

# Making sense of DNA methylation data

Peter Hickey



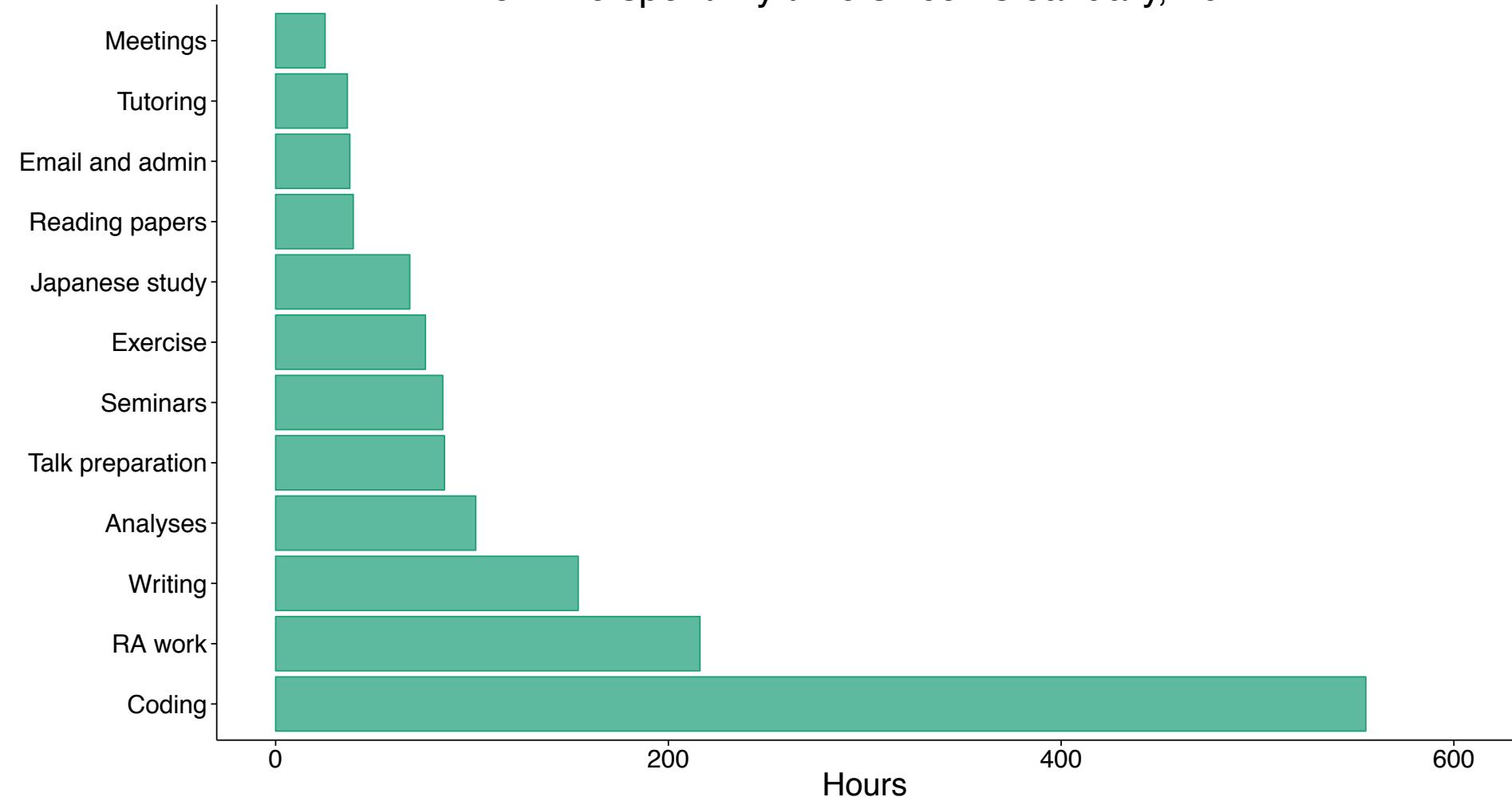
@PeteHaitch

15 September 2014

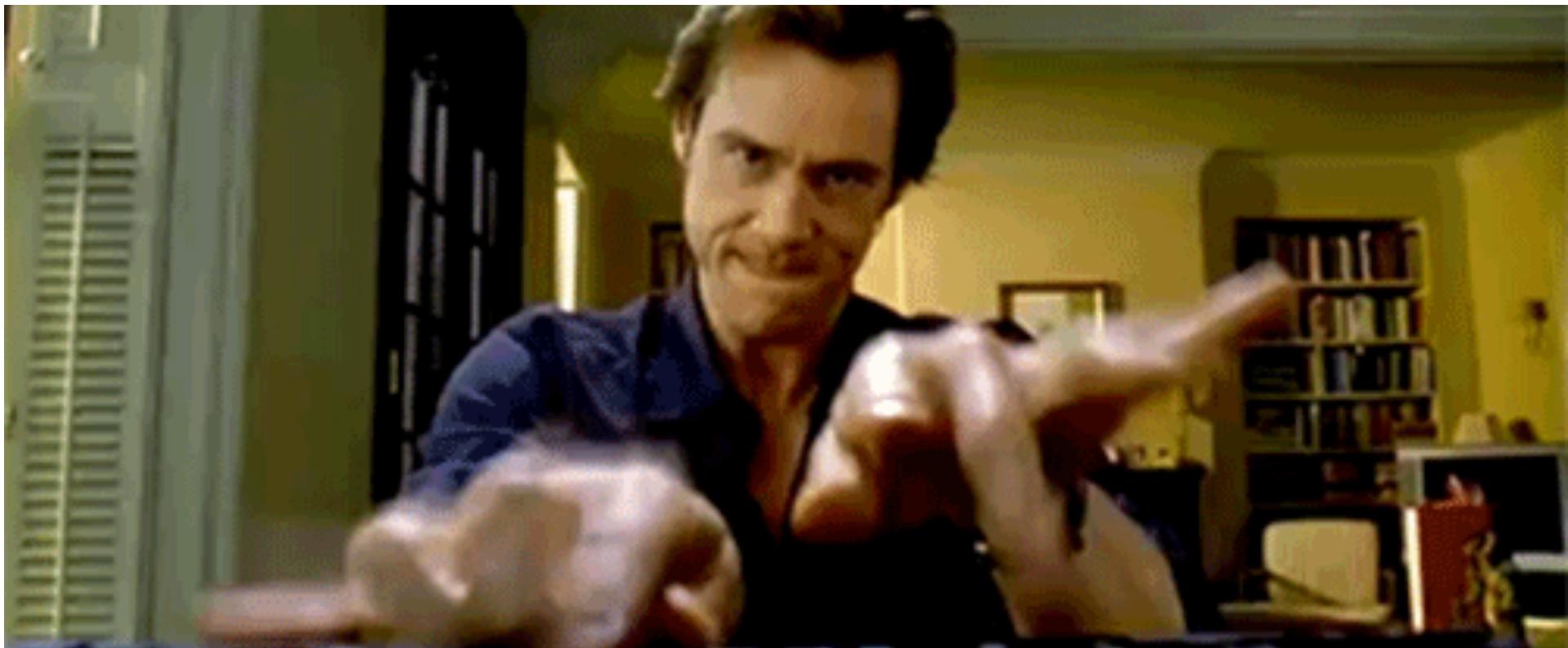
How I spend my time

**“Oh, you’re a statistician...”**

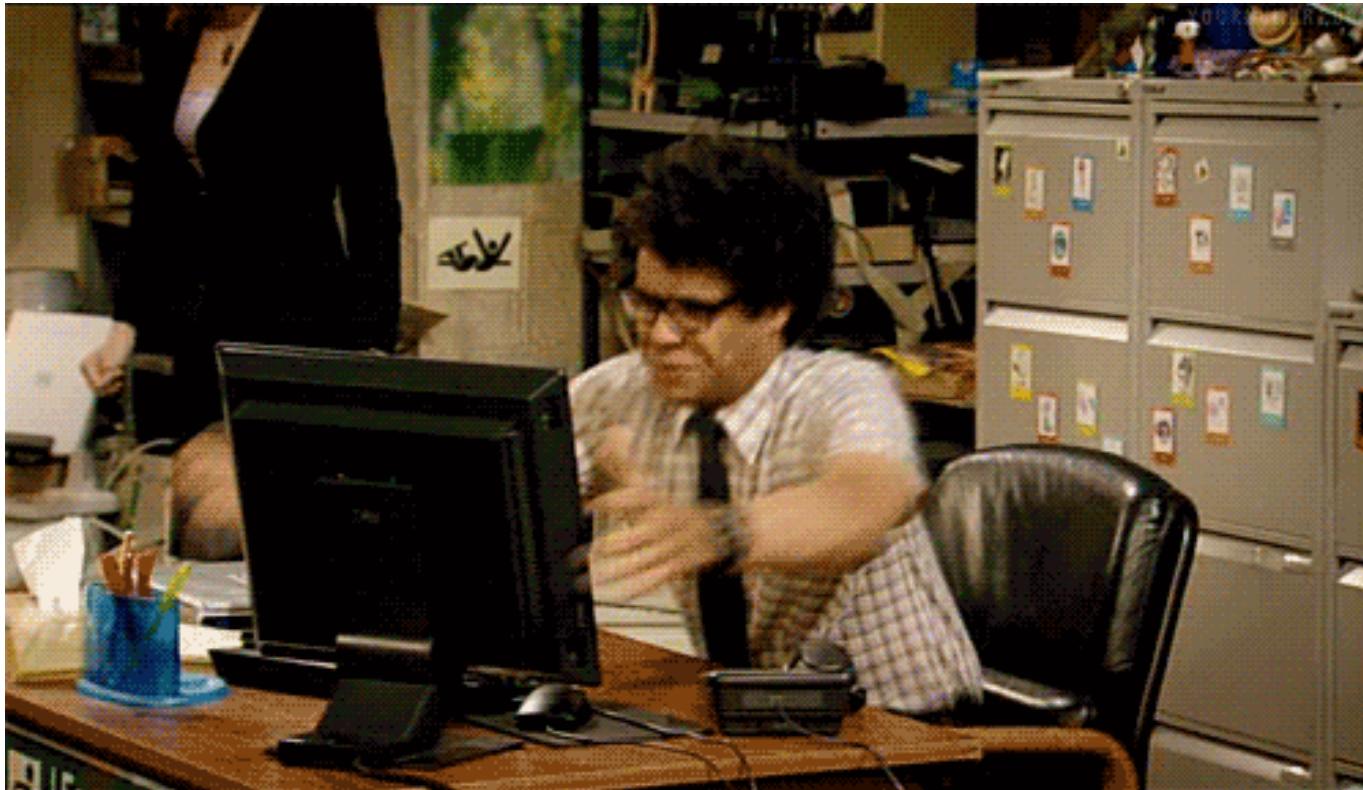
# How I've spent my time since 13 January, 2014



# How I spend my time



# How I spend my time



# Project aim

*How to analyse whole-genome bisulfite-sequencing experiments to learn about DNA methylation?*

