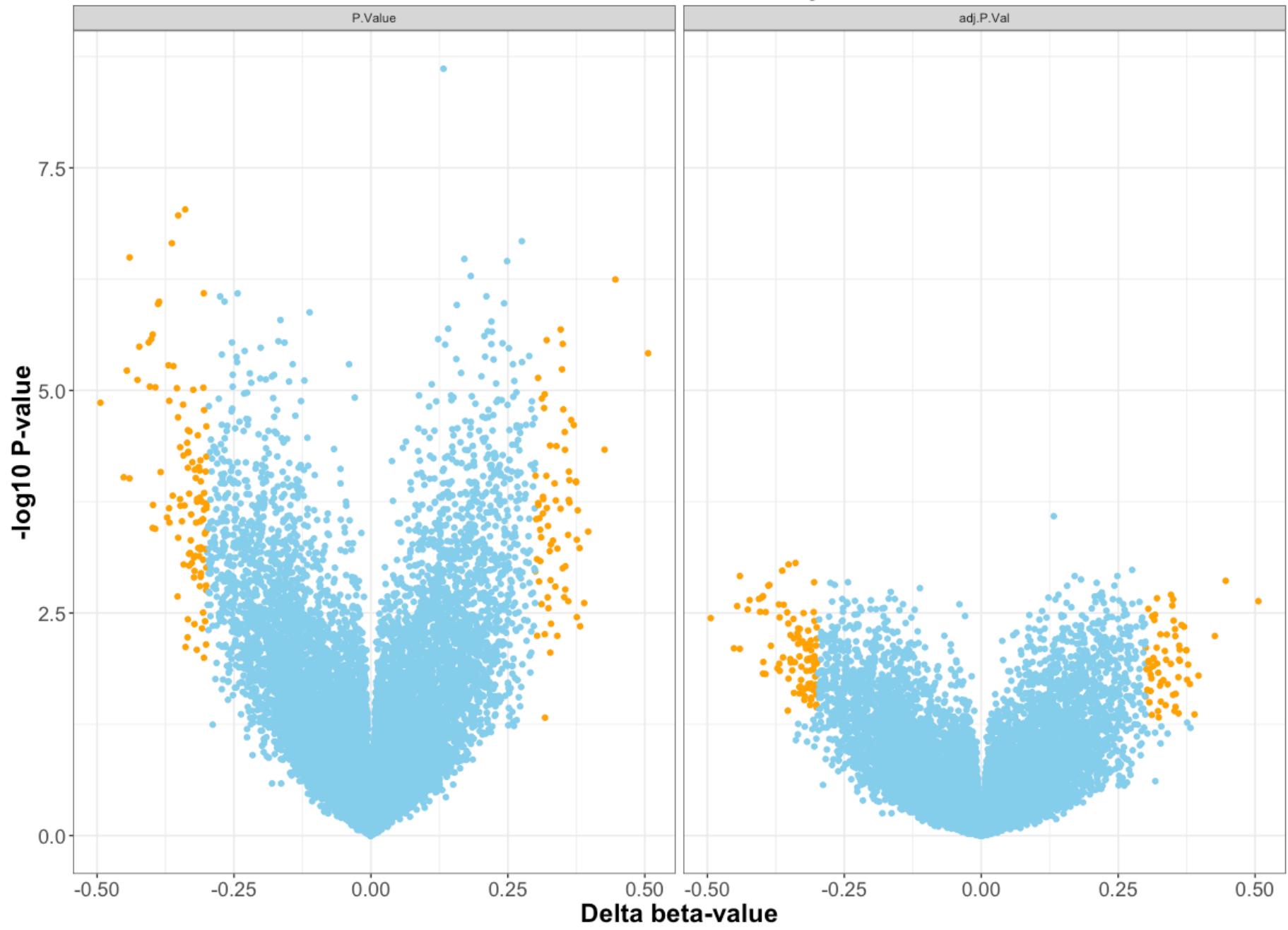
$n = 450,000$



20,000 permutations of H0 = True



P-values before and after adjustment

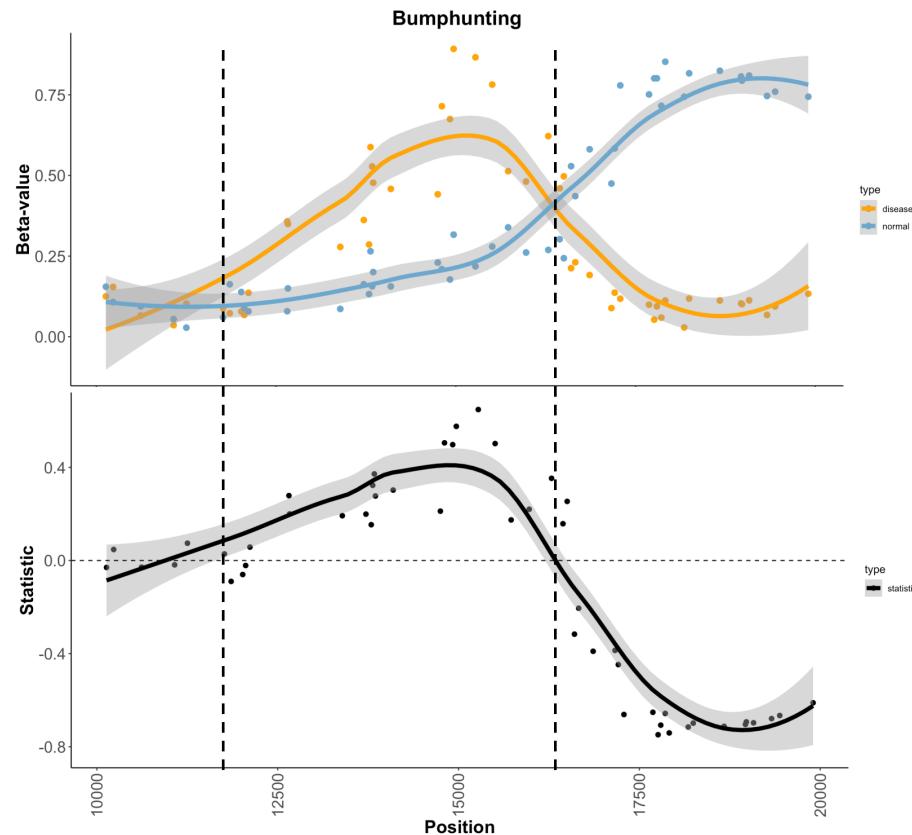


Differentially methylated regions (DMRs)

- Identify regions of differential methylation rather than a CpG site
 - Change across a region is biologically more relevant
 - Less features being tested
- Calculate differential methylation probabilities across continuous sites over a region

Software:

- minfi::bumphunter
- DMRcate
- CHAMP::probe lasso
- comb-p
- BS-seq



Choose assay for analysing DNA methylation



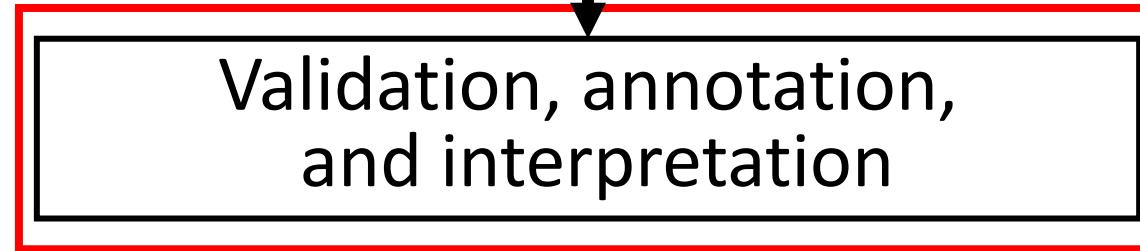
Pre-processing and quality control

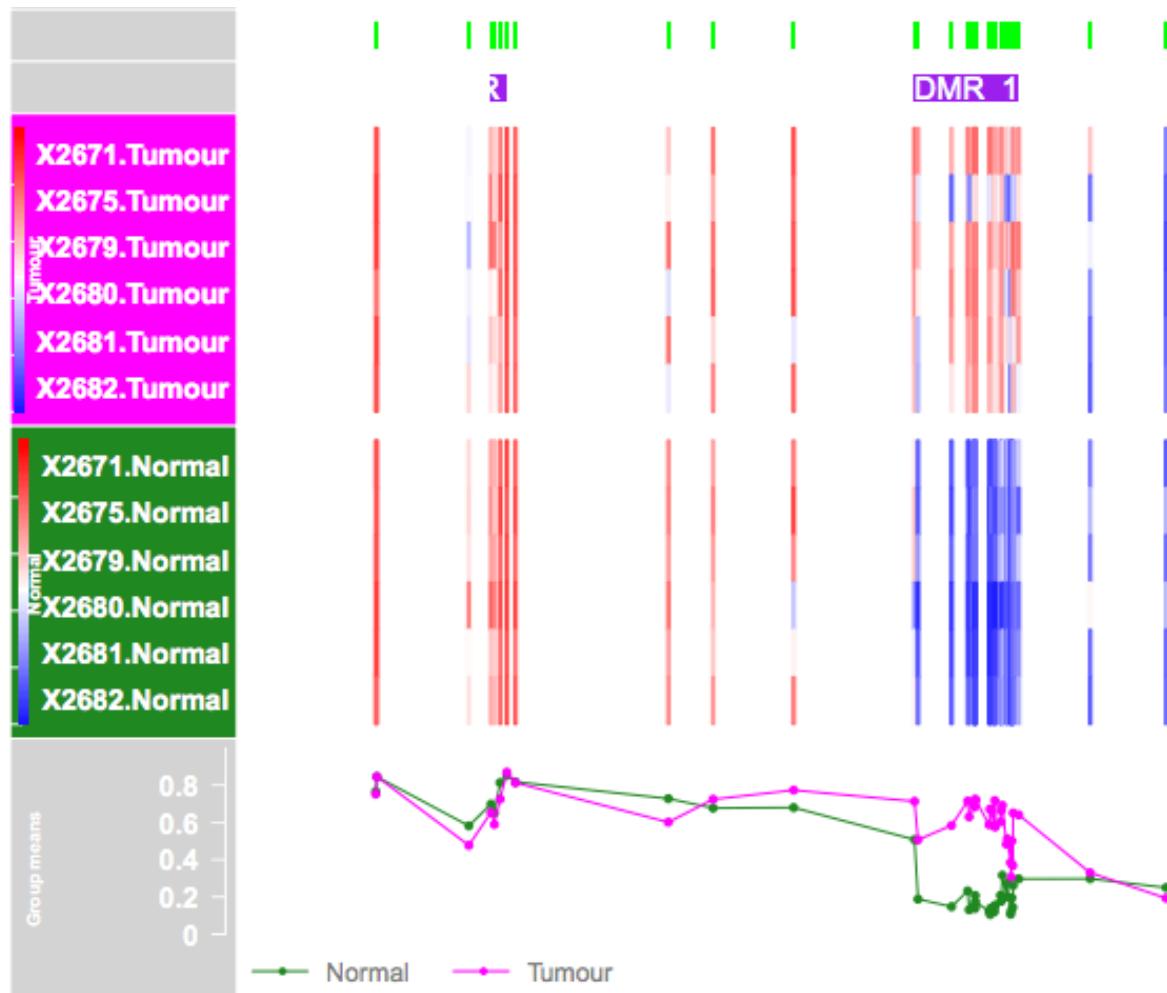


Statistical analysis and data visualization



Validation, annotation,
and interpretation

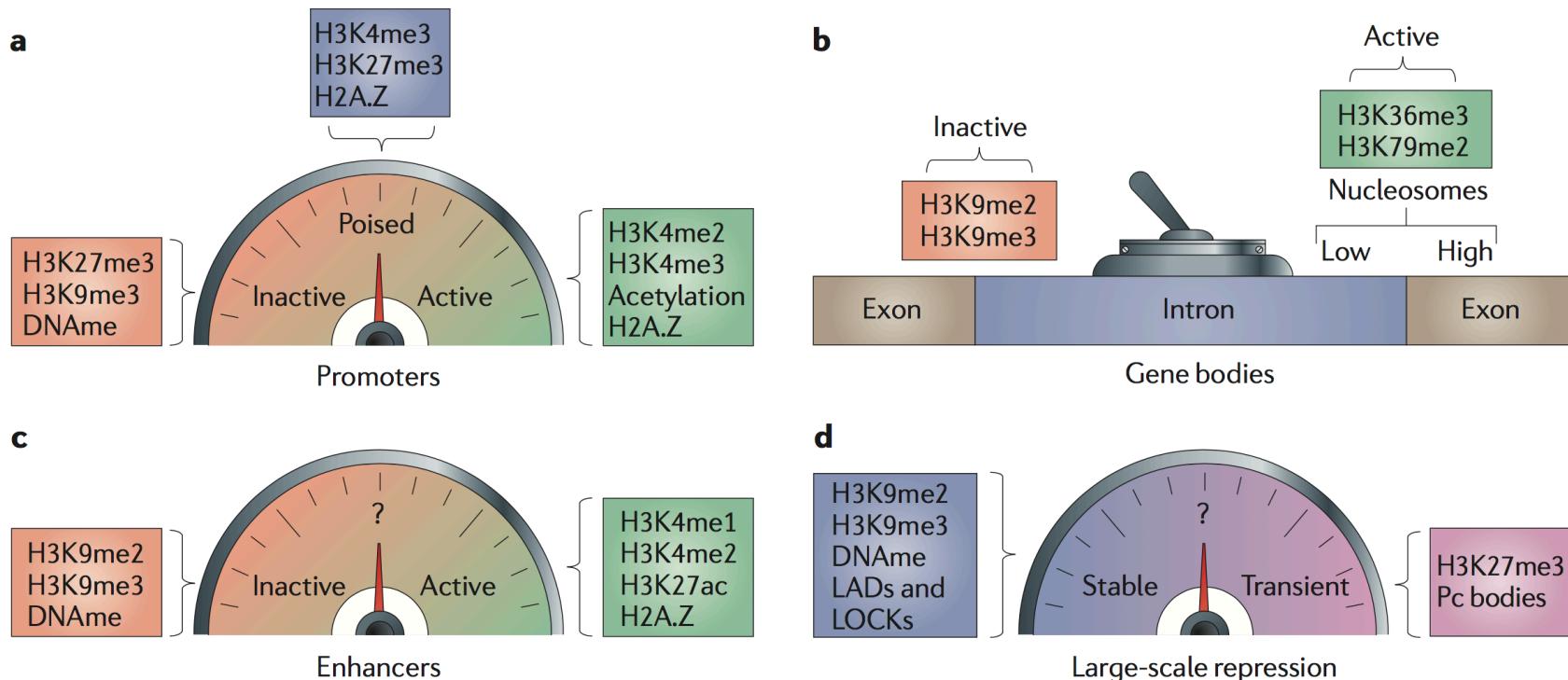




Histone modifications

The histone code hypothesis

- Combinatorial histone post-translational modifications (PTMs) encrypt the recruitment of different histone effectors
- These combinations define specific chromatin states
- These states are heritable

