

# **DNA methylation**

## study design, bench work, and statistical approaches

Amanda Lea  
September 22, 2016

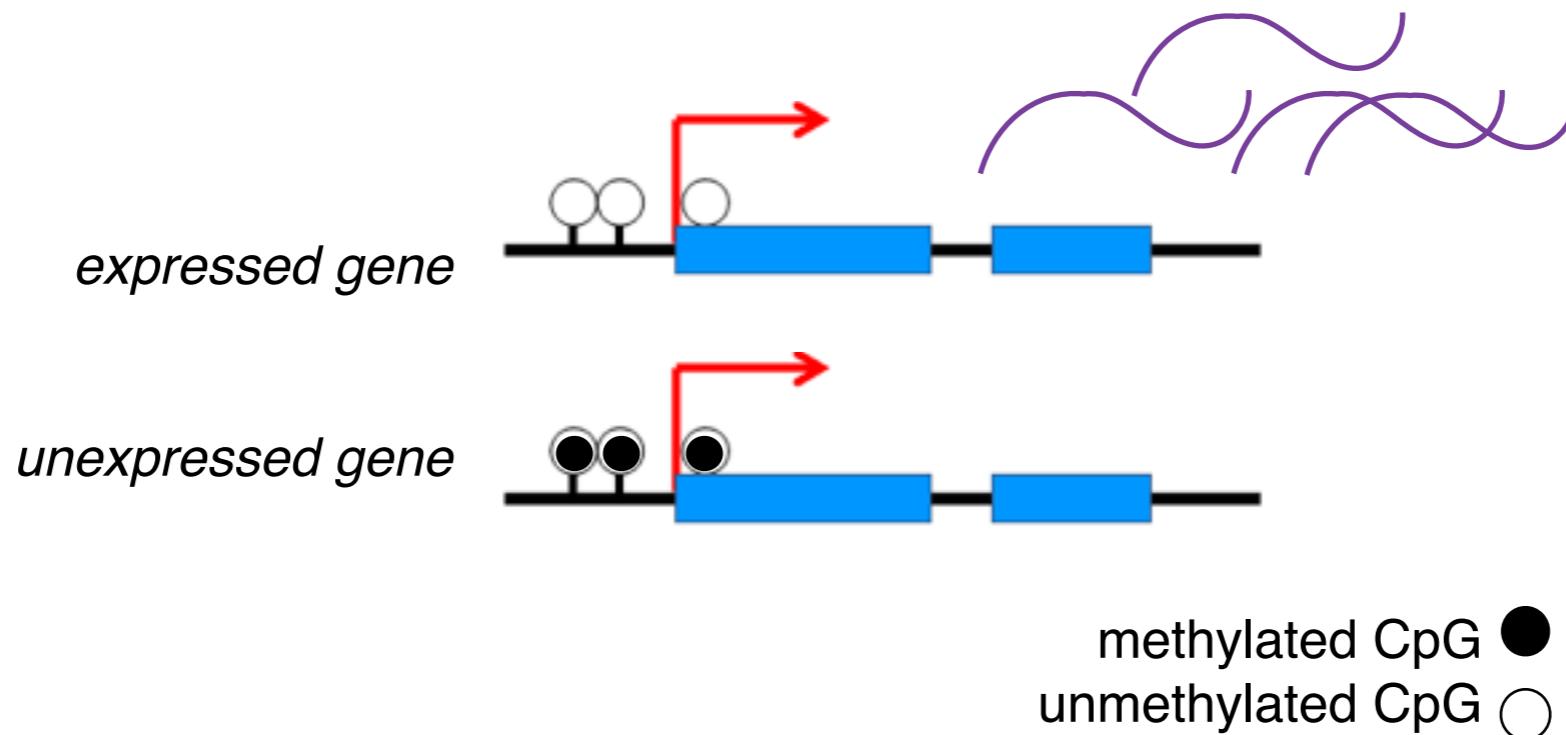


## Some things we'll talk about:

- What is DNA methylation? Why do we care about it?
- How can we measure genome-wide DNA methylation (in a cost effective way, in non-model organisms)?
- What do (sequencing-based) DNA methylation data look like? What are some particular challenges in working with these data?
- What models have been proposed to analyze differential methylation (and what are their relative shortcomings)?

# DNA methylation

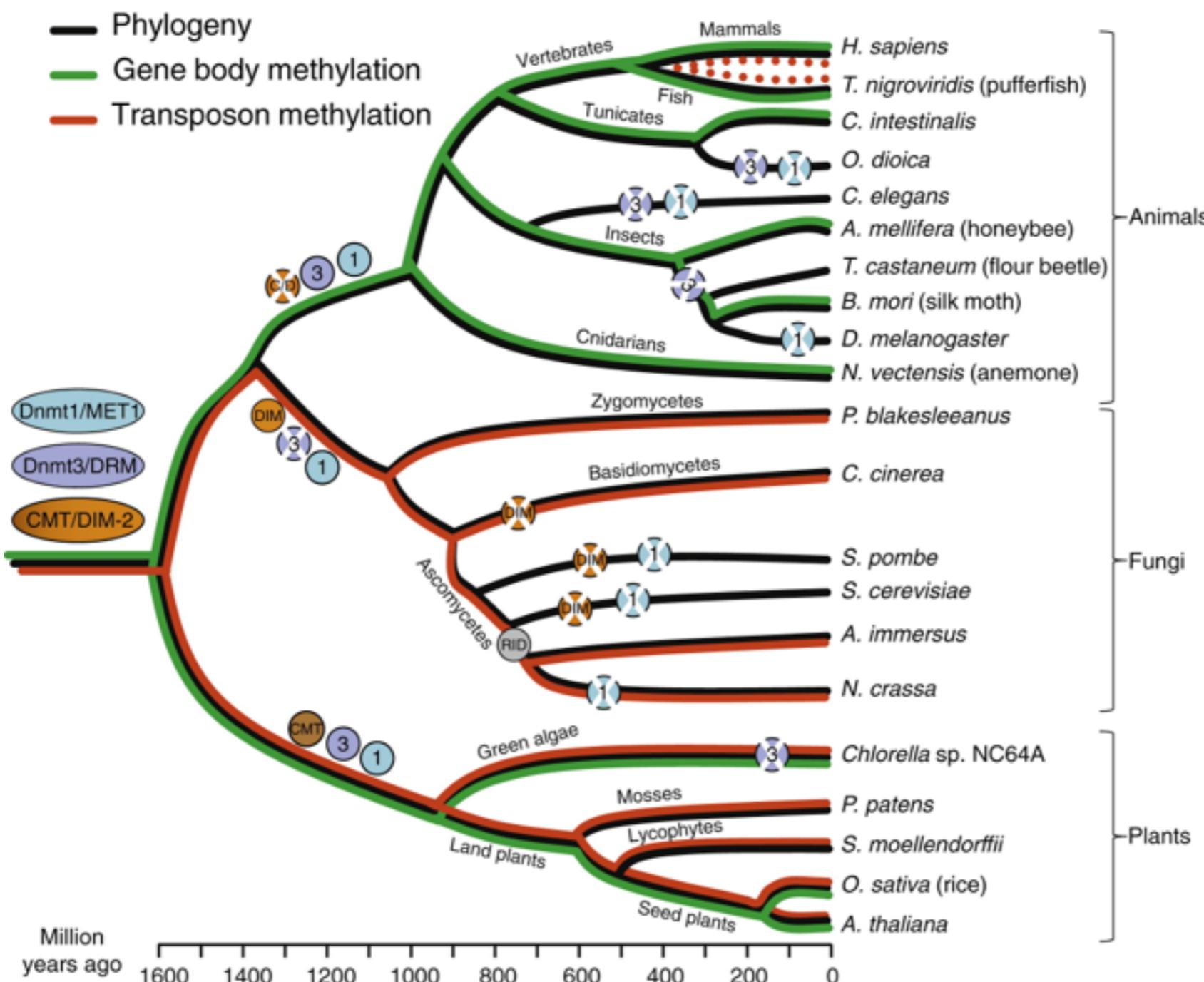
- one of many epigenetic gene regulatory mechanisms
- covalent addition of methyl groups to cytosine bases
- occurs at ‘CpG sites’ in vertebrates
- enriched in/near genes and regulatory regions
- key functions = X inactivation, suppression of transposable elements, imprinted genes



To review, differences in DNA methylation are associated with:

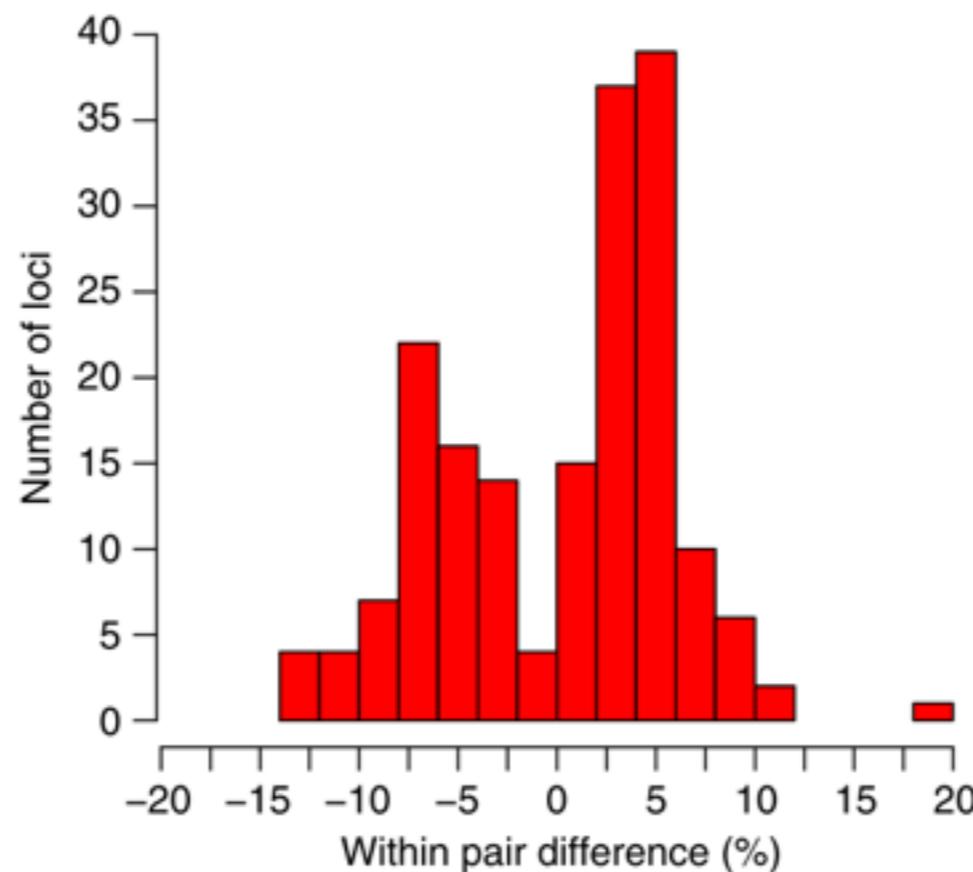
To review, differences in DNA methylation are associated with:

- evolutionary history, species differences



To review, differences in DNA methylation are associated with:

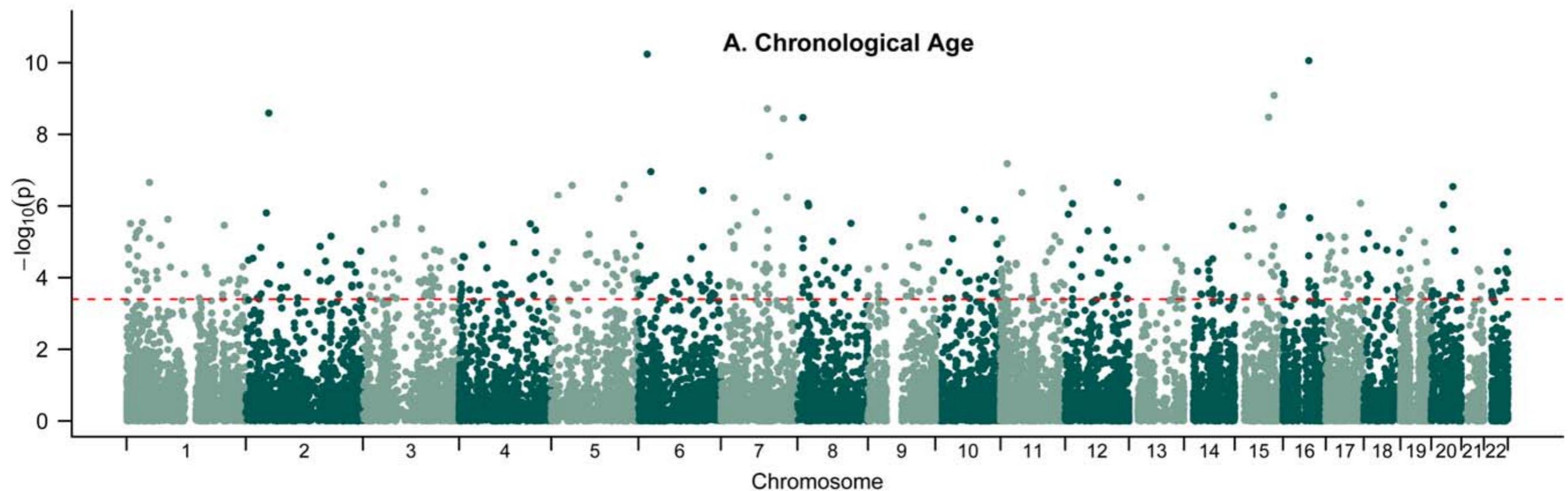
- evolutionary history, species differences
- environmental exposures



**Figure 2 | The average within-pair difference for the 181 regions associated with prenatal famine exposure after correction for multiple testing.** A histogram for the average within pair difference (%) between the exposed and unexposed siblings. A positive number reflects relative higher DNA methylation levels in the exposed.

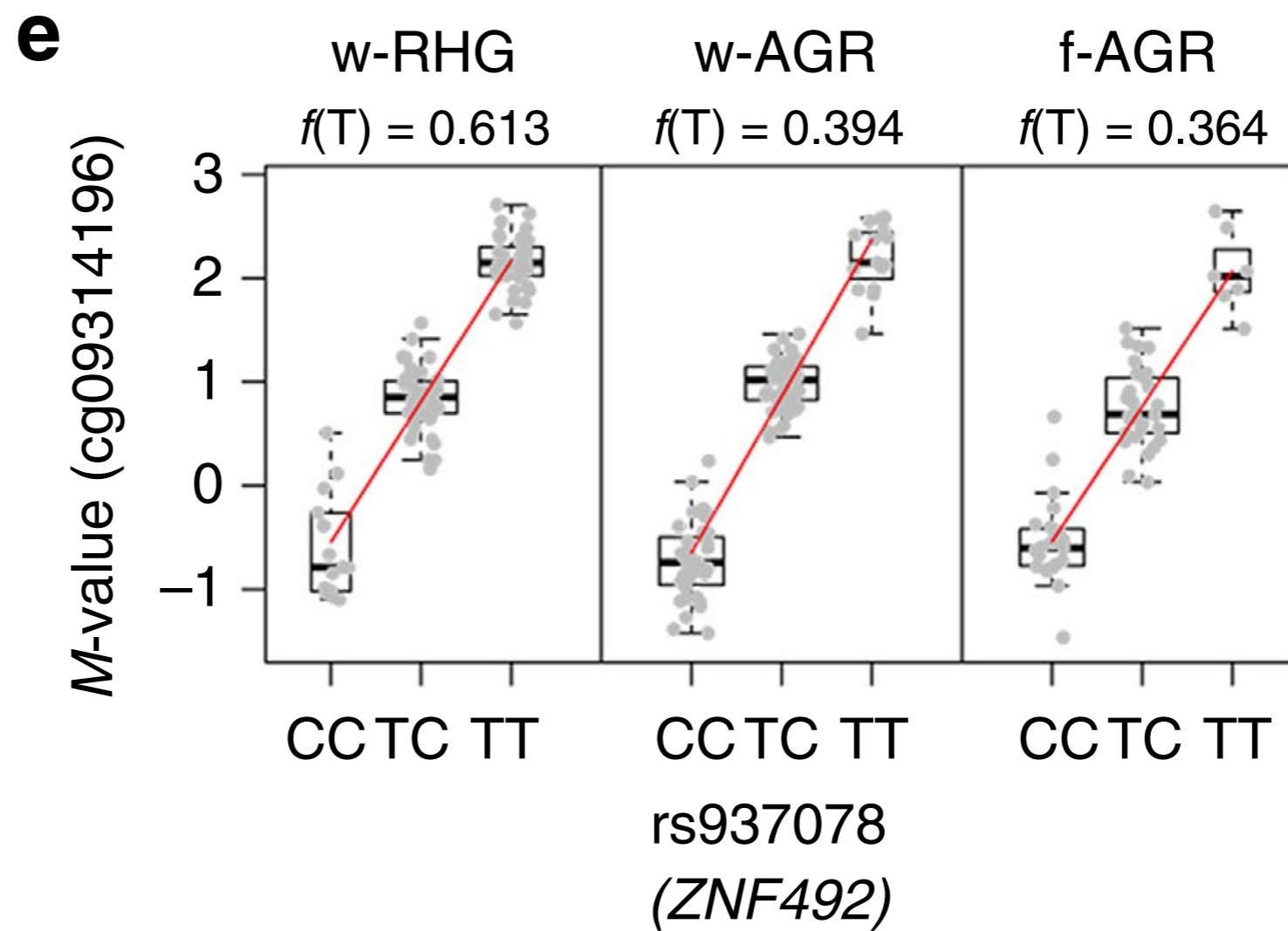
To review, differences in DNA methylation are associated with:

- evolutionary history, species differences
- environmental exposures
- age (also life history or developmental stage)



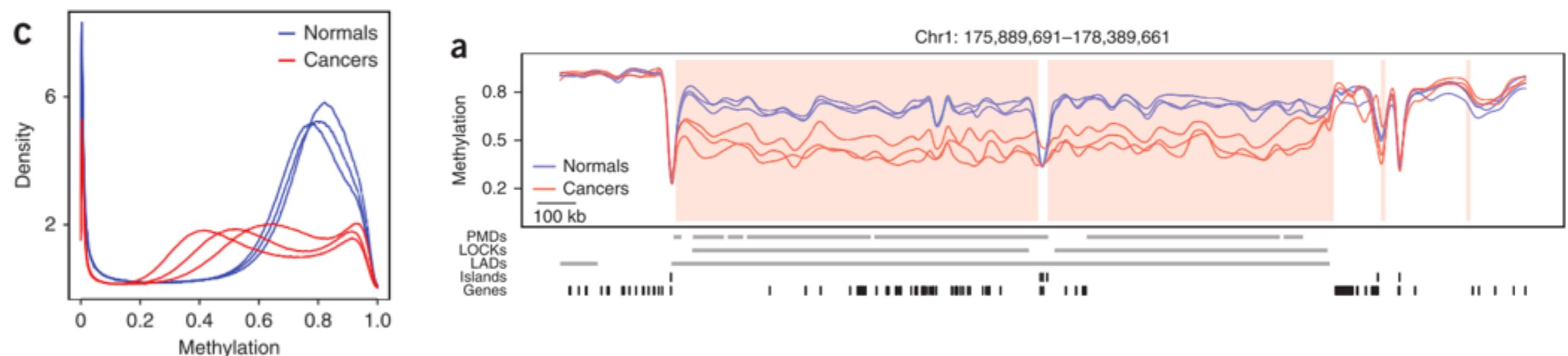
To review, differences in DNA methylation are associated with:

- evolutionary history, species differences
- environmental exposures
- age (also life history or developmental stage)
- genetic variation



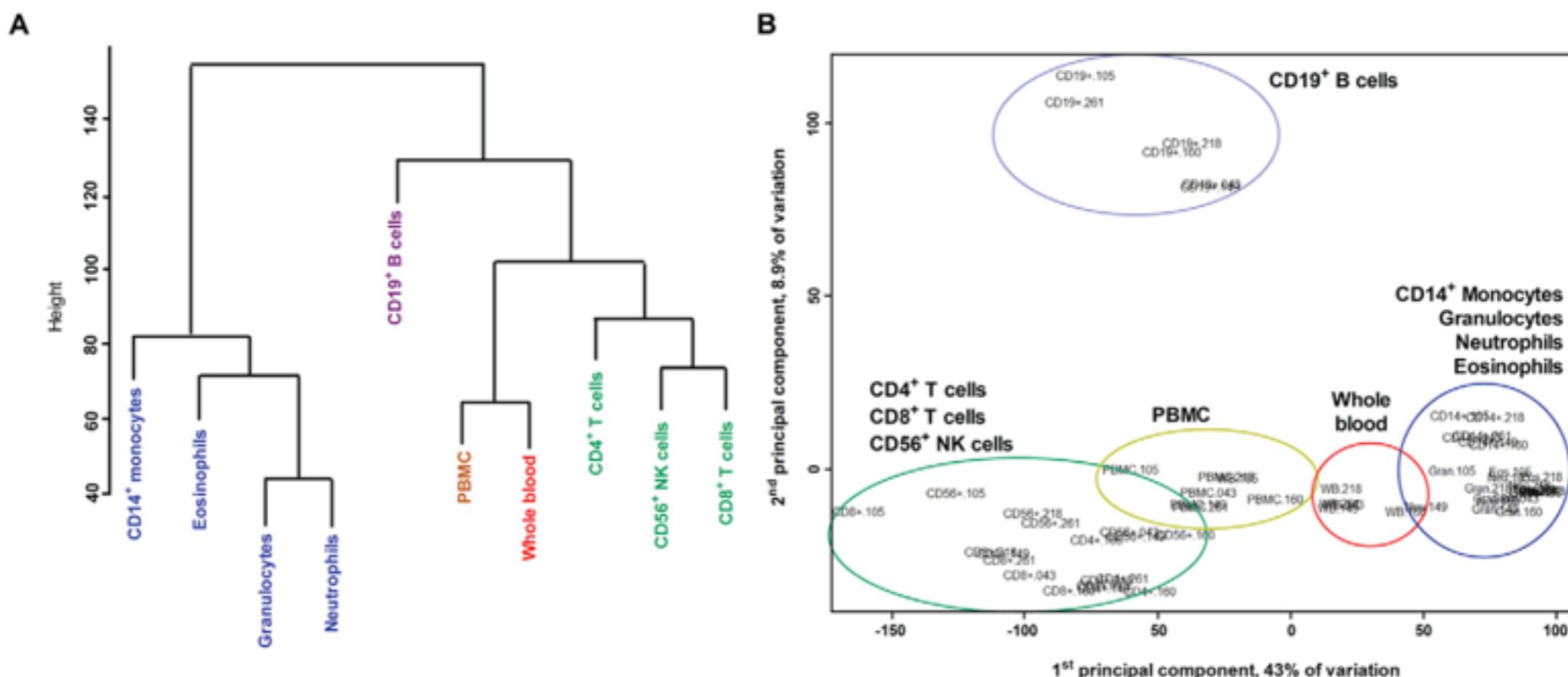
To review, differences in DNA methylation are associated with:

- evolutionary history, species differences
- environmental exposures
- age (also life history or developmental stage)
- genetic variation
- disease status



To review, differences in DNA methylation are associated with:

- evolutionary history, species differences
- environmental exposures
- age (also life history or developmental stage)
- genetic variation
- disease status
- cell type identity/composition



To review, differences in DNA methylation are associated with:

- evolutionary history, species differences
- environmental exposures
- age (also life history or developmental stage)
- genetic variation
- disease status
- cell type identity/composition

***Interesting!***

***Also, important to account for these factors in your study!***

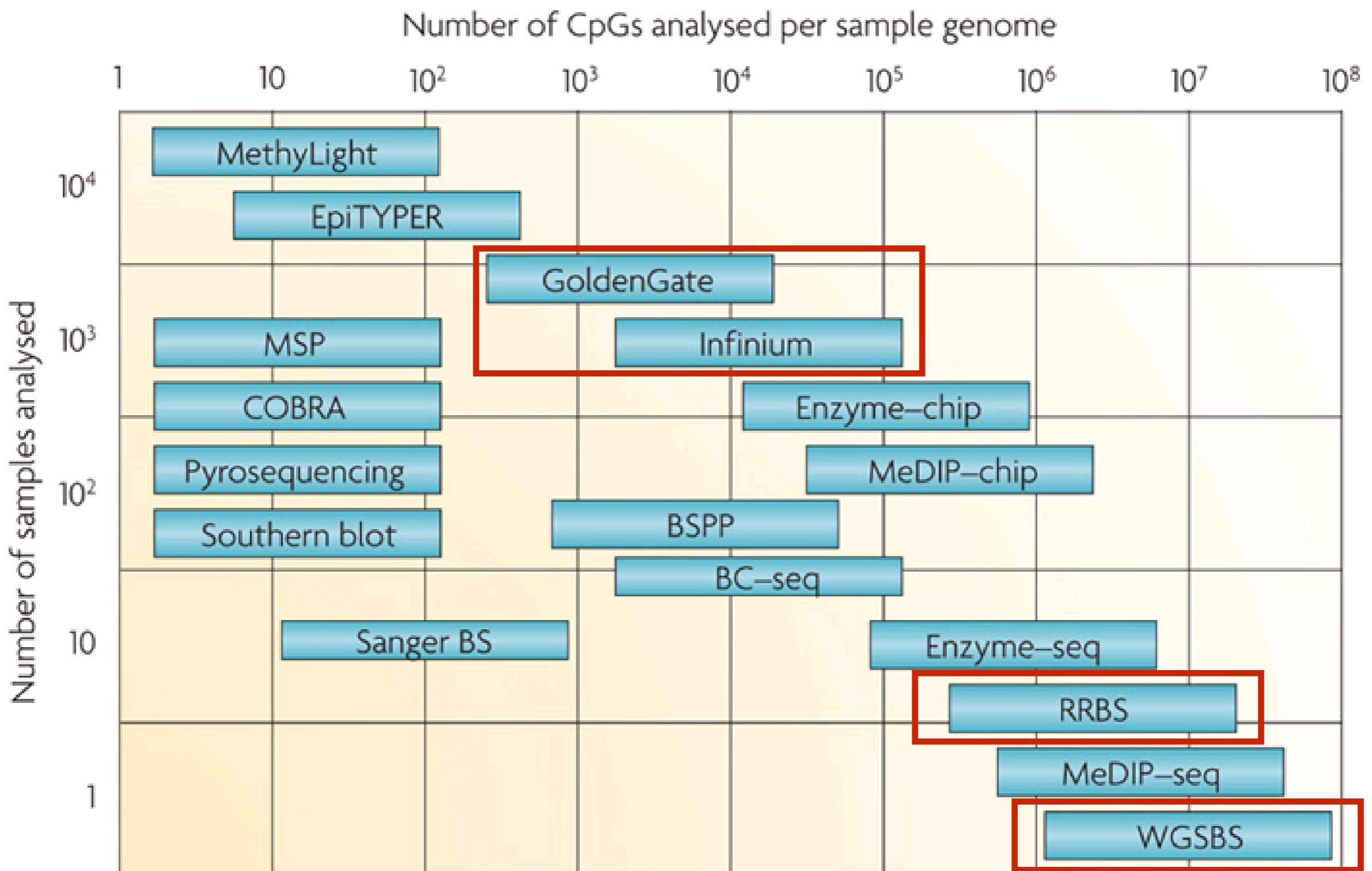
We're talking about this kind of DNA methylation study:

1. Get biological samples
2. Extract DNA
3. Measure DNA methylation levels at many loci using a high-throughput approach
4. Test for associations between predictor of interest and DNA methylation levels at each measured site (controlling for other stuff)
5. Correct for multiple hypothesis testing (e.g., FDR)
6. Make some inference based on which genes/regions exhibit differential methylation

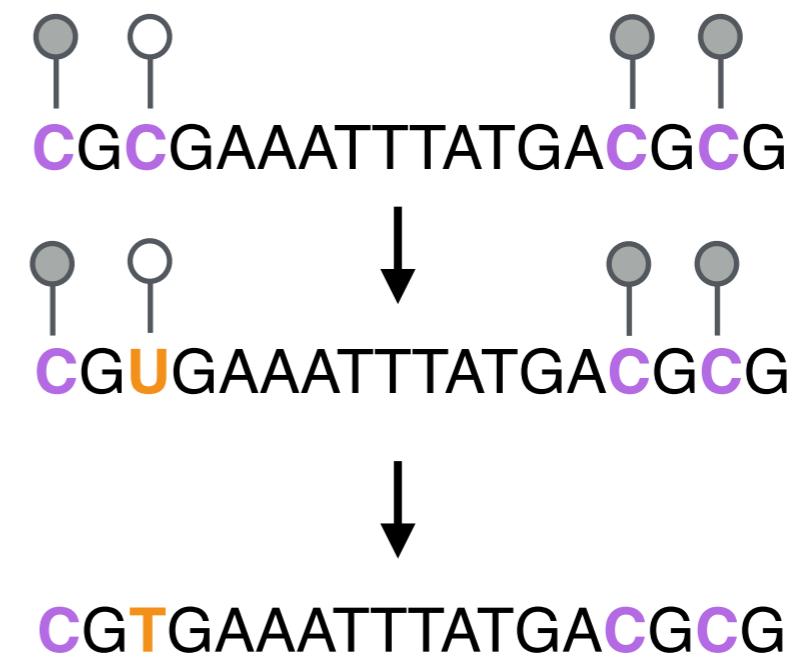
We're talking about this kind of DNA methylation study:

1. Get biological samples
2. Extract DNA
3. Measure DNA methylation levels at many loci using a high-throughput approach
4. Test for associations between predictor of interest and DNA methylation levels at each measured site (controlling for other stuff)
5. Correct for multiple hypothesis testing (e.g., FDR)
6. Make some inference based on which genes/regions exhibit differential methylation