# Exploratory data analyses



# Developing methods



# Implementing methods



# Collaboration



**5mCs and one assay**

ACGCGAAACGTTCTATCGG  
TGCCTTGCAAGATAGCC

ACGCGAAACGTTCTATCGG  
TGC CGCTTTGCAAGATAGCC

m

m

m

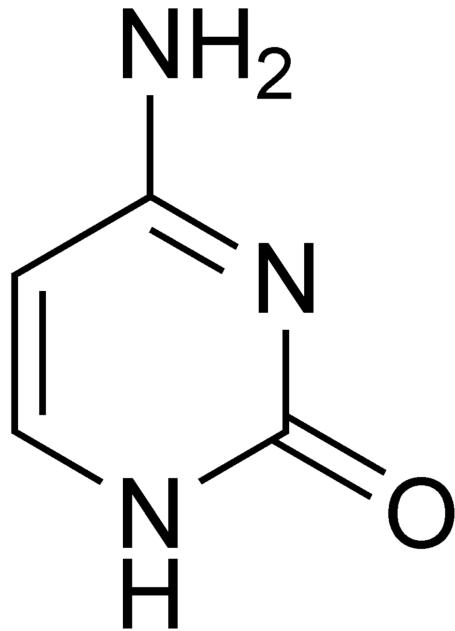
m

ACGCGAAACGTTCTATCGG  
TGCGCTTGTCAAGATAGCC

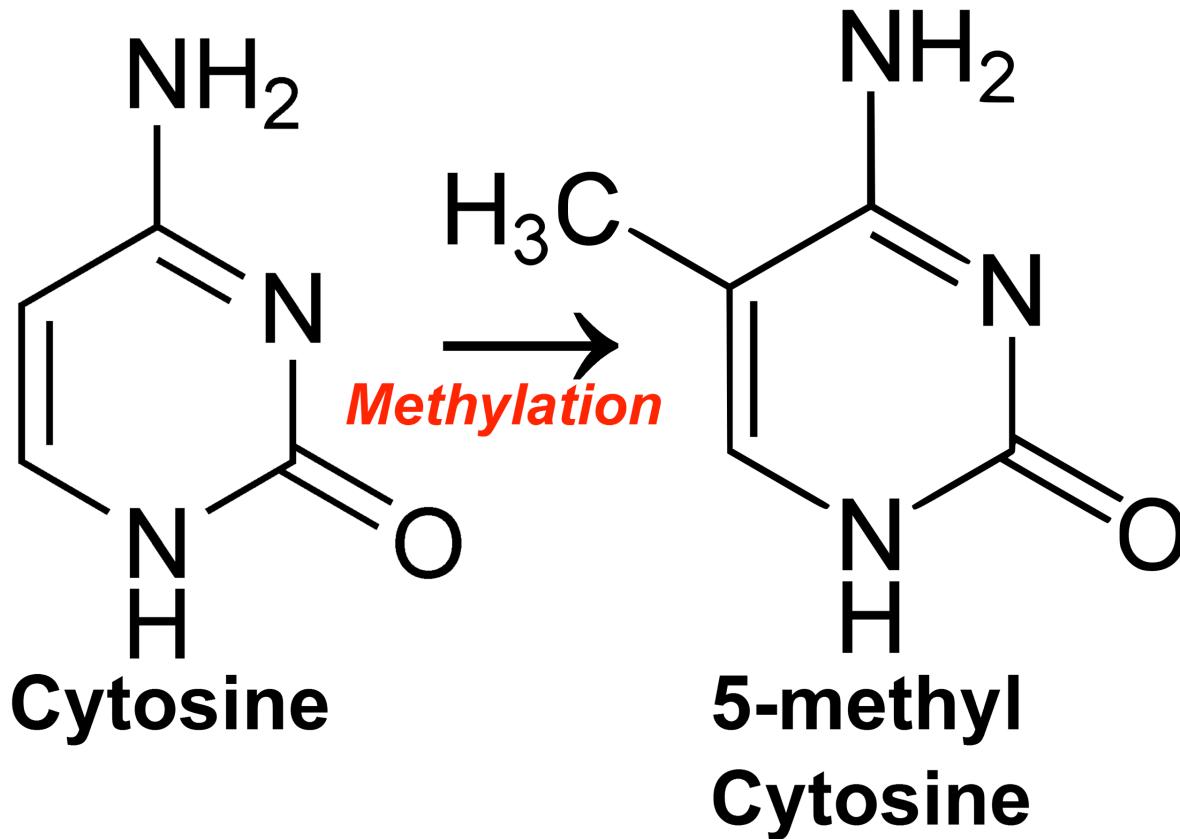
m

m

m



**Cytosine**



ACGCGAAACGTTCTATCGG  
TGC CGCTTTGCAAGATAGCC

m

m

m

m

ACGCGAAACGTTCTATCGG  
TGCGCTTGTCAAGATAGCC

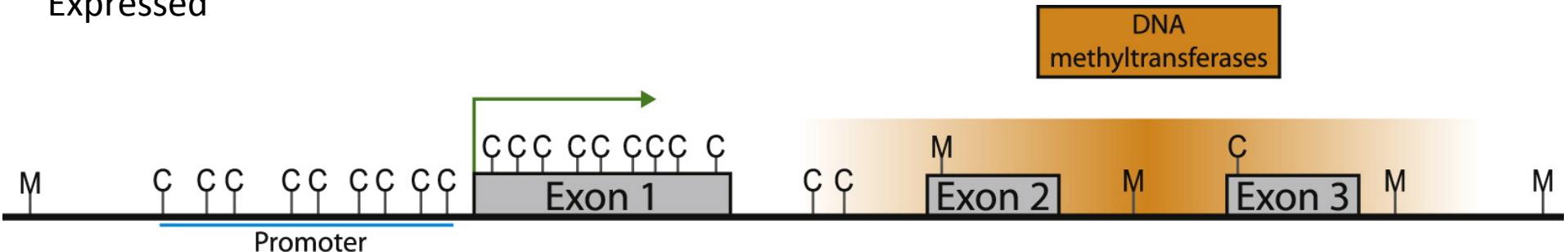
m

m

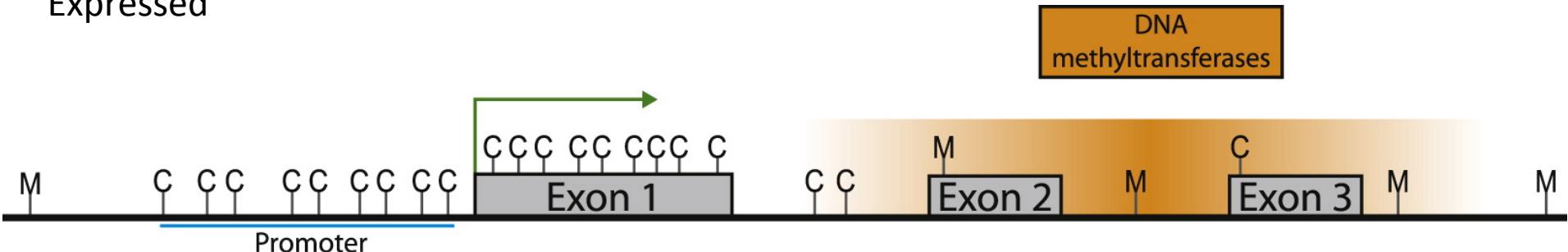
m



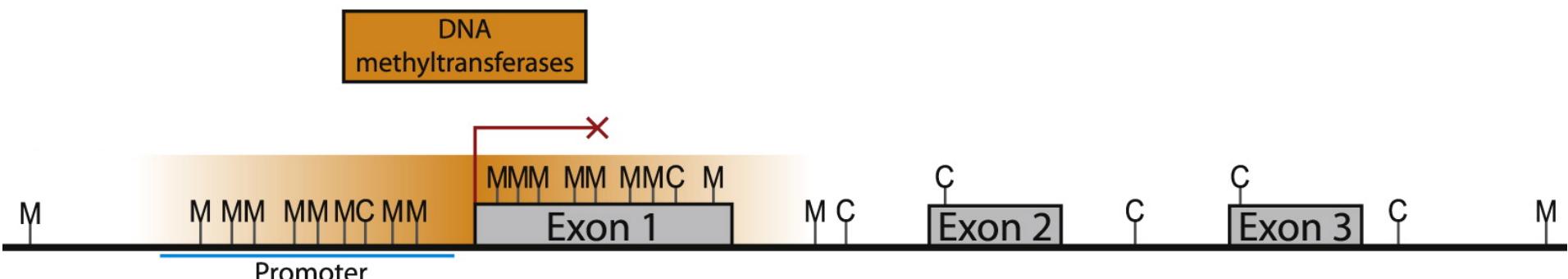
Expressed



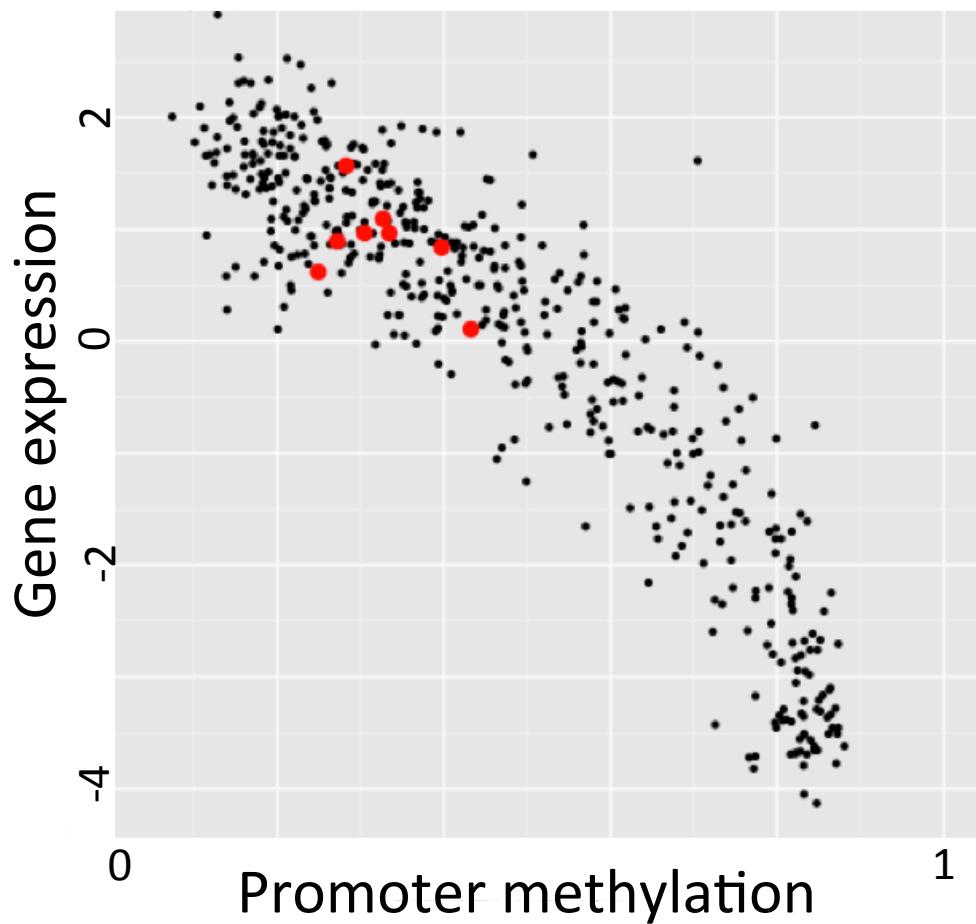
## Expressed



Not expressed

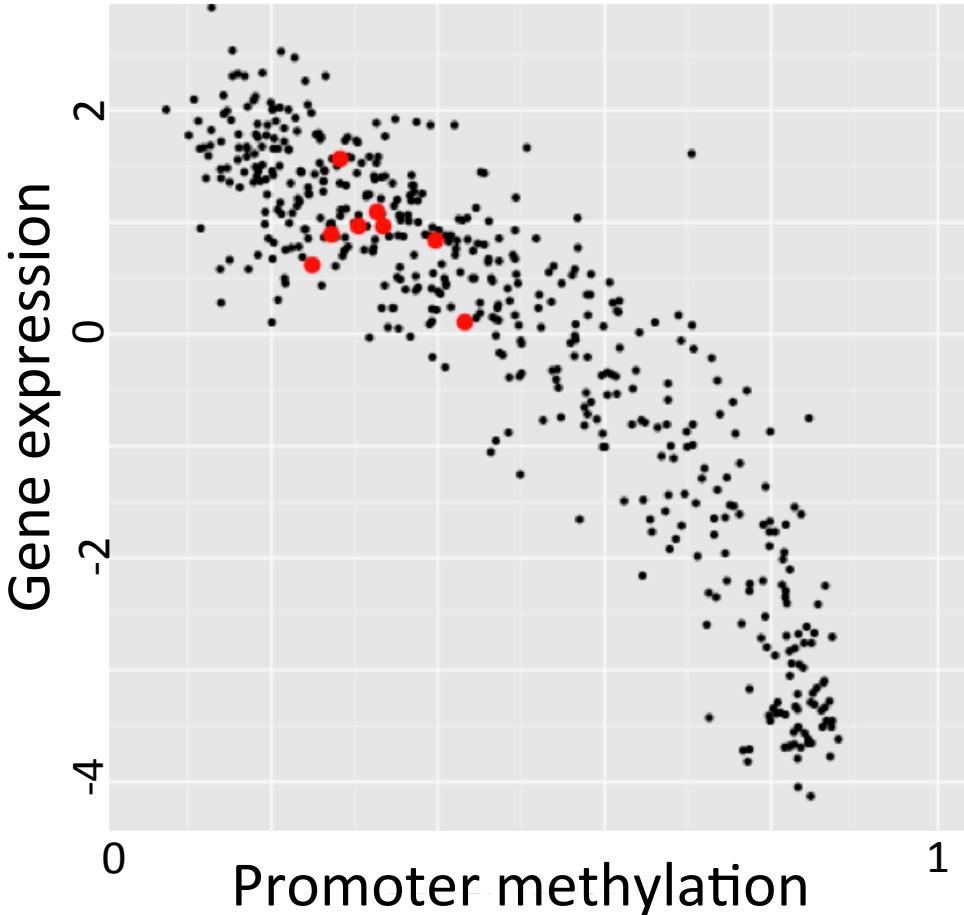


*AMT*



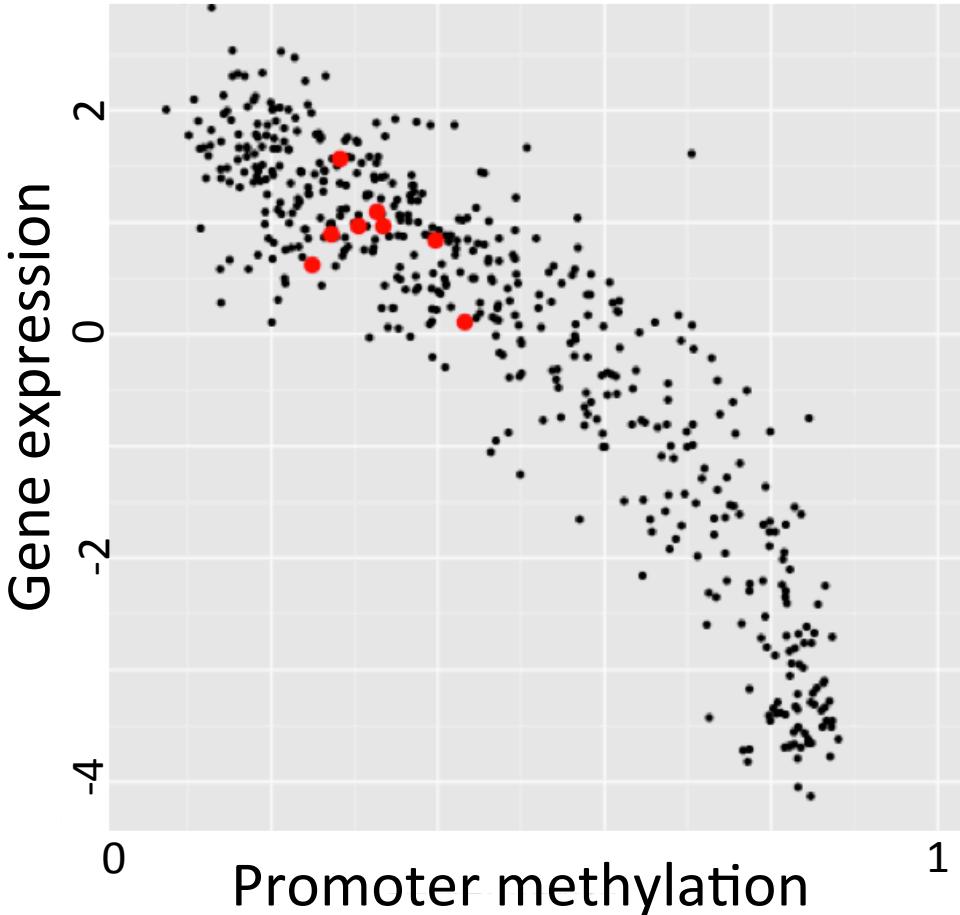
Cancer Genome Atlas Research Network. "Integrated genomic analyses of ovarian carcinoma." Nature 474.7353 (2011): 609-615.

*AMT*

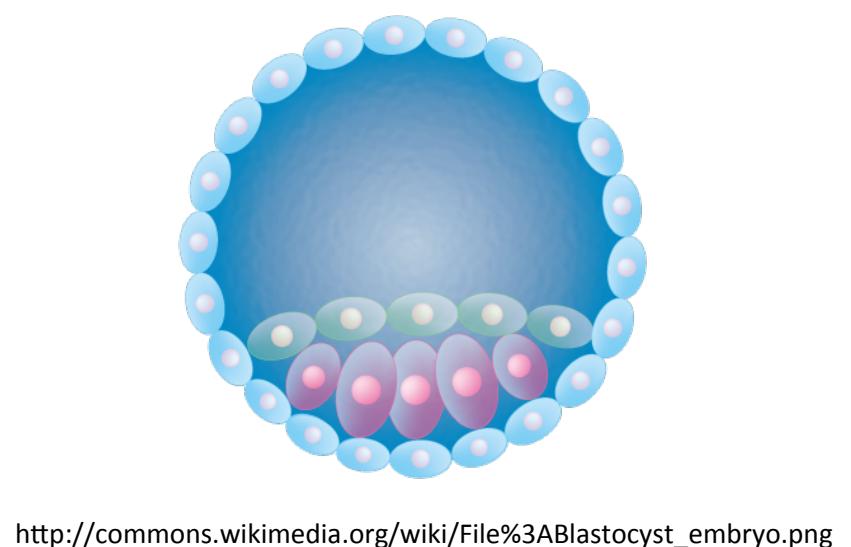


Cancer Genome Atlas Research Network. "Integrated genomic analyses of ovarian carcinoma." Nature 474.7353 (2011): 609-615.

# *AMT*



Cancer Genome Atlas Research Network. "Integrated genomic analyses of ovarian carcinoma." *Nature* 474.7353 (2011): 609-615.