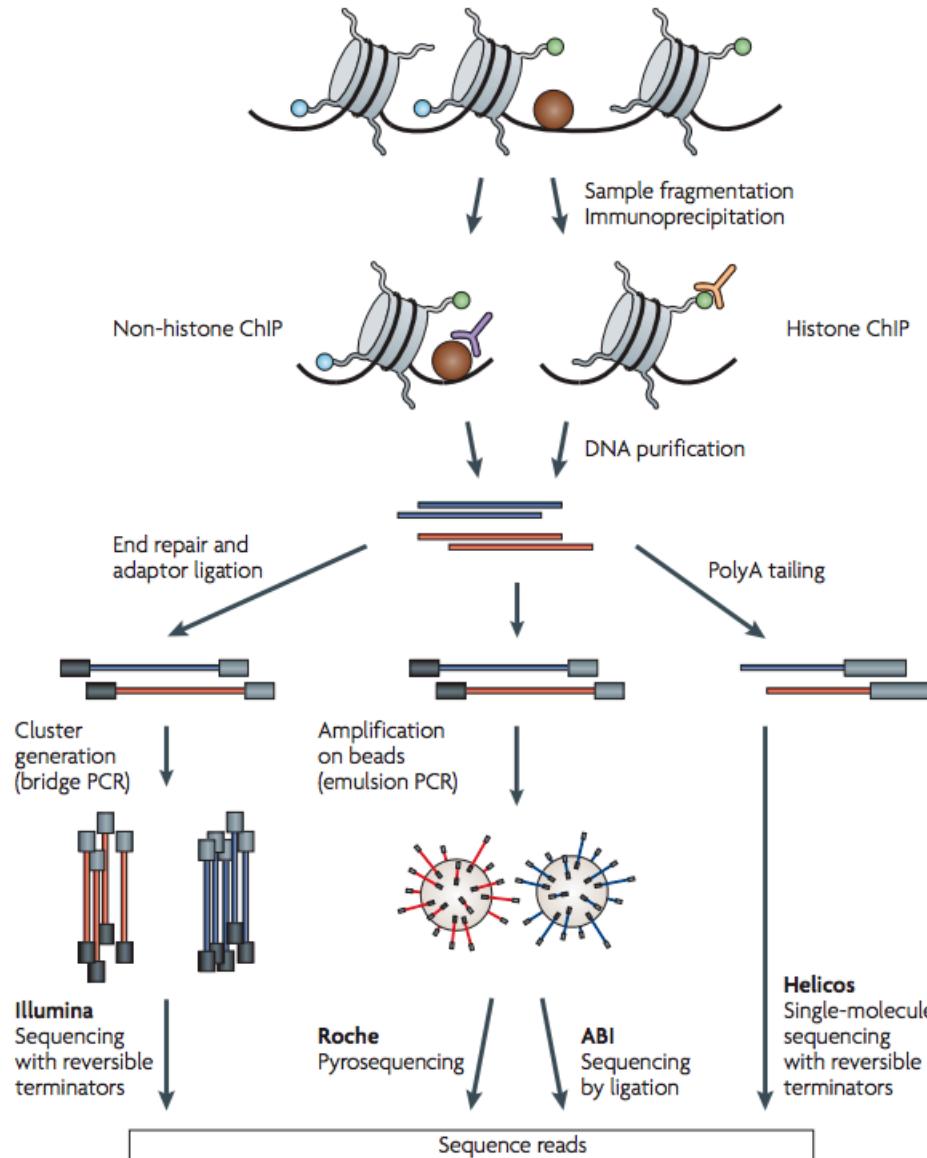


Strahl, B. D., & Allis, C. D. (2000). The language of covalent histone modifications. *Nature*, 403(6765), 41

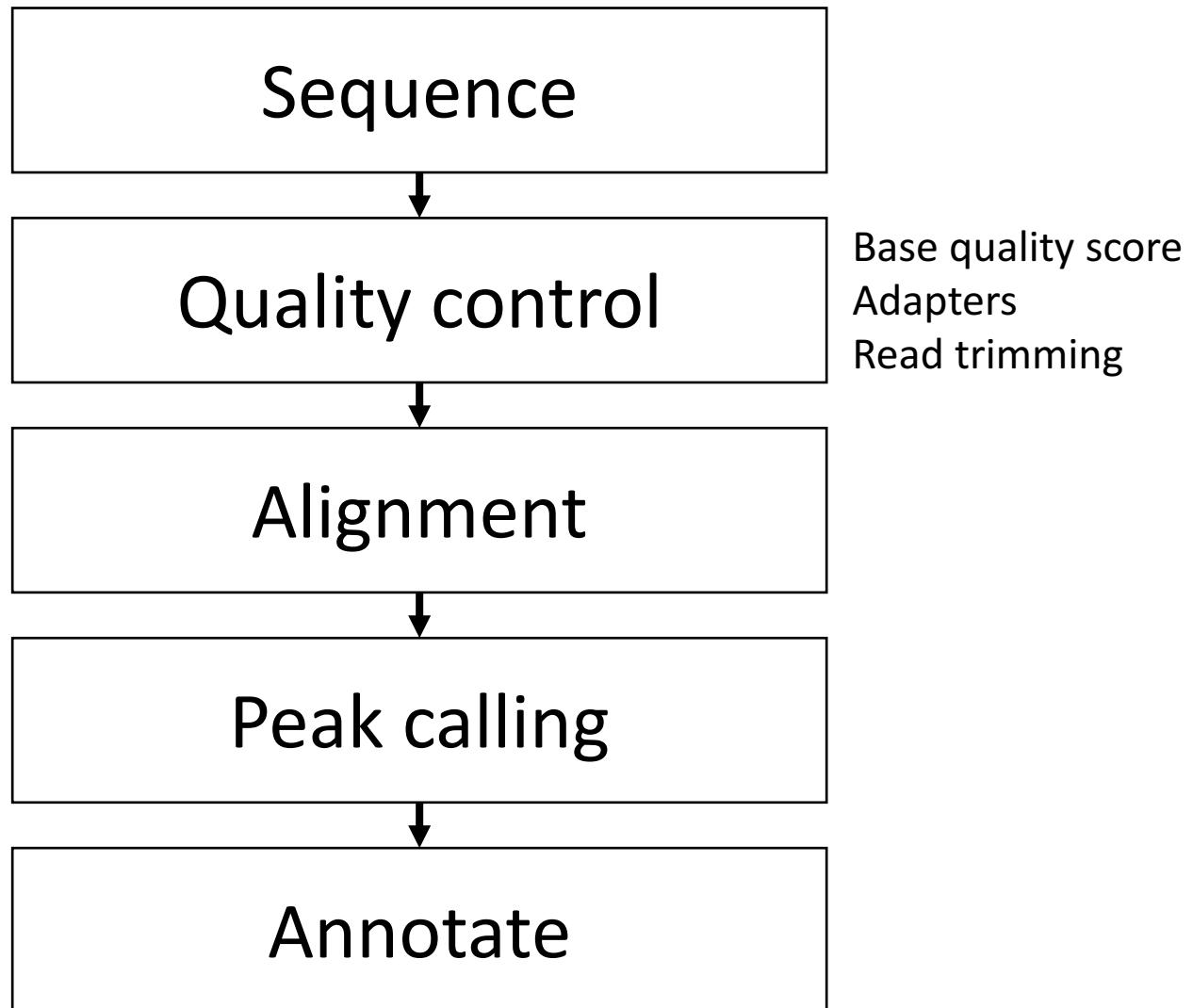
Turner, B. M. (2000). Histone acetylation and an epigenetic code. *Bioessays*, 22(9), 836-845

Zhou, V. W., Goren, A., & Bernstein, B. E. (2011). Charting histone modifications and the functional organization of mammalian genomes. *Nature Reviews Genetics*, 12(1), 7