# Ways to measure DNA methylation levels



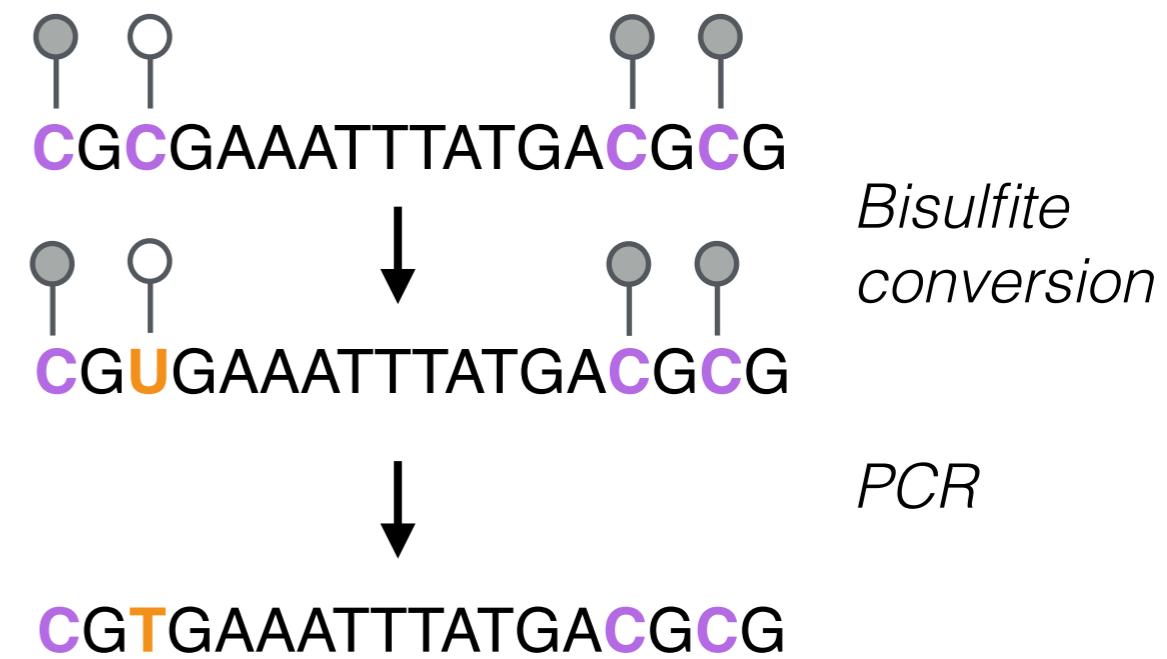
# Bisulfite sequencing



*Bisulfite  
conversion*

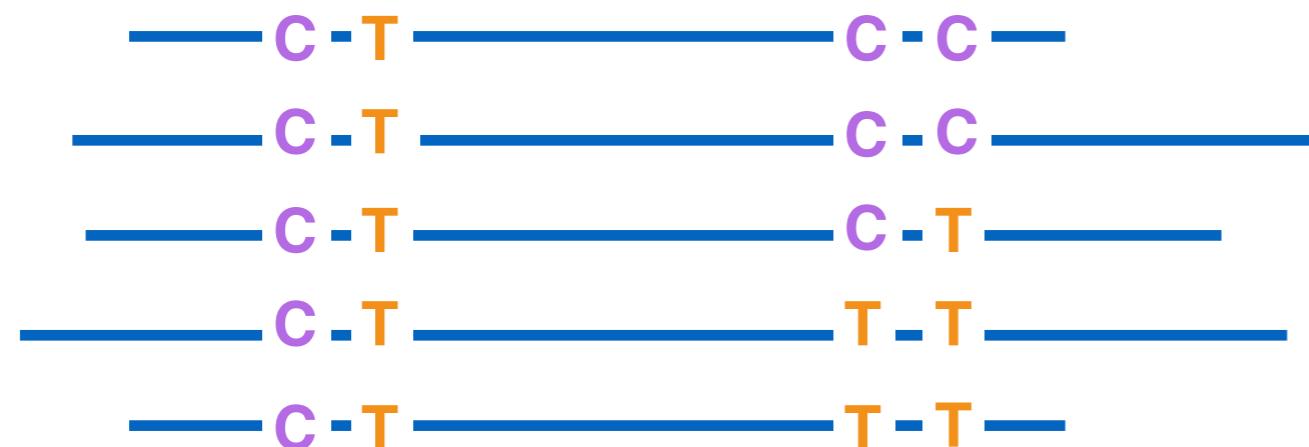
*PCR*

# Bisulfite sequencing



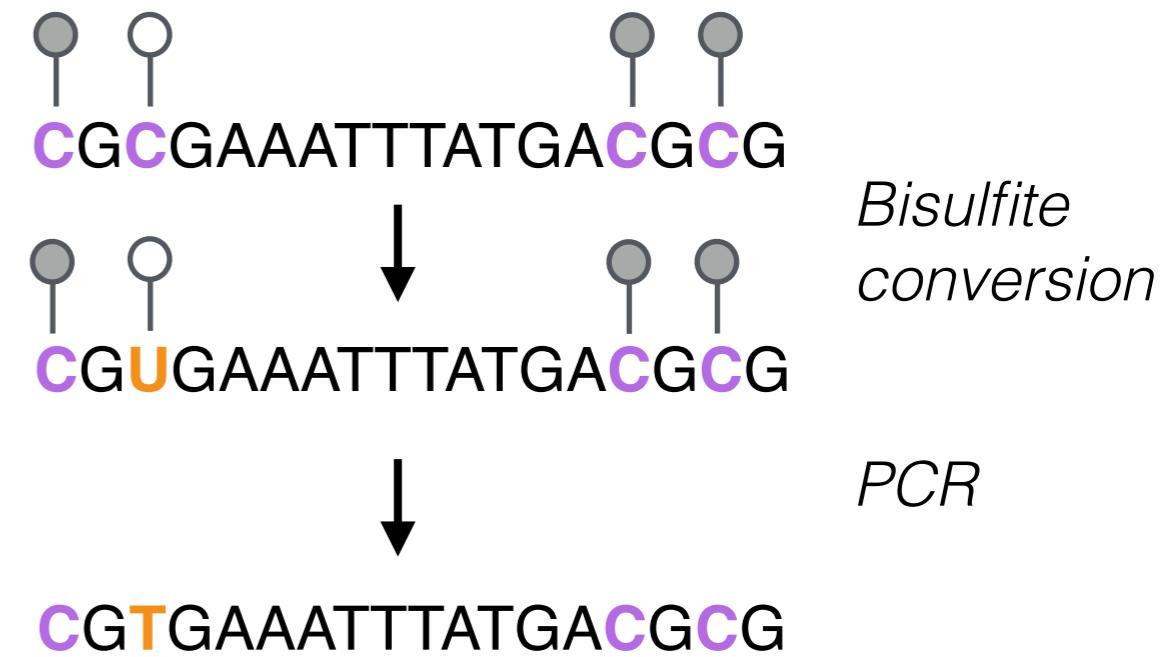
chr1

AAAATTACGCGAAATTATGACGCGAAAATTA



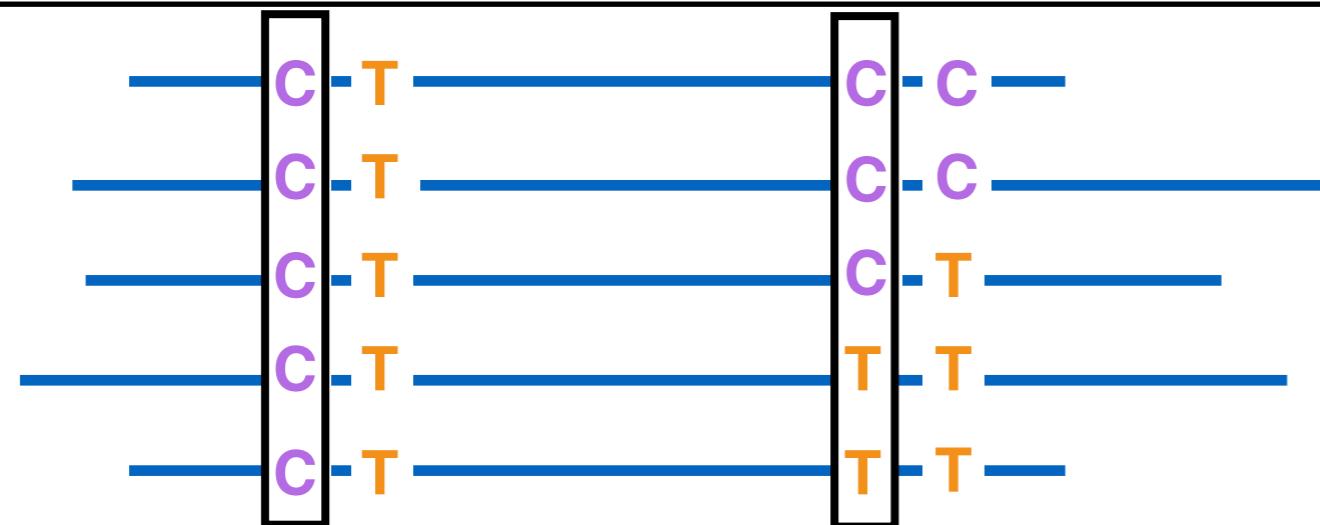
Reads from individual 1

# Bisulfite sequencing



chr1

AAAATTACGCGAAATTATGACGCGAAAATTA



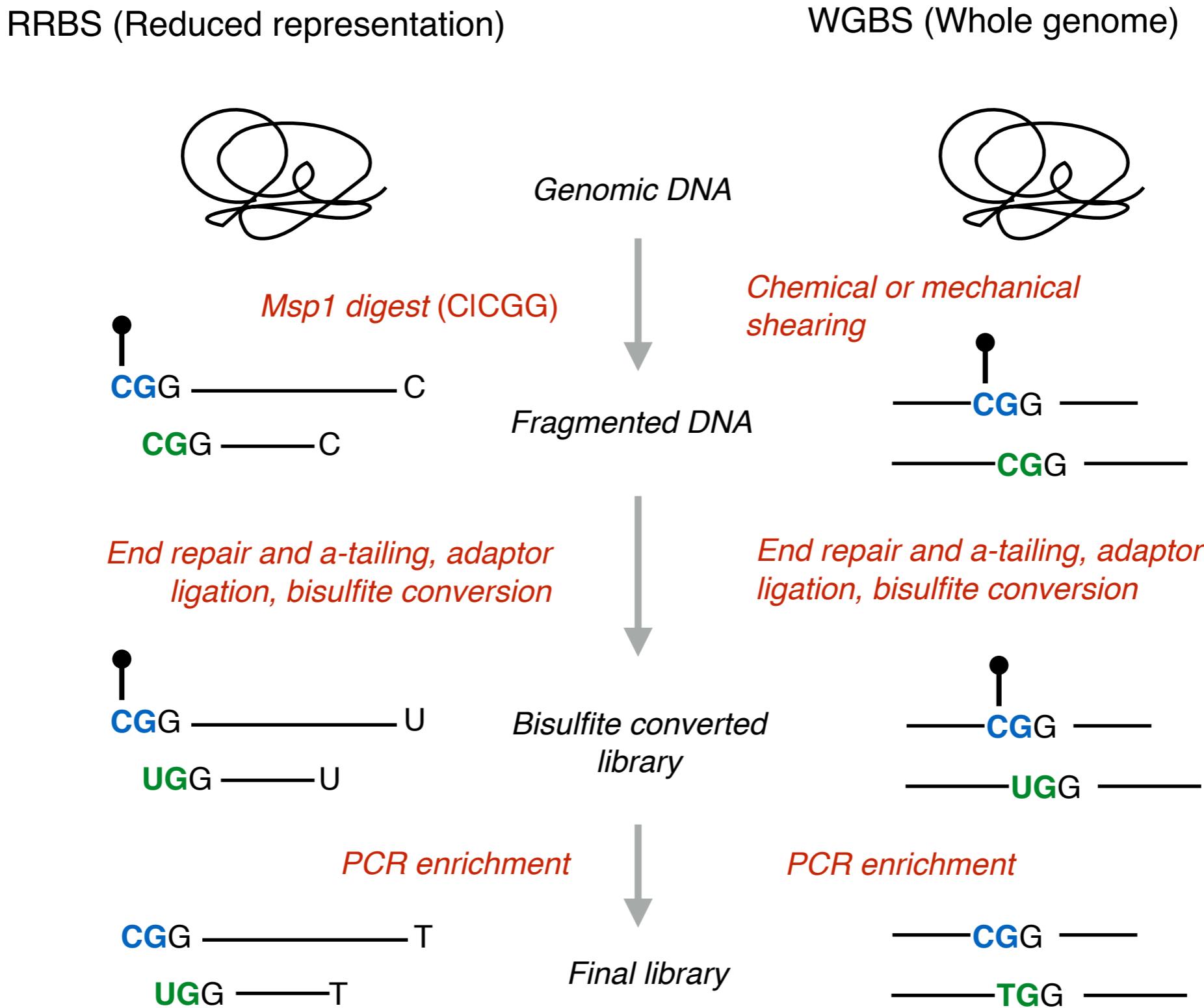
Reads from individual 1

5/5 = 100% methylated

3/5 = 60% methylated

DNA methylation level = C / (C+T)

# Bisulfite sequencing

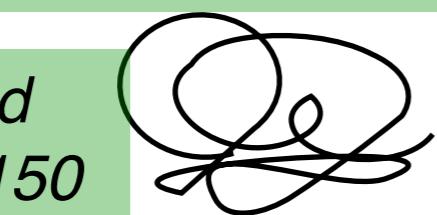


# Bisulfite sequencing

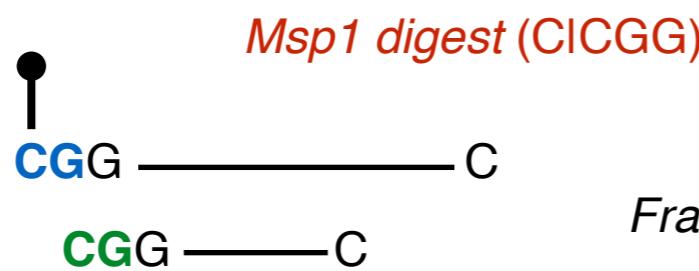
RRBS (Reduced representation)  
*a few mil CpG sites*

*for a primate-sized genome, cost <\$150*

*need <200ng of DNA*

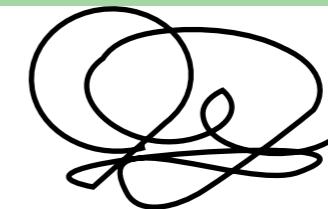


*Genomic DNA*

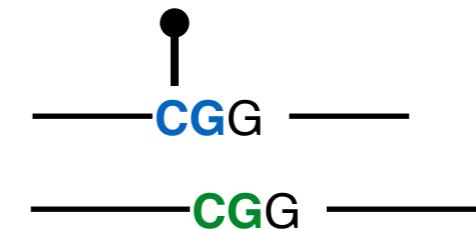


*Fragmented DNA*

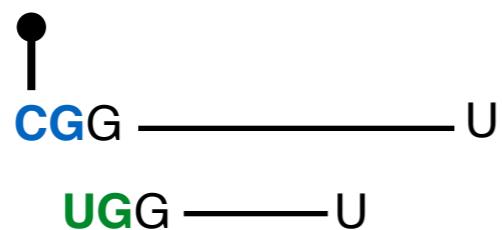
WGBS (Whole genome)  
*~28 mil CpG sites*



*Chemical or mechanical shearing*

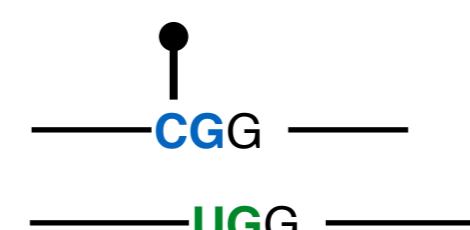


*End repair and a-tailing, adaptor ligation, bisulfite conversion*

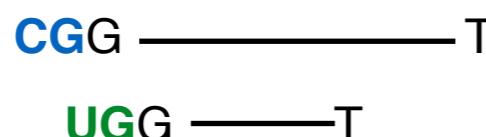


*Bisulfite converted library*

*End repair and a-tailing, adaptor ligation, bisulfite conversion*

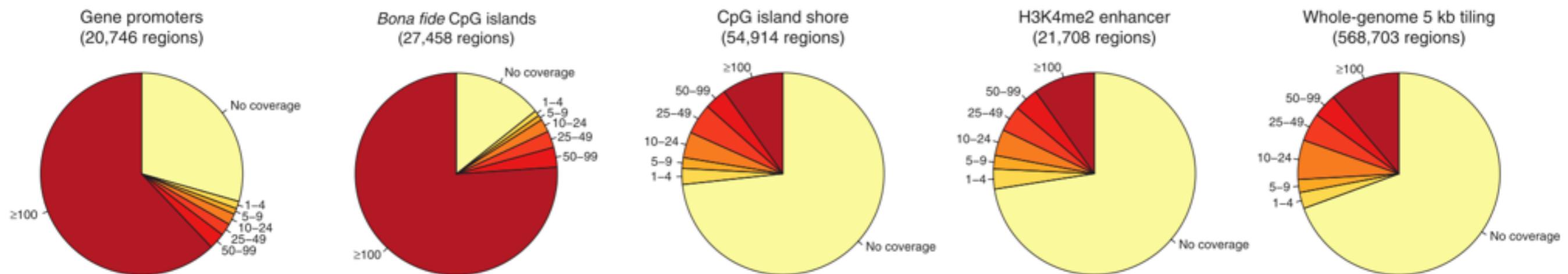


*PCR enrichment*

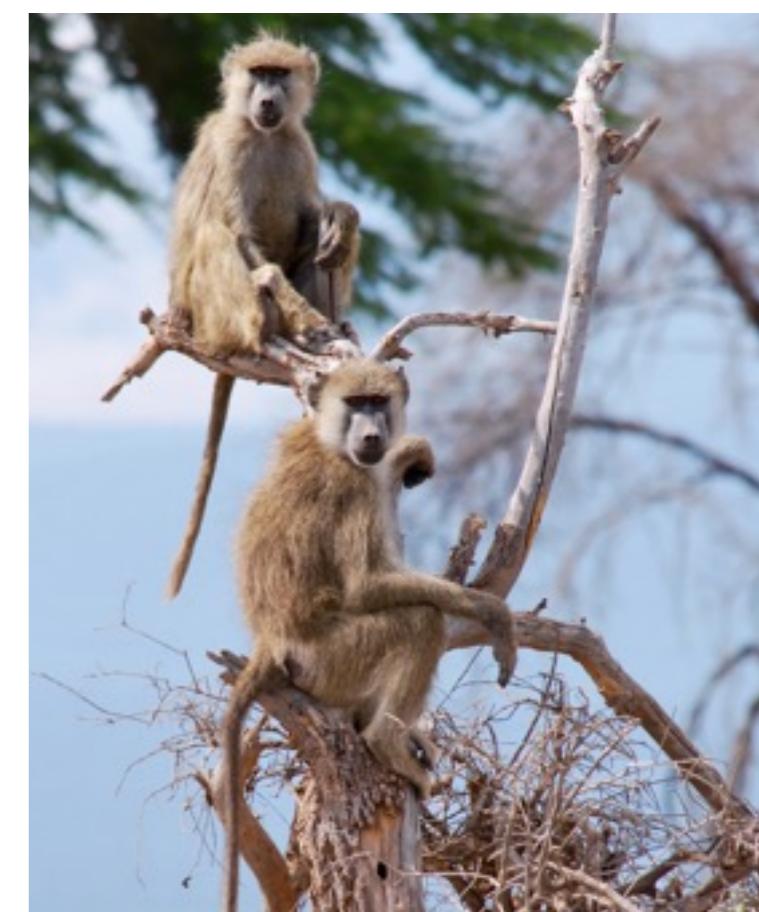
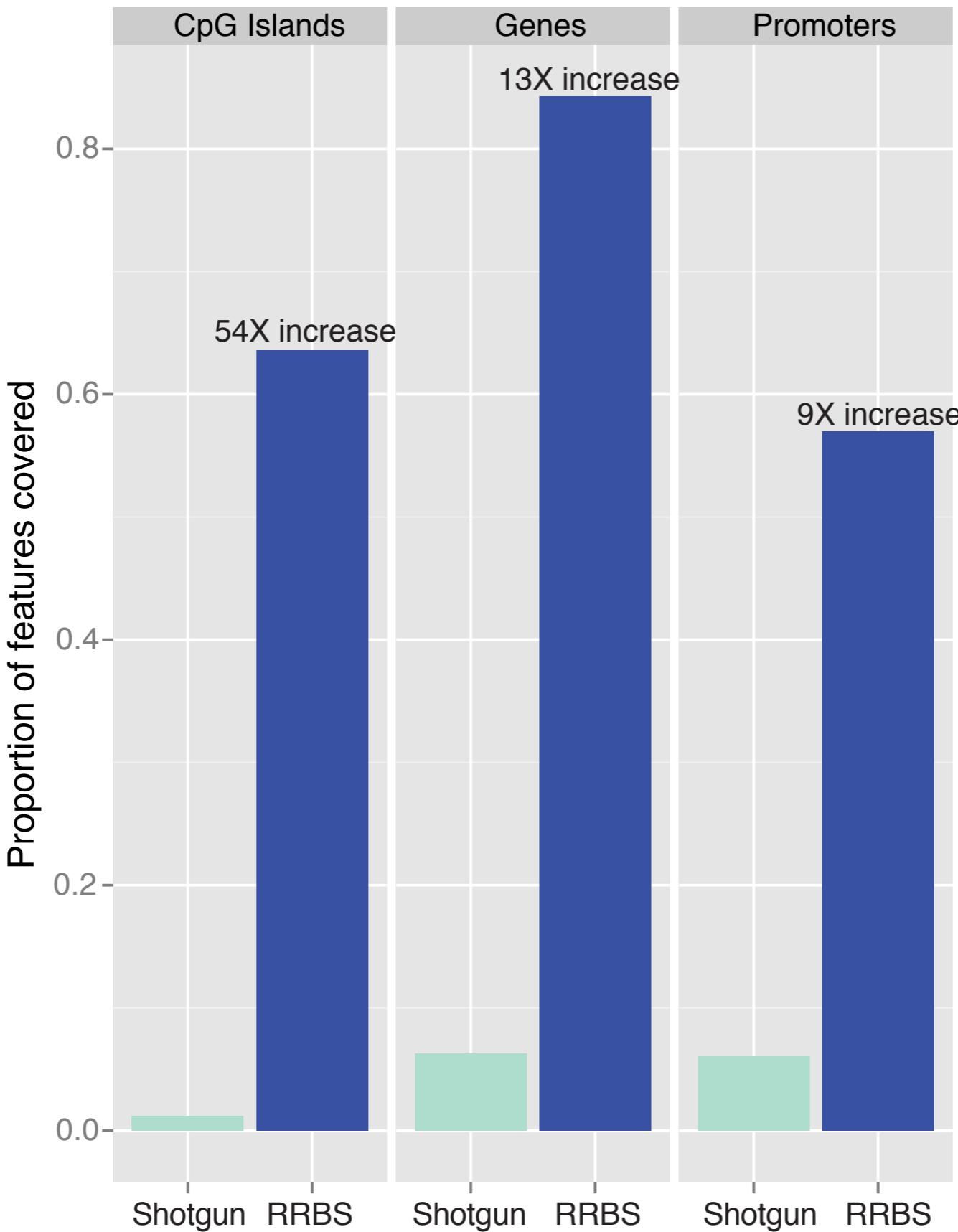


*Final library*

# RRBS enriches for functionally important regions of the genome



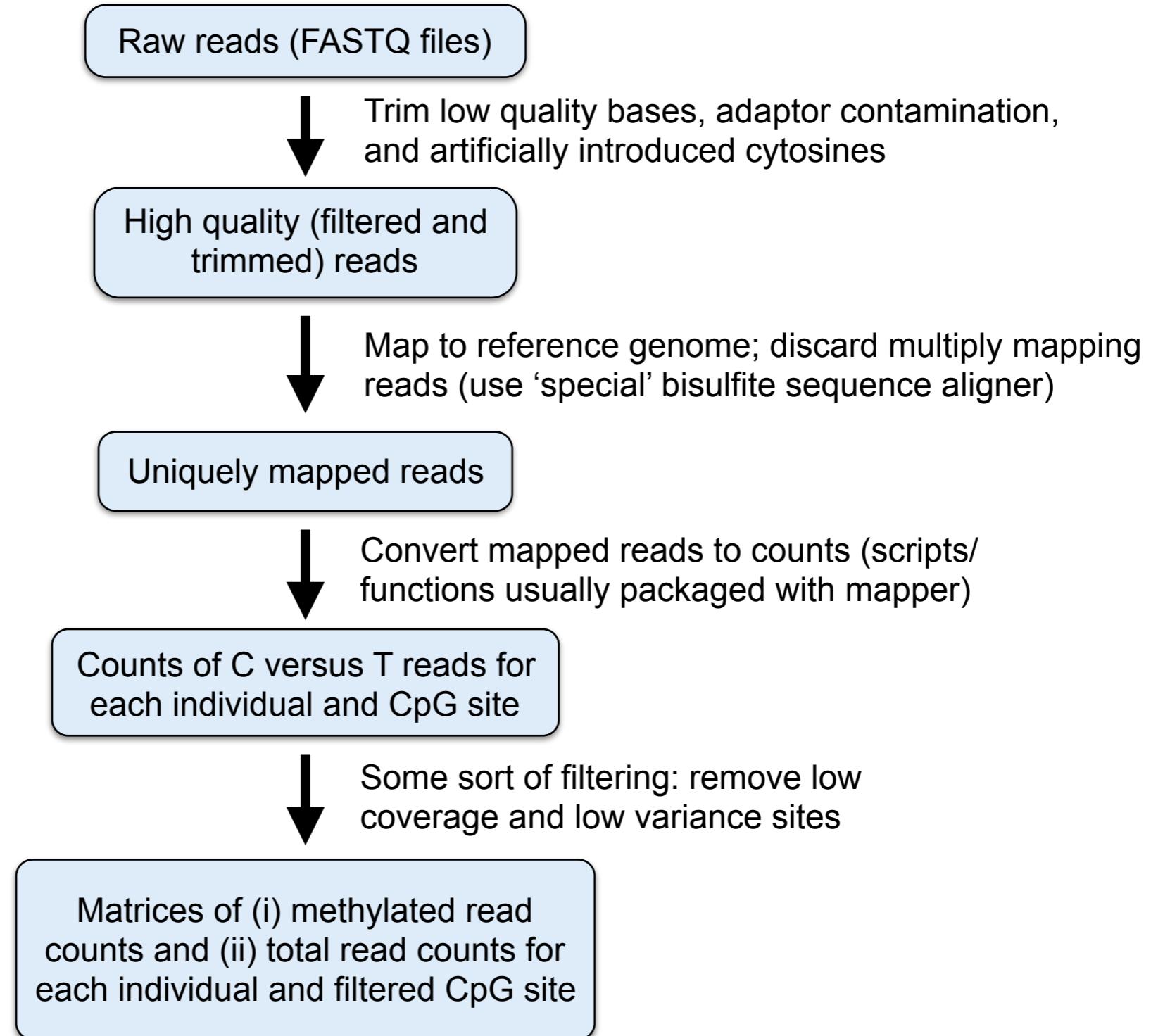
# RRBS enriches for functionally important regions of the genome



# Bisulfite sequencing data

# Bisulfite sequencing data

- First - process, map, and filter your data



Note, there are special mappers for bisulfite sequencing data

**BMC Bioinformatics**



Methodology article

Open Access

## **BSMAP: whole genome bisulfite sequence MAPping program**

Yuanxin Xi and Wei Li\*

**BIOINFORMATICS APPLICATIONS NOTE**

Vol. 27 no. 11 2011, pages 1571–1572  
doi:10.1093/bioinformatics/btr167

Sequence analysis

Advance Access publication April 14, 2011

## **Bismark: a flexible aligner and methylation caller for Bisulfite-Seq applications**

Felix Krueger\* and Simon R. Andrews

Cell Reports  
**Resource**

## **Differential DNA Methylation Analysis without a Reference Genome**

Johanna Klughammer,<sup>1</sup> Paul Datlinger,<sup>1</sup> Dieter Printz,<sup>2</sup> Nathan C. Sheffield,<sup>1</sup> Matthias Ertl,<sup>1</sup> Johanna Hadler,<sup>1</sup> Gerhard Fritsch,<sup>2</sup> and Christoph Bock<sup>1,3,4</sup>

**SOFTWARE**

Open Access

## **BS-Seeker2: a versatile aligning pipeline for bisulfite sequencing data**

Weilong Guo<sup>1,2</sup>, Petko Fiziev<sup>3</sup>, Weihong Yan<sup>4</sup>, Shawn Cokus<sup>2</sup>, Xueguang Sun<sup>5</sup>, Michael Q Zhang<sup>1,6</sup>, Pao-Yang Chen<sup>7\*</sup> and Matteo Pellegrini<sup>2,8\*</sup>

## Bisulfite sequencing data

- Now you have a count matrix!