m

m

m

m

ACGCGAAACGTTCTATCGG

TGCGCTTTCGAAGATAGCC

m

m

m

m

m

m

m

ACGCGAAACGTTCTATCGG

ACGCGAAACGTTCTATCGG

ACGCGAAACGTTCTATCGG

+

PCR amplification

=

ACGCGAAACGTTCTATCGG

+

PCR amplification

=

ACGCGAAACGTTCTATCGG

ACGCGAAACGTTCTATCGG

ACGCGAAACGTTCTATCGG

+

Sodium bisulfite

=

ACGCGAAACGTTCTATCGG

+

Sodium bisulfite

=

ACGUGAAACGTTCTATCGG

ACGUGAAACGTTCTATCGG

ACGUGAAACGTTCTATCGG

+

PCR amplification

=

ACGUGAAACGTTCTATCGG

+

PCR amplification

=

ACGTGAAACGTTCTATCGG

ACGUGAAACGTTCTATCGG

+

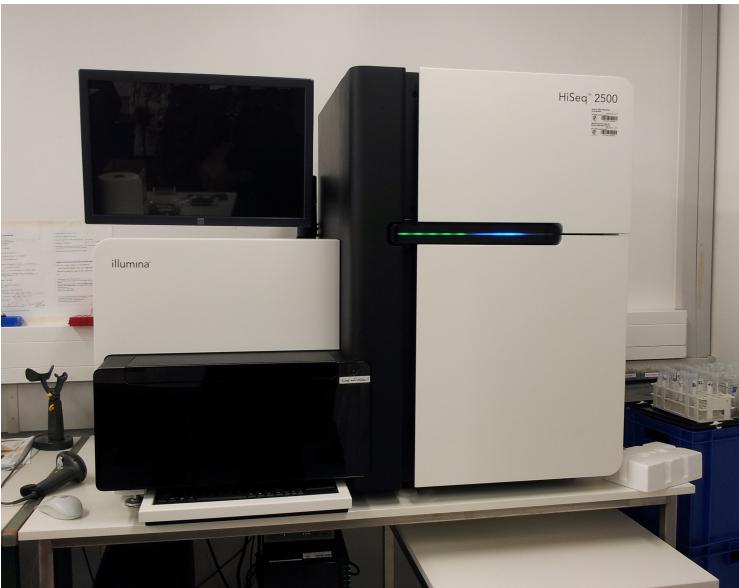
PCR amplification

=

ACGTGAAACGTTCTATCGG

# Bisulfite treatment of DNA

+



=

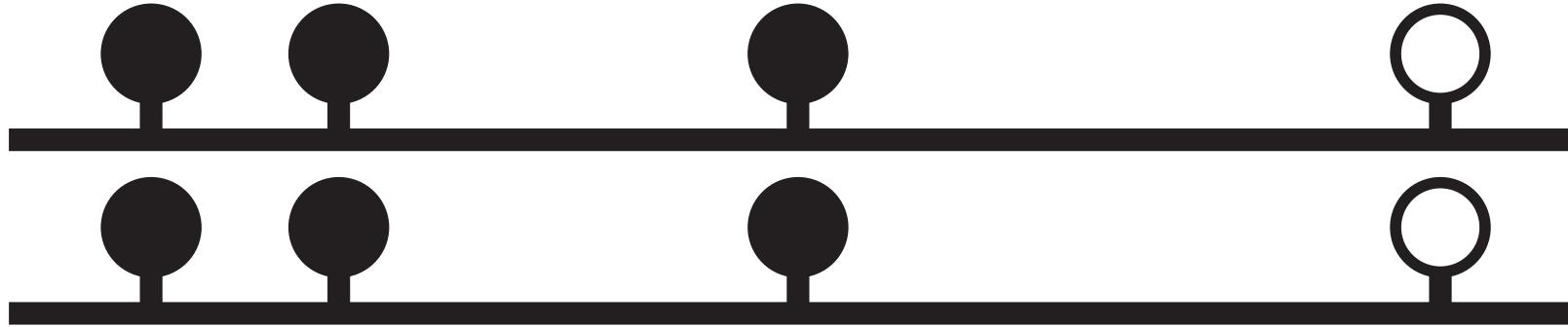