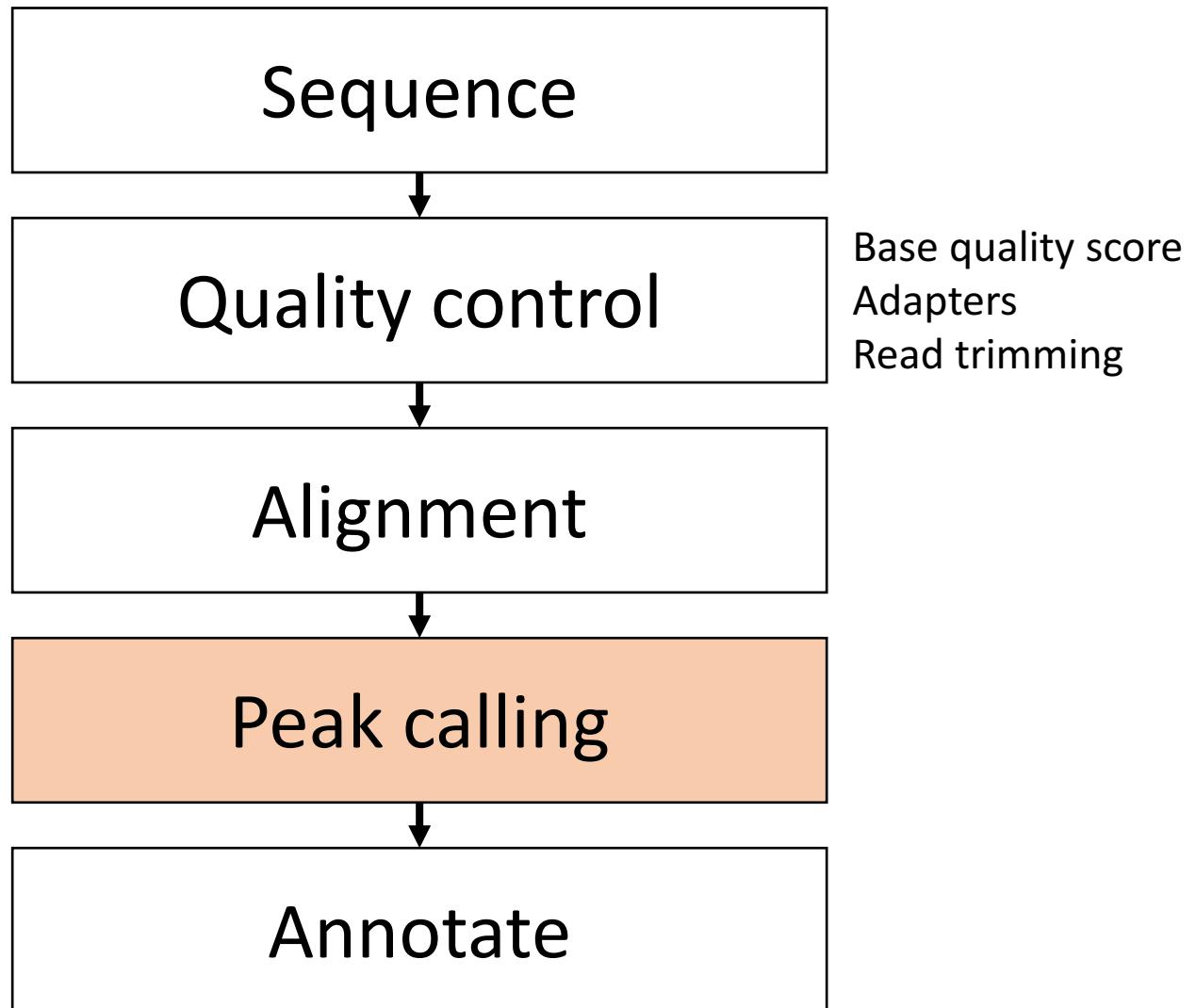
ChIP-Seq: Assay



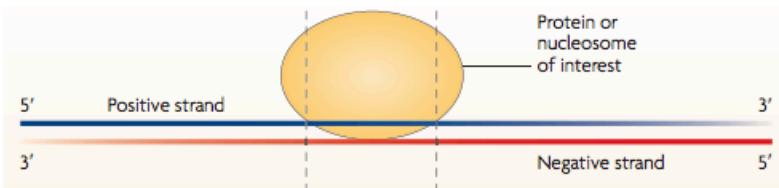
ChIP-Seq: Analysis pipeline



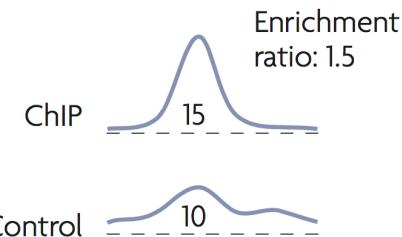
ChIP-Seq: Analysis pipeline



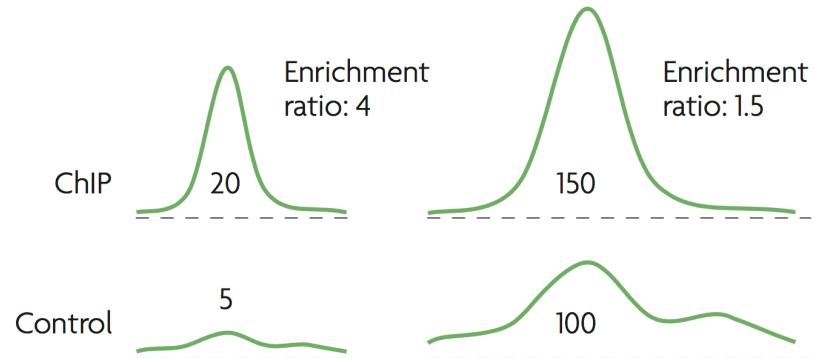
Goal: Identify enrichment peaks



Ba Not statistically significant

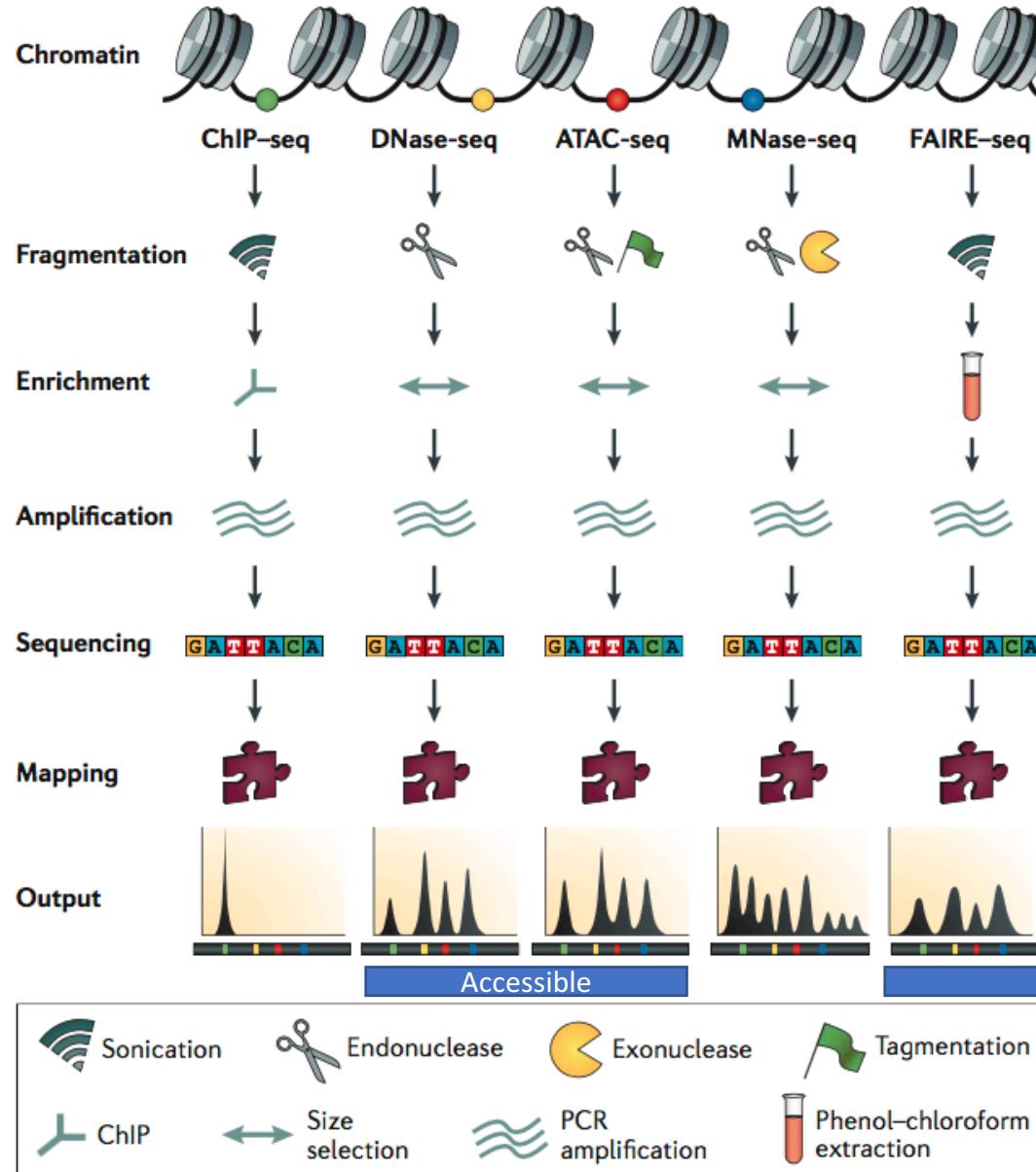


Bb Statistically significant

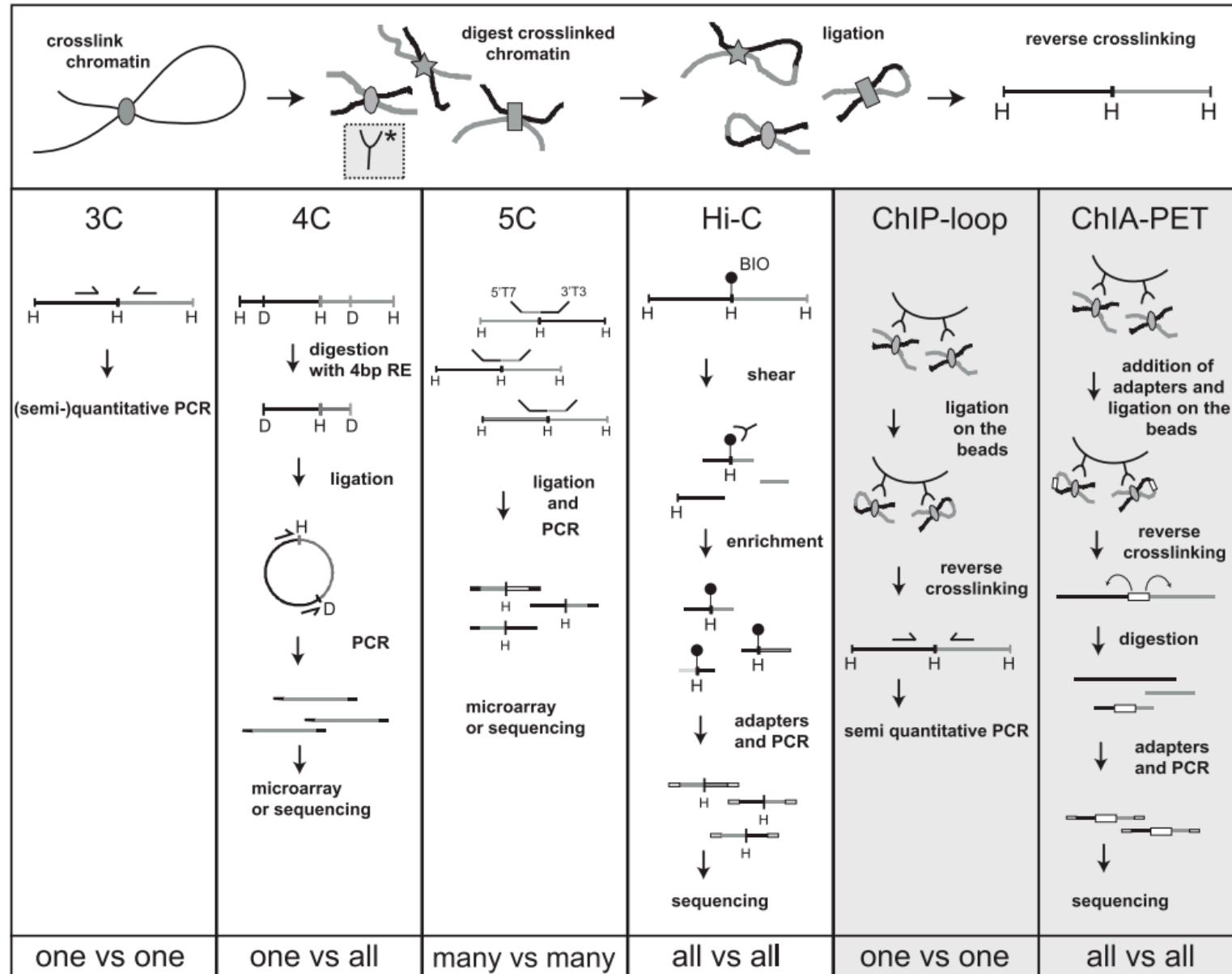


Chromatin features

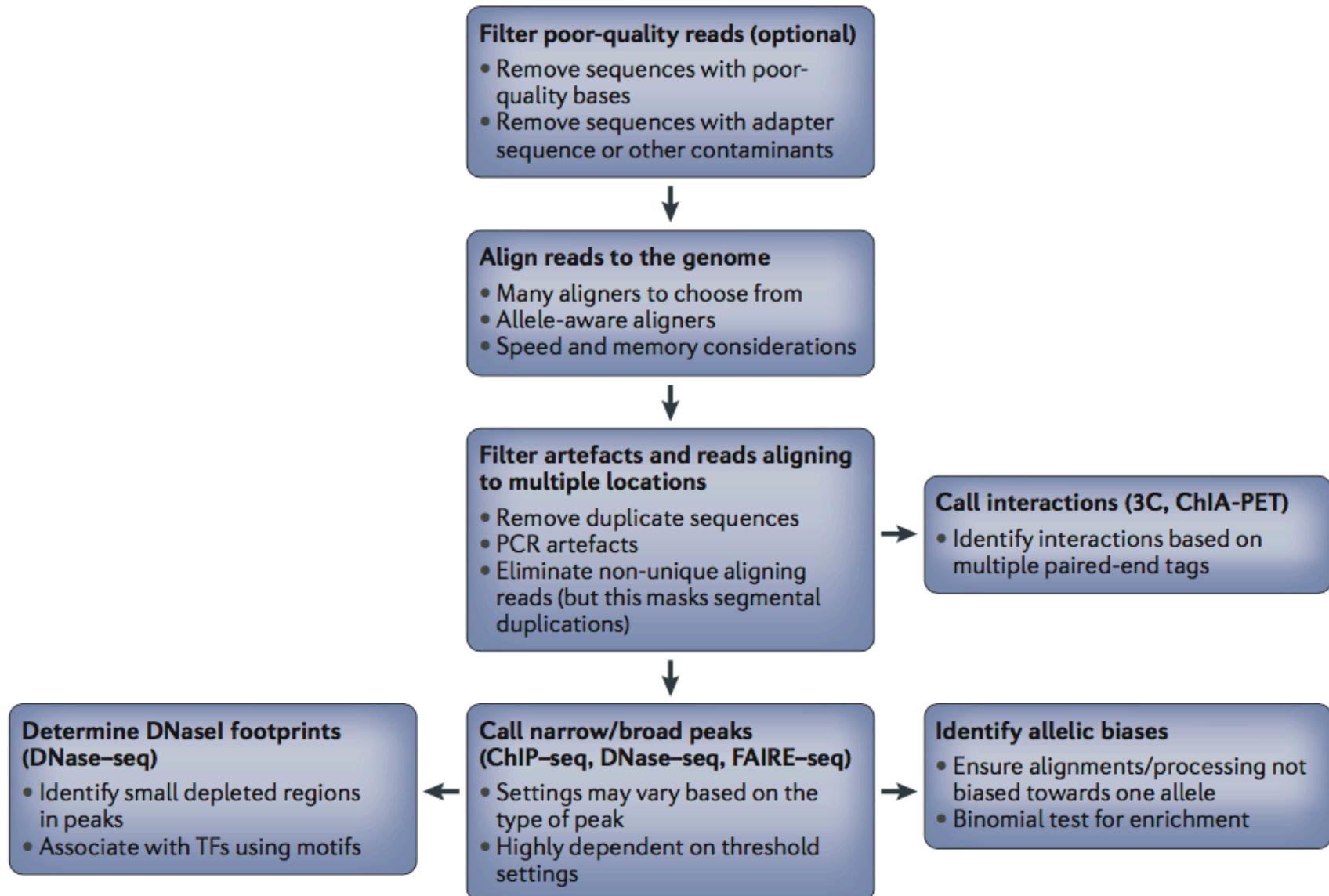
Chromatin features: accessibility



Chromatin features: long range interactions



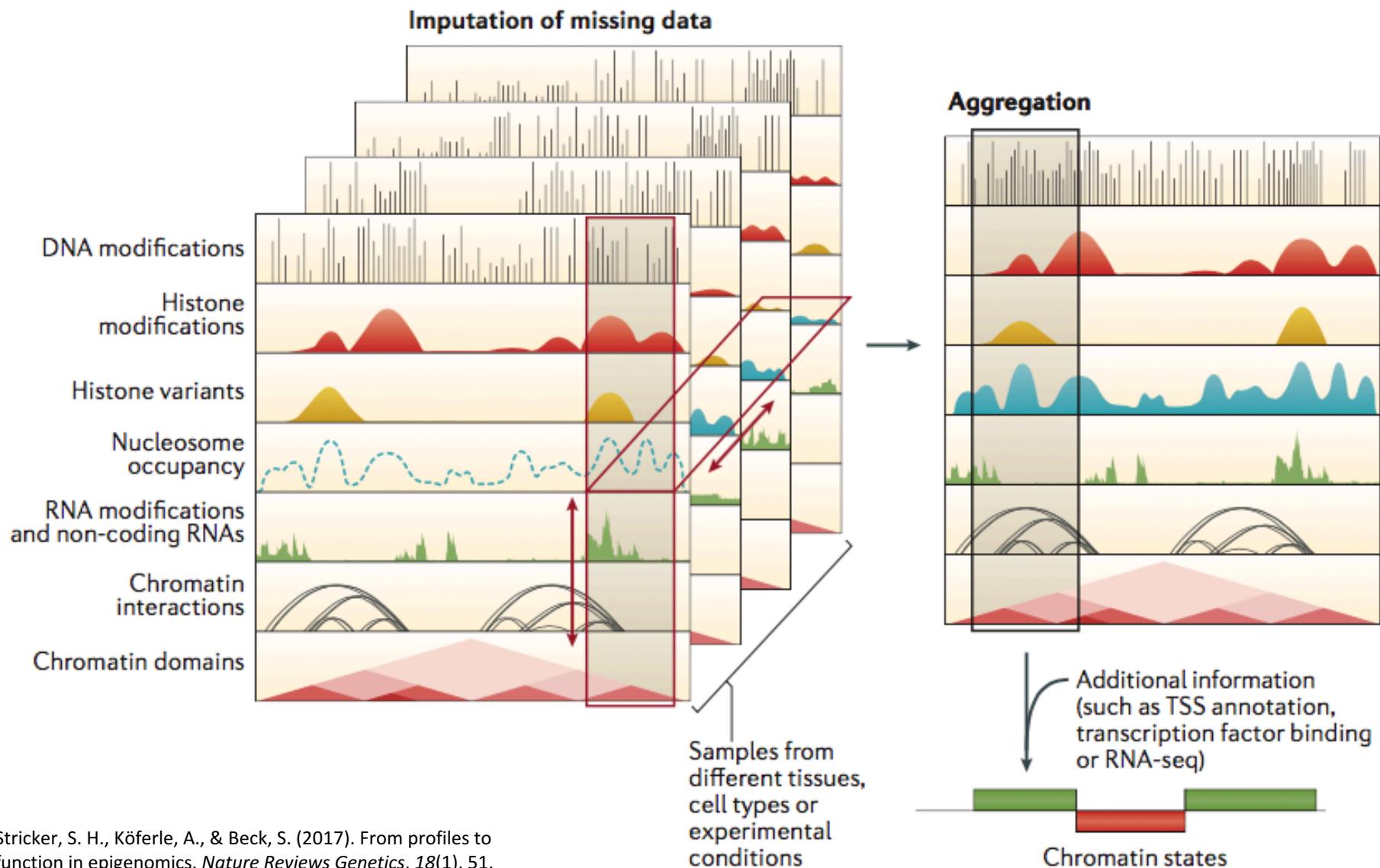
Post-sequencing bioinformatics pipeline



Integration

Putting it all together

Imputation and integration to infer functional states

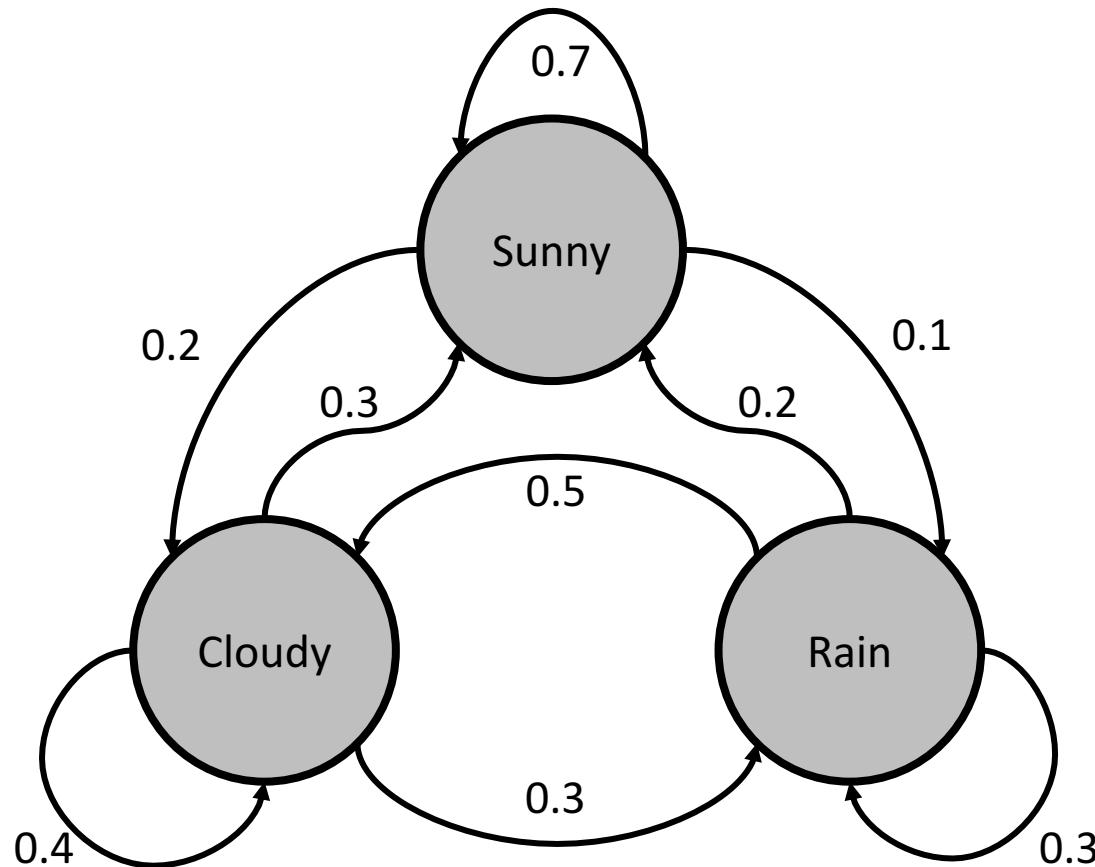