Data matrix for  $n$  individuals measured at  $q$  loci  
(methylated reads / total reads)

	ID_1	ID_2	...	ID_n
loc_1	2/7	20/70	...	40/140
loc_2	2/3	40/60	...	40/140
loc_3	5/5	50/50		100/100
...	...	...	...	...
loc_q	0/0	1/3	...	4/6

# Bisulfite sequencing data

Data matrix for  $n$  individuals measured at  $q$  loci  
(methylated reads / total reads)

	ID_1	ID_2	...	ID_n
loc_1	2/7	20/70	...	40/140
loc_2	2/3	40/60	...	40/140
loc_3	5/5	50/50		100/100
...	...	...	...	...
loc_q	0/0	1/3	...	4/6

Biological variable = a proportion or level (ranging from 0-1)

Data = counts

# Bisulfite sequencing data

Data matrix for  $n$  individuals measured at  $q$  loci  
(methylated reads / total reads)

	ID_1	ID_2	...	ID_n
loc_1	2/7	20/70	...	40/140
loc_2	2/3	40/60	...	40/140
loc_3	5/5	50/50		100/100
...	...	...	...	...
loc_q	0/0	1/3	...	4/6

Missing data

# Bisulfite sequencing data

Data matrix for  $n$  individuals measured at  $q$  loci  
(methylated reads / total reads)

	ID_1	ID_2	...	ID_n
loc_1	2/7	20/70	...	40/140
loc_2	2/3	40/60	...	40/140
loc_3	5/5	50/50		100/100
...	...	...	...	...
loc_q	0/0	1/3	...	4/6

Sites with little variance in methylation levels

# Bisulfite sequencing data

Data matrix for  $n$  individuals measured at  $q$  loci  
(methylated reads / total reads)

	ID_1	ID_2	...	ID_n
loc_1	2/7	20/70	...	40/140
loc_2	2/3	40/60	...	40/140
loc_3	5/5	50/50		100/100
...	...	...	...	...
loc_q	0/0	1/3	...	4/6

Variation in coverage  
across **individuals** (by  
orders of magnitude)

# Bisulfite sequencing data

Data matrix for  $n$  individuals measured at  $q$  loci  
(methylated reads / total reads)

	ID_1	ID_2	...	ID_n
loc_1	2/7	20/70	...	40/140
loc_2	2/3	40/60	...	40/140
loc_3	5/5	50/50		100/100
...	...	...	...	...
loc_q	0/0	1/3	...	4/6

Variation in coverage across **sites** (by orders of magnitude)

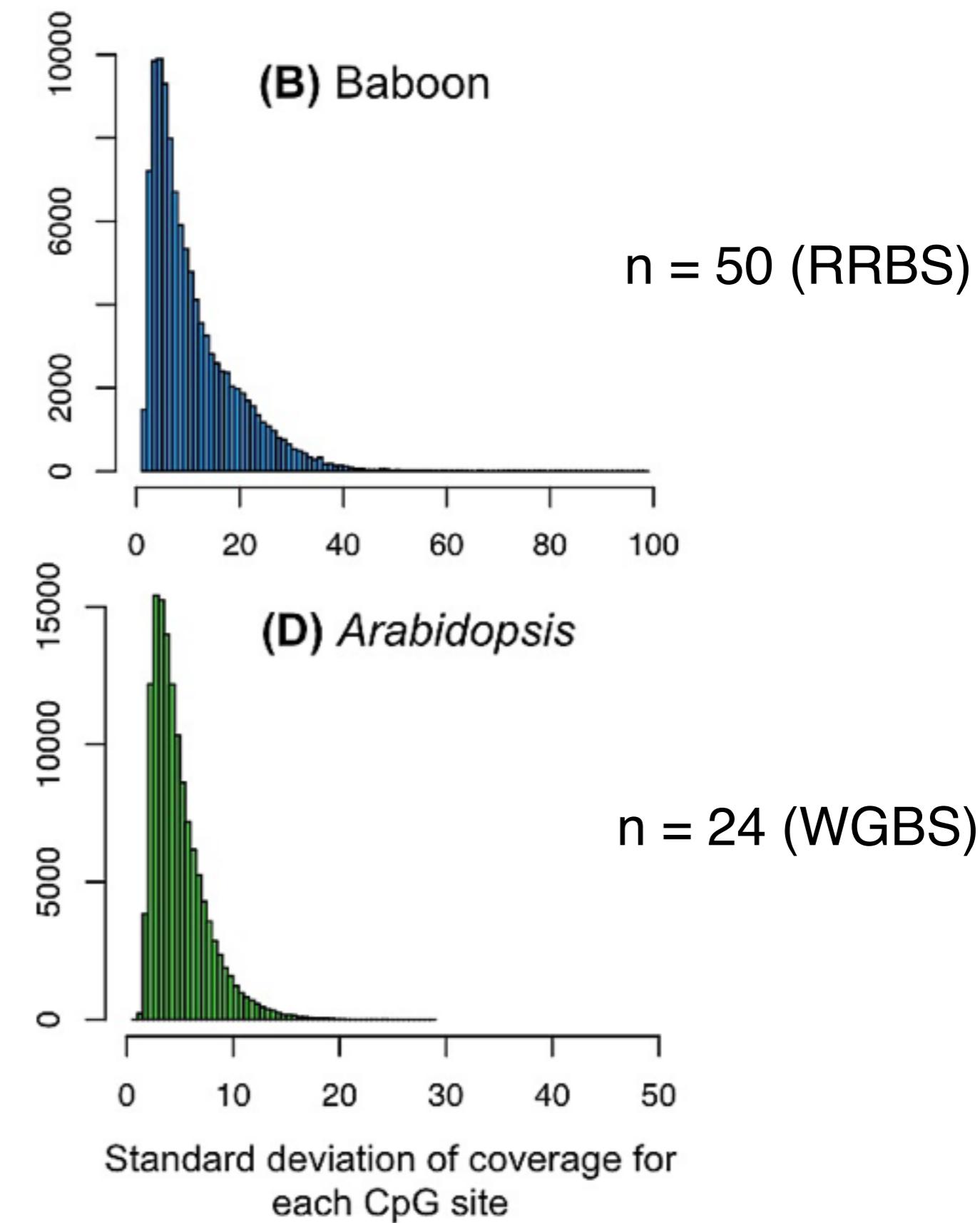
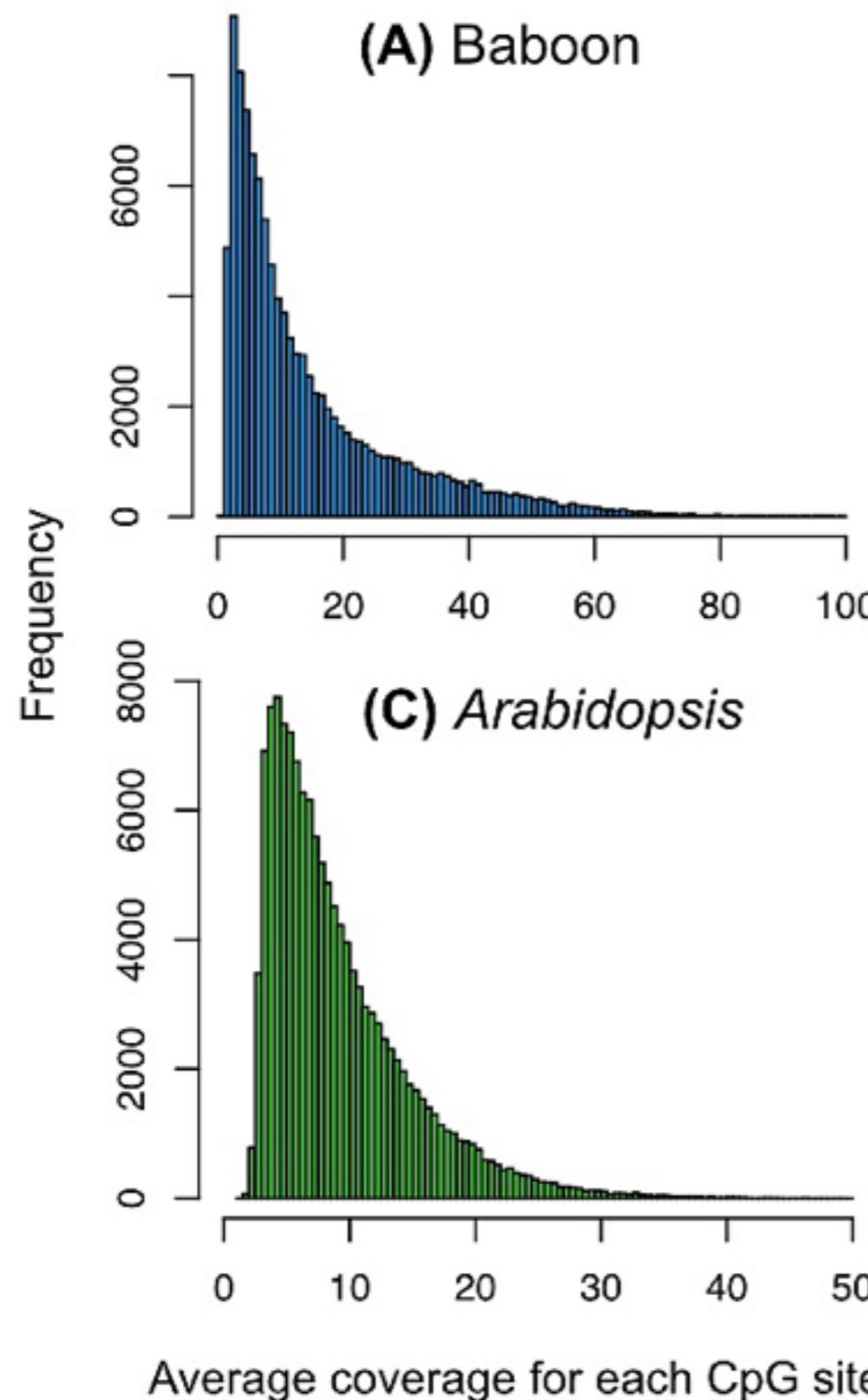
# Bisulfite sequencing data

Data matrix for  $n$  individuals measured at  $q$  loci  
(methylated reads / total reads)

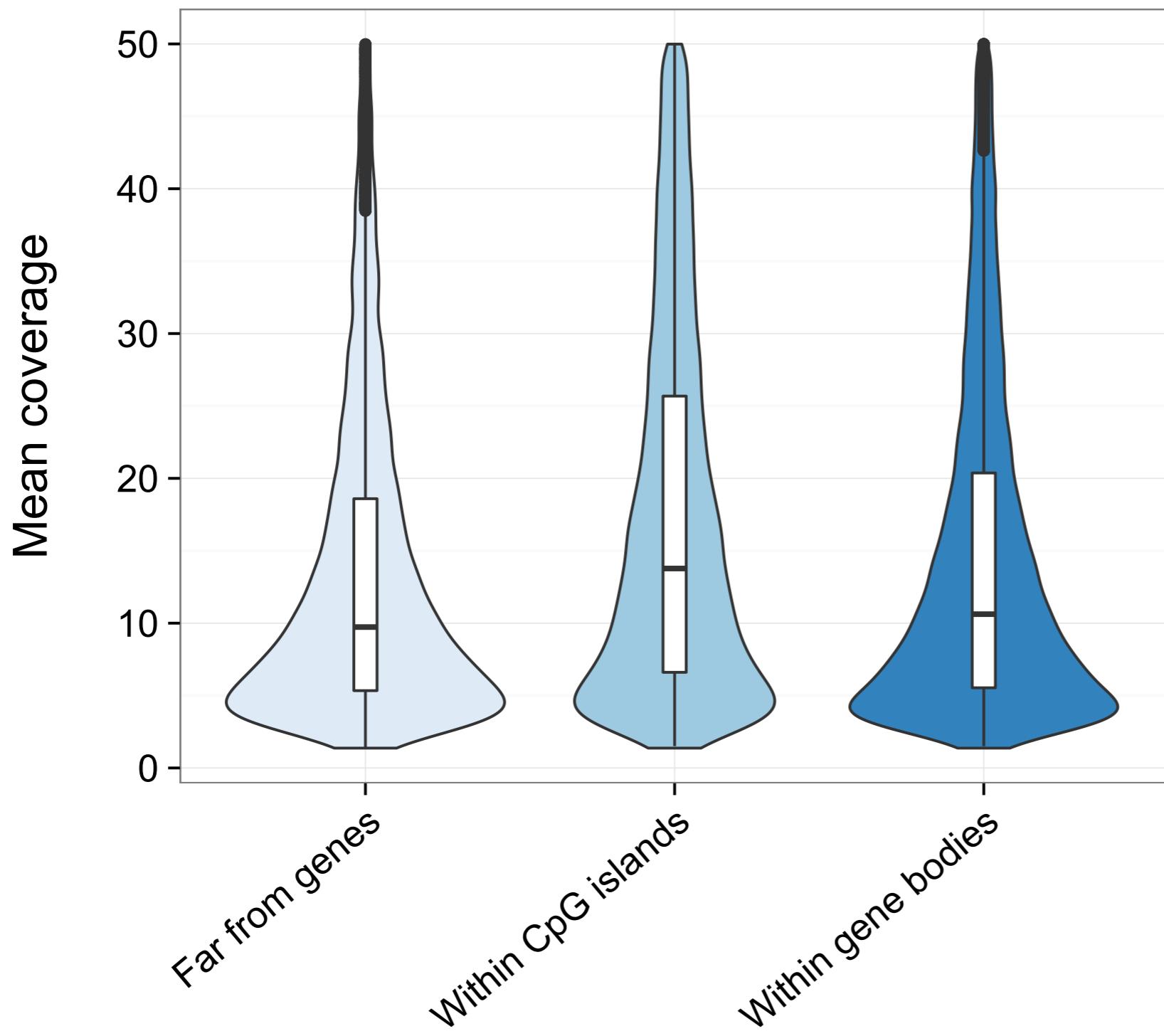
	ID_1	ID_2	...	ID_n
loc_1	2/7	20/70	...	40/140
loc_2	2/3	40/60	...	40/140
loc_3	5/5	50/50		100/100
...	...	...	...	...
loc_q	0/0	1/3	...	4/6

Unlike RNA-seq, all of  
this total read count  
variation is technical

# Bisulfite sequencing data



# Bisulfite sequencing data



# Bisulfite sequencing data

	ID_1	ID_2	...	ID_n
loc_1	2/7	20/70	...	40/140
loc_2	2/3	40/60	...	40/140
loc_3	5/5	50/50		100/100
...	...	...	...	...
loc_q	0/0	1/3	...	4/6

transform?