# Whole-genome bisulfite-sequencing

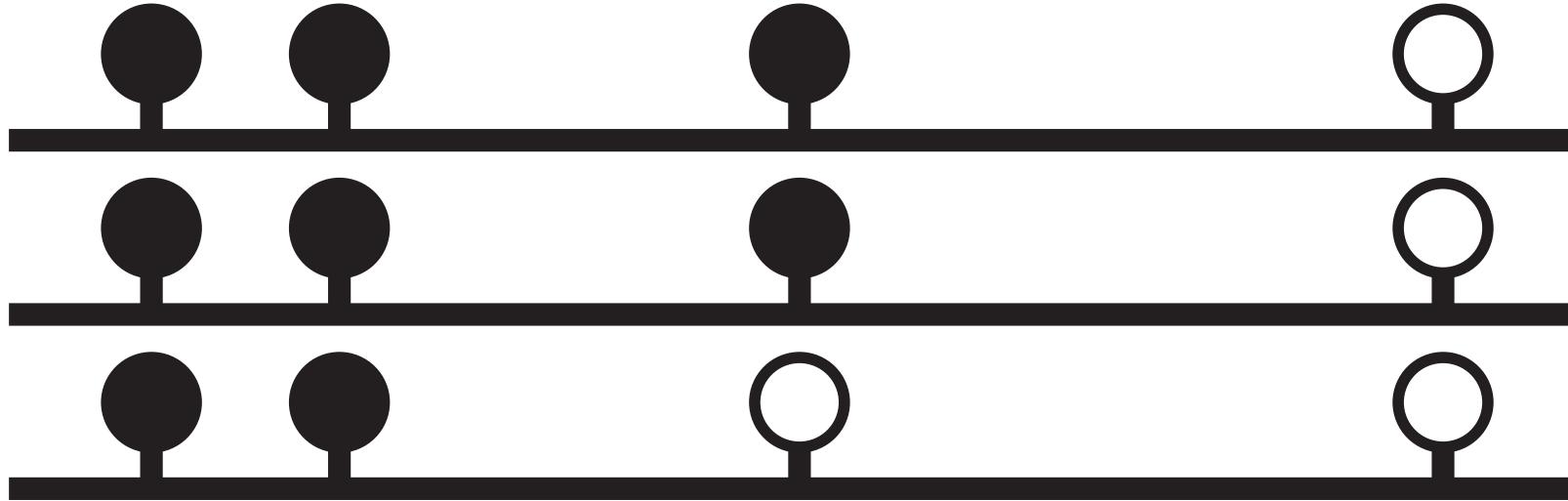
ACGCGAAACGTTCTATCG

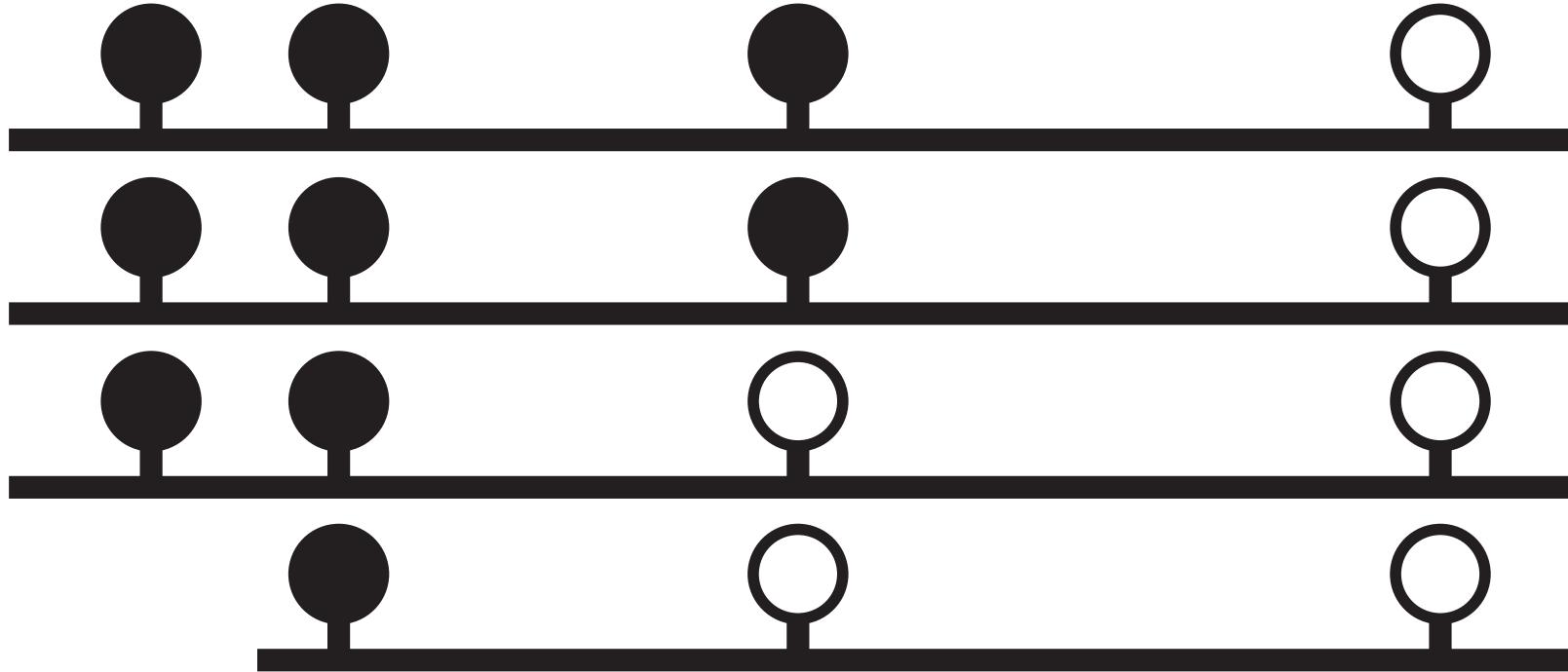
ACGCGAAACGTTCTATCG

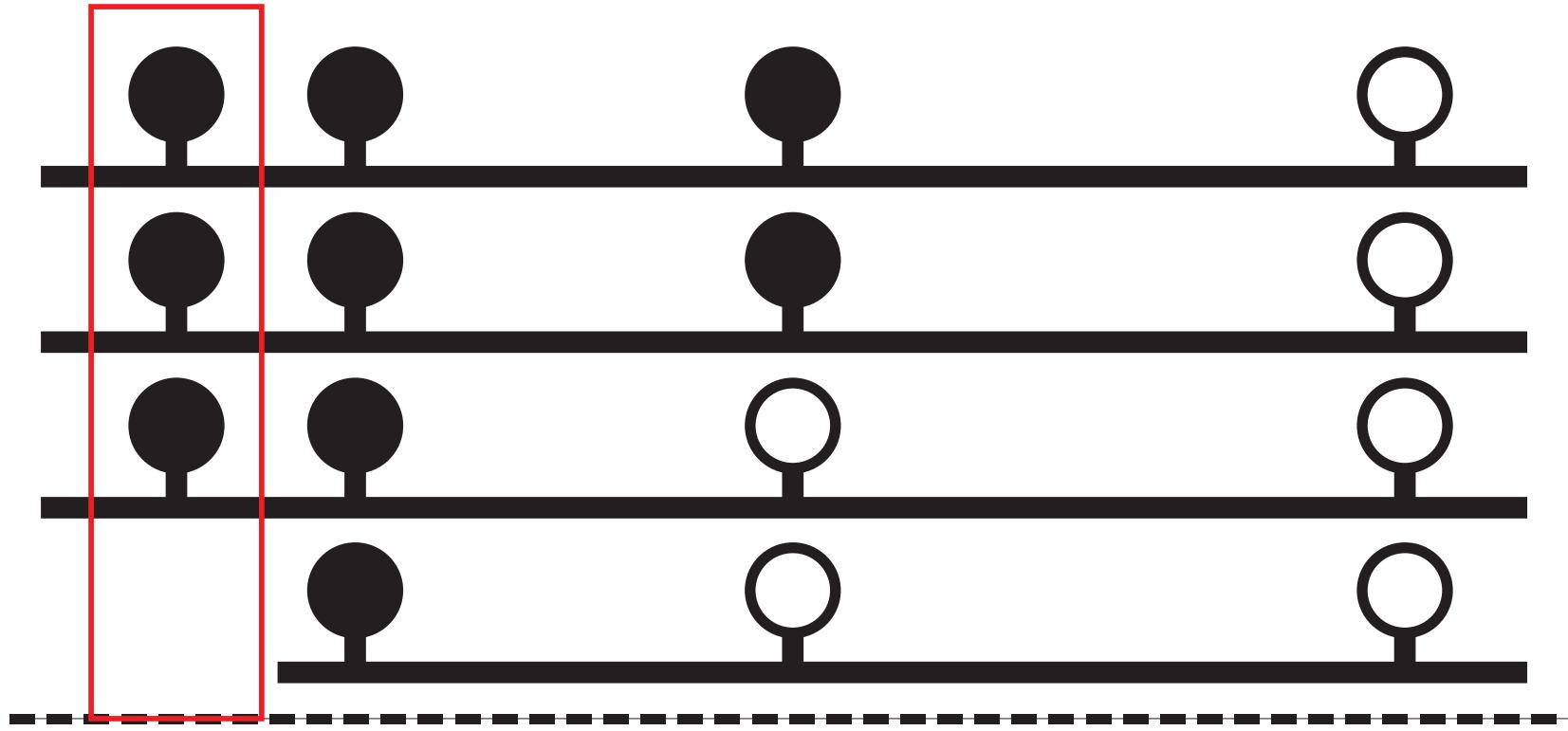




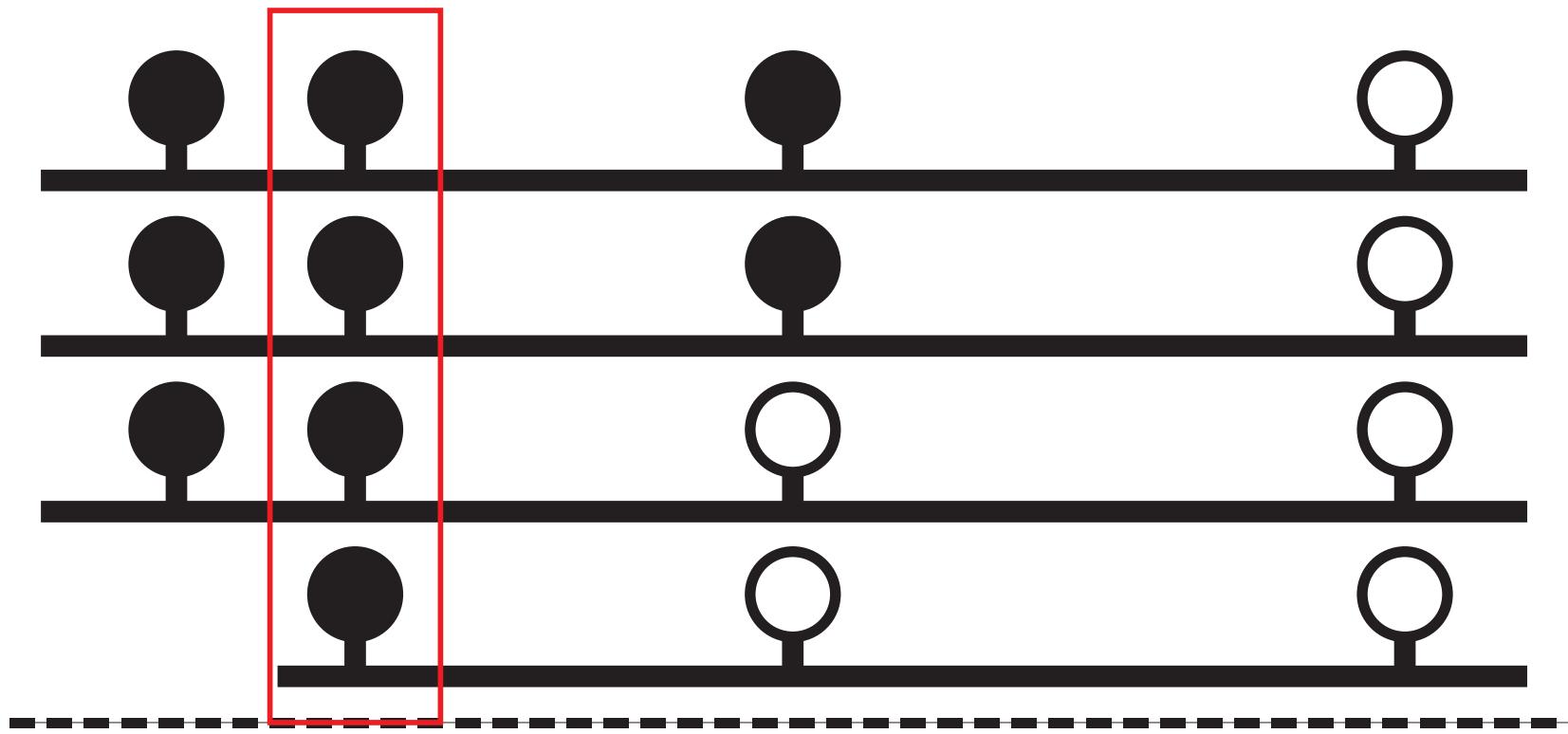




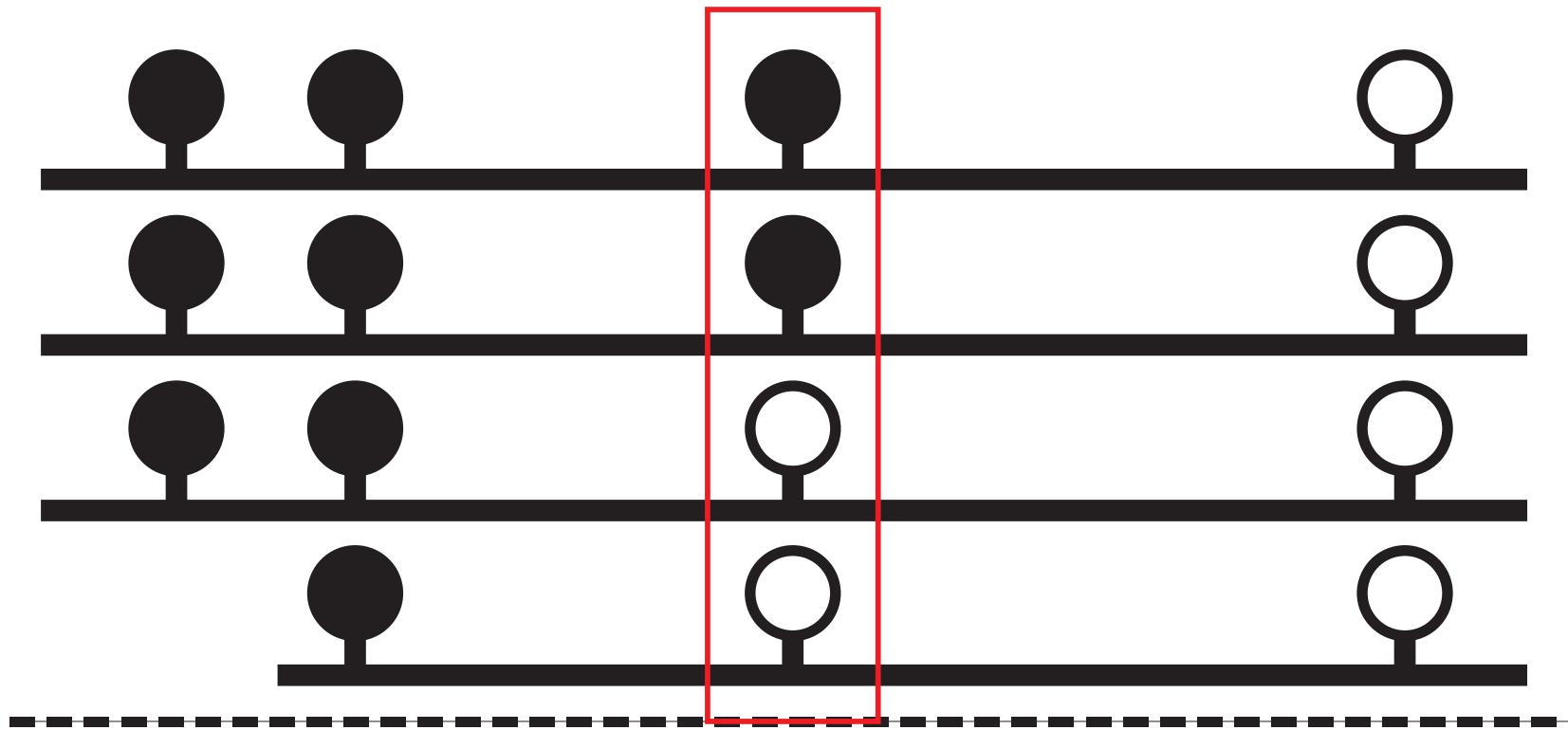




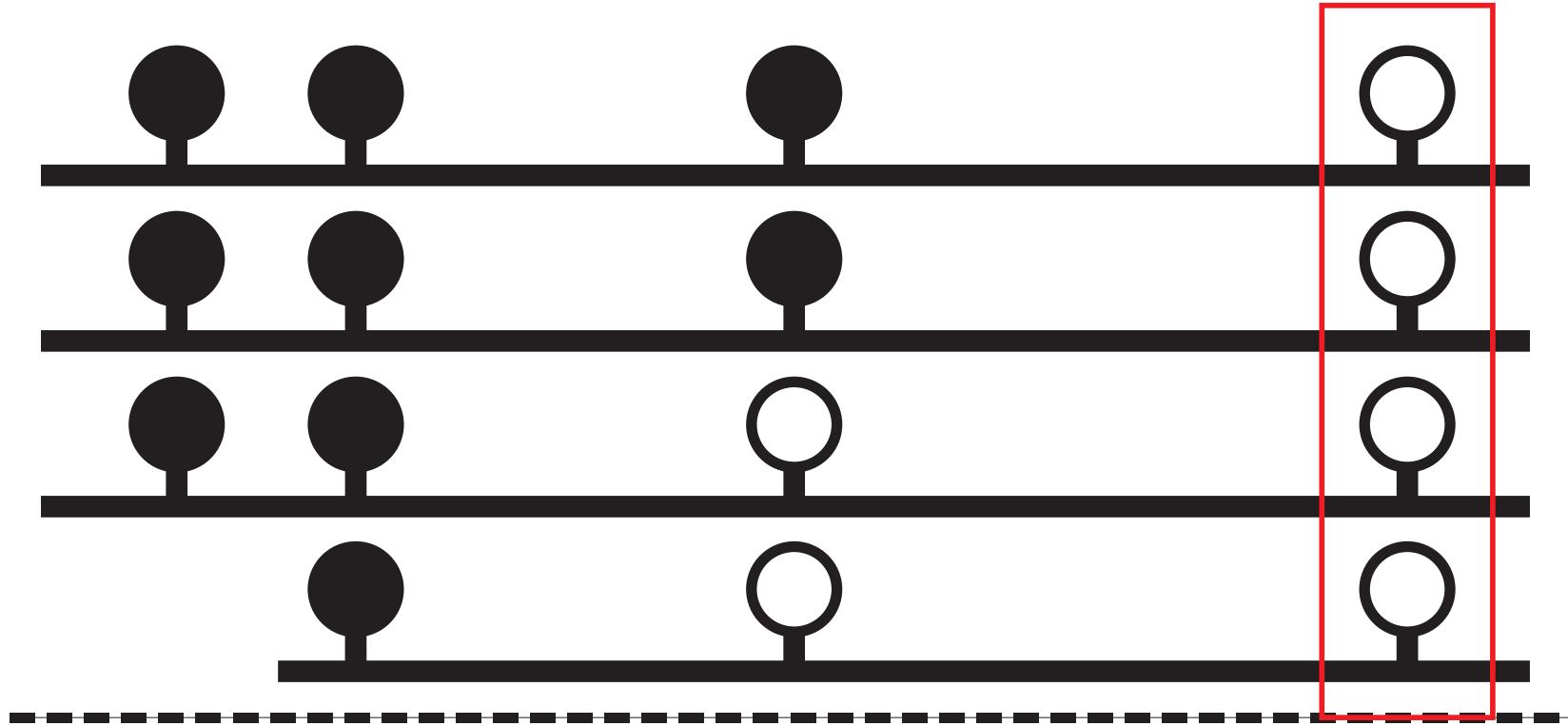
$$\beta_1 = 3/3$$



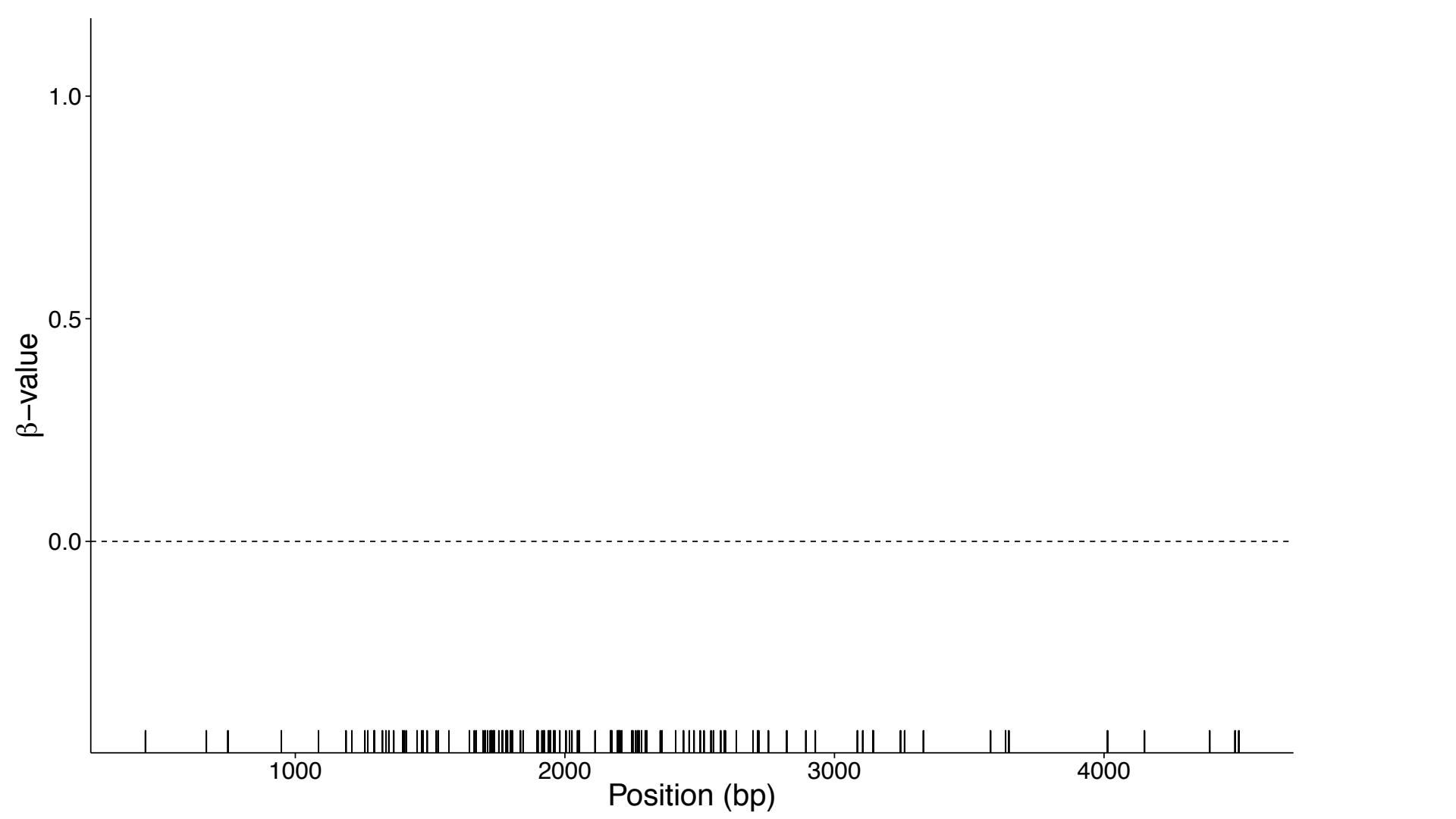
$$\beta_2 = 4/4$$

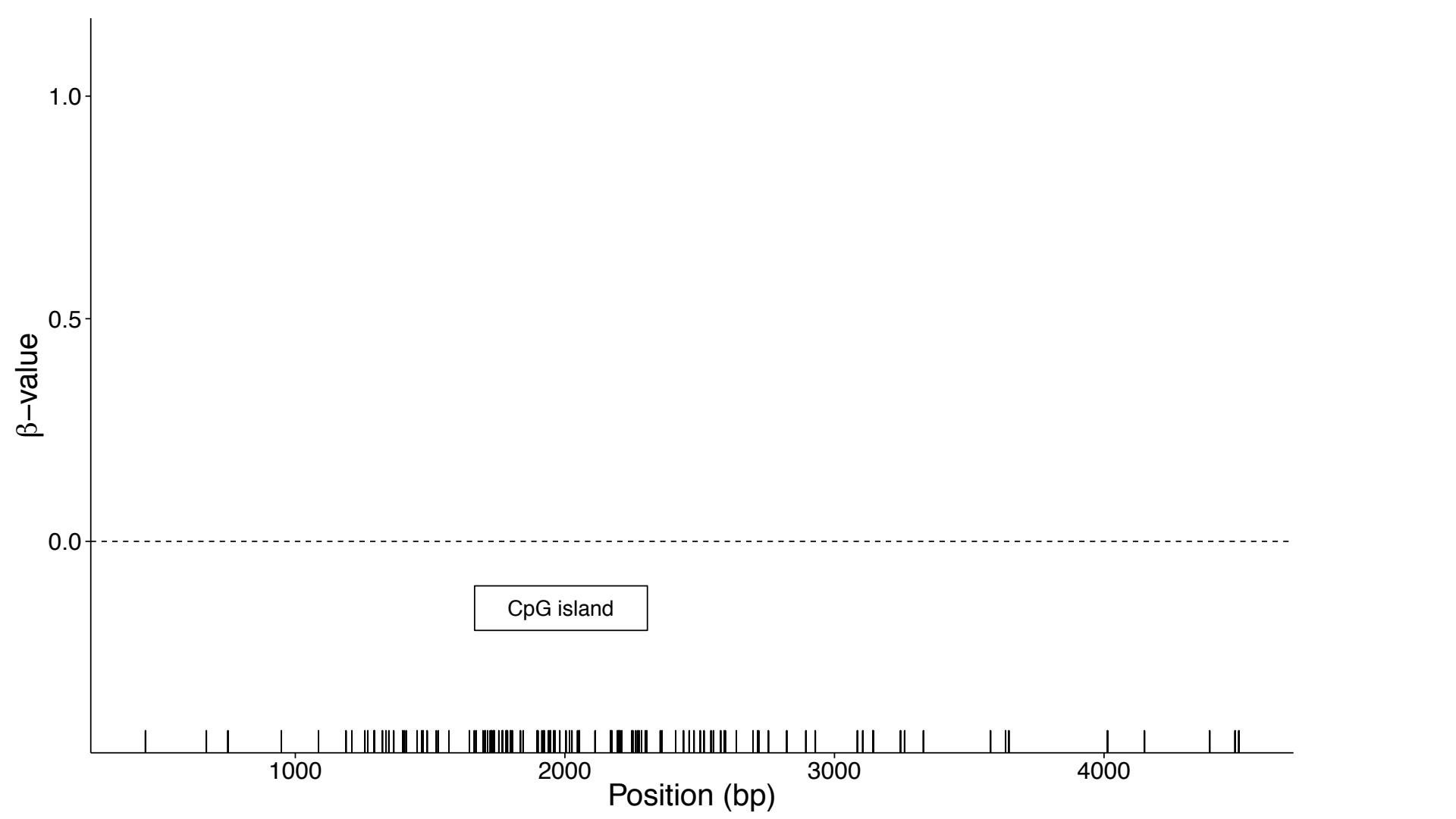


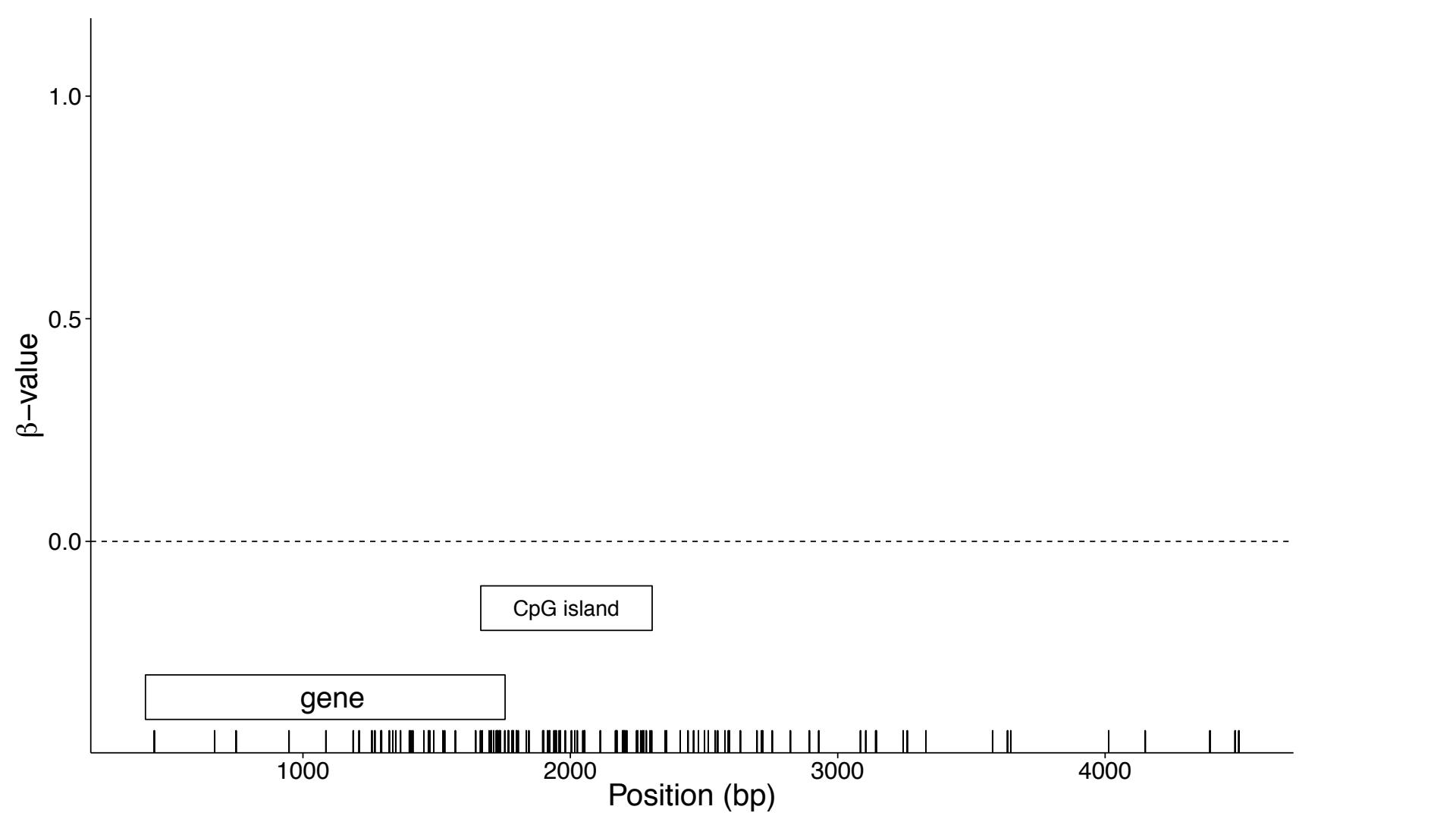
$$\beta_3 = 2/4$$

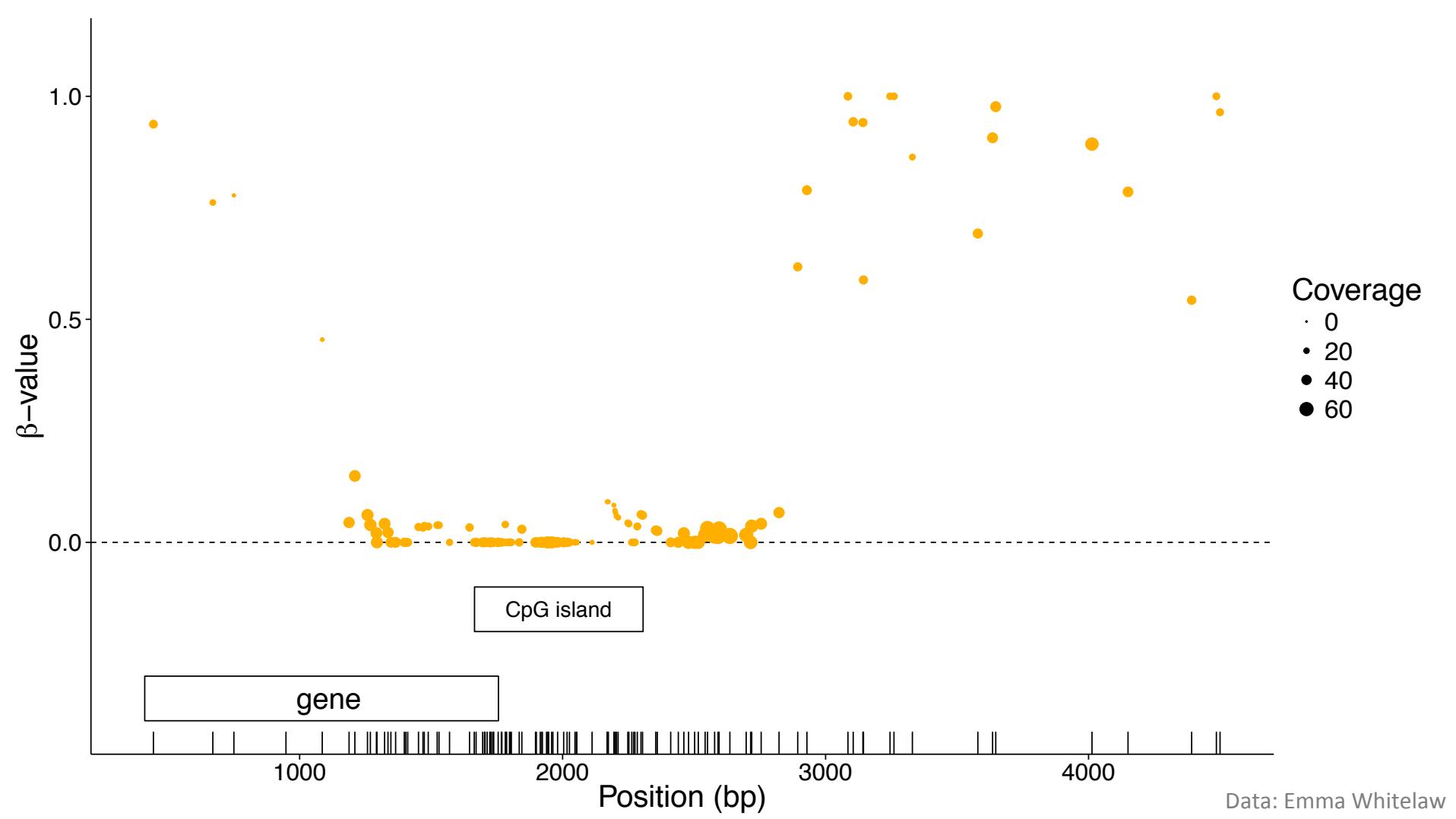


$$\beta_4 = 0/4$$





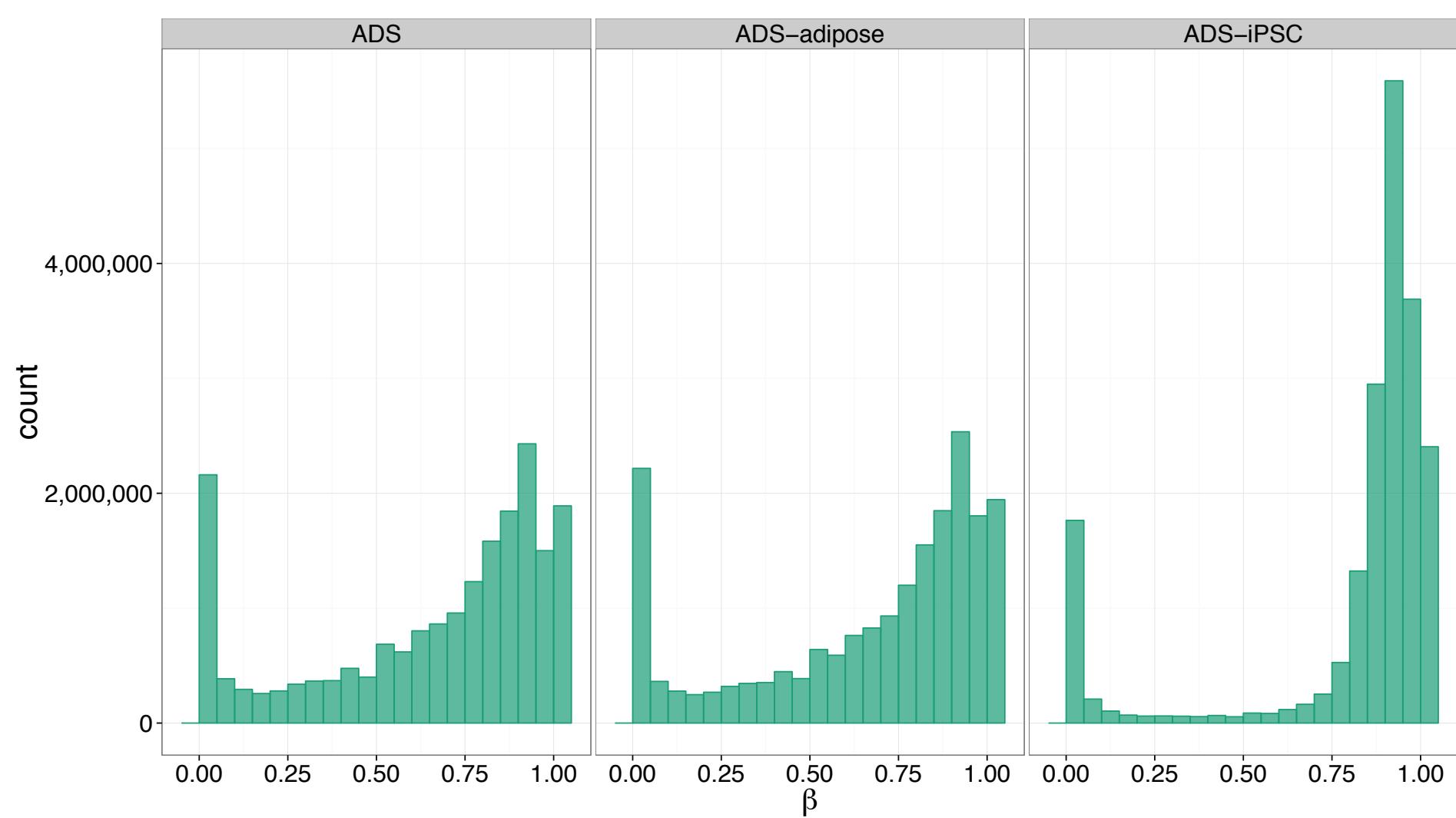


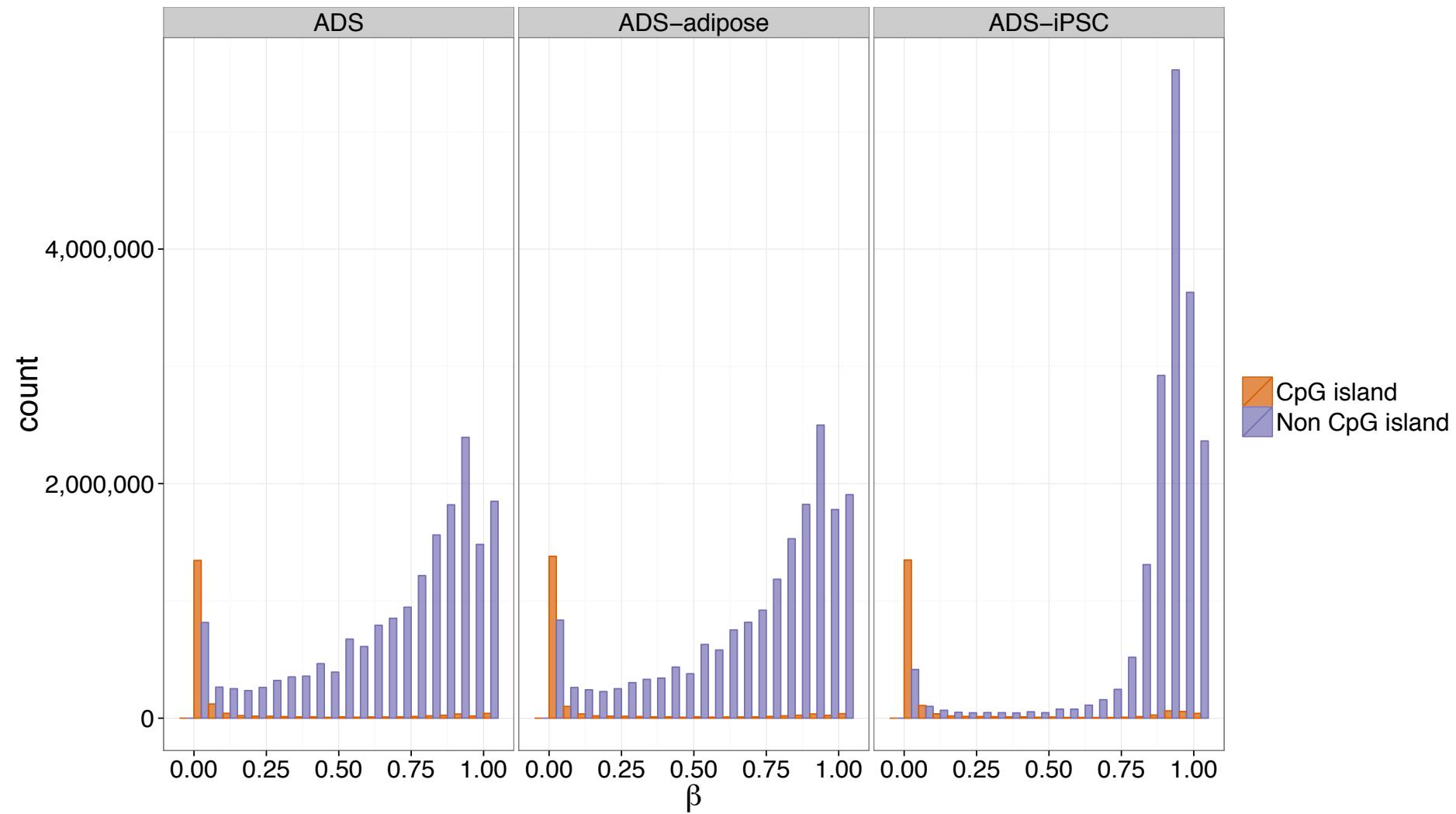