ChromHMM: Inferring chromatin states

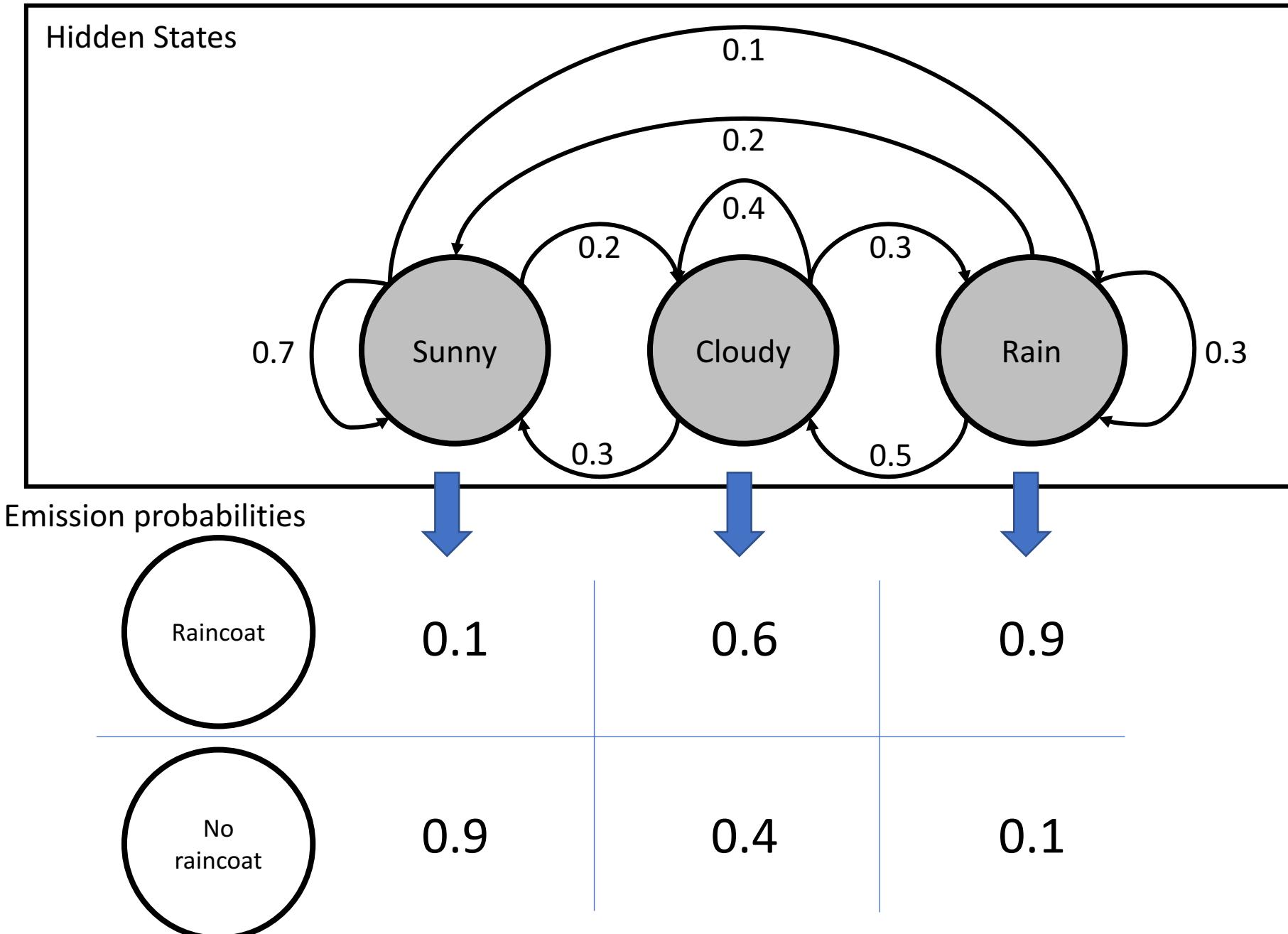
- Purpose: identify segments of the genome as biologically relevant functional elements
- ChromHMM uses a multivariate hidden Markov model (HMM) to infer chromatin states
- Chromatin states are derived from combinations of epigenetics marks (histone modifications, chromatin features)
- Chromatin states include:
 - Active/Weak
 - Promoter, enhancer, insulator
 - Polycomb

Simple Hidden Markov Model

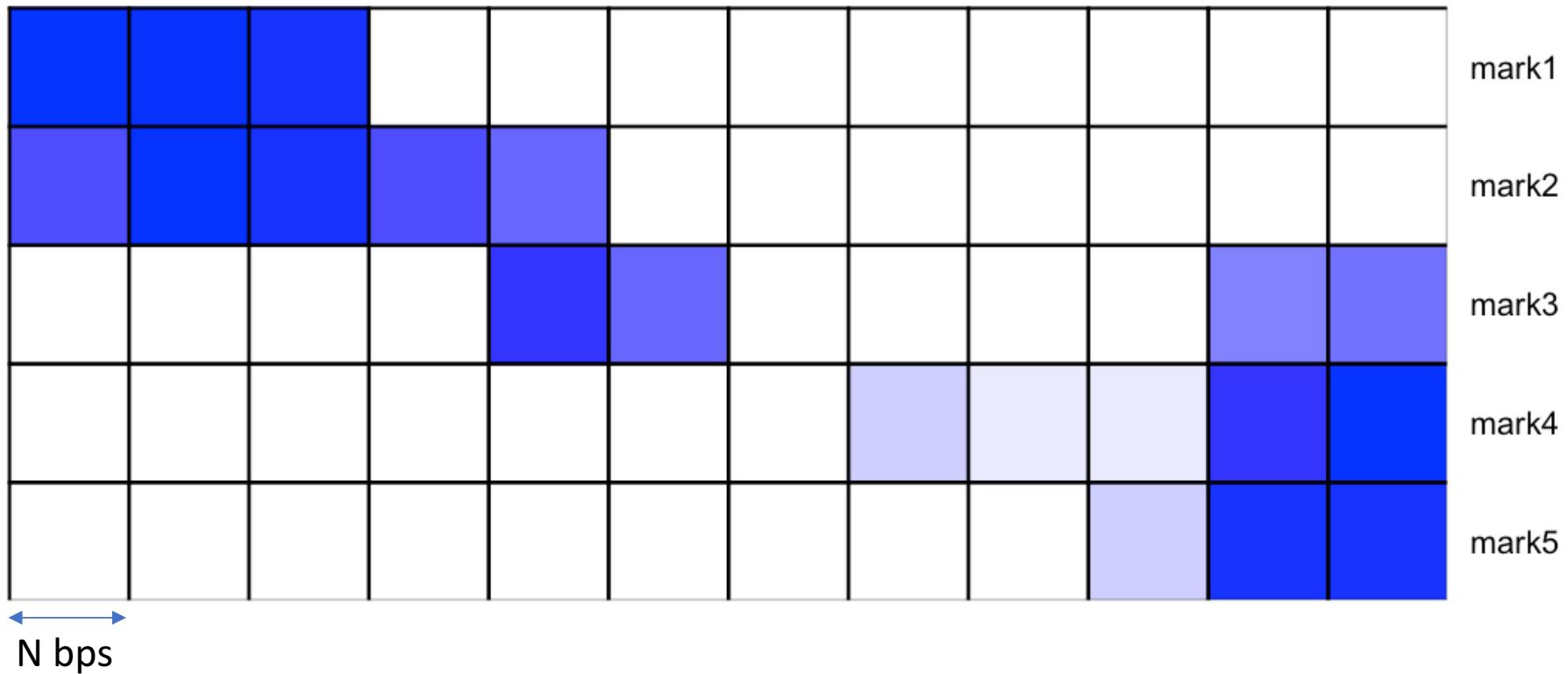
Markov Chain for weather



Simple Hidden Markov Model



Observation

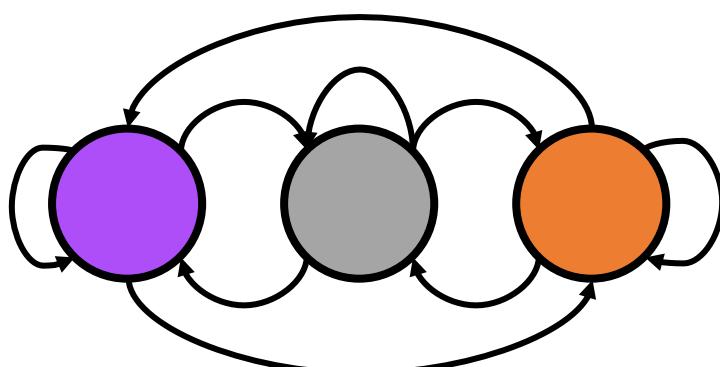


Observation

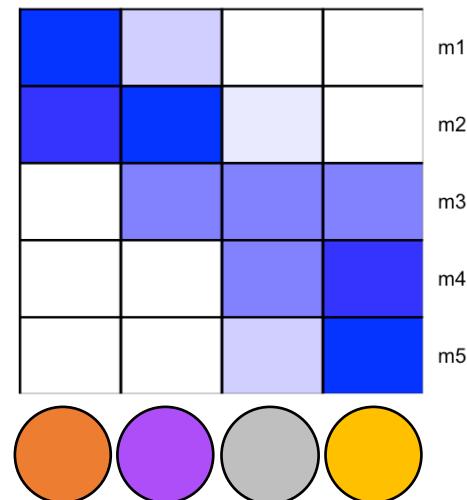
mark1												
mark2												
mark3												
mark4												
mark5												