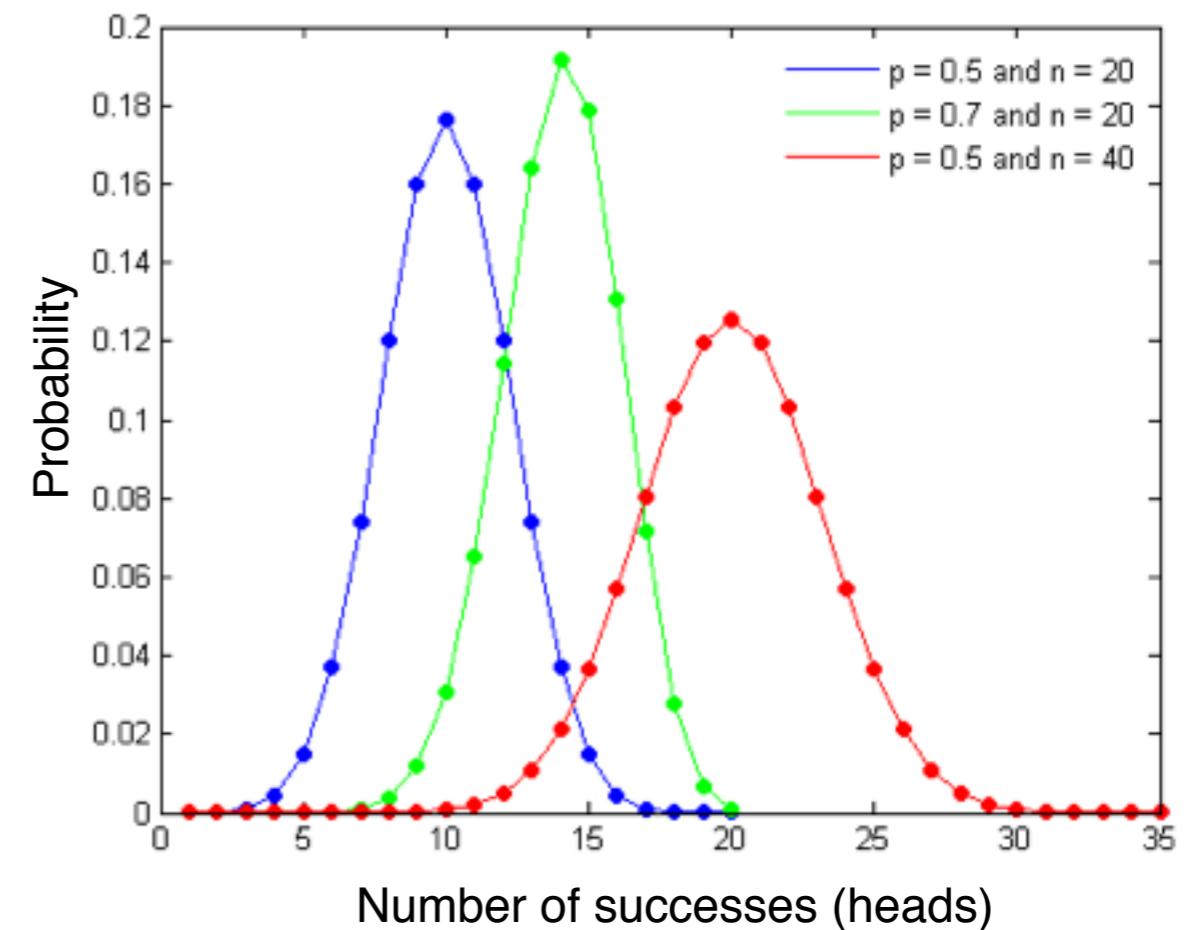
	ID_1	ID_2	...	ID_n
loc_1	0.28	0.28	...	0.29
loc_2	0.67	0.67	...	0.29
loc_3	1	1		1
...	...	...	...	...
loc_q	NA	0.33	...	0.67

# Methods for analyzing bisulfite sequencing data

Programs	Method	Model counts?
Many	t-test or Wilcoxon rank-sum test	No (model proportion)
Many	Fisher's exact test	Yes
Many	Linear regression	No (model normalized proportion)
Many	Logistic/Binomial regression	Yes
<b>DSS</b> (Feng et al. 2014), <b>MOABS</b> (Sun et al. 2014), <b>RADMeth</b> (Dolzhenko & Smith 2014)	Based on a beta binomial regression	Yes

## Beta binomial models

- $m_i \sim \text{Bin}(n_i, p_i)$ , where  $p_i$  is the unknown methylation level of the site ( $m_i$  and  $n_i$  are the observed counts)



## Beta binomial models

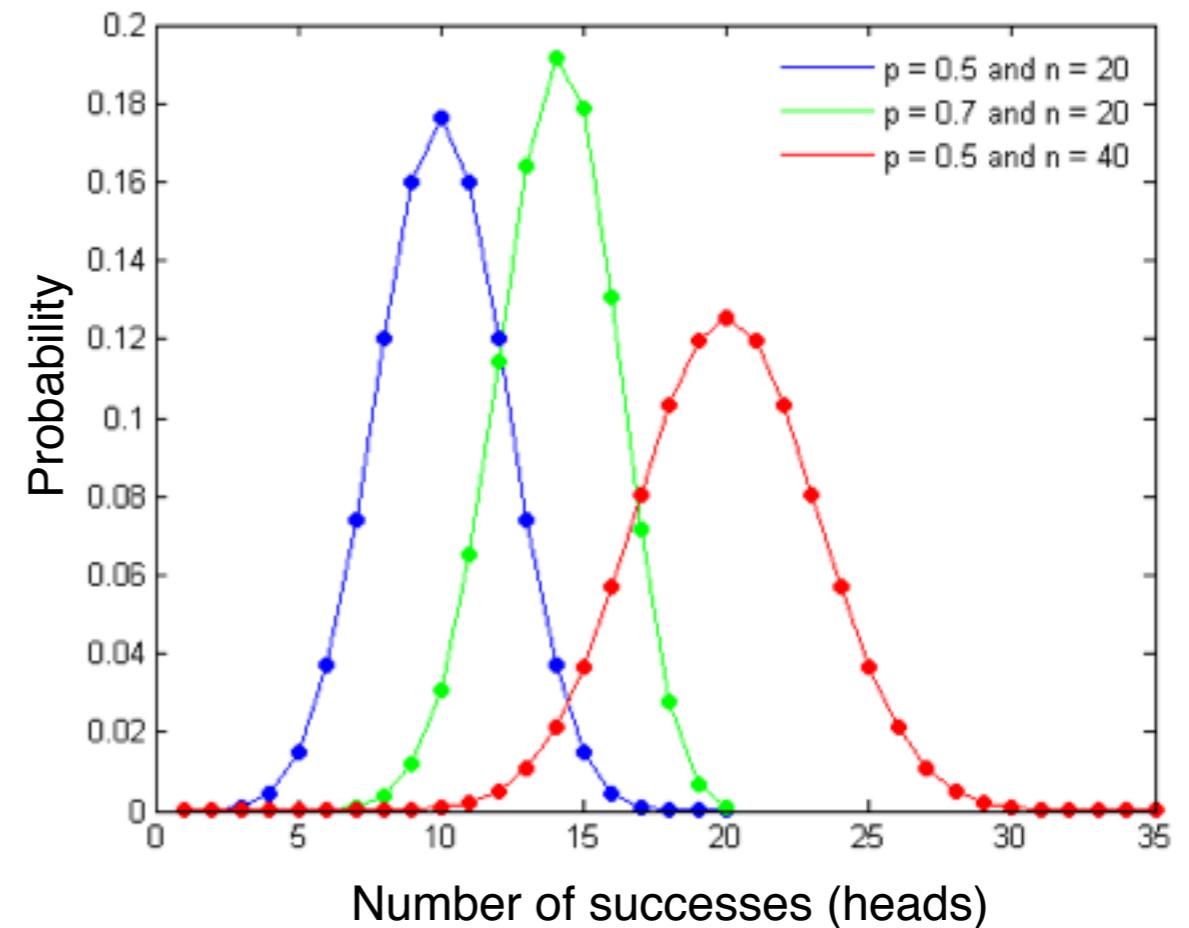
- $m_i \sim \text{Bin}(n_i, p_i)$ , where  $p_i$  is the unknown methylation level of the site ( $m_i$  and  $n_i$  are the observed counts)

- The *binomial* assumes:

$$E(m_i) = n_i * p_i$$

$$\text{Var}(m_i) = n_i * p_i * (1 - p_i)$$

- But counts from bisulfite-sequencing data are over dispersed



# Beta binomial models

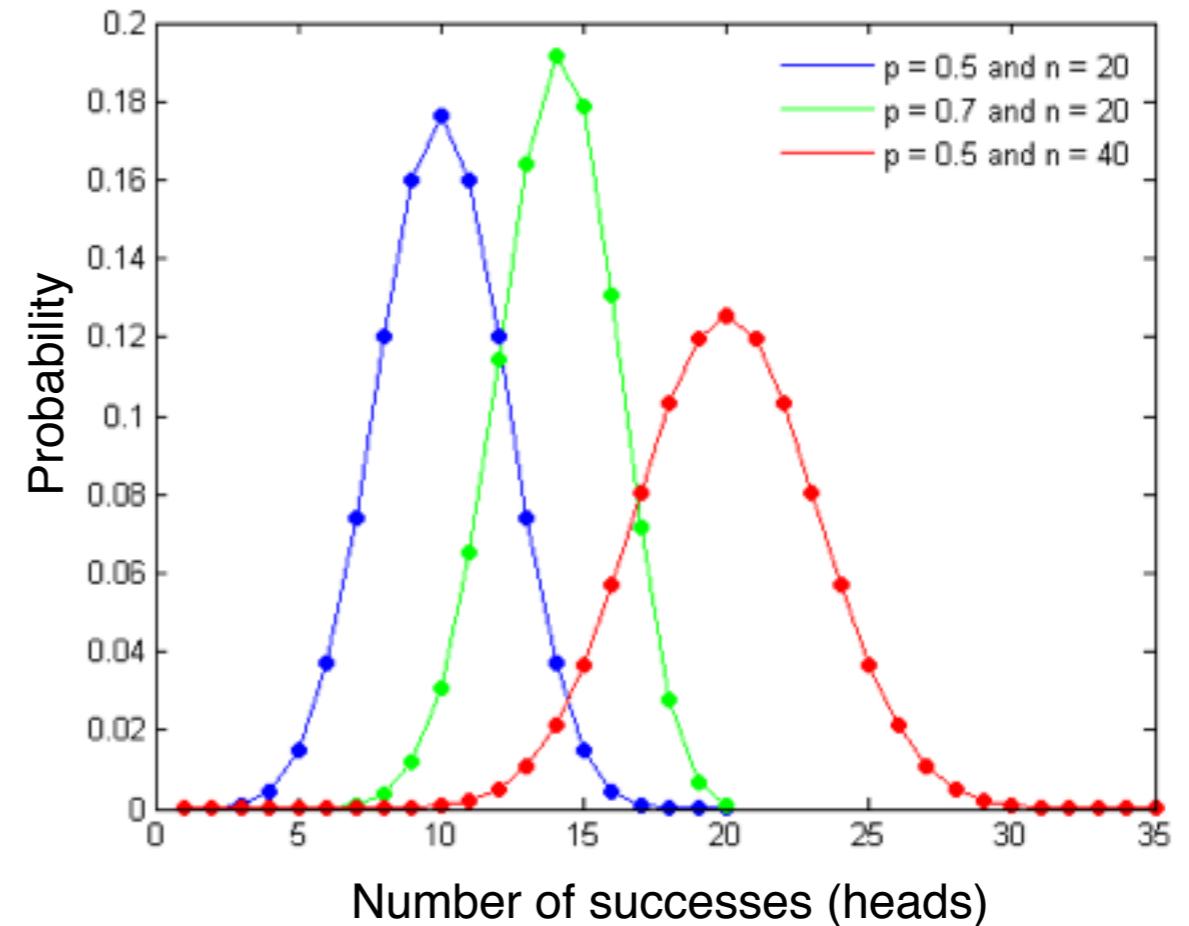
- $m_i \sim \text{Bin}(n_i, p_i)$ , where  $p_i$  is the unknown methylation level of the site ( $m_i$  and  $n_i$  are the observed counts)

- The *binomial* assumes:

$$E(m_i) = n_i * p_i$$

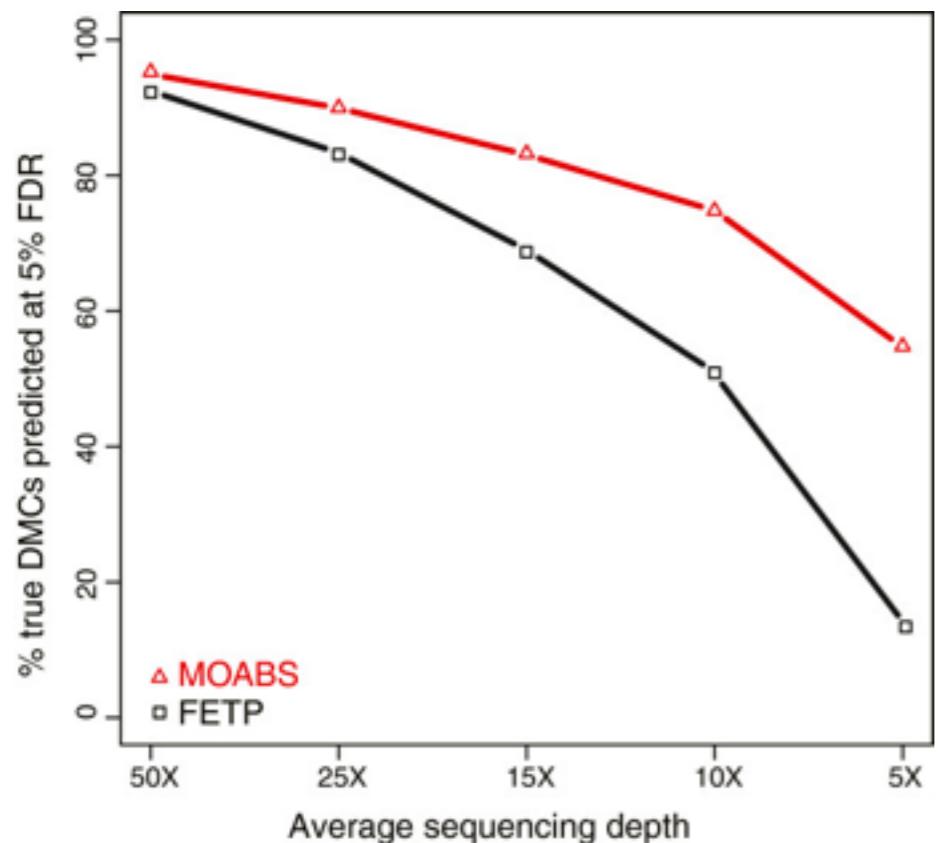
$$\text{Var}(m_i) = n_i * p_i * (1 - p_i)$$

- But counts from bisulfite-sequencing data are over dispersed
- The *beta-binomial* assumes  $p_i$  is a random variable and follows a beta distribution (accounts for overdispersion)
- Not surprisingly, *beta-binomial* models outperform non count-based methods, as well as binomial models for analyzing DNA methylation datasets