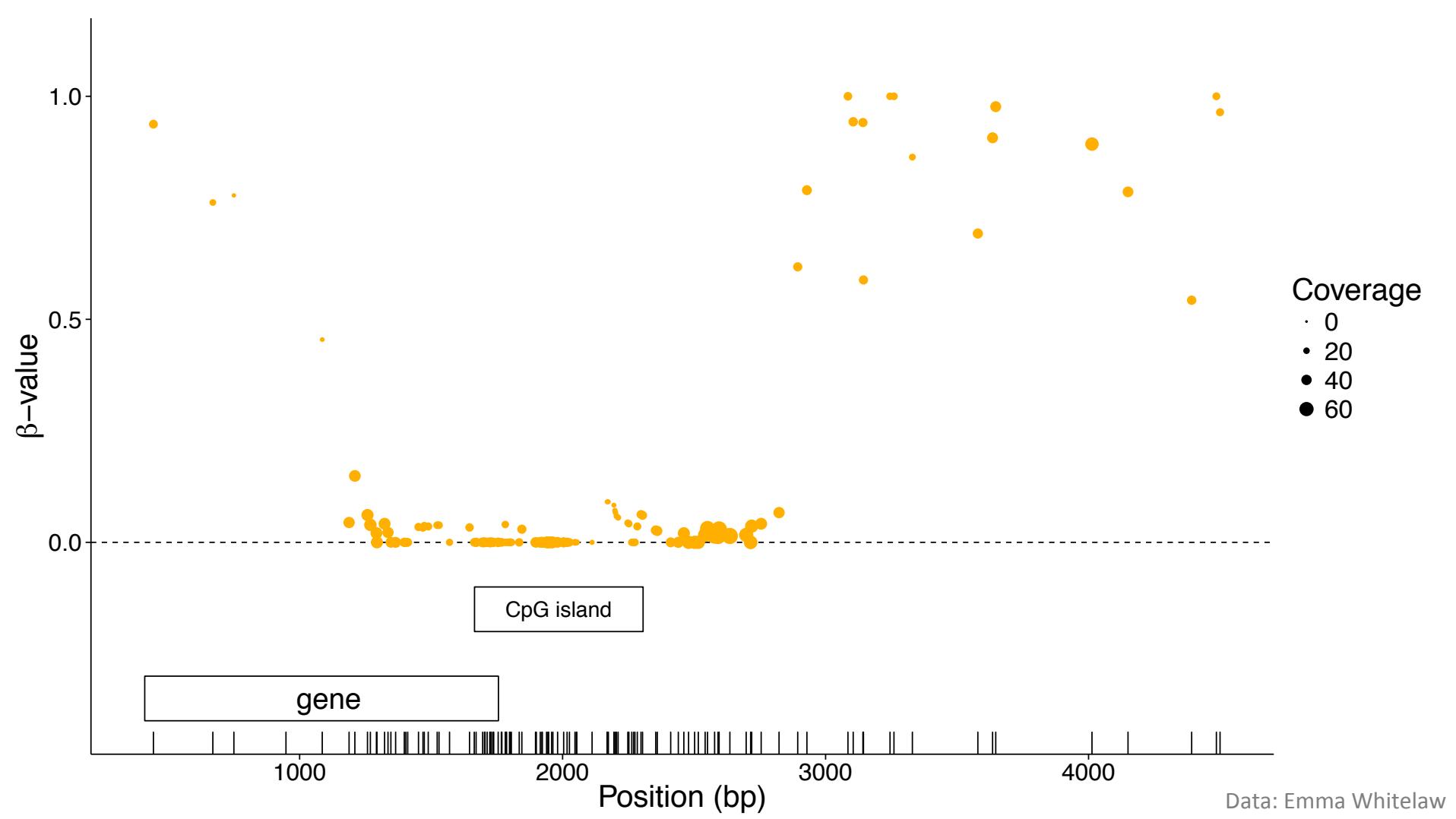
# Lister data

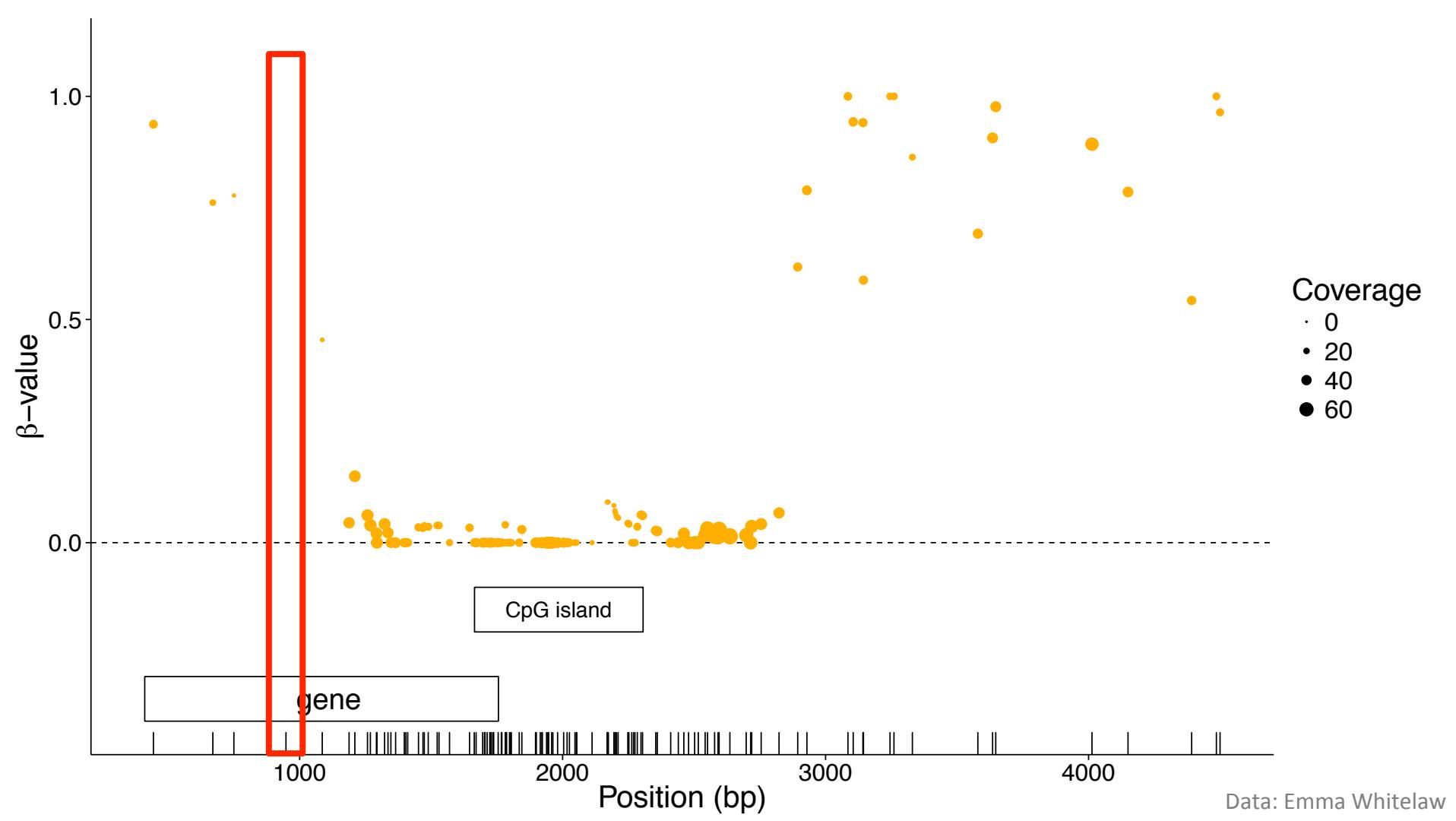
	<i>ADS</i>	<i>ADS-adipose</i>	<i>ADS-iPSC</i>
<b>Organism</b>	Human (female)	Human (female)	Human (female)
<b>Cell type</b>	Somatic	Somatic	Induced pluripotent stem cell (iPSC)
<b>Description</b>	Adipose	Adipocytes derived from <i>ADS</i>	iPSC line derived from <i>ADS</i>
<b>Sequencing</b>	75 bp paired-end	75 bp paired-end	75 bp paired-end
<b>Average coverage</b>	23×	24×	26×

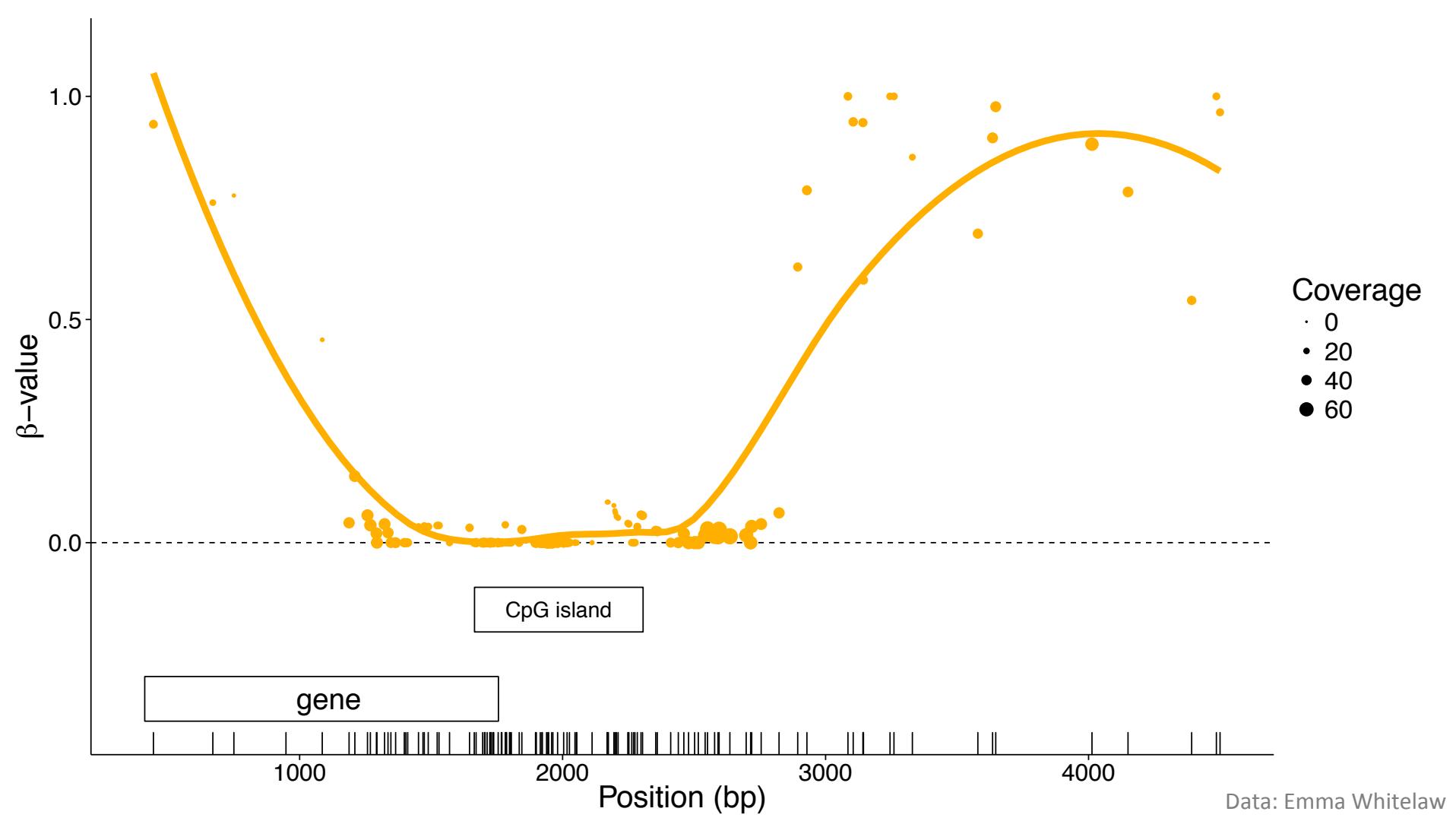
Lister, Ryan, et al. "Hotspots of aberrant epigenomic reprogramming in human induced pluripotent stem cells." Nature 471.7336 (2011): 68-73.