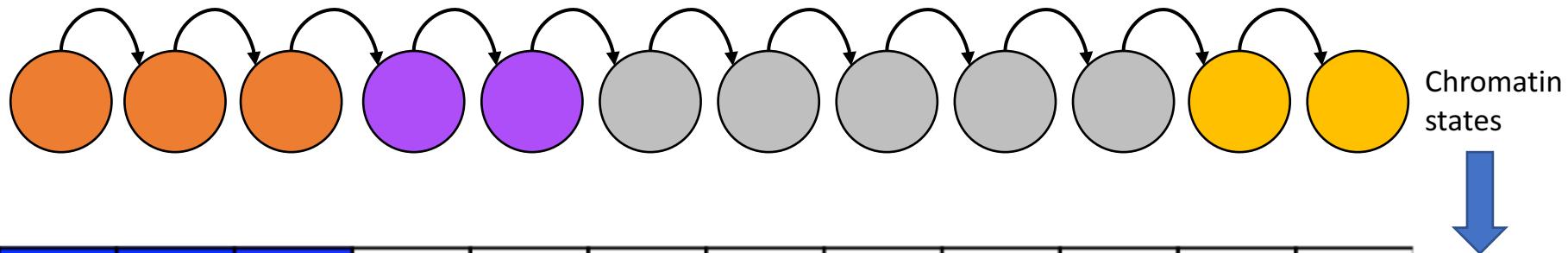
Learn from the data

1. Transition probabilities

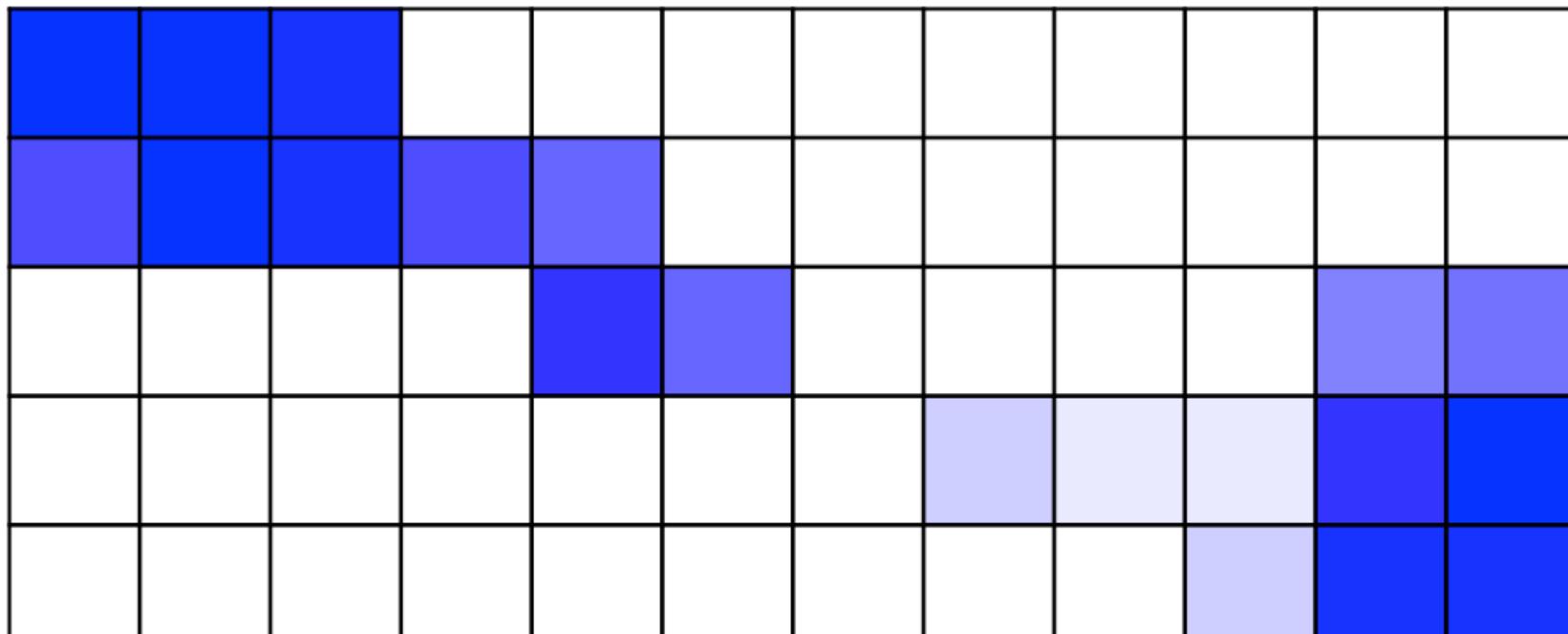


2. Emission probabilities





Chromatin
states



Genomic Annotations

TSS

Gene body

Enhancer

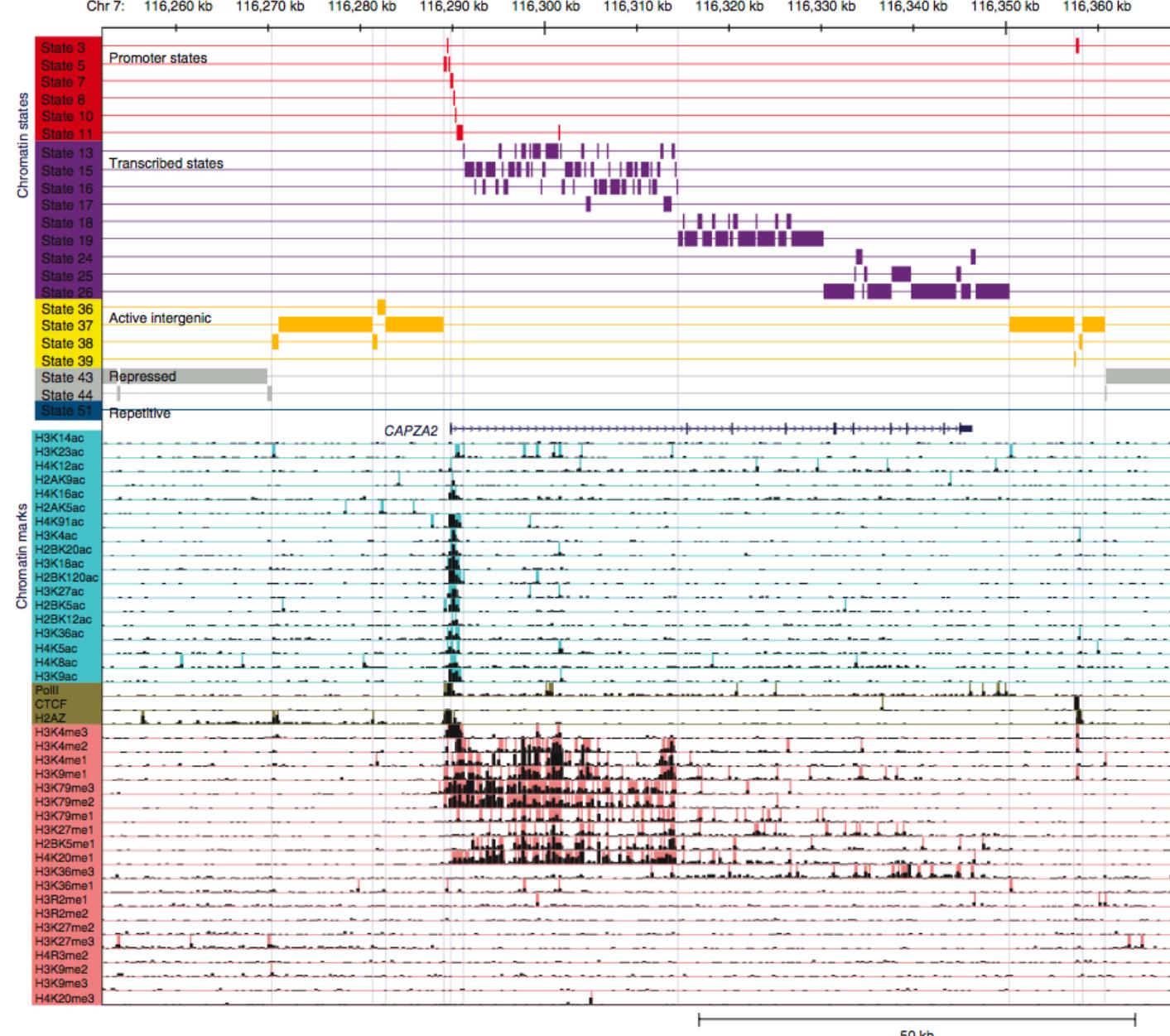
Chromatin states

Promoters

Transcribed states

Active intergenic

Repressed



Integrative genomics

Case 02:

Integrating genomics, epigenomics, and experimental biology to identify causal allele and regulatory mechanism in obesity



ESTABLISHED IN 1812

SEPTEMBER 3, 2015

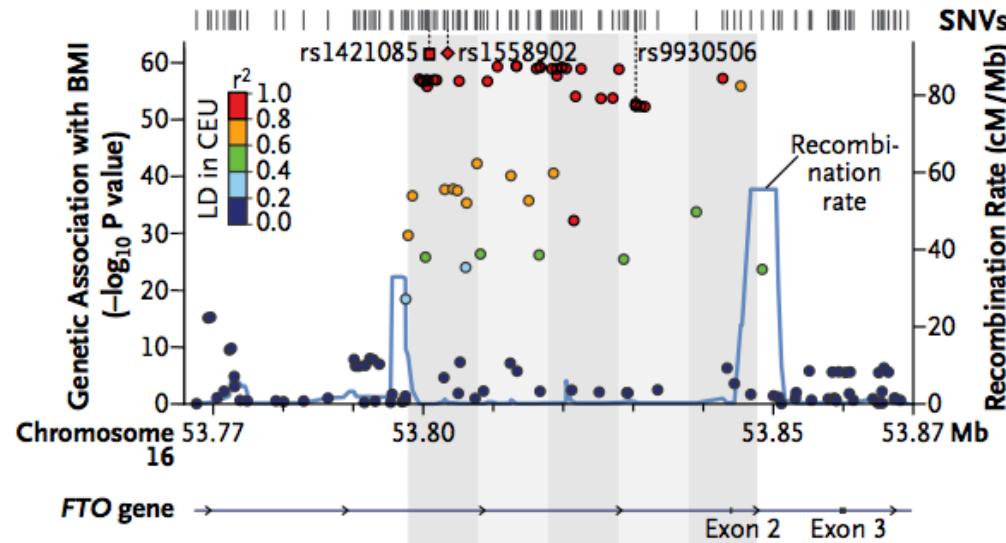
VOL. 373 NO. 10

FTO Obesity Variant Circuitry and Adipocyte Browning in Humans

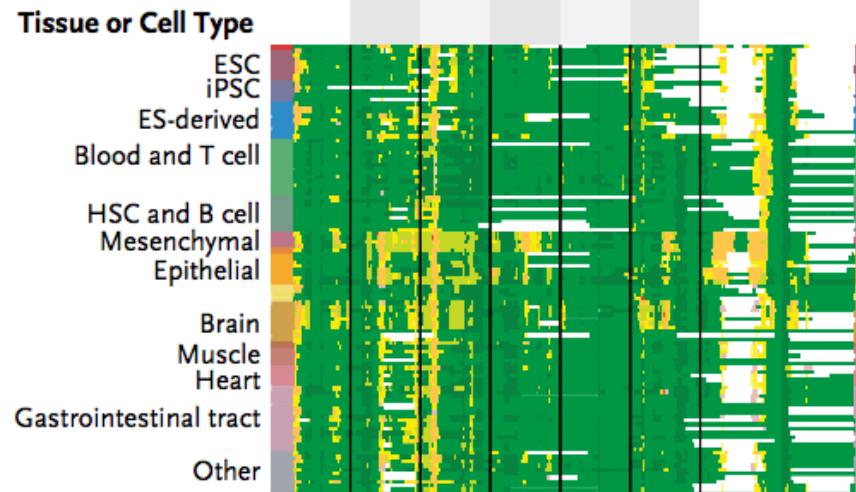
Melina Claussnitzer, Ph.D., Simon N. Dankel, Ph.D., Kyoung-Han Kim, Ph.D., Gerald Quon, Ph.D.,
Wouter Meuleman, Ph.D., Christine Haugen, M.Sc., Viktoria Glunk, M.Sc., Isabel S. Sousa, M.Sc.,
Jacqueline L. Beaudry, Ph.D., Vijitha Puvilindran, B.Sc., Nezar A. Abdennur, M.Sc., Jannel Liu, B.Sc.,
Per-Arne Svensson, Ph.D., Yi-Hsiang Hsu, Ph.D., Daniel J. Drucker, M.D., Gunnar Mellgren, M.D., Ph.D.,
Chi-Chung Hui, Ph.D., Hans Hauner, M.D., and Manolis Kellis, Ph.D.