# Beta binomial models

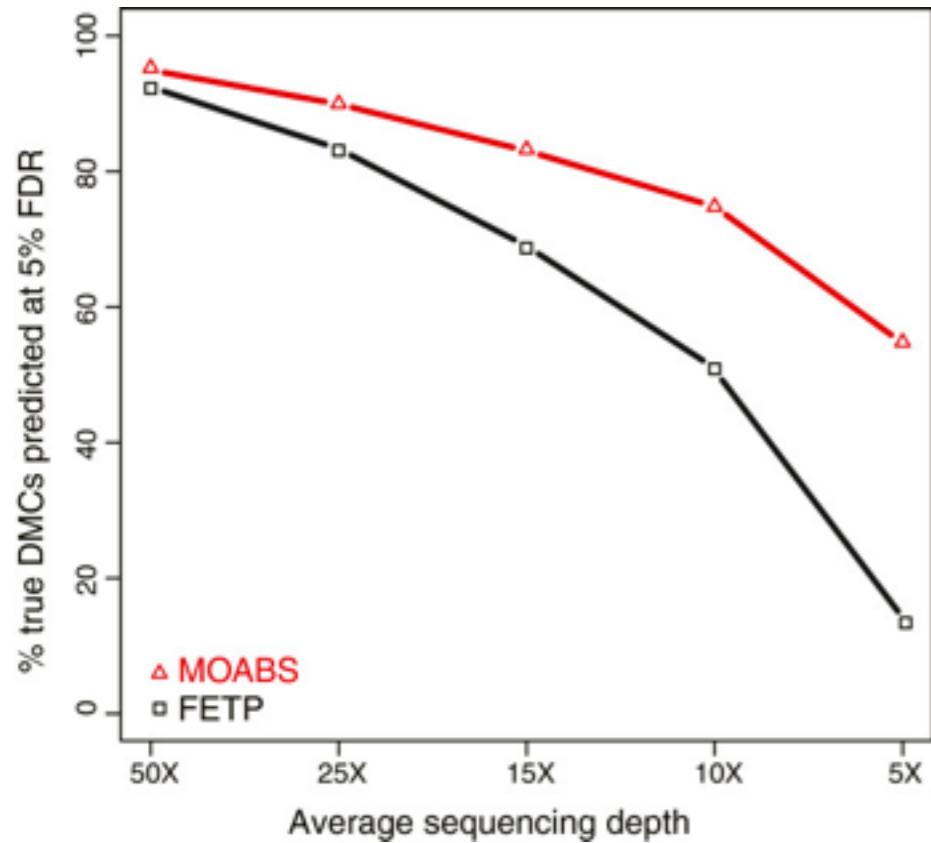
Comparison between 2 samples



MOABS: Sun et al. 2014, Genome Biology

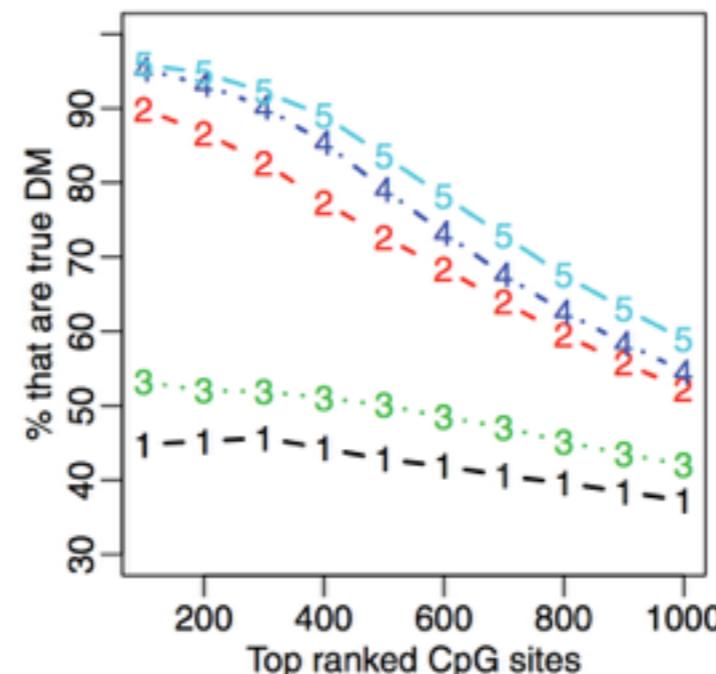
# Beta binomial models

Comparison between 2 samples

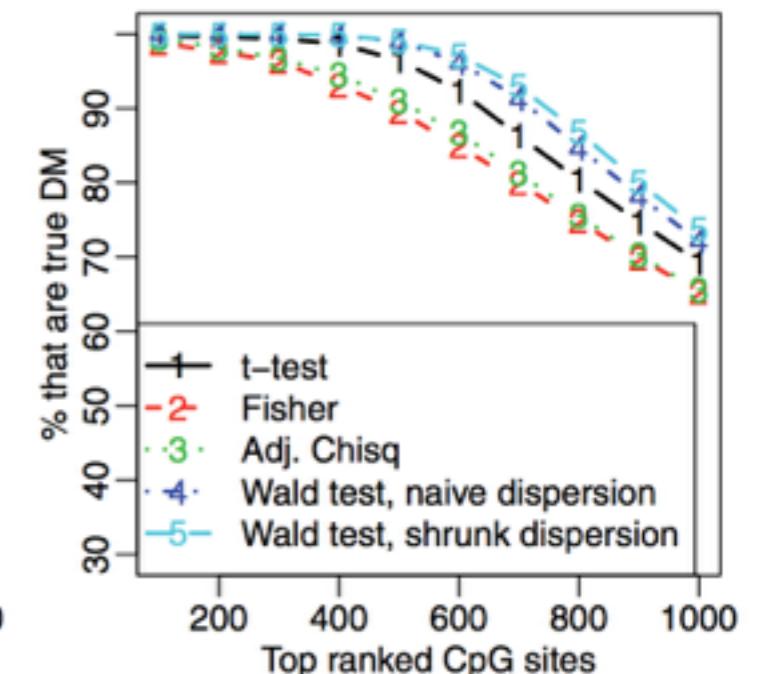


MOABS: Sun et al. 2014, Genome Biology

Comparison between 2 groups

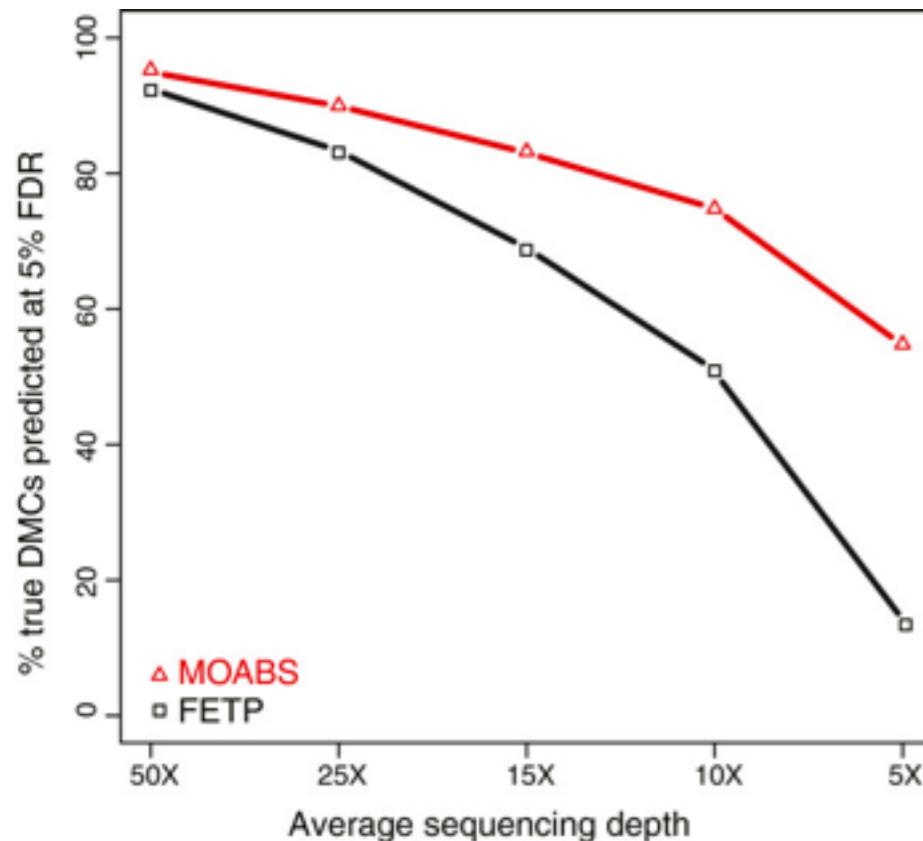


DSS: Feng et al. 2014, Nucleic Acids Research



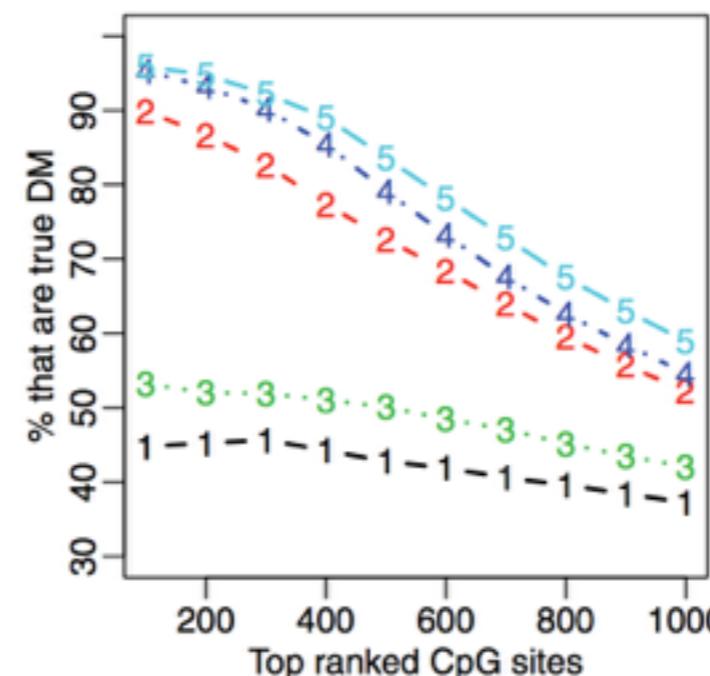
# Beta binomial models

Comparison between 2 samples

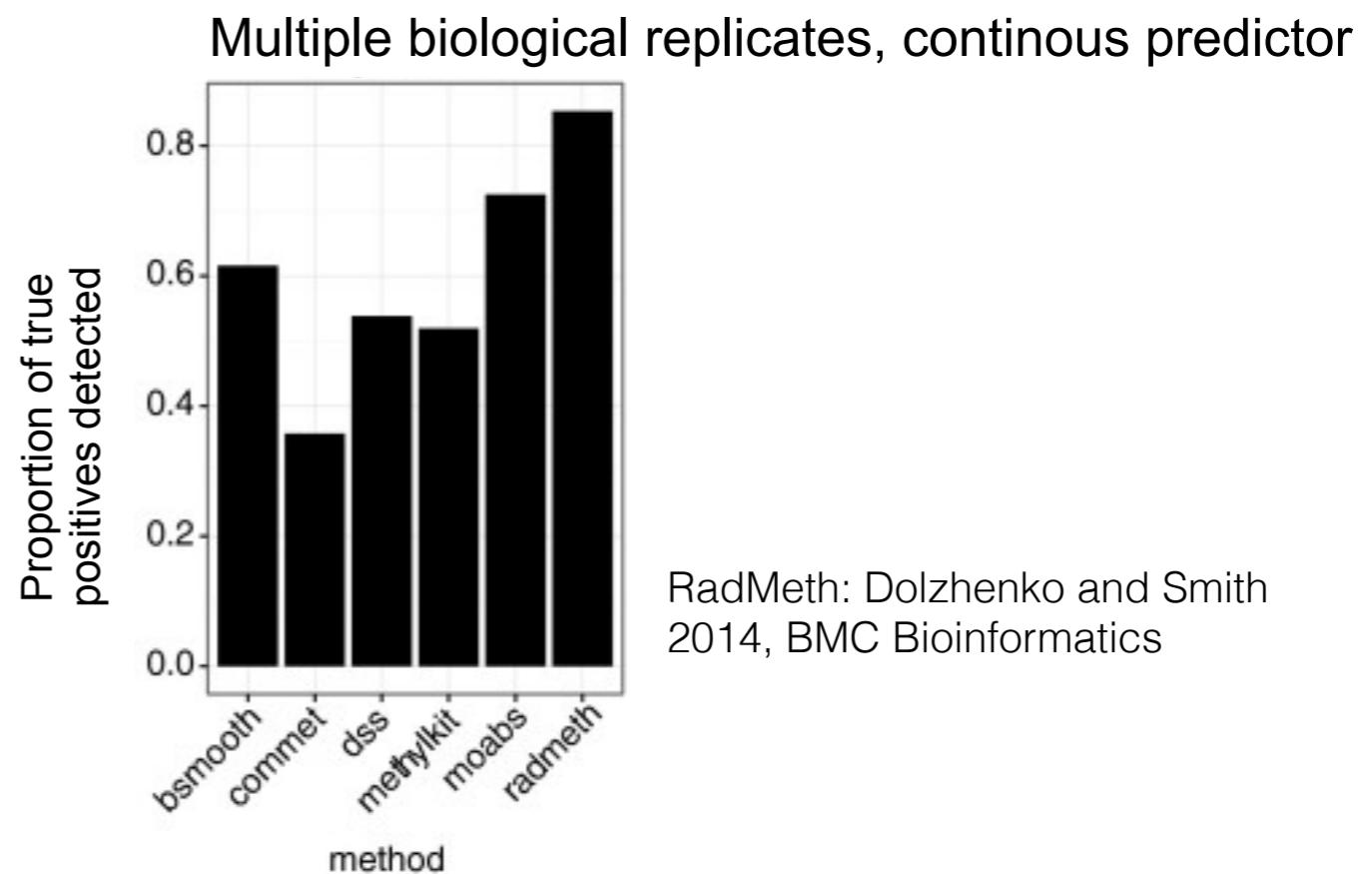


MOABS: Sun et al. 2014, Genome Biology

Comparison between 2 groups



DSS: Feng et al. 2014, Nucleic Acids Research



RadMeth: Dolzhenko and Smith  
2014, BMC Bioinformatics

## Beta binomial models

These tools are an improvement, but are relatively inflexible and few allow researchers to model additional covariates (e.g.,):

- technical/batch effects effects
- age
- sex
- cell type heterogeneity

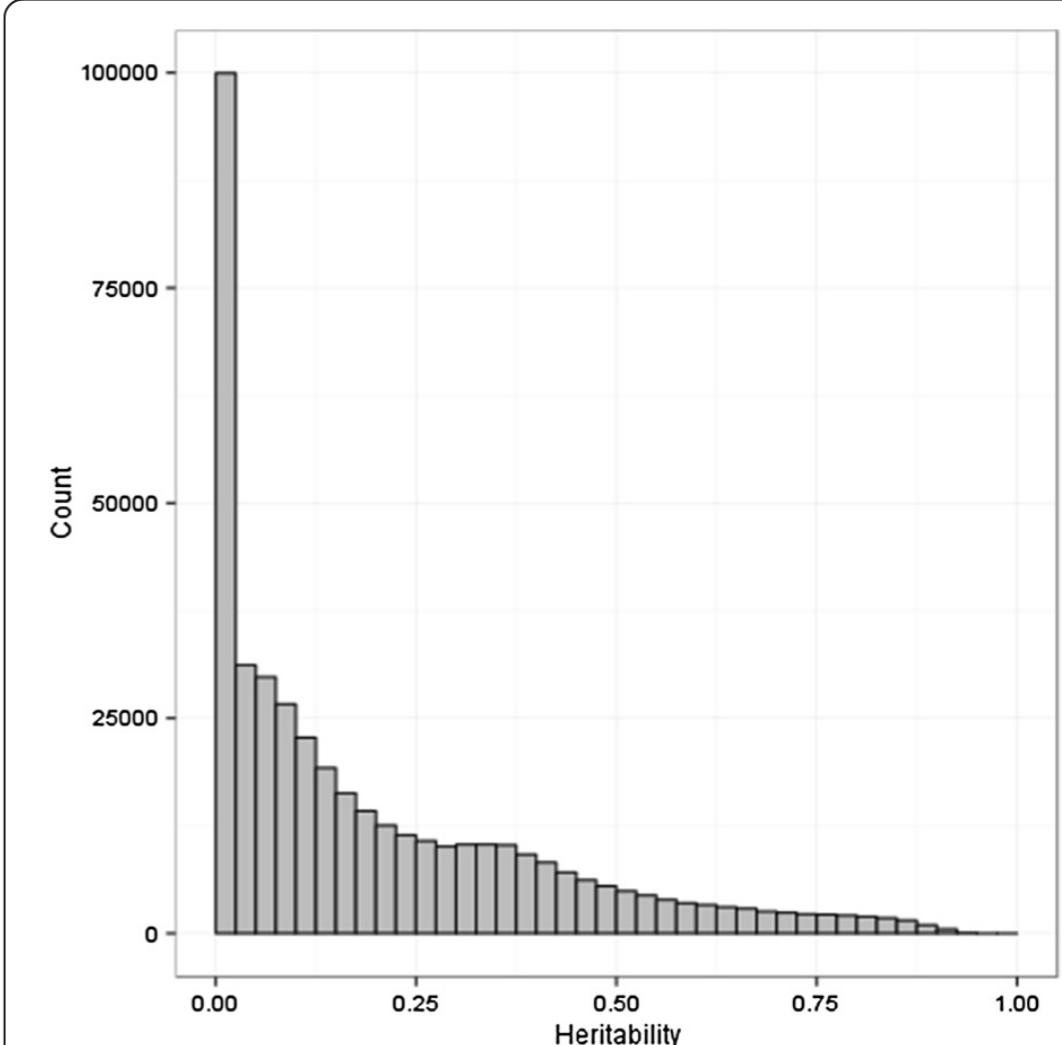
# Beta binomial models

These tools are an improvement, but are relatively inflexible and few allow researchers to model additional covariates (e.g.,):

- technical/batch effects effects
- age
- sex
- cell type heterogeneity
- also, genetic effects!

**Table 1 Average correlation across all probes of normalised methylation measurements between relative pairs**

Relationship	Pairs (n)	Correlation	Expected <sup>a</sup>
MZ twins	67	0.200	$h^2$
DZ twins	111	0.109	$h^2/2$
Siblings	262 <sup>b</sup>	0.090	$h^2/2$
Parent-Offspring	362 <sup>b</sup>	0.089	$h^2/2$
Mother-Offspring	190	0.097	$h^2/2$
Father-Offspring	172	0.085	$h^2/2$
Parent-Parent	58	0.023	0
Unrelated	187,331 <sup>b</sup>	-0.002	0



**Figure 1 Distribution of heritability estimates for DNA methylation levels.** The average genetic heritability estimate is 0.199. A zero estimate for genetic heritability was observed in 17.1% of cases indicating that genetic heritability results in transgenerational inheritance of DNA methylation for at least 65.8% of probes.

# Methods for analyzing bisulfite sequencing data

Programs	Method	Model counts?	Control for relatedness/structure?	Control for covariates?
Many	t-test or Wilcoxon rank-sum test	No	No	No
Many	Fisher's exact test	Yes	No	No
Many	Linear regression	No	No	Yes
Many	Logistic/Binomial regression	Yes	No	Yes
DSS (Feng et al. 2014), MOABS (Sun et al. 2014), RADMeth (Dolzhenko & Smith 2014)	Based on a beta binomial model	Yes	No	No/Yes

# Methods for analyzing bisulfite sequencing data

Programs	Method	Model counts?	Control for relatedness/structure?	Control for covariates?
Many	t-test or Wilcoxon rank-sum test	No	No	No
Many	Fisher's exact test	Yes	No	No
Many	Linear regression	No	No	Yes
Many	Logistic/Binomial regression	Yes	No	Yes
<b>DSS</b> (Feng et al. 2014), <b>MOABS</b> (Sun et al. 2014), <b>RADMeth</b> (Dolzhenko & Smith 2014)	Based on a beta binomial regression	Yes	No	No/Yes
<b>EMMA/EMMAX</b> (Kang et al. 2010), <b>GEMMA</b> (Zhou & Stephens 2014)	Based on a linear mixed effects model	No	<b>Yes</b>	Yes

# The Linear Mixed Model

$$\mathbf{y} = \mu + \mathbf{x}\beta + \mathbf{u} + \boldsymbol{\epsilon},$$

$$\mathbf{u} \sim \text{MVN}_n(0, \sigma_g^2 \mathbf{K}), \quad \boldsymbol{\epsilon} \sim \text{MVN}_n(0, \sigma_e^2 \mathbf{I}_{n \times n}),$$

- ▶  $\mathbf{y}$  is a  $n$ -vector of phenotypes,
- ▶  $\mathbf{x}$  is a  $n$ -vector of genotypes for the SNP,
- ▶  $\mathbf{u}$  is a  $n$ -vector of random effects,
- ▶  $\boldsymbol{\epsilon}$  is a  $n$ -vector of residual errors,
- ▶  $\mathbf{K}$  is a known  $n \times n$  relatedness matrix, estimated either from a pedigree or genotypes.

$$\sigma_g^2 = \sigma^2 h^2 \text{ and } \sigma_e^2 = \sigma^2 (1-h)^2$$

# Methods for analyzing bisulfite sequencing data

Programs	Method	Model counts?	Control for relatedness/structure?	Control for covariates?
Many	t-test or Wilcoxon rank-sum test	No	No	No
Many	Fisher's exact test	Yes	No	No
Many	Linear regression	No	No	Yes
Many	Logistic/Binomial regression	Yes	No	Yes
<b>DSS</b> (Feng et al. 2014), <b>MOABS</b> (Sun et al. 2014), <b>RADMeth</b> (Dolzhenko & Smith 2014)	Based on a beta binomial regression	Yes	No	No/Yes
<b>EMMA/EMMAX</b> (Kang et al. 2010), <b>GEMMA</b> (Zhou & Stephens 2014)	Based on a linear mixed effects model	<b>No</b>	Yes	Yes

# Methods for analyzing bisulfite sequencing data

Programs	Method	Model counts?	Control for relatedness/structure?	Control for covariates?
Many	t-test or Wilcoxon rank-sum test	No	No	No
Many	Fisher's exact test	Yes	No	No
Many	Linear regression	No	No	Yes
Many	Logistic/Binomial regression	Yes	No	Yes
<b>DSS</b> (Feng et al. 2014), <b>MOABS</b> (Sun et al. 2014), <b>RADMeth</b> (Dolzhenko & Smith 2014)	Based on a beta binomial regression	Yes	No	No/Yes
<b>EMMA/EMMAX</b> (Kang et al. 2010), <b>GEMMA</b> (Zhou & Stephens 2014)	Based on a linear mixed effects model	No	Yes	Yes
<b>MACAU</b> (Lea et al. 2015)	Based on a binomial mixed effects model	Yes	Yes	Yes

# The Linear Mixed Model

$$\mathbf{y} = \mu + \mathbf{x}\beta + \mathbf{u} + \boldsymbol{\epsilon},$$

$$\mathbf{u} \sim \text{MVN}_n(0, \sigma_g^2 \mathbf{K}), \quad \boldsymbol{\epsilon} \sim \text{MVN}_n(0, \sigma_e^2 \mathbf{I}_{n \times n}),$$

- ▶  $\mathbf{y}$  is a  $n$ -vector of phenotypes,
- ▶  $\mathbf{x}$  is a  $n$ -vector of genotypes for the SNP,
- ▶  $\mathbf{u}$  is a  $n$ -vector of random effects,
- ▶  $\boldsymbol{\epsilon}$  is a  $n$ -vector of residual errors,
- ▶  $\mathbf{K}$  is a known  $n \times n$  relatedness matrix, estimated either from a pedigree or genotypes.