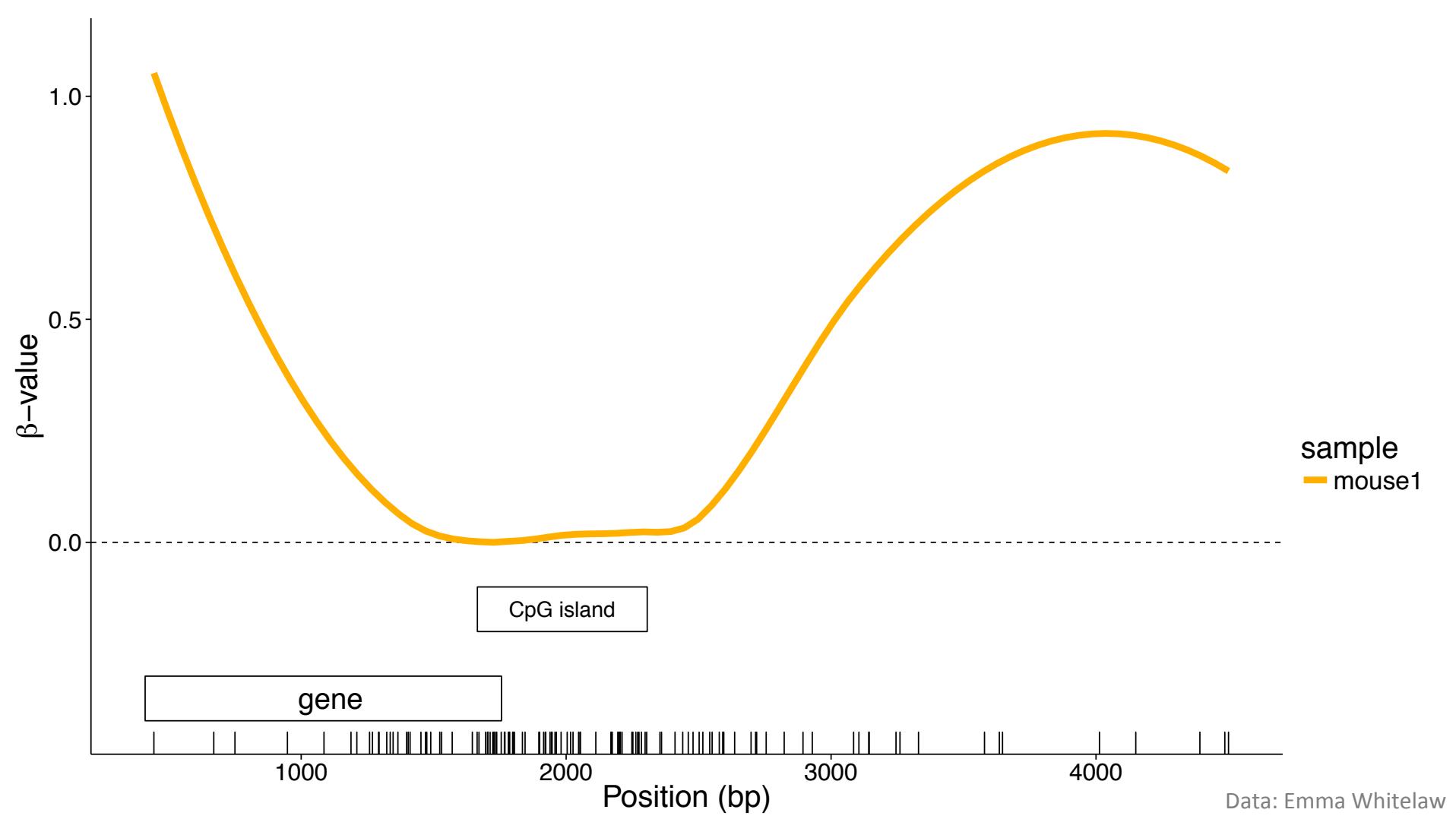


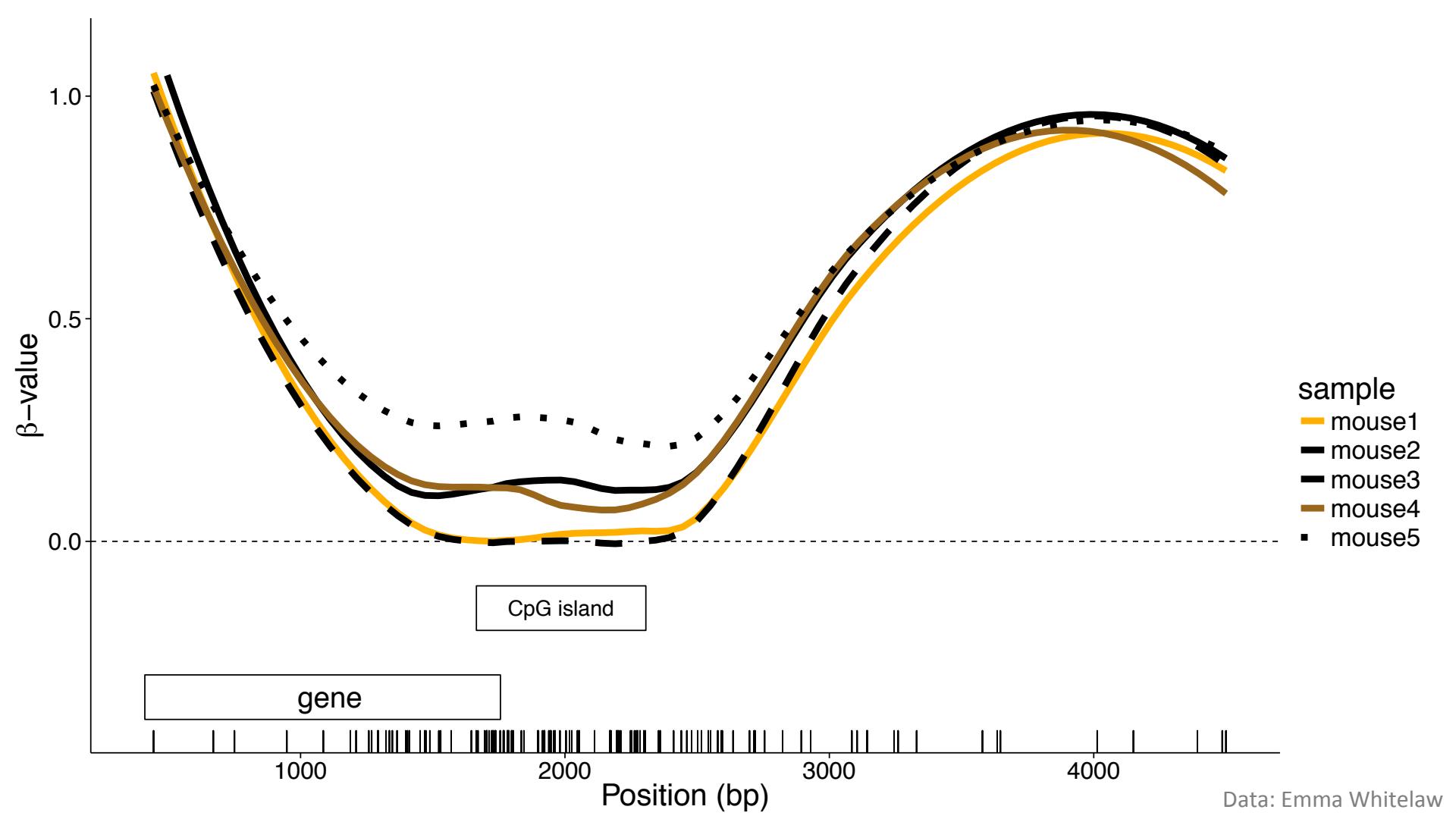


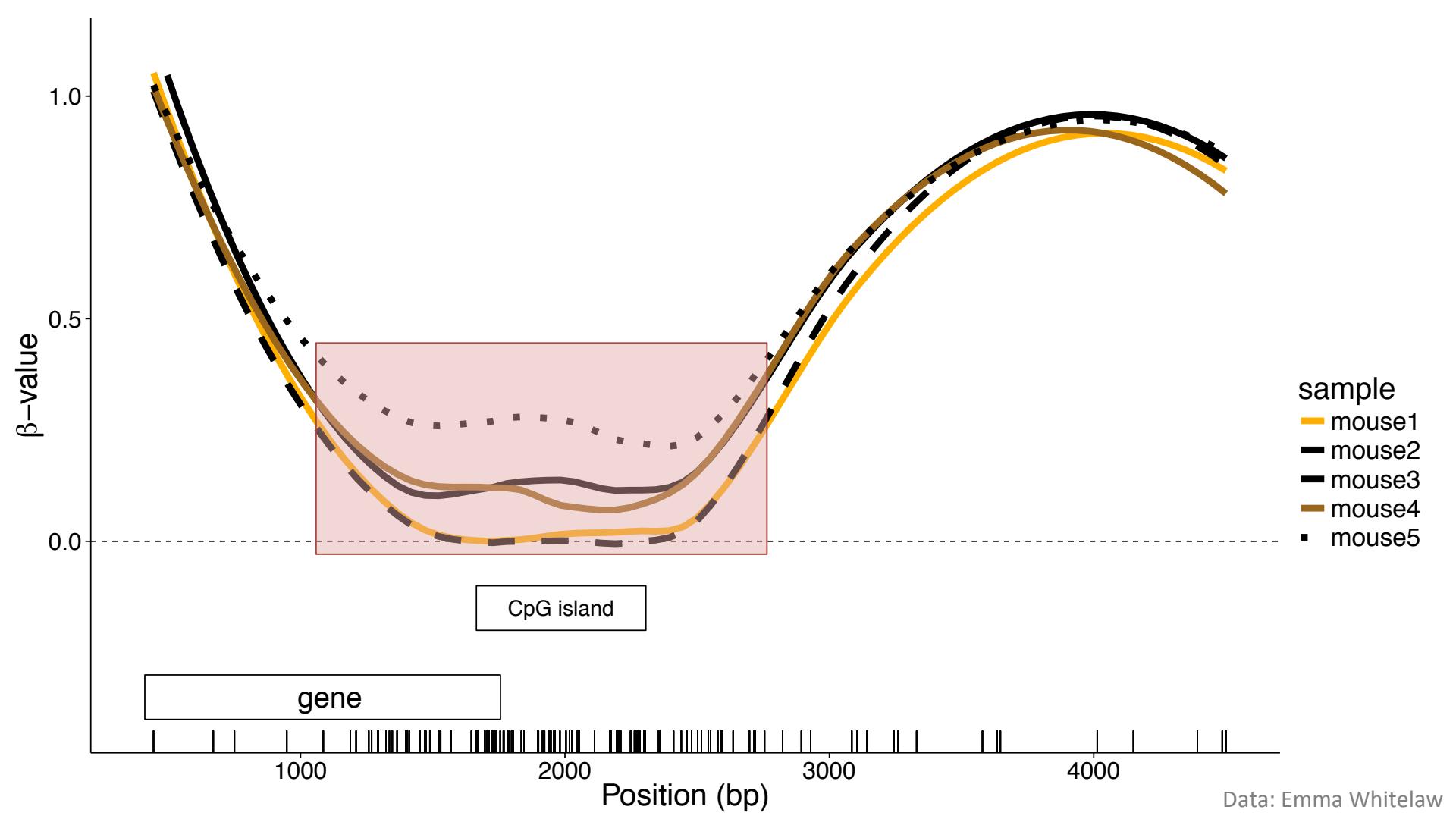












# **Co-methylation**

# Co-methylation = co-occurrence

*“The presence of methylation over a stretch of neighboring CpG positions”*

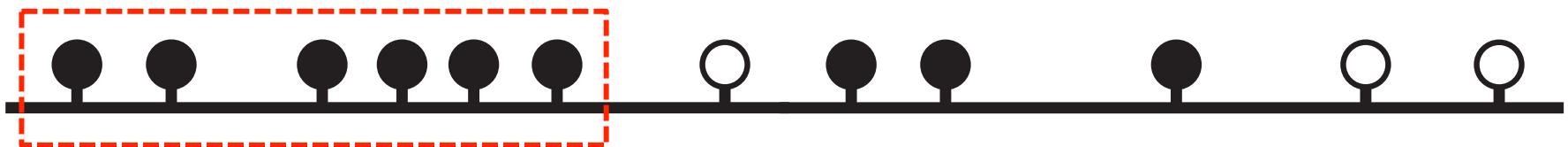
# Co-methylation = co-occurrence

*“The presence of methylation over a stretch of neighboring CpG positions”*



# Co-methylation = co-occurrence

*“The presence of methylation over a stretch of neighboring CpG positions”*



# Co-methylation = correlation

*“The relationship between the degree of methylation over distance”*

# Co-methylation = correlation

“*The relationship between the degree of methylation over distance*”

# Co-methylation = correlation

“*The relationship between the degree of methylation over distance*”

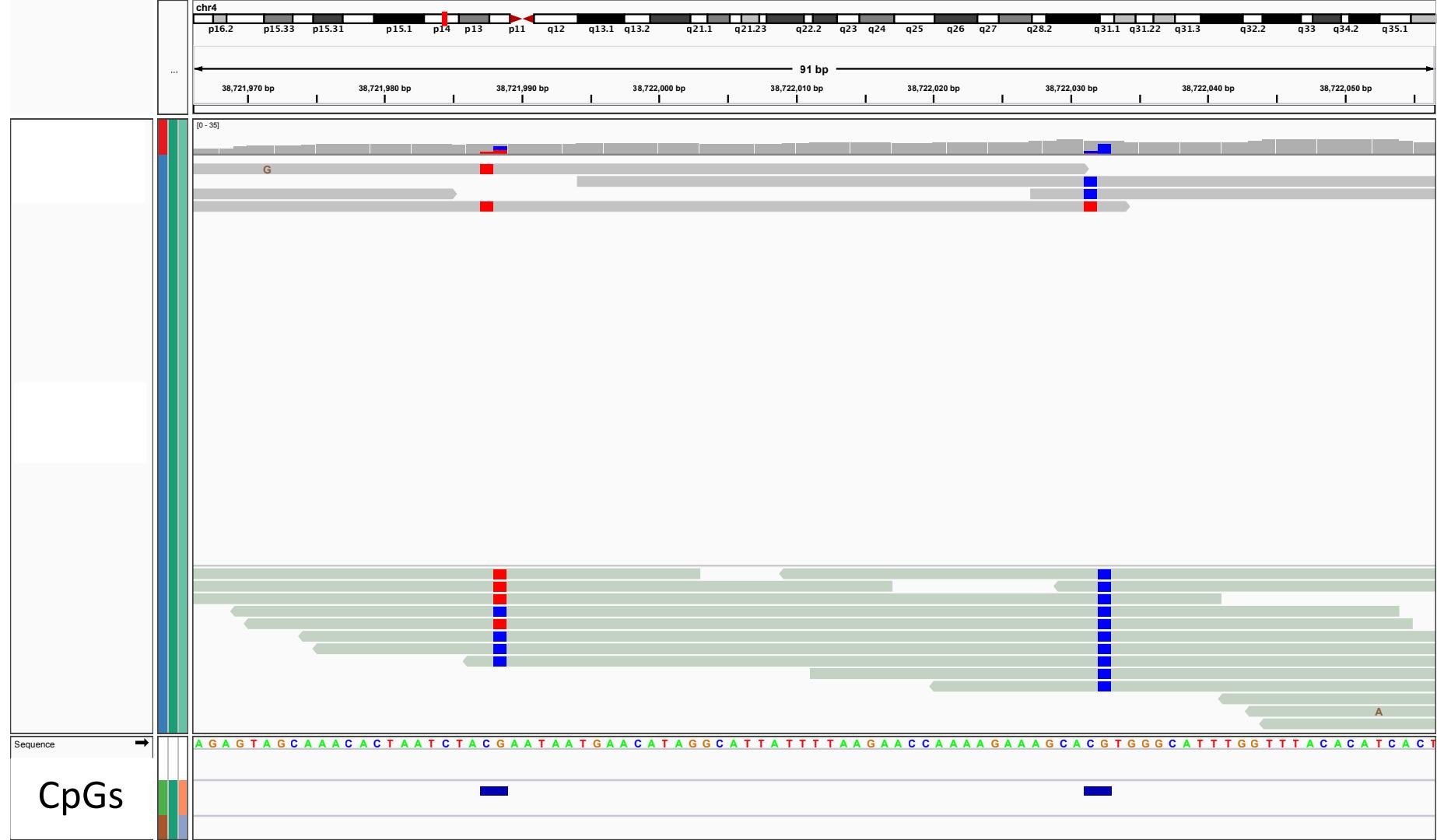
1. Within-fragment co-methylation
2. Correlation of  $\beta$ -values

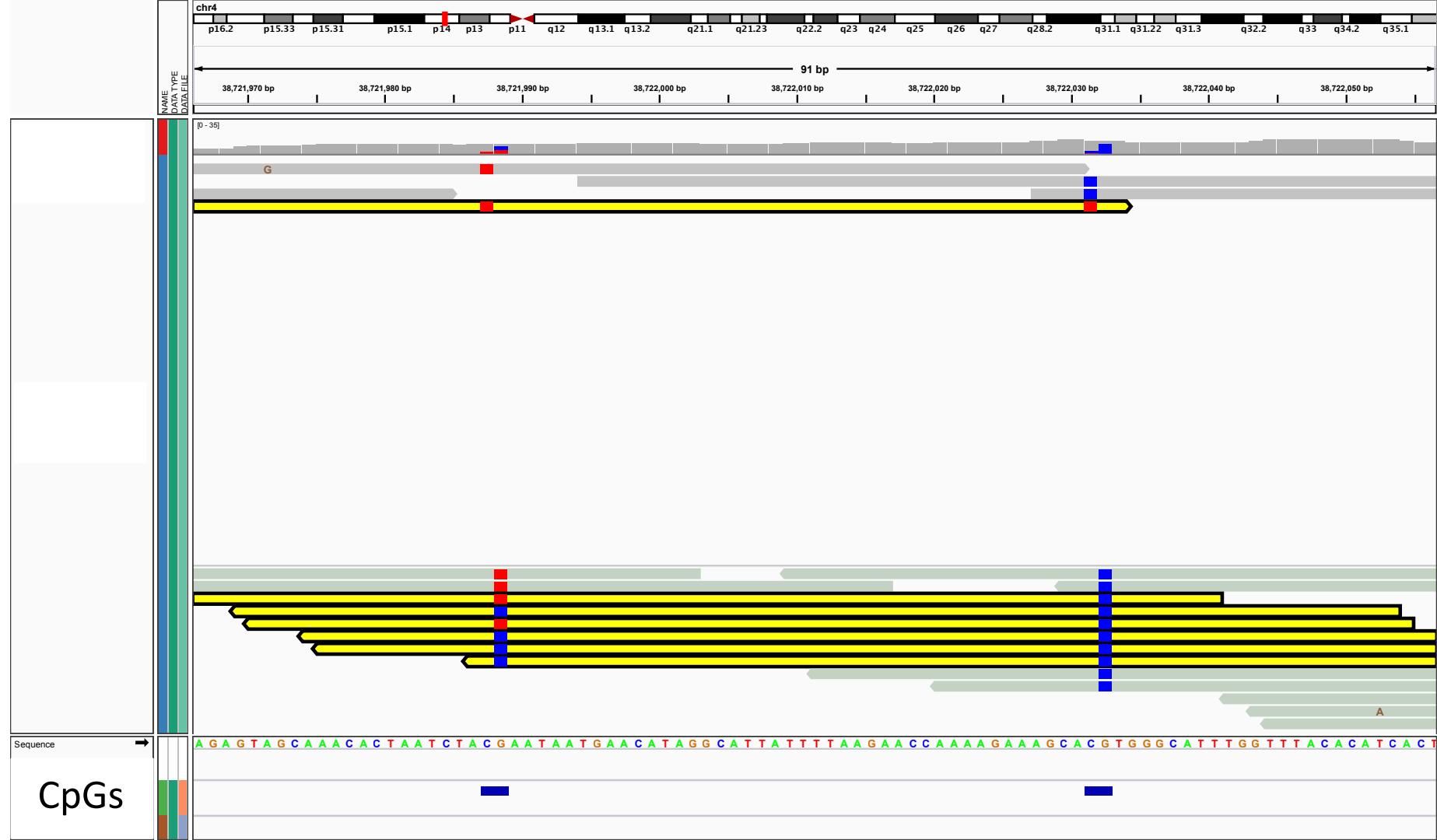
# Co-methylation = correlation

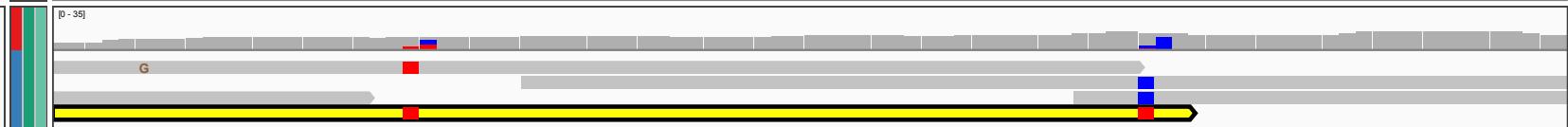
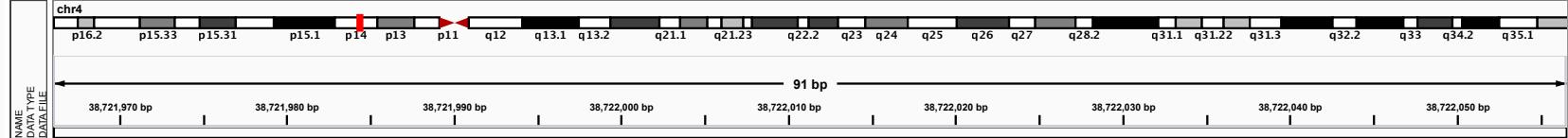
*“The relationship between the degree of methylation over distance”*

## 1. Within-fragment co-methylation

2. Correlation of  $\beta$ -values

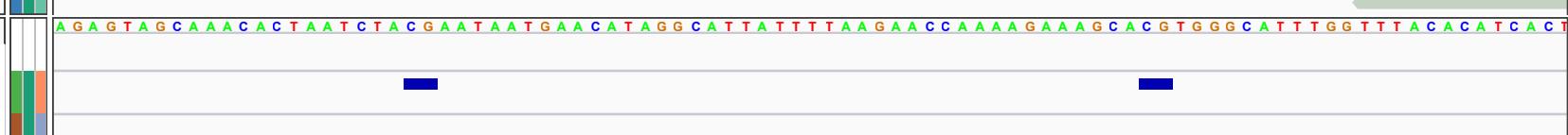
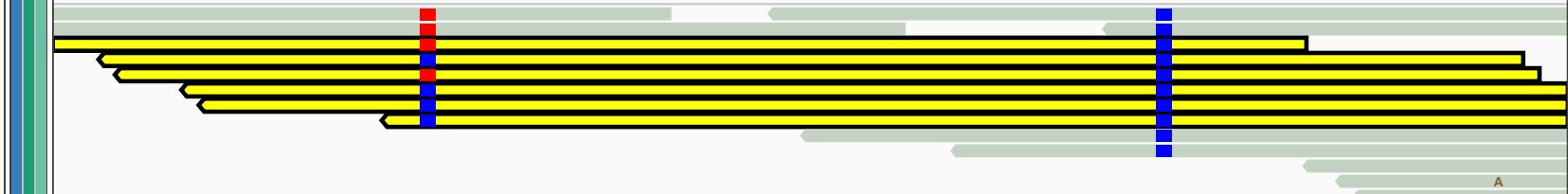