GWAS locus prioritization using epigenome map of human cells

A

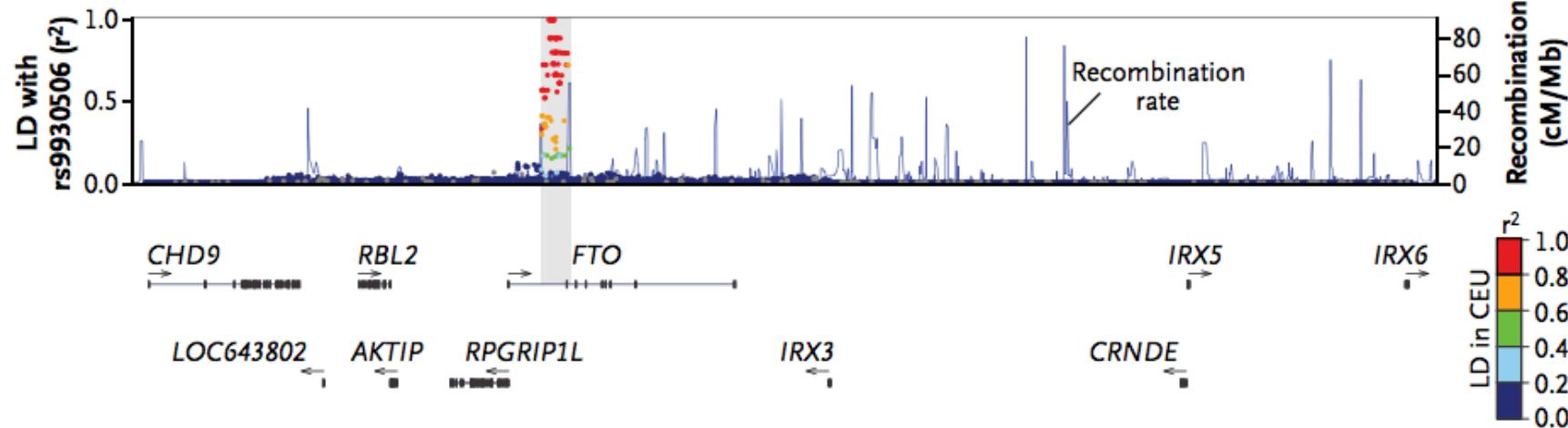


B

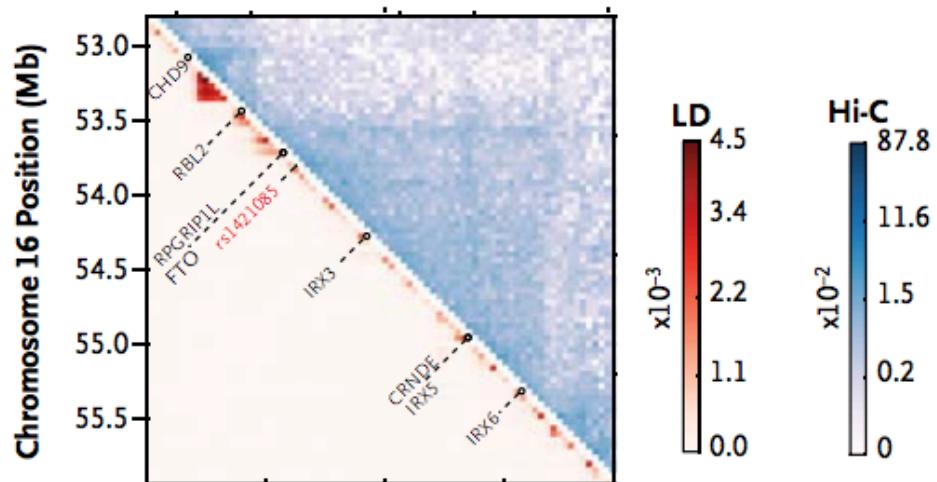


Hi-C data reveals 2-Mb topologically associated domain (TAD)

A

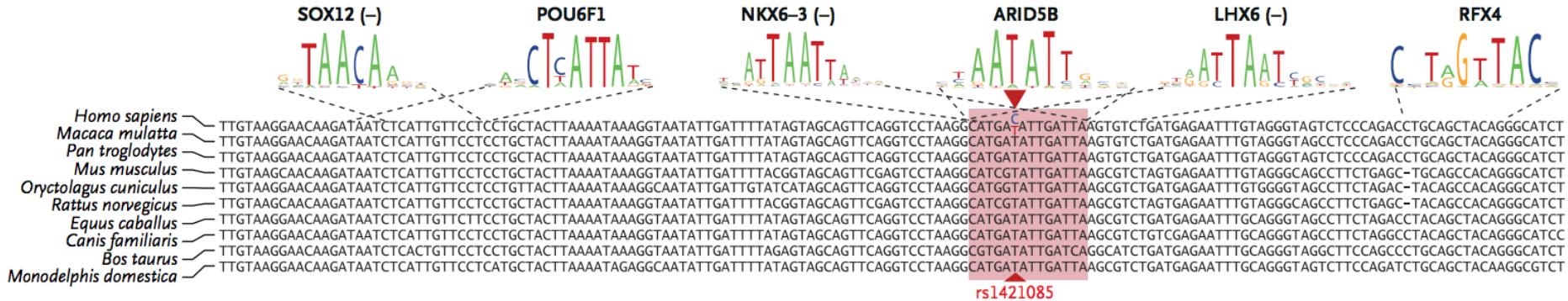


B



Identification of an SNV in linkage disequilibrium with candidate risk allele which disrupts an ARID5B repressor motif

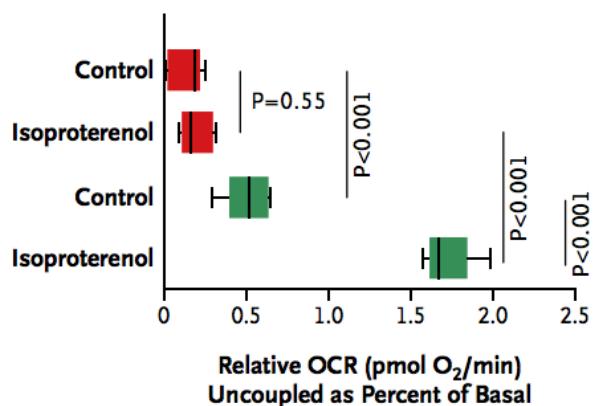
A



CRISPR-CAS9 rescue of homozygous risk allele restores oxygen consumption in adipocytes from carriers

rs1421085 (risk background)

- CC risk allele
- CC→TT rescue (CRISPR–Cas9 editing)



Limitations and considerations

Enrichment methods

Bisulfite-based methods

- Incomplete bisulfite conversion
 - Leads to increased methylation signal
 - Identification
 - Control probes in microarray
 - Methylation of non-CpG sites
 - Methylation-specific PCR of non-CpG sites

Affinity-based methods

- Highly dependent on antibody used
 - Include input control and other appropriate controls

Measurement

Array

- Sequence specific bias within reference genome
- Neighboring SNPs may interfere with binding

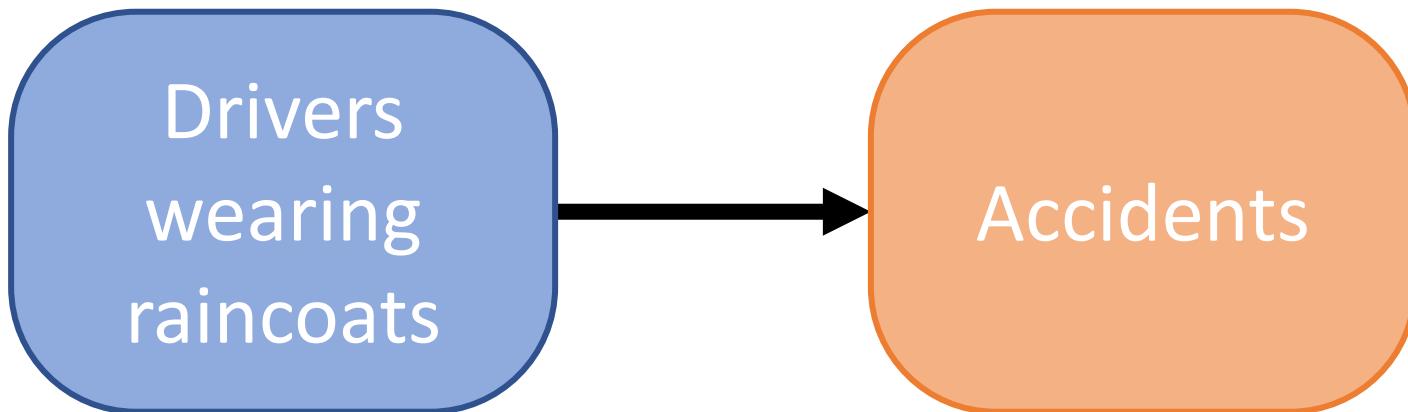
Sequencing

- Alignment
 - Repetitive regions are difficult to align
 - WGBS aligners have different biases depending on reference genome
 - Operator can also change parameters such as mismatch penalty
- Uneven depth across genome, may bias peak calling
 - Include input DNA for capture-based methods

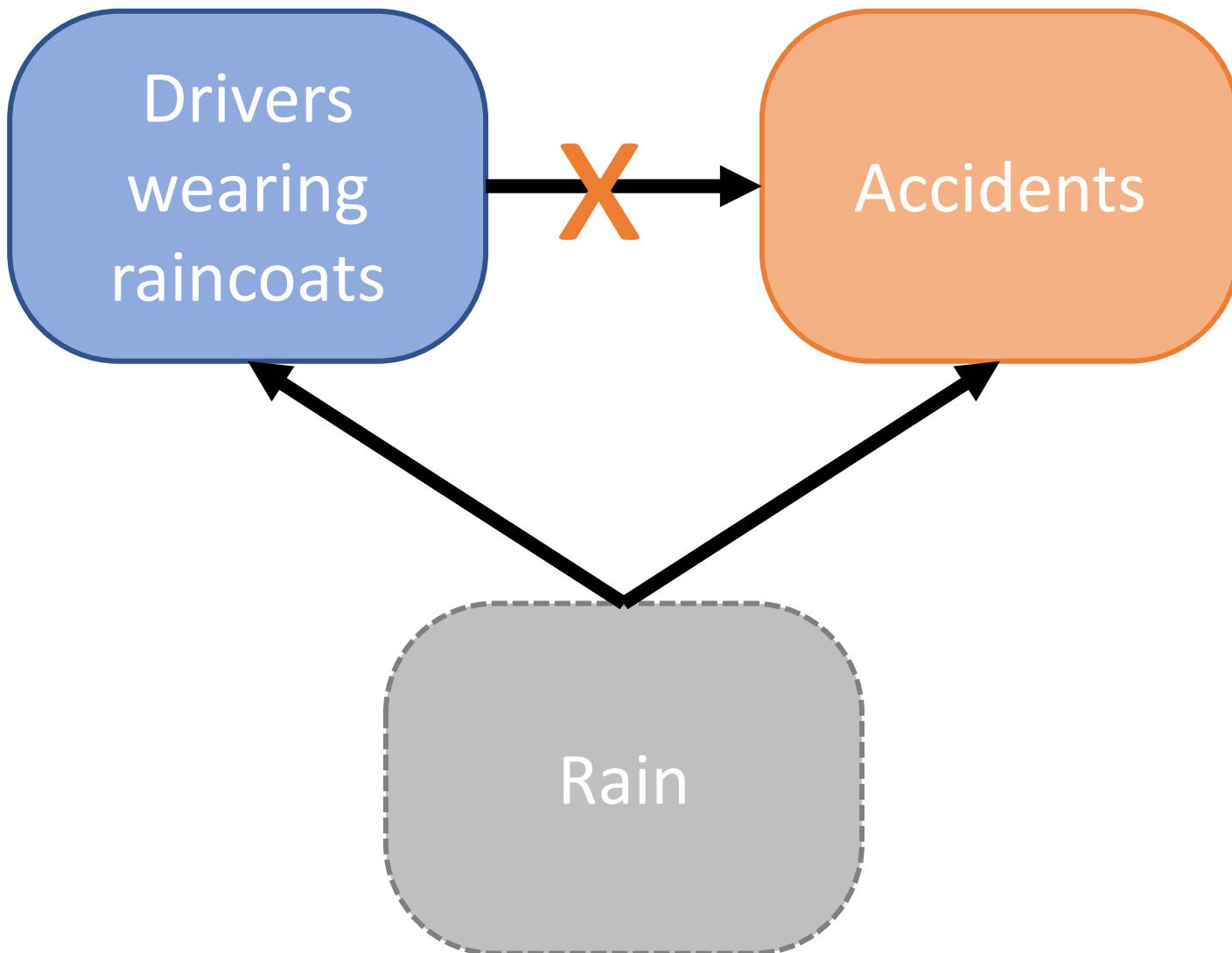
General statistical & design considerations

- Correlation does not mean causation
 - Experiments are required to establish causality
- Balanced study
- Replication using an independent cohort is important
 - Especially for biomarker or clustering studies
 - Training and validation cohorts should be representative of each other
- Orthogonal validation of sequencing or microarray results
- Accounting for potential sources of variation