$$\sigma_g^2 = \sigma^2 h^2 \text{ and } \sigma_e^2 = \sigma^2 (1-h)^2$$

# A binomial mixed effects model

methylated counts  $\sim \text{Bin}(\text{total counts}, p)$

$$\log(p/(1-p)) = \mu + \mathbf{x}\beta + \mathbf{u} + \epsilon,$$

$$\mathbf{u} \sim \text{MVN}_n(0, \sigma_g^2 \mathbf{K}), \quad \epsilon \sim \text{MVN}_n(0, \sigma_e^2 \mathbf{I}_{n \times n}),$$

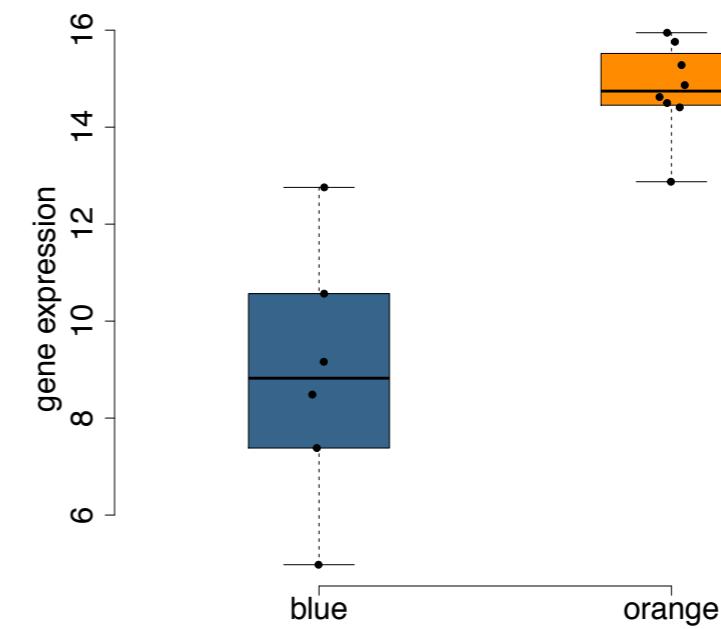
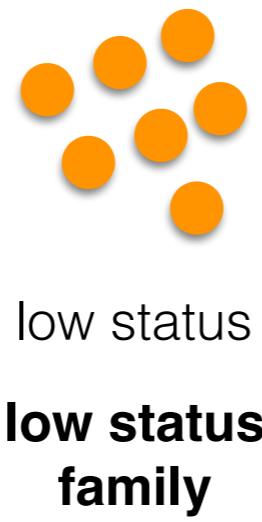
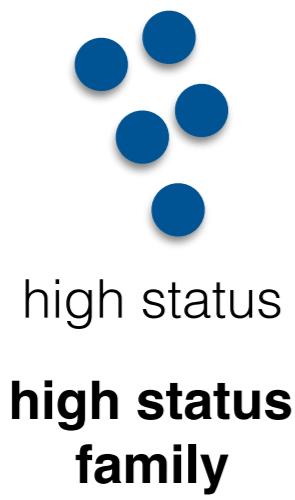
- ▶  $\mathbf{y}$  is a  $n$ -vector of phenotypes,
- ▶  $\mathbf{x}$  is a  $n$ -vector of genotypes for the SNP,
- ▶  $\mathbf{u}$  is a  $n$ -vector of random effects,
- ▶  $\epsilon$  is a  $n$ -vector of residual errors,
- ▶  $\mathbf{K}$  is a known  $n \times n$  relatedness matrix, estimated either from a pedigree or genotypes.

$$\sigma_g^2 = \sigma^2 h^2 \text{ and } \sigma_e^2 = \sigma^2 (1-h)^2$$



Xiang Zhou  
University of Michigan

Correlation ( $R^2$ ) between population structure and predictor variable

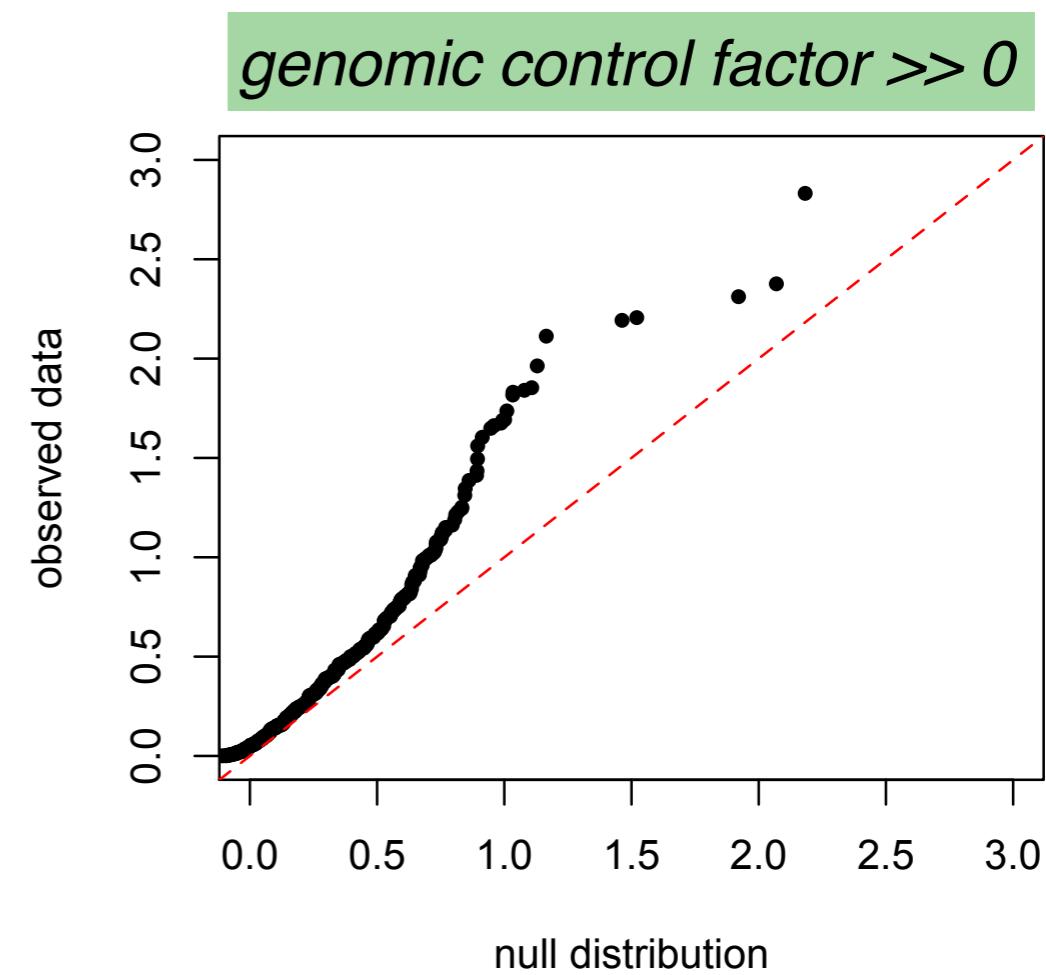
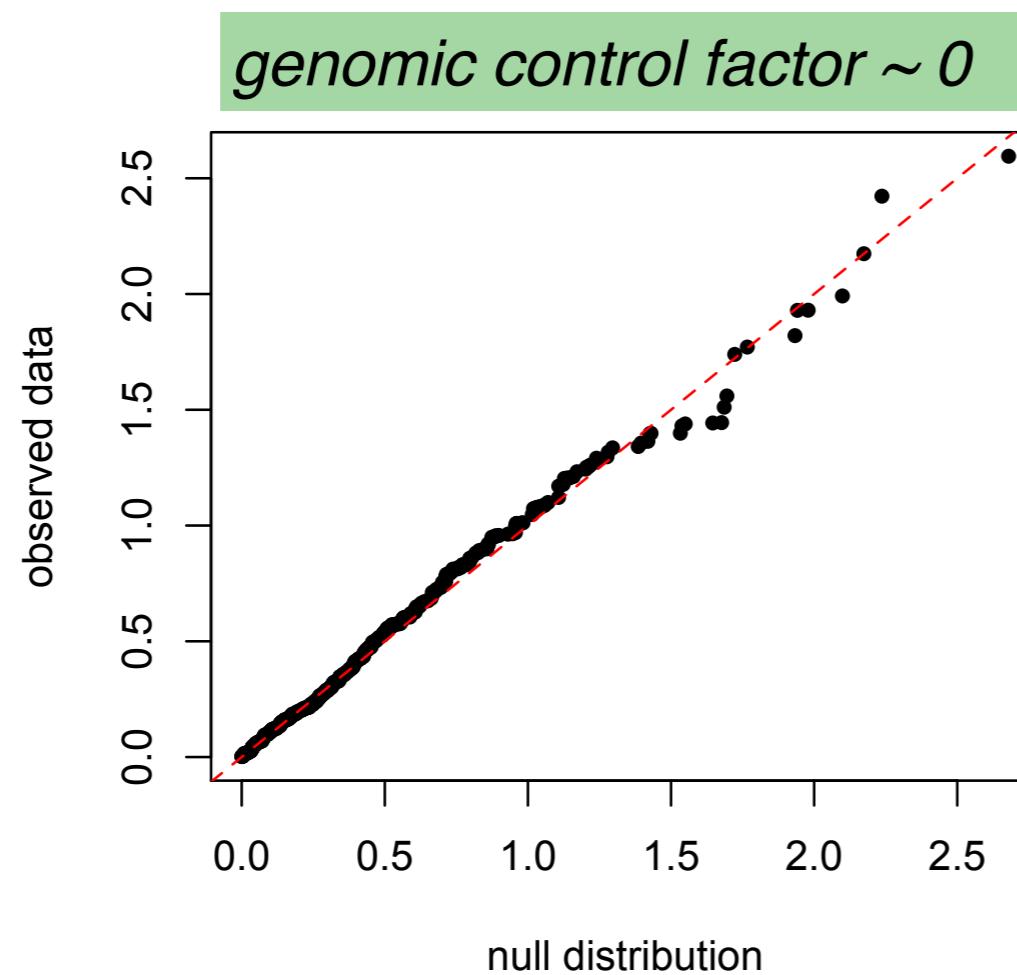


Correlation ( $R^2$ ) between population structure and predictor variable

Genomic control factor ( $\lambda$ )

Correlation ( $R^2$ ) between population structure and predictor variable

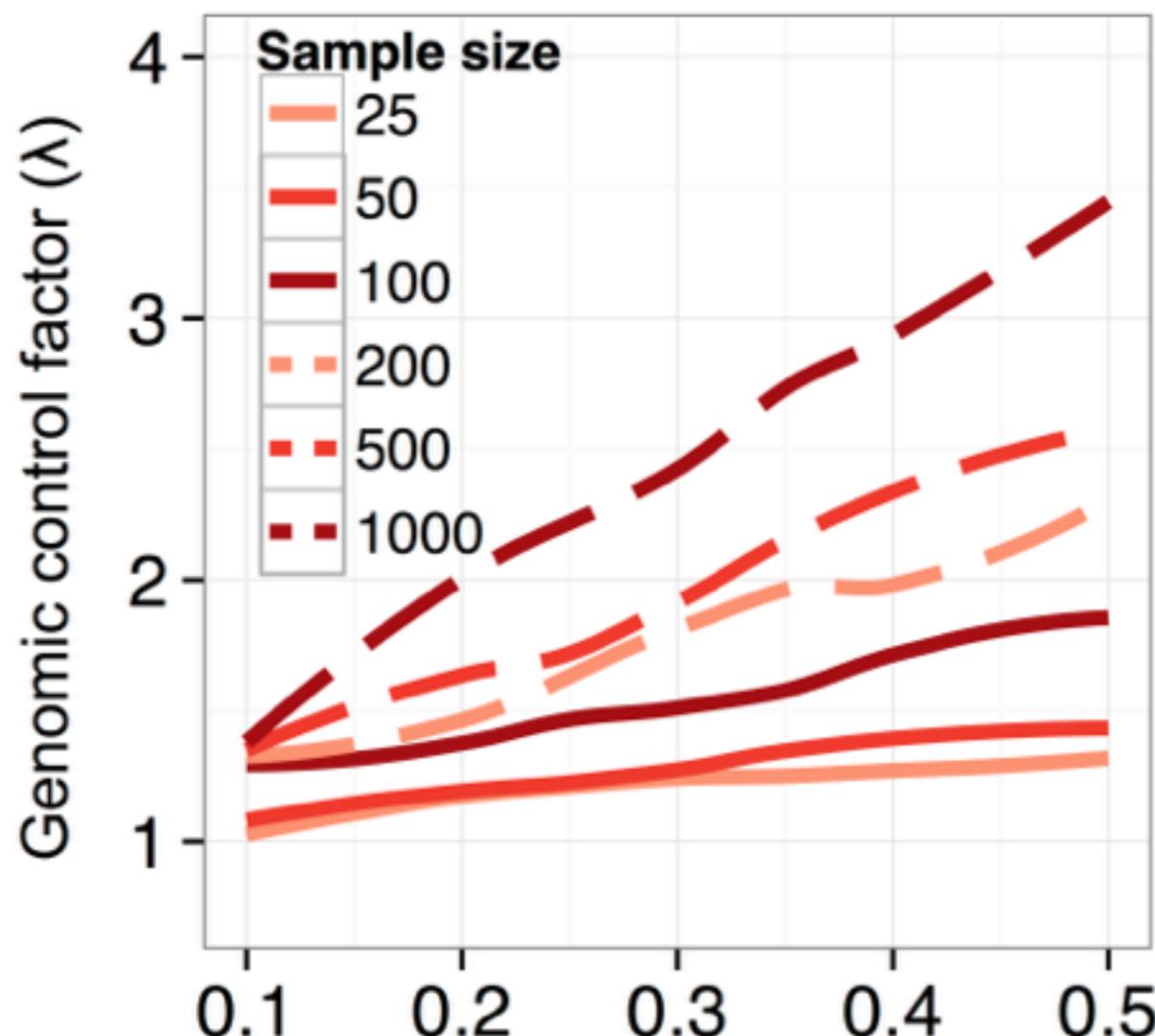
Genomic control factor ( $\lambda$ )



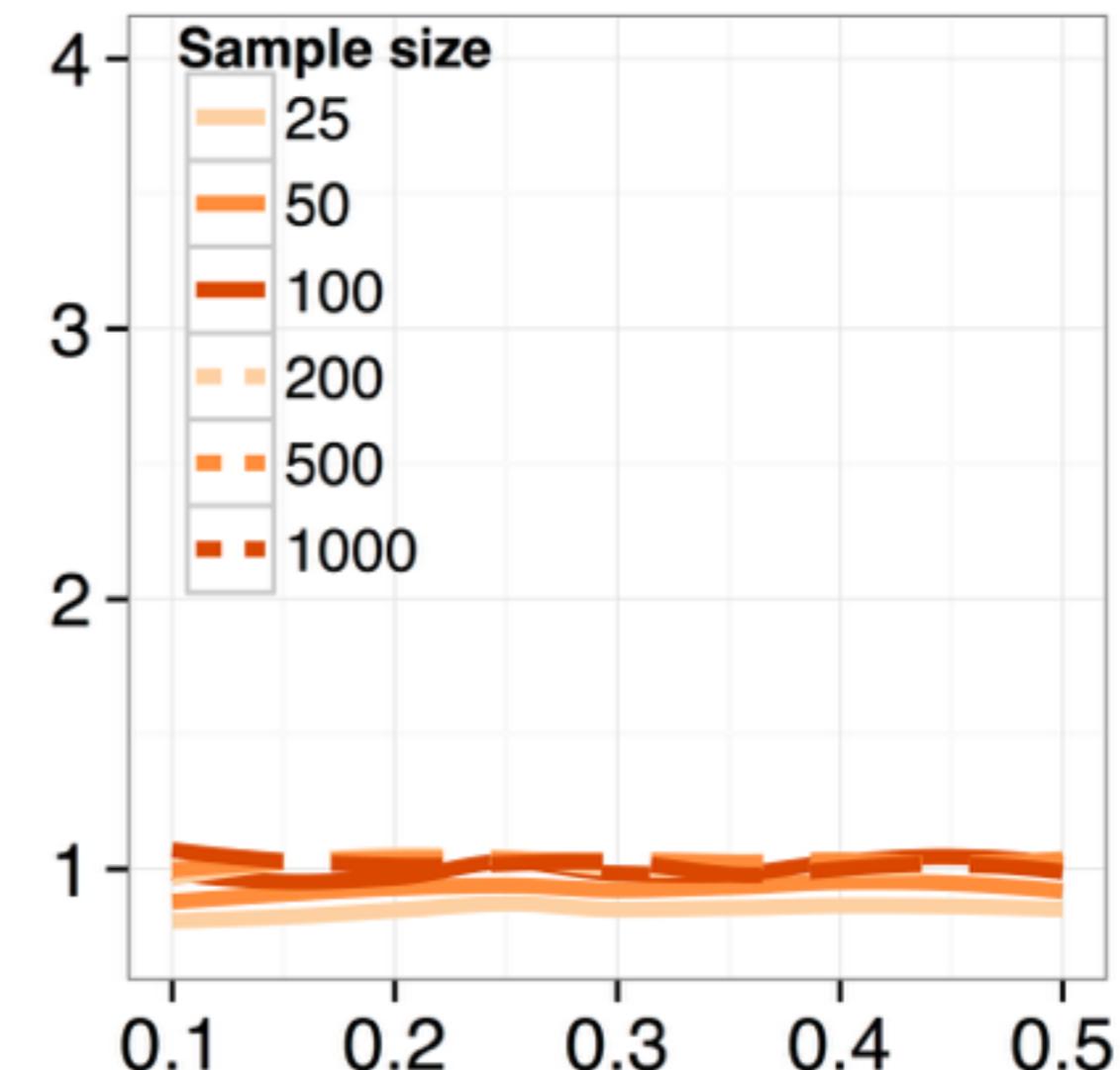
Correlation ( $R^2$ ) between population structure and predictor variable

Binomial mixed effects models minimize false positives when genetic structure is confounded with the predictor of interest (beta binomial models do not)