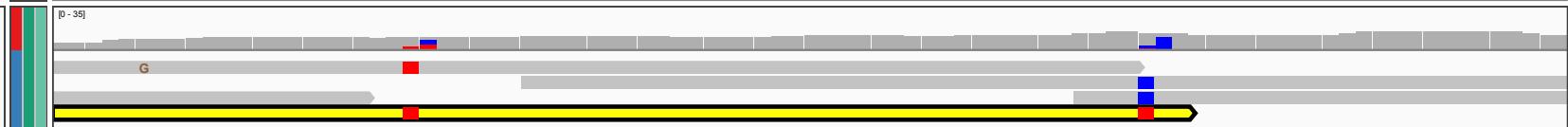
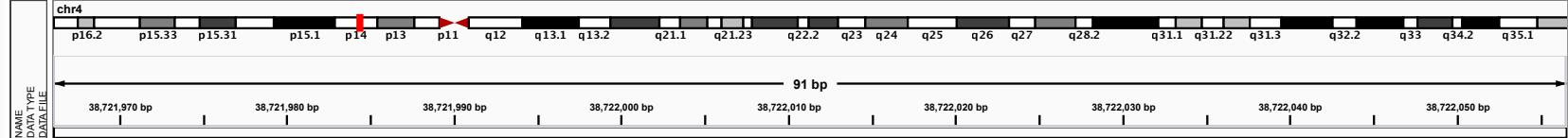




### *Second CpG*

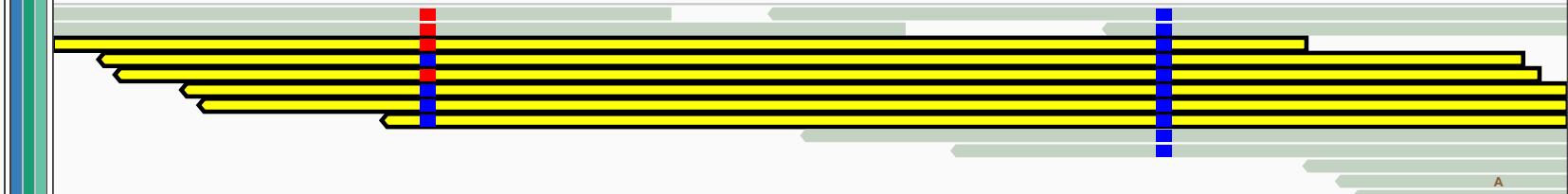
	Methylated	Unmethylated	Total
<i>First CpG</i>	Methylated	1	2
	Unmethylated	0	4
<b>Total</b>	1	6	7





	Methylated	Unmethylated	Total	
First CpG	Methylated	1	2	3
	Unmethylated	0	4	4
Total	1	6	7	

$$\text{log-odds ratio} = \log_2 \left( \frac{1.5 \times 4.5}{2.5 \times 0.5} \right) = 2.4$$



Sequence

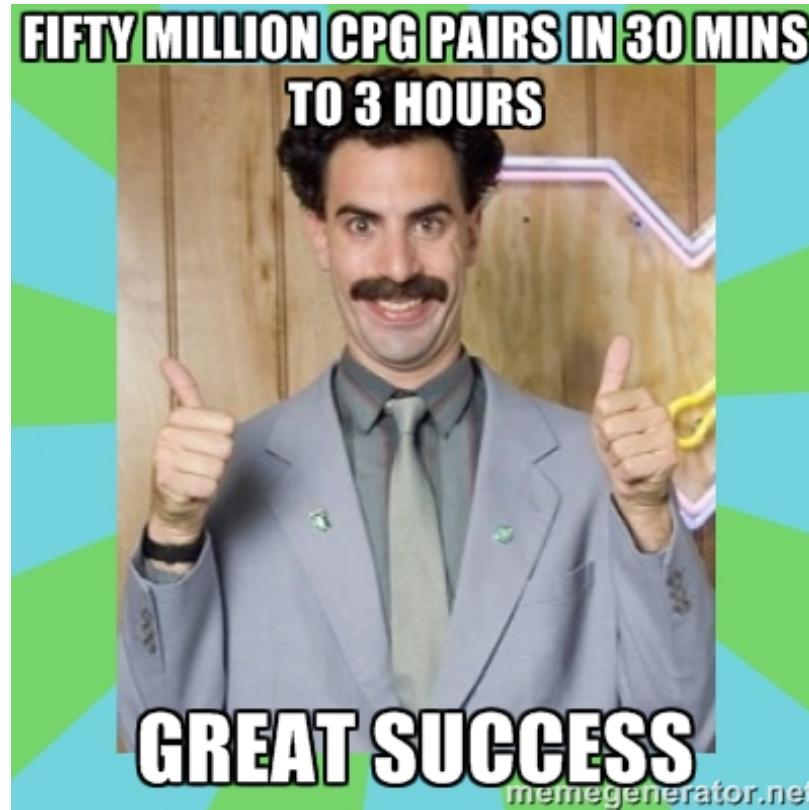
CpGs

Do this 50 million times per sample

# Do this 50 million times per sample

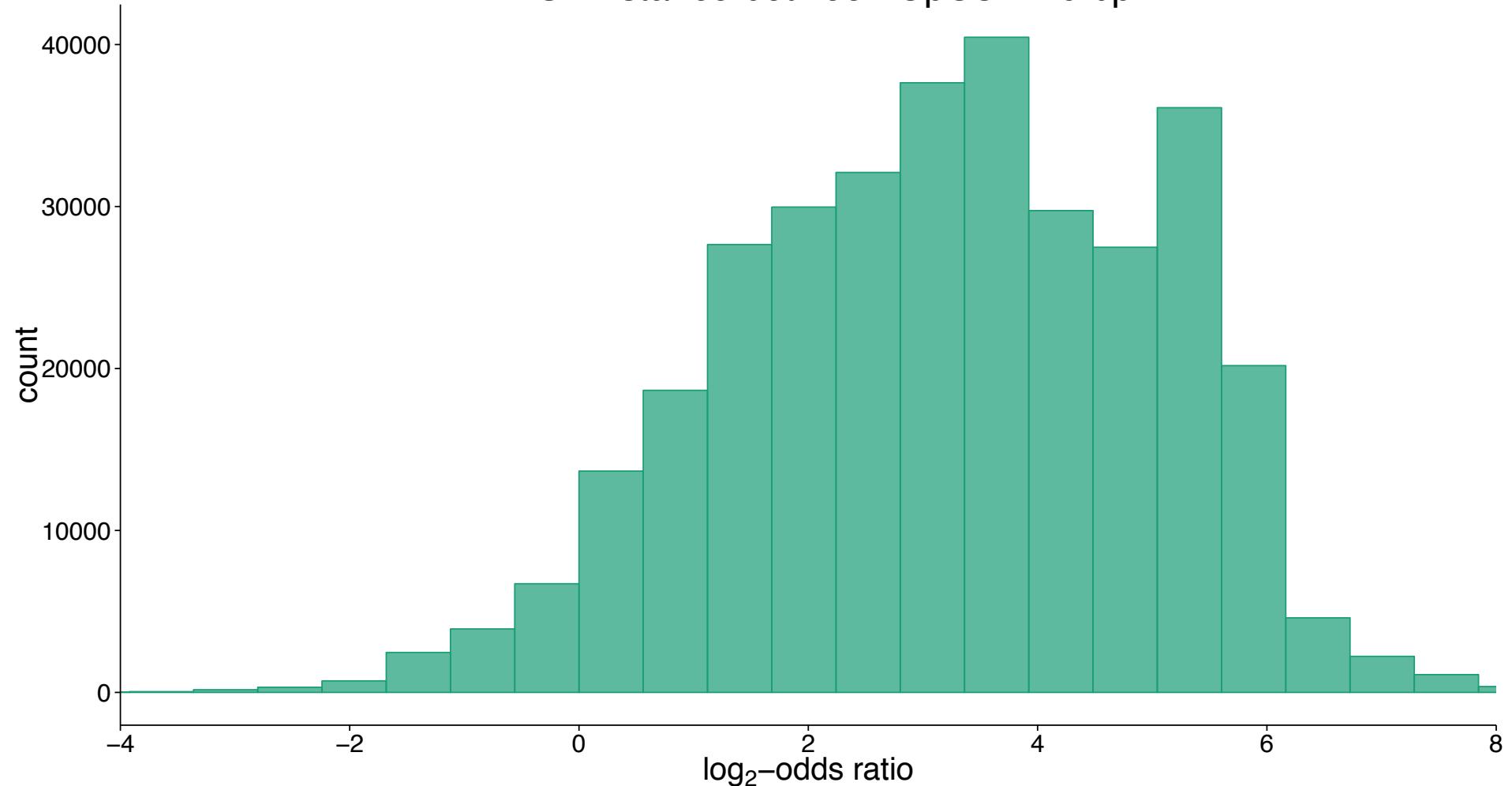


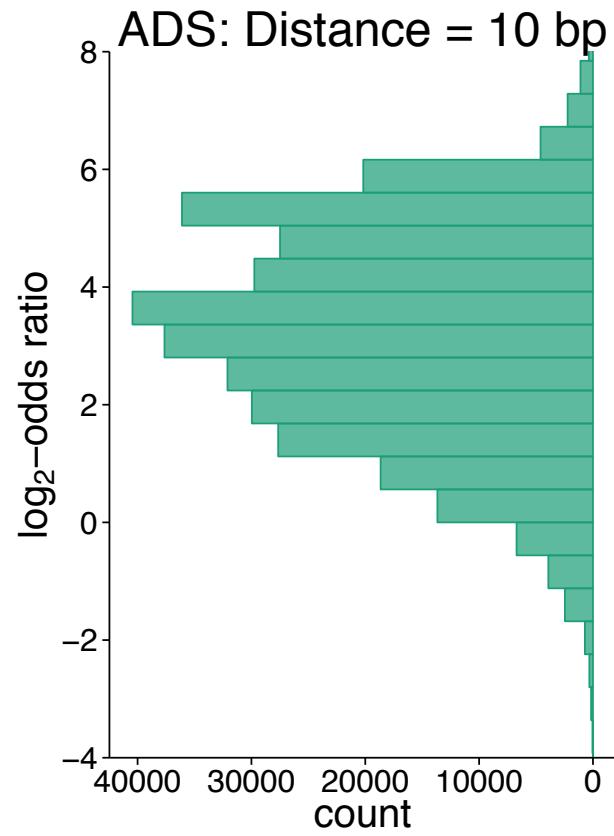
# methtuple

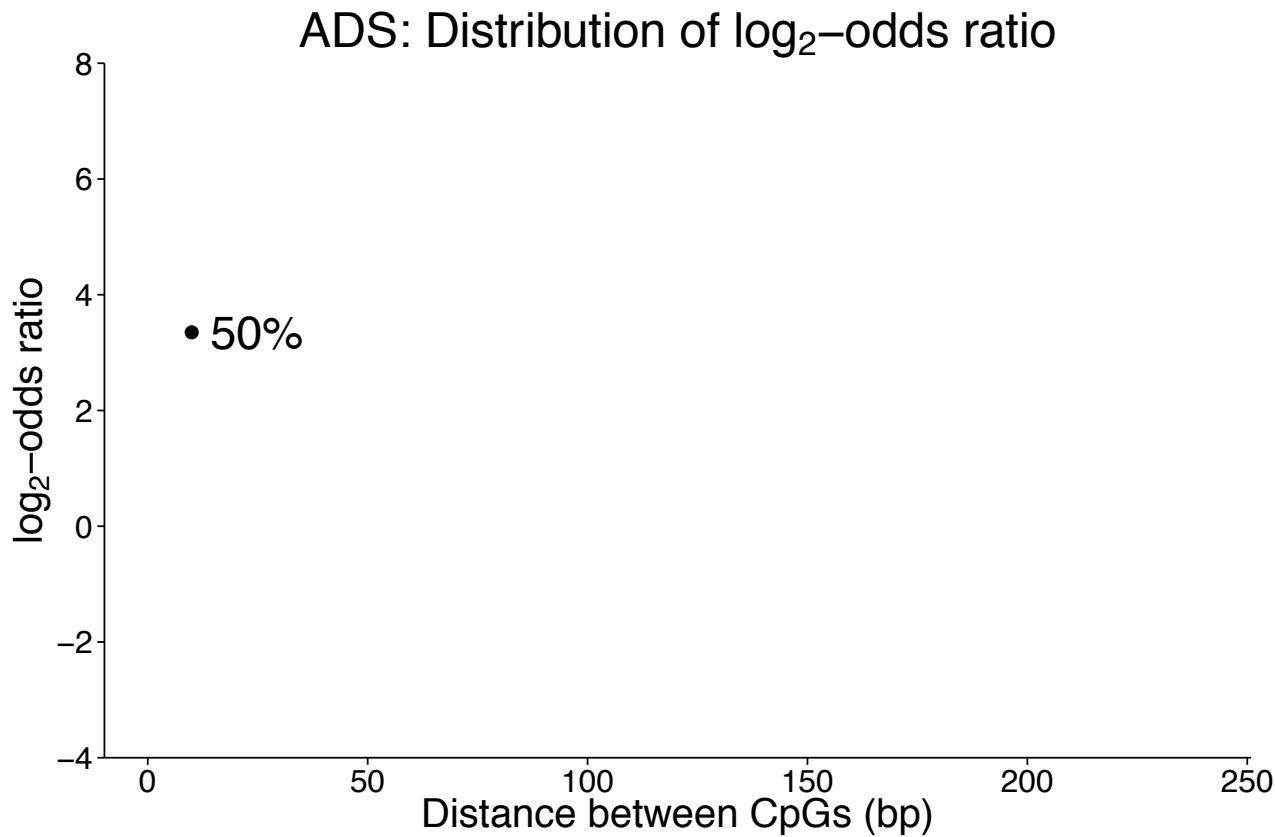
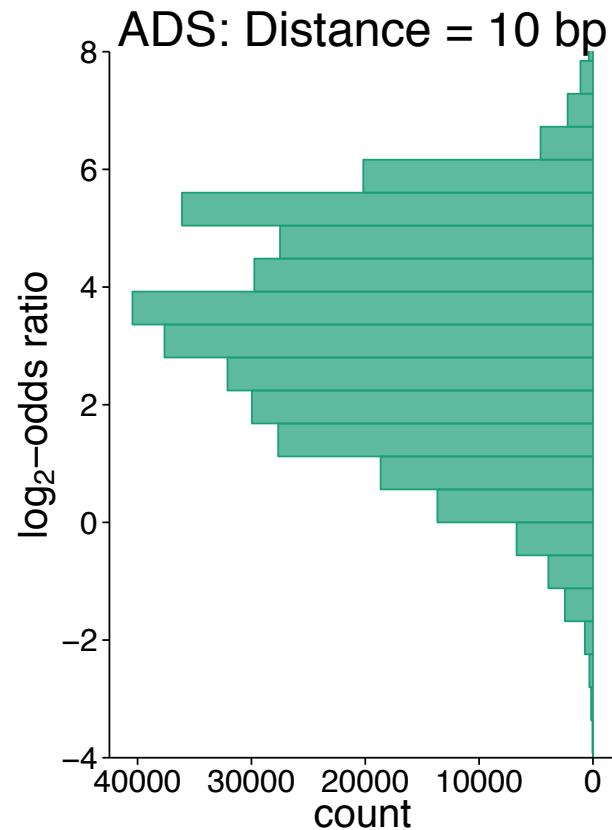


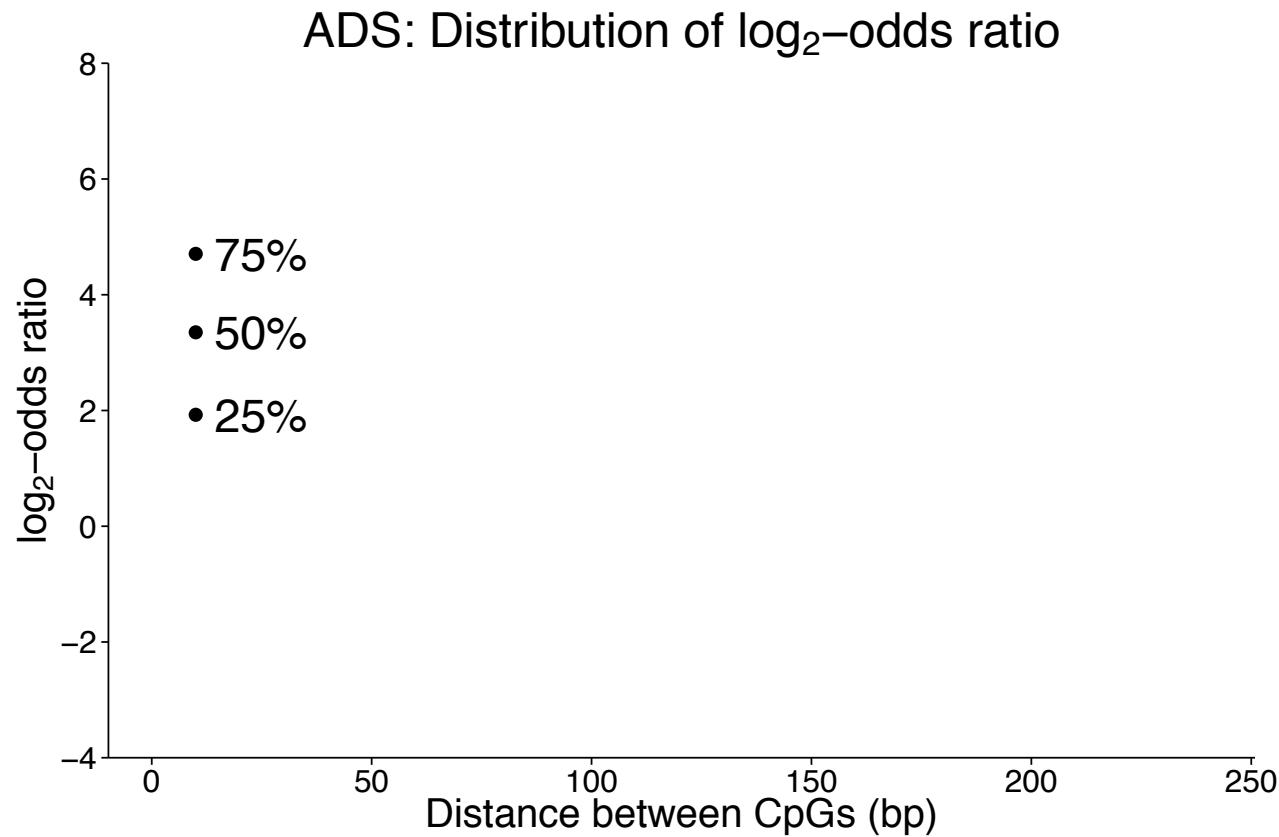
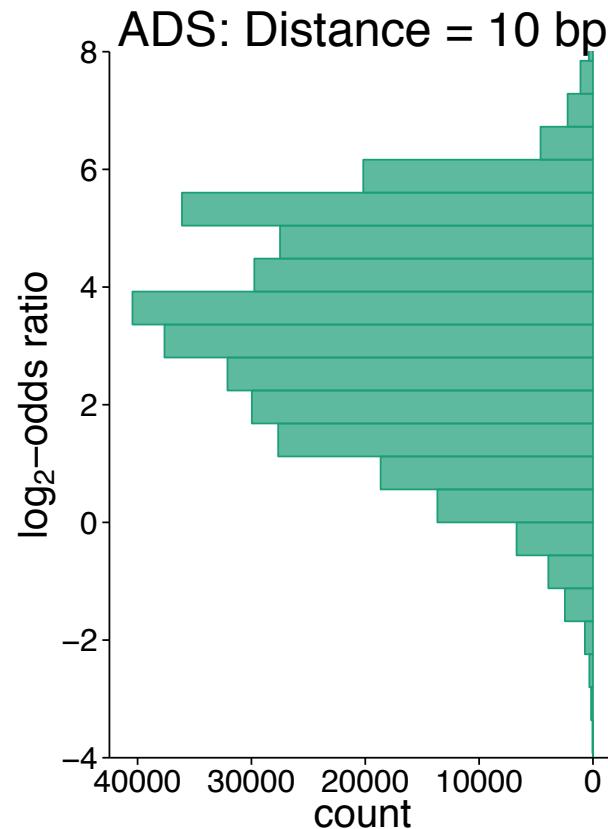
[www.github.com/PeteHaitch/methtuple](https://www.github.com/PeteHaitch/methtuple)