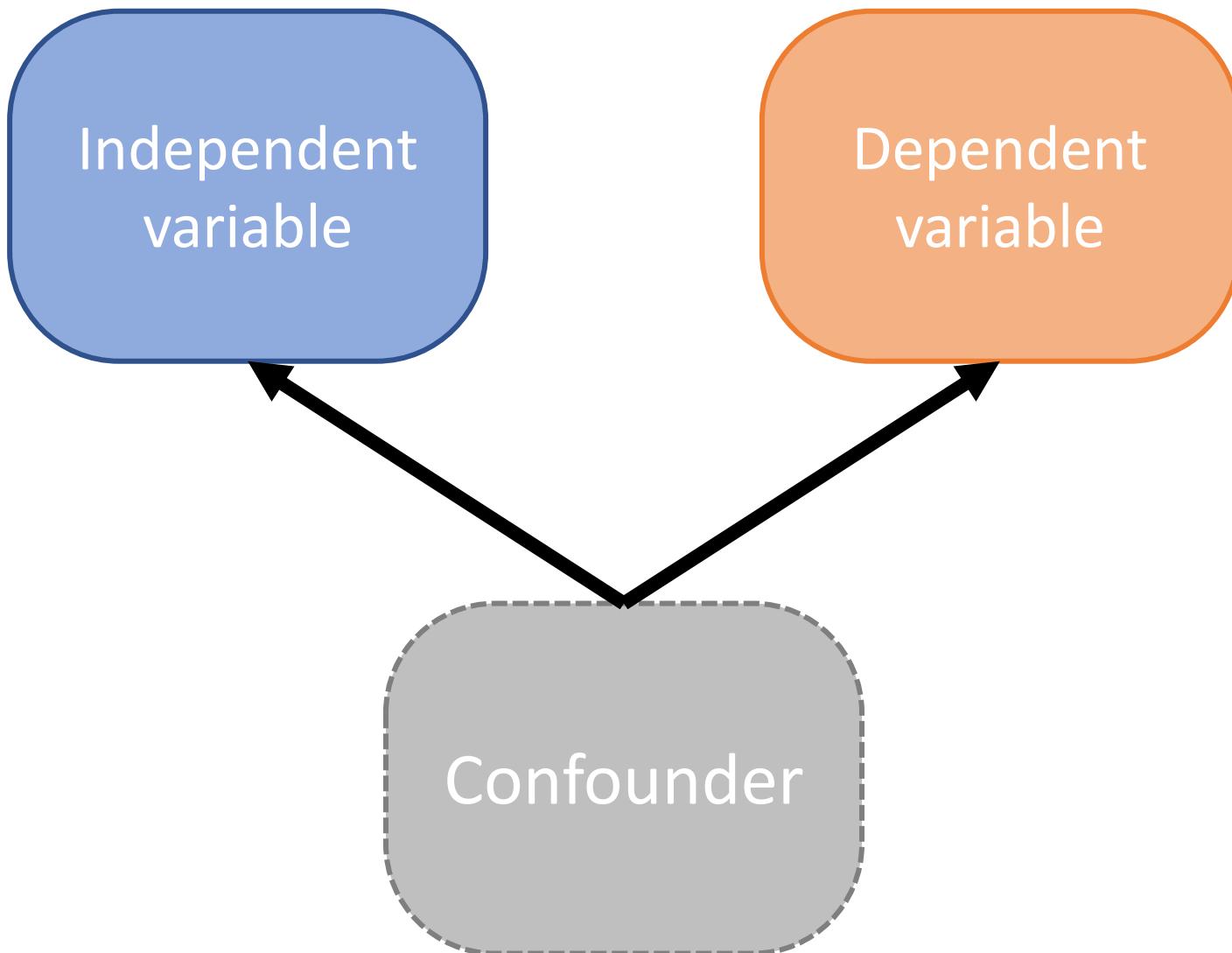
Confounding



Confounding

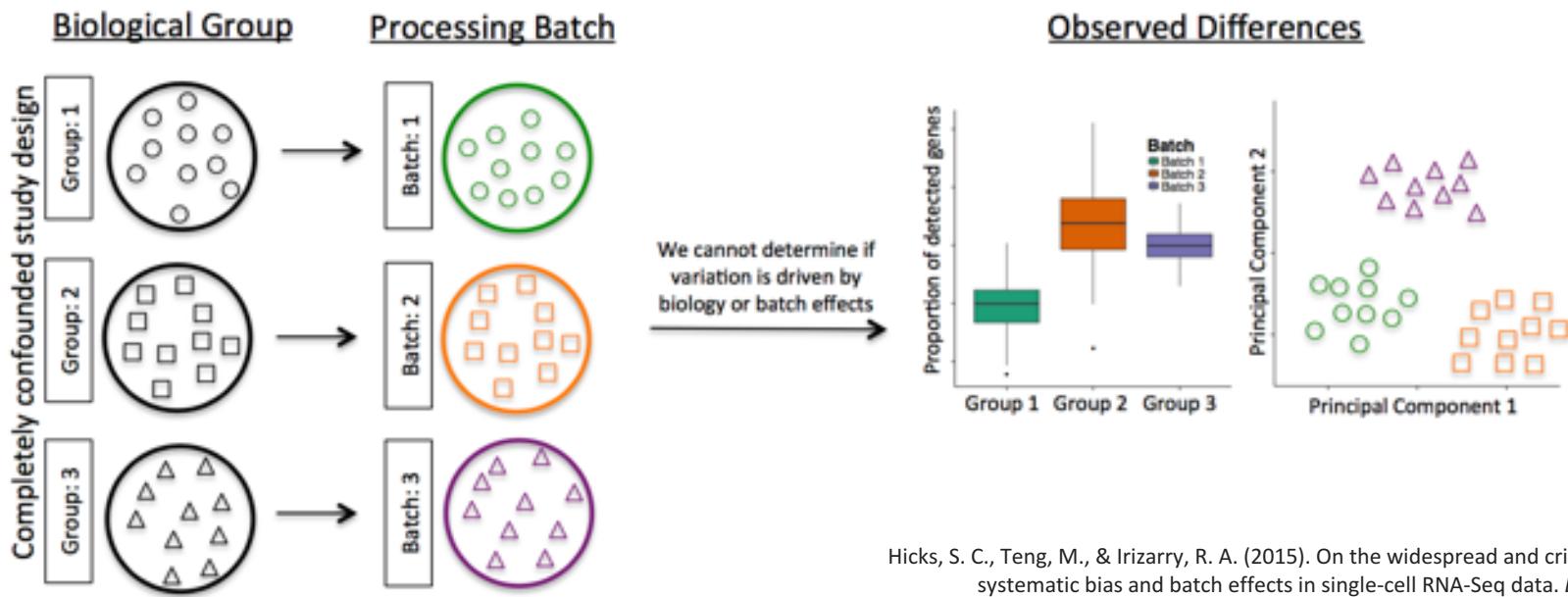


Confounding



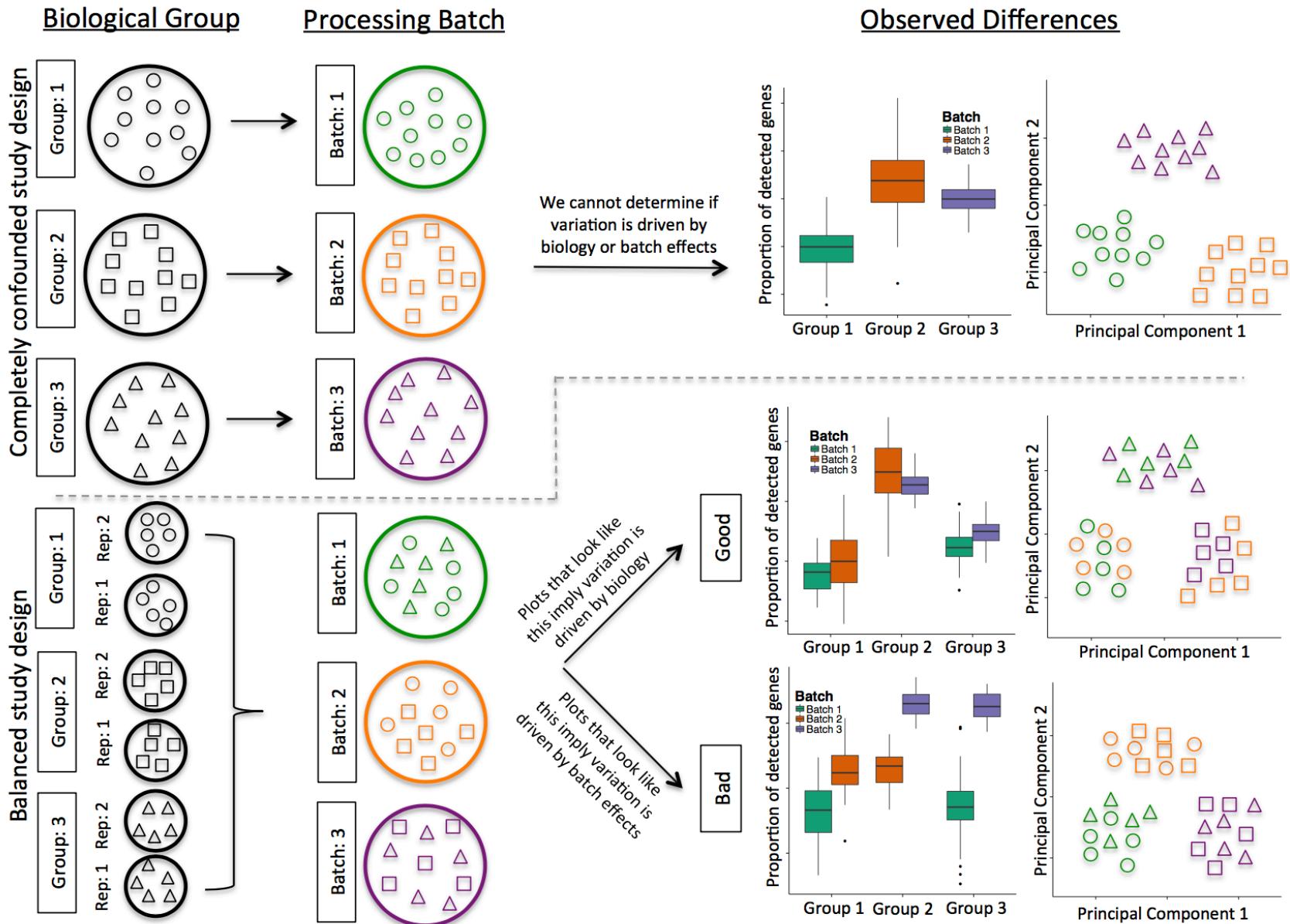
Batch as a confounder

Disease Date
Treatment Operator
Handedness Plate



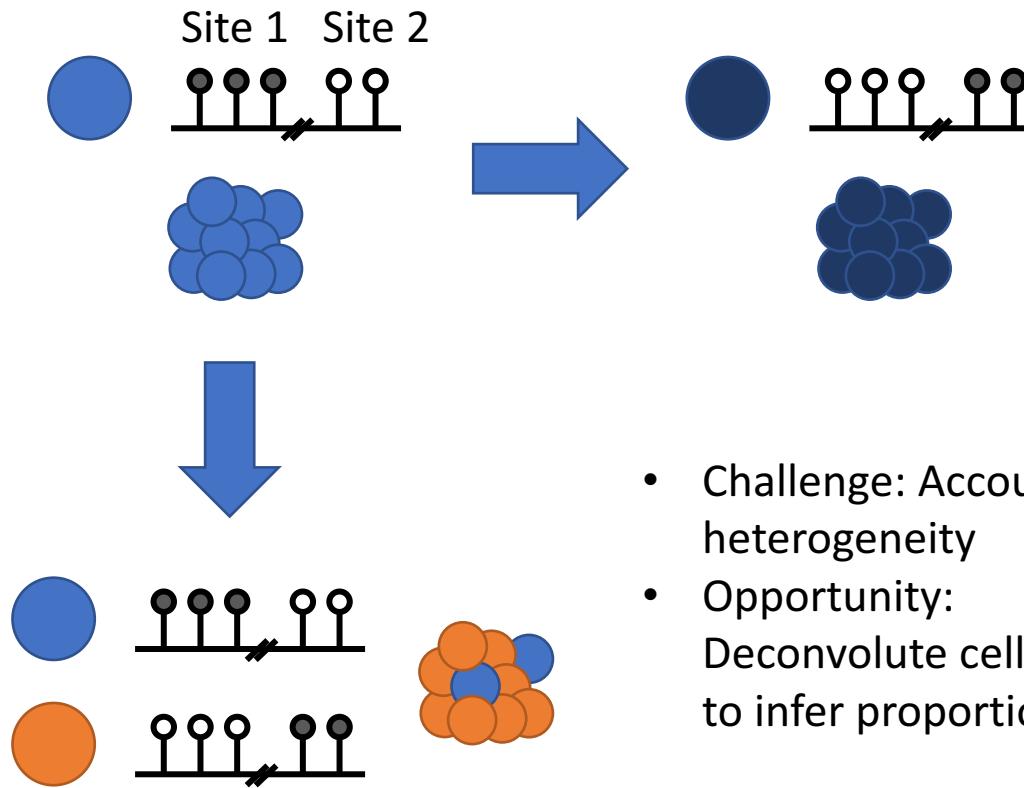
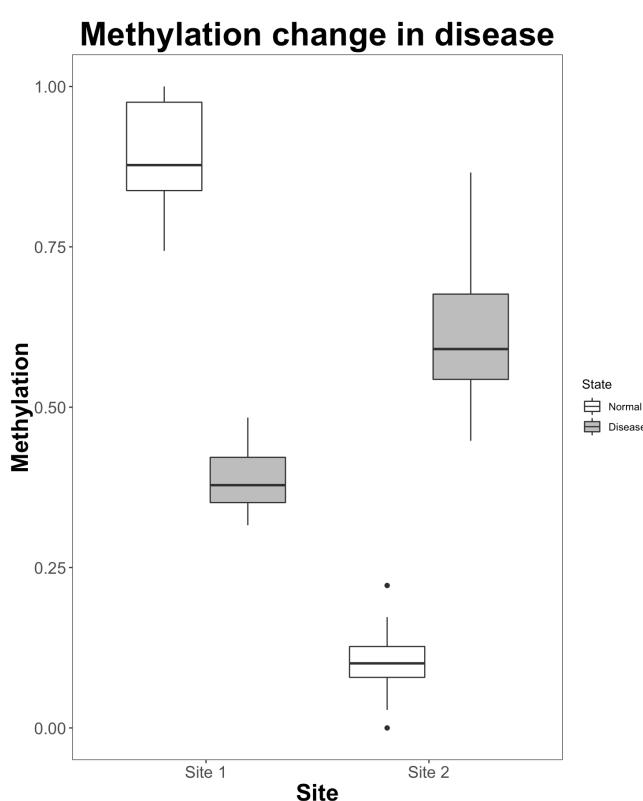
- Study design: randomization or blocked design
- Batch correction: estimate signals contributed from known batches and account in model

Batch as a confounder



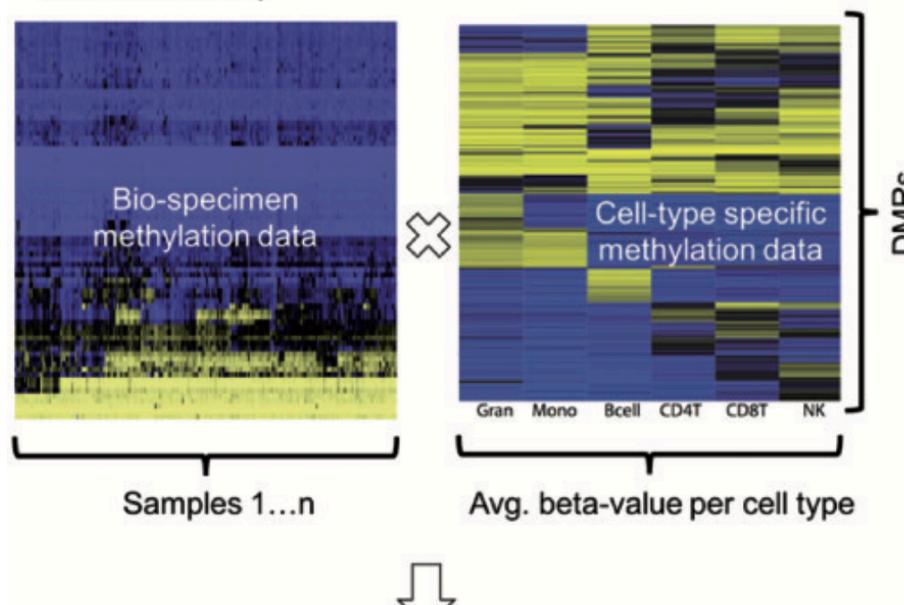
Heterogeneity

- Cell methylation exists in 3 states:
 - Unmethylated
 - Hemi-methylated
 - Methylated
- NOTE: Measurement includes a 4th state; missing data
- How are we measuring methylation as a continuous value from 0 – 100%?
- We are measuring a heterogeneous admixture of cells



- Challenge: Account for heterogeneity
- Opportunity: Deconvolute cell types to infer proportion

1. Reference-based cell type deconvolution using immune cell DMRs (e.g. Houseman method)



The resulting matrix shows the estimated immune cell type proportions for each sample. The columns represent different cell types: Gran, Mono, B-cell, CD4T, CD8T, and NK. The rows represent individual samples labeled Gran₁, Mono₁, B-cell₁, CD4T₁, CD8T₁, NK₁; Gran₂, Mono₂, B-cell₂, CD4T₂, CD8T₂, NK₂; ..., ..., ..., ..., ..., ...; and Gran_n, Mono_n, B-cell_n, CD4T_n, CD8T_n, NK_n.