**A**

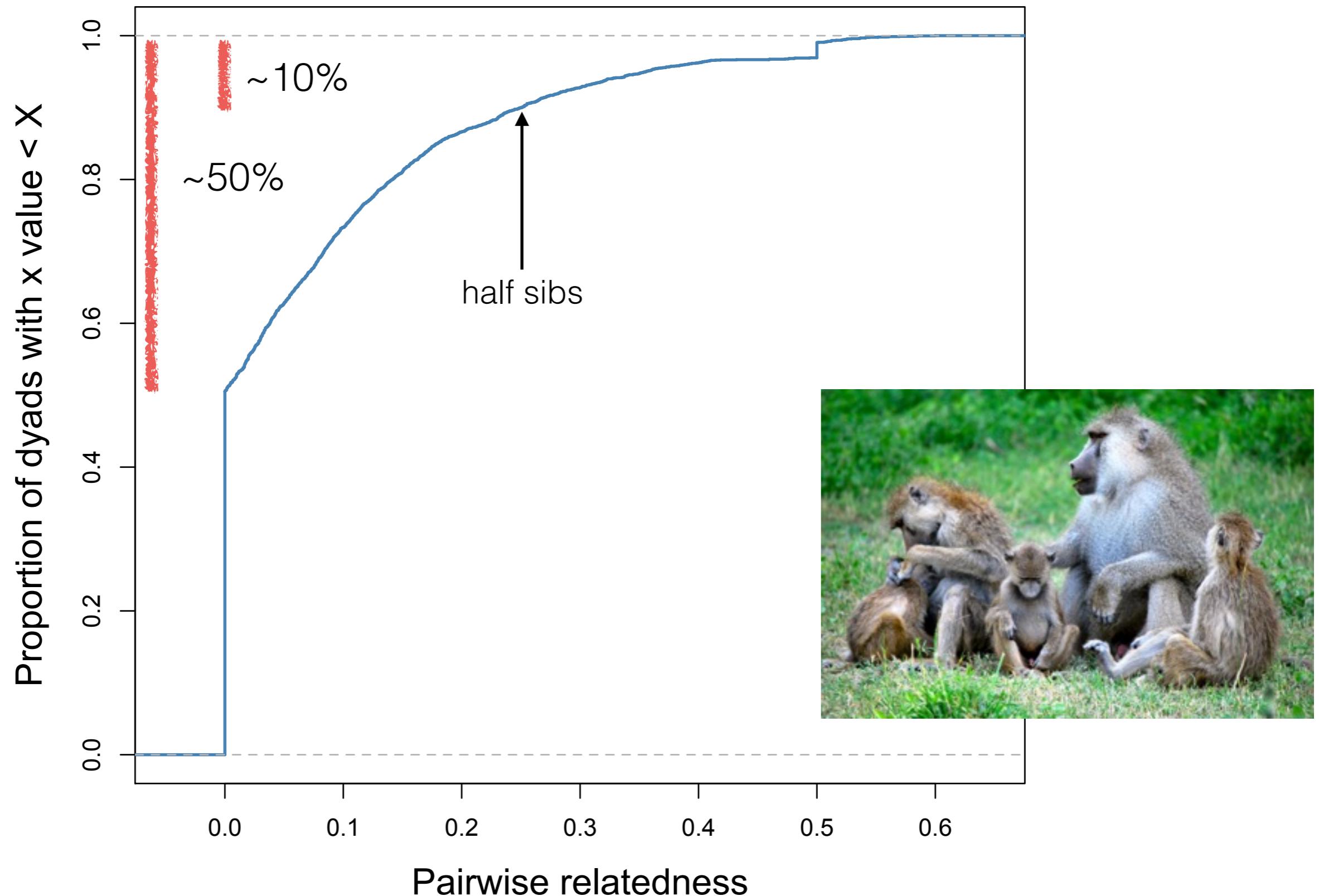


**B**



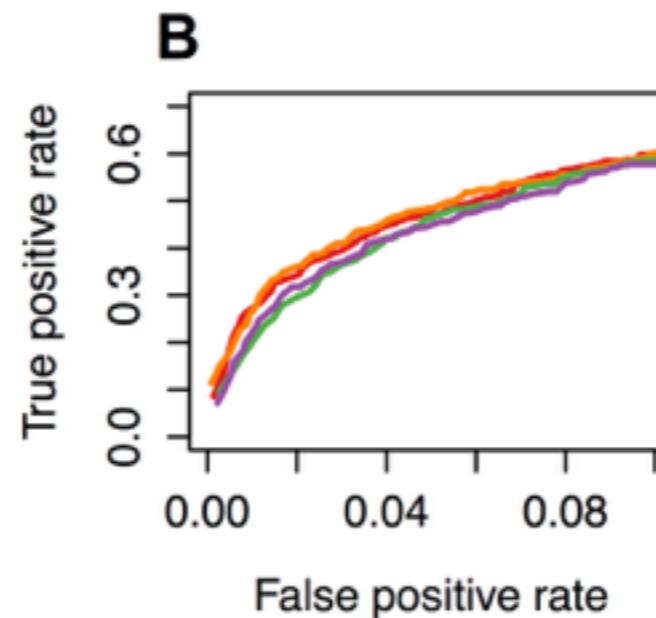
Correlation ( $R^2$ ) between population structure and predictor variable

What about samples that contain relatives (but there's no cofound)?

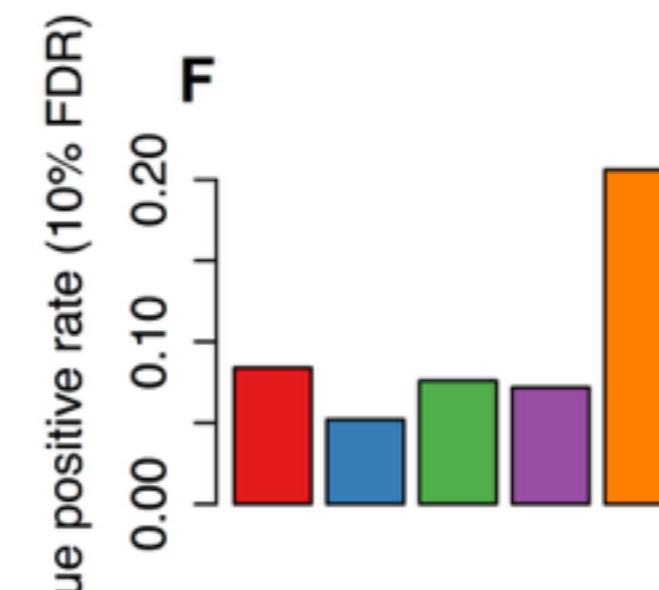
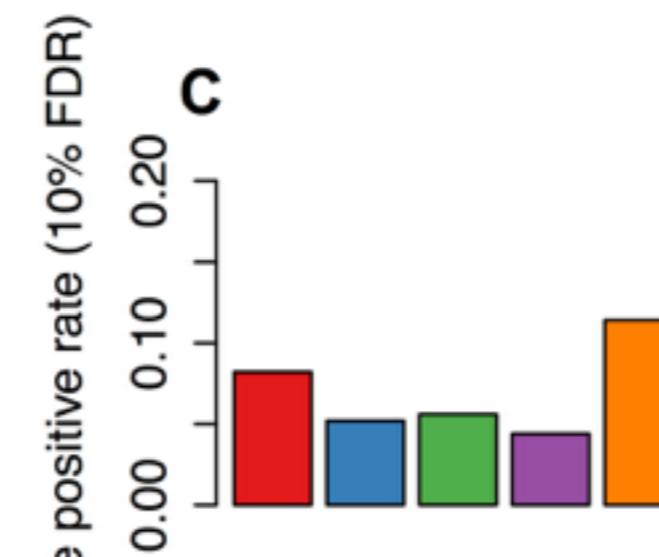
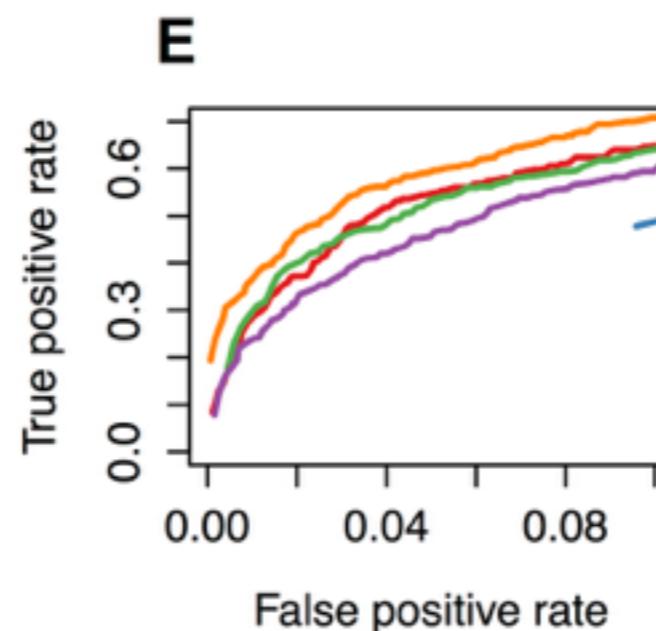


Binomial mixed effects models improve power when the sample contains related individuals (relative to beta binomial model)

$h^2=0.3$



$h^2=0.6$



■ Beta-binomial ■ Binomial ■ GEMMA ■ Linear ■ MACAU

To wrap up, some key points:

- DNA methylation is affected by many factors (age, cell type composition, genetic variation, environmental conditions) that we want to study but should also account for
- Bisulfite sequencing (RRBS or WGBS) is the method of choice for those that work on non-model organisms
- Bisulfite sequencing data is count data, with a lot of variation/unique coverage properties
- Count-based models generally provide better power (than non count-based models)
- Controlling for relatedness is important and can (i) reduce false positive associations and (ii) increase power to detect true effects

Some important things we didn't talk about:

## Some important things we didn't talk about:

- Batch effects (they exist in DNA methylation data too!)
  - Bisulfite conversion rate, sequencing depth, library prep, etc.
  - Same approaches to identify them apply, e.g., PCA

## Some important things we didn't talk about:

- Batch effects (they exist in DNA methylation data too)!
  - Bisulfite conversion rate, sequencing depth, library prep, etc.
  - Some approaches to identify them apply, e.g., PCA
- False discovery rates and empirical nulls (they apply here too)!

## Some important things we didn't talk about:

- Batch effects (they exist in DNA methylation data too)!
  - Bisulfite conversion rate, sequencing depth, library prep, etc.
  - Some approaches to identify them apply, e.g., PCA
- False discovery rates and empirical nulls (they apply here too)!
- Cell type heterogeneity
  - Measure it!
  - Deconvolution
  - Understand magnitude of bias using methylation data from pure cell types

## Some important things we didn't talk about:

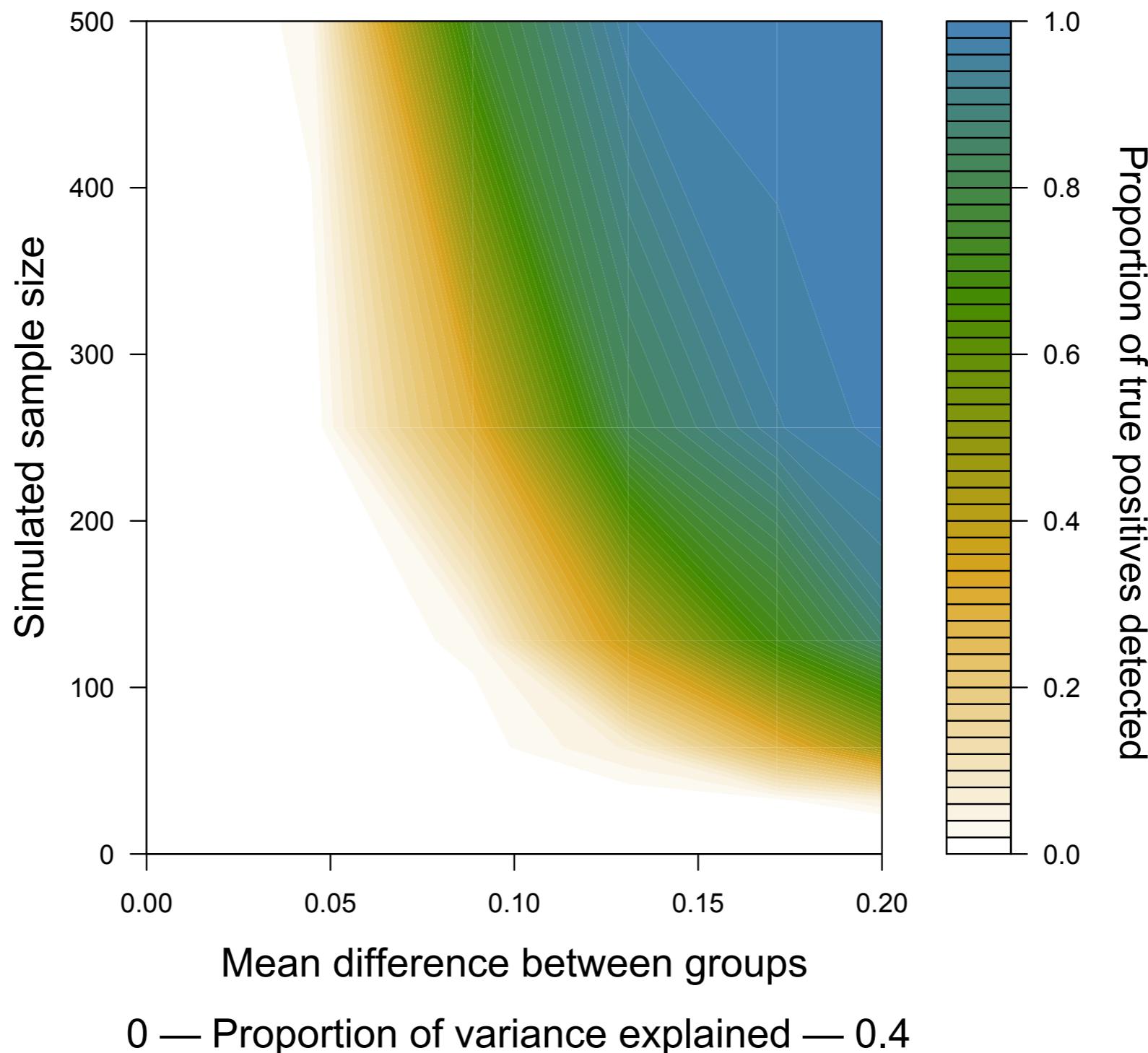
- Batch effects (they exist in DNA methylation data too!)
  - Bisulfite conversion rate, sequencing depth, library prep, etc.
  - Some approaches to identify them apply, e.g., PCA
- False discovery rates and empirical nulls (they apply here too)!
- Cell type heterogeneity
  - Measure it!
  - Deconvolution
  - Understand magnitude of bias using methylation data from pure cell types
- Other kinds of analyses, e.g.,
  - Genotyping from bisulfite-sequencing data
  - meQTL
  - Identifying differentially methylated regions

## Some important things we didn't talk about:

- Power and sample size
- More biological replicates is always good (if you're on a budget, think about using RRBS and/or reducing your sequencing depth)

## Some important things we didn't talk about:

- Power and sample size
- More biological replicates is always good (if you're on a budget, think about using RRBS and/or reducing your sequencing depth)



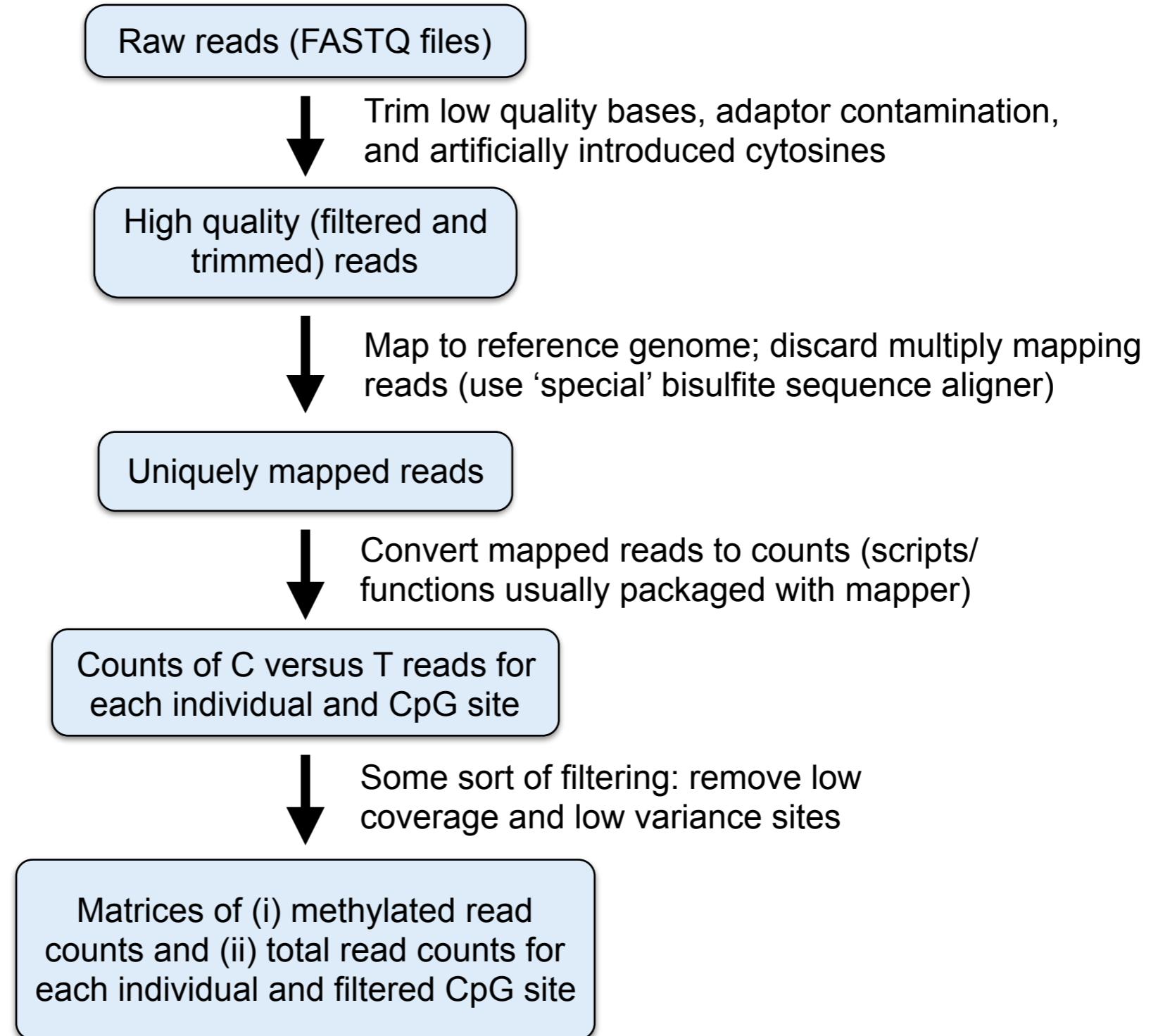
# **DNA methylation**

## hands on data analysis

Amanda Lea  
September 22, 2016

# Bisulfite sequencing data

- First - process, map, and filter your data



# Bisulfite sequencing data

- Now you have a count matrix!

Data matrix for  $n$  individuals measured at  $q$  loci  
(methylated reads / total reads)

	ID_1	ID_2	...	ID_n
loc_1	2/7	20/70	...	40/140
loc_2	2/3	40/60	...	40/140
loc_3	5/5	50/50		100/100
...	...	...	...	...
loc_q	0/0	1/3	...	4/6

# An RRBS dataset



FYI - you can use these simulations to do power analyses

- Is your sample size sufficient to detect your expected effect size?