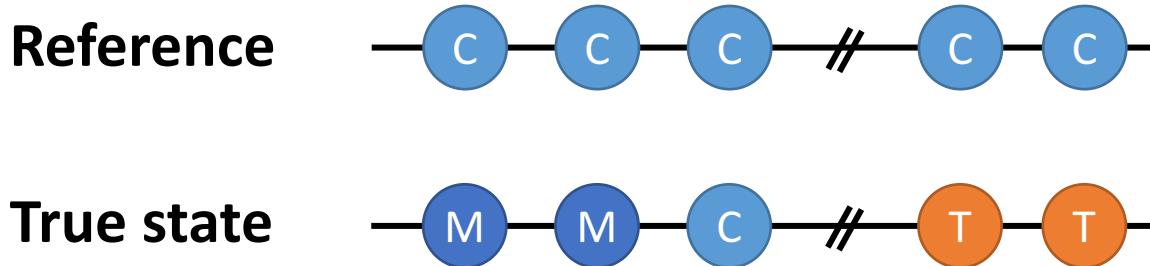
Samples 1...n		Gran	Mono	B-cell	CD4T	CD8T	NK
Gran ₁	Mono ₁	B-cell ₁	CD4T ₁	CD8T ₁	NK ₁		
Gran ₂	Mono ₂	B-cell ₂	CD4T ₂	CD8T ₂	NK ₂		
...
Gran _n	Mono _n	B-cell _n	CD4T _n	CD8T _n	NK _n		

Immune proportion estimates for samples 1...n

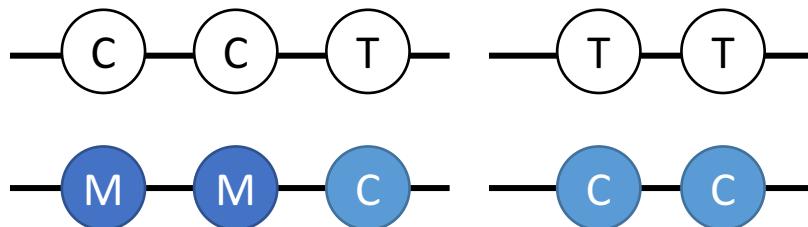
Emerging technologies in epigenomics.

Nanopore sequencing

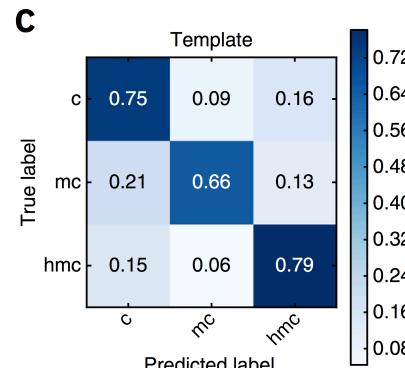
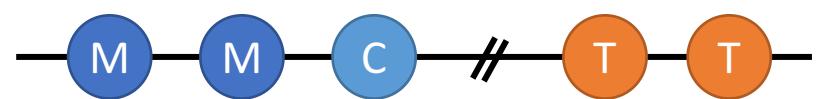
Allows for the detection of 5mC and 5hmC without bisulfite treatment and distinguishes T's from unmethylated C's.



Bisulfite sequencing

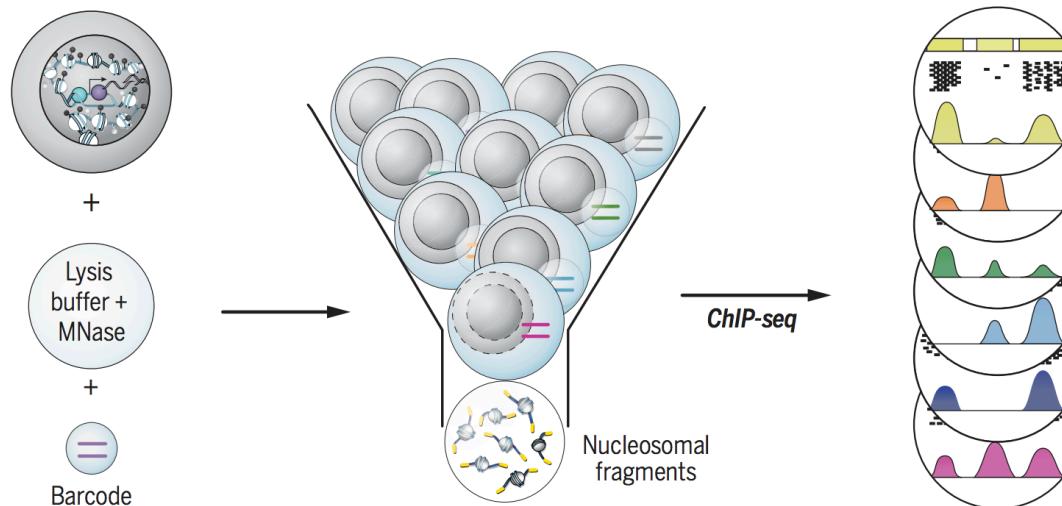


Nanopore sequencing

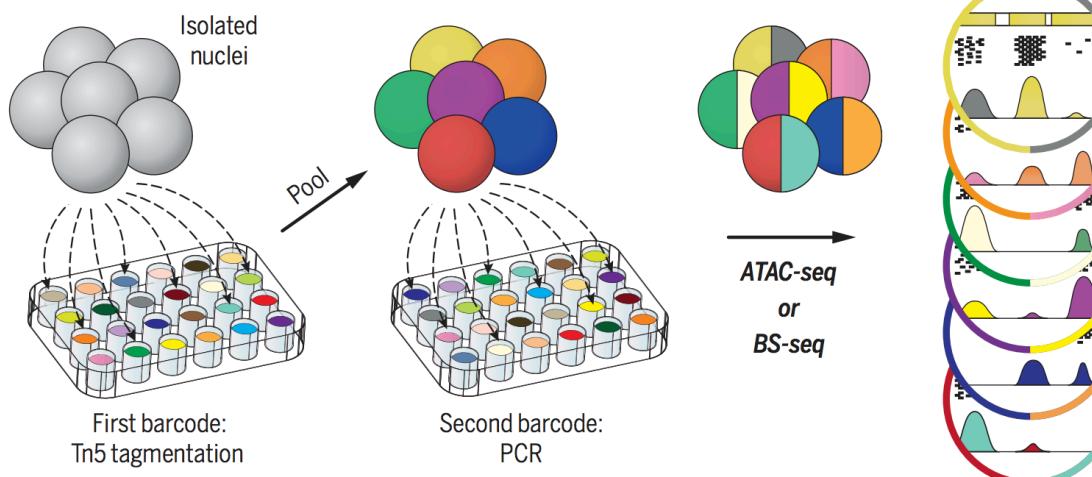


Single cell epigenomics

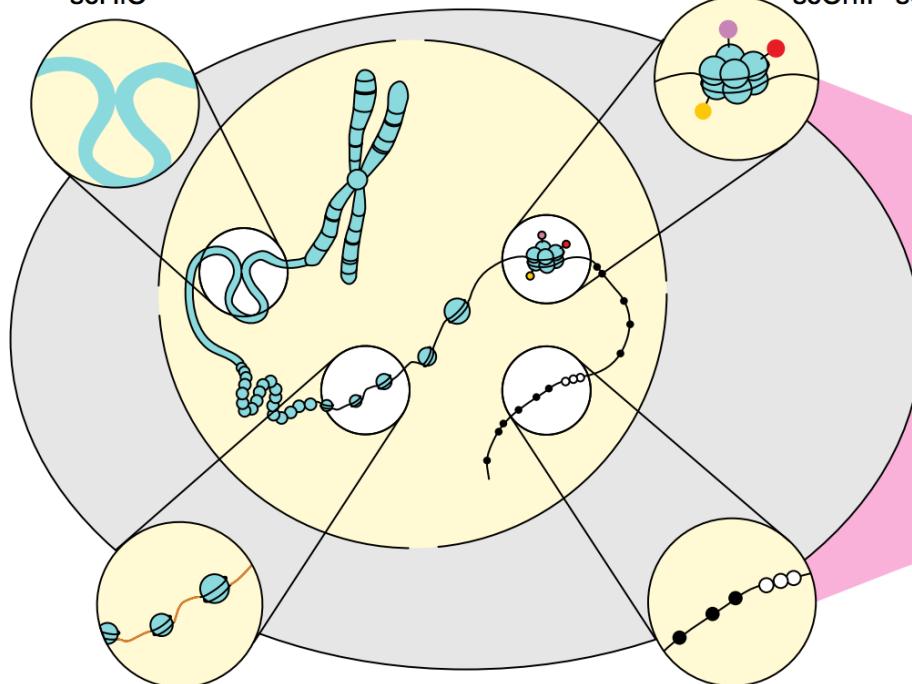
Droplet barcoding



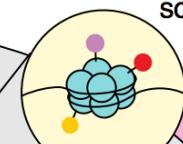
Combinatorial barcoding



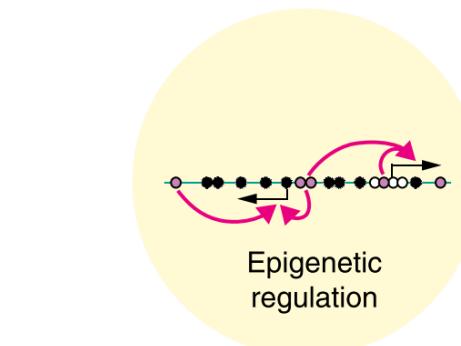
Chromosome conformation
scHiC



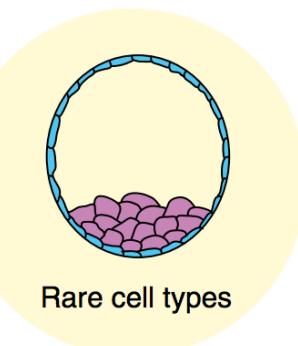
Histone modifications
scChIP-seq



DNA accessibility
scATAC-seq, scDNase-seq



DNA modifications
scRRBS, scBS-seq



Transcriptomics

Epigenomics

Genomics

Integrated multi-Omics

Thank you.

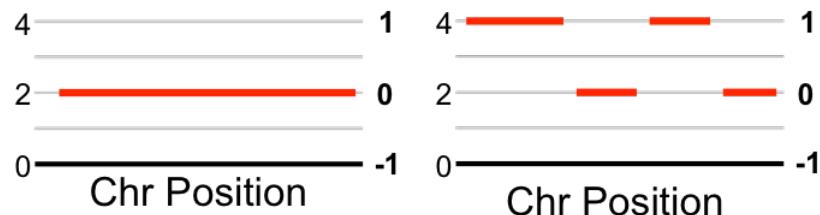
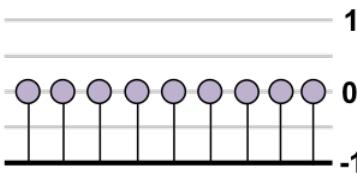
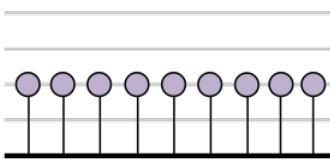
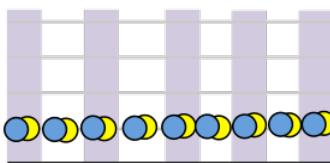
QUESTIONS?

APPENDIX

SNP Array

Normal Sample

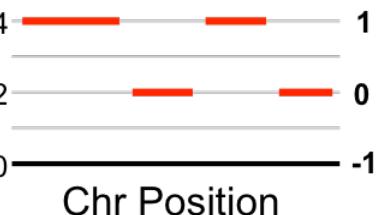
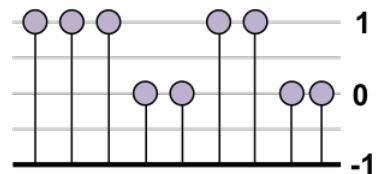
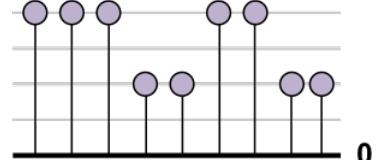
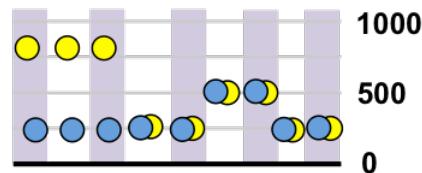
- Allele A = A_i
- Allele B = B_i



SNP Array

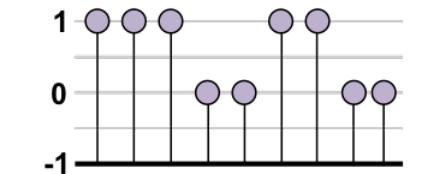
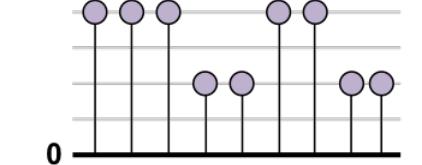
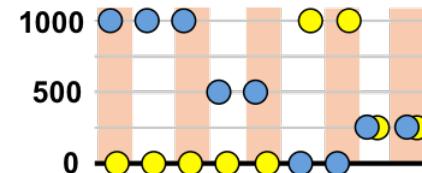
Tumor Sample

- Allele A = A_i
- Allele B = B_i



Methylation Array

- Unmethylated, T = A_i
- Methylated, C = B_i



Raw signal intensity:
 A_i, B_i

↓
 Observed copy number:
 $R_i = A_i + B_i$

↓
 Log R Ratio:
 $LRR = \log(R_i/R')$

↓
 Segmented LRR