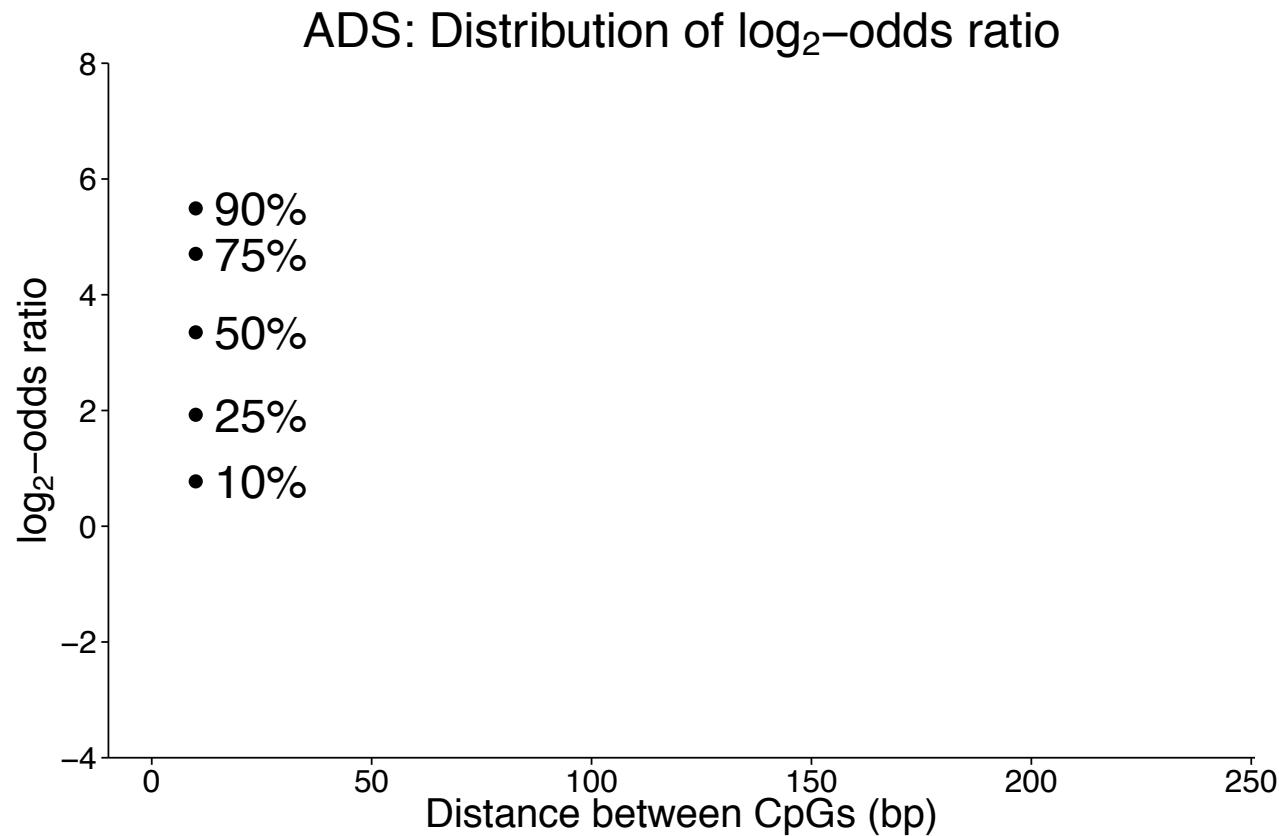
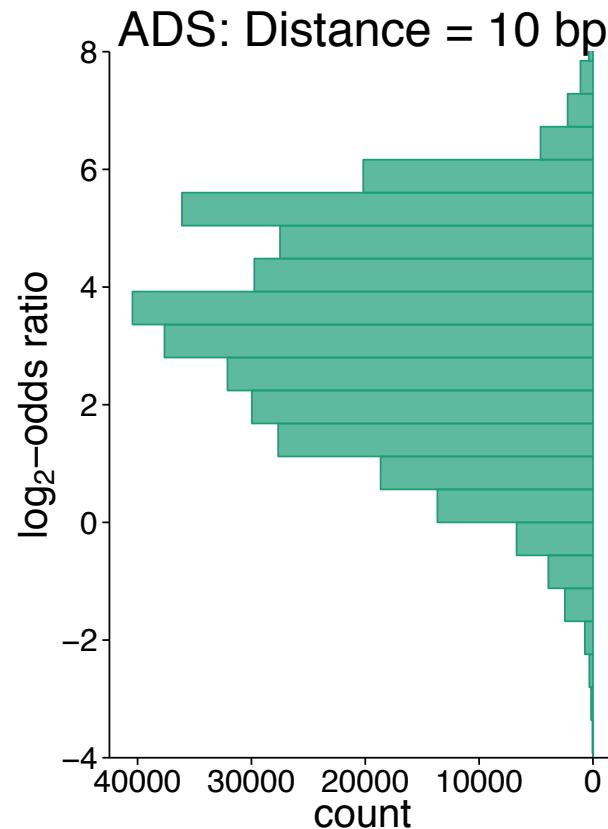
# ADS: Distance between CpGs = 10 bp



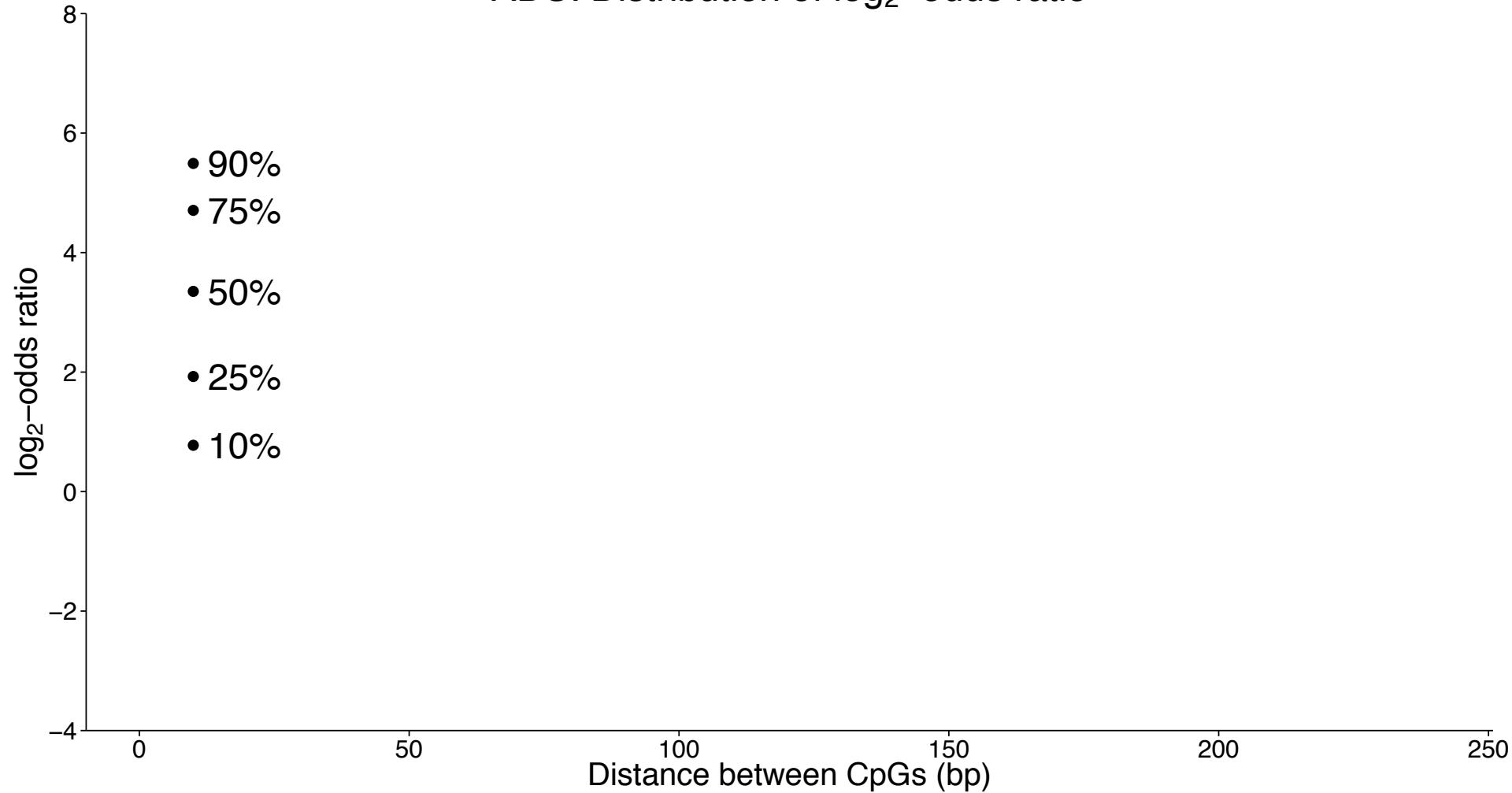








# ADS: Distribution of $\log_2$ -odds ratio



# ADS: Distribution of $\log_2$ -odds ratio

$\log_2$ -odds ratio

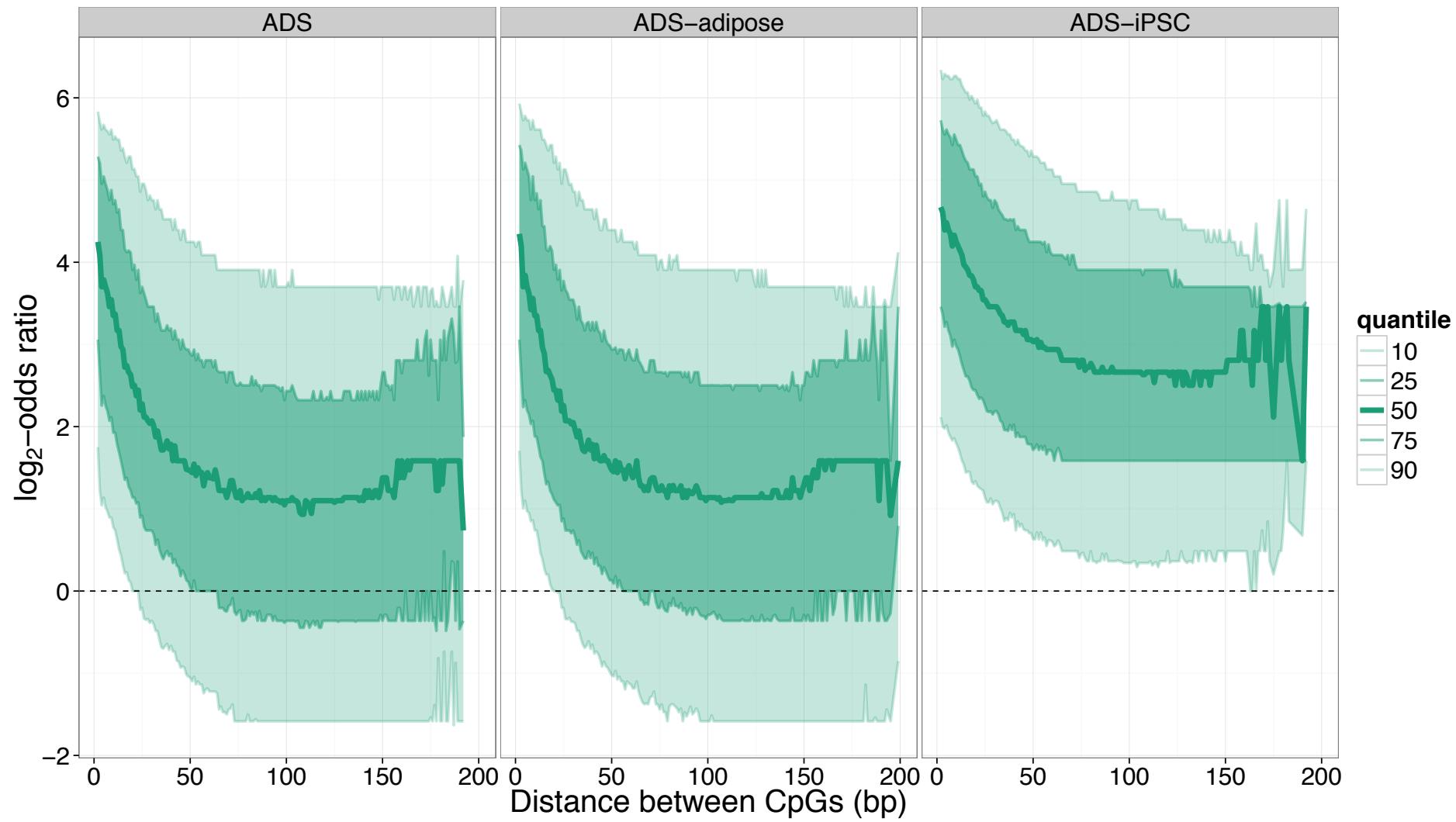
8  
6  
4  
2  
0  
-2  
-4

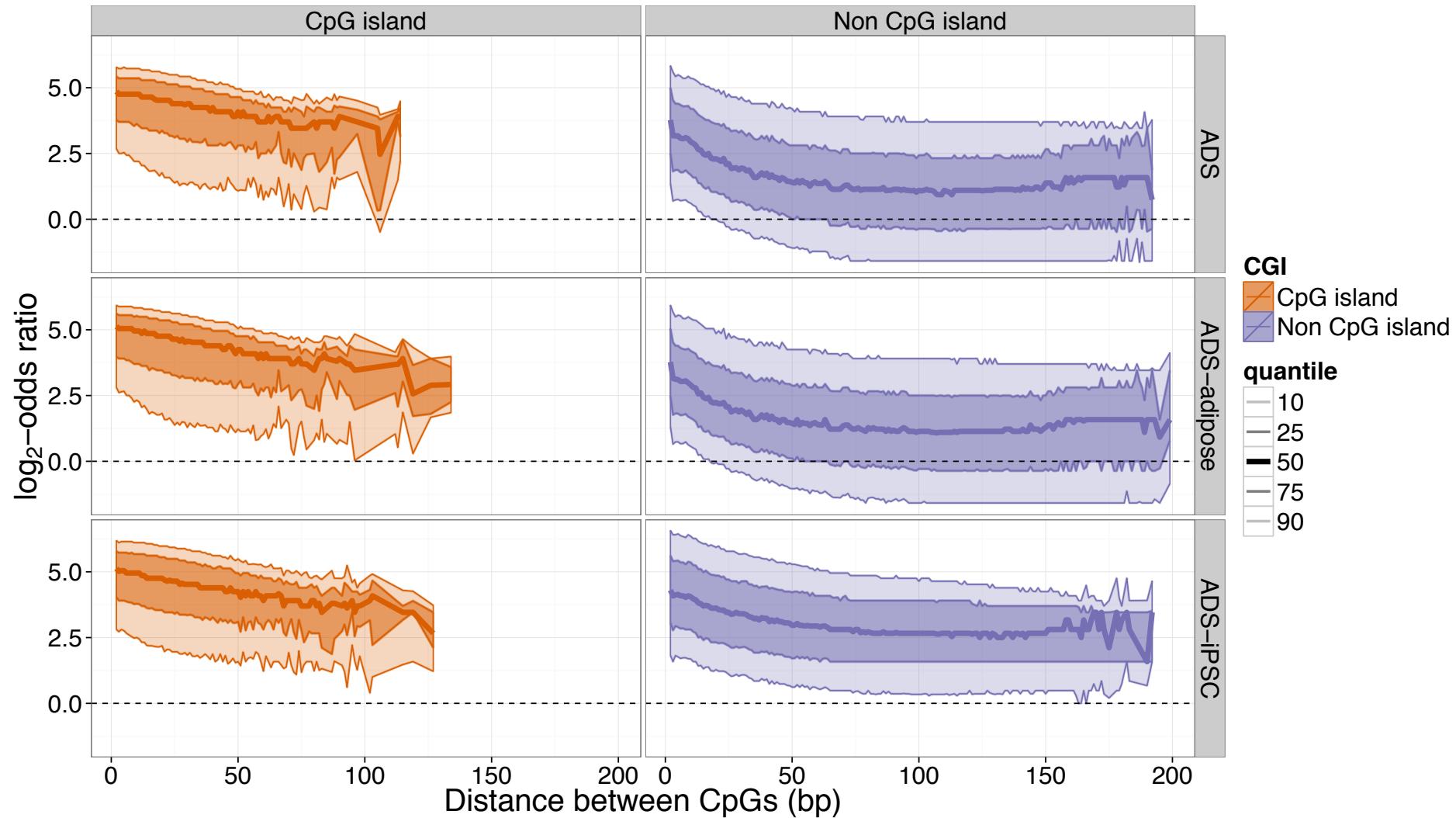
- 90%
- 75%
- 50%
- 25%
- 10%

quantile  
10  
25  
50  
75  
90

0 50 100 150 200

Distance between CpGs (bp)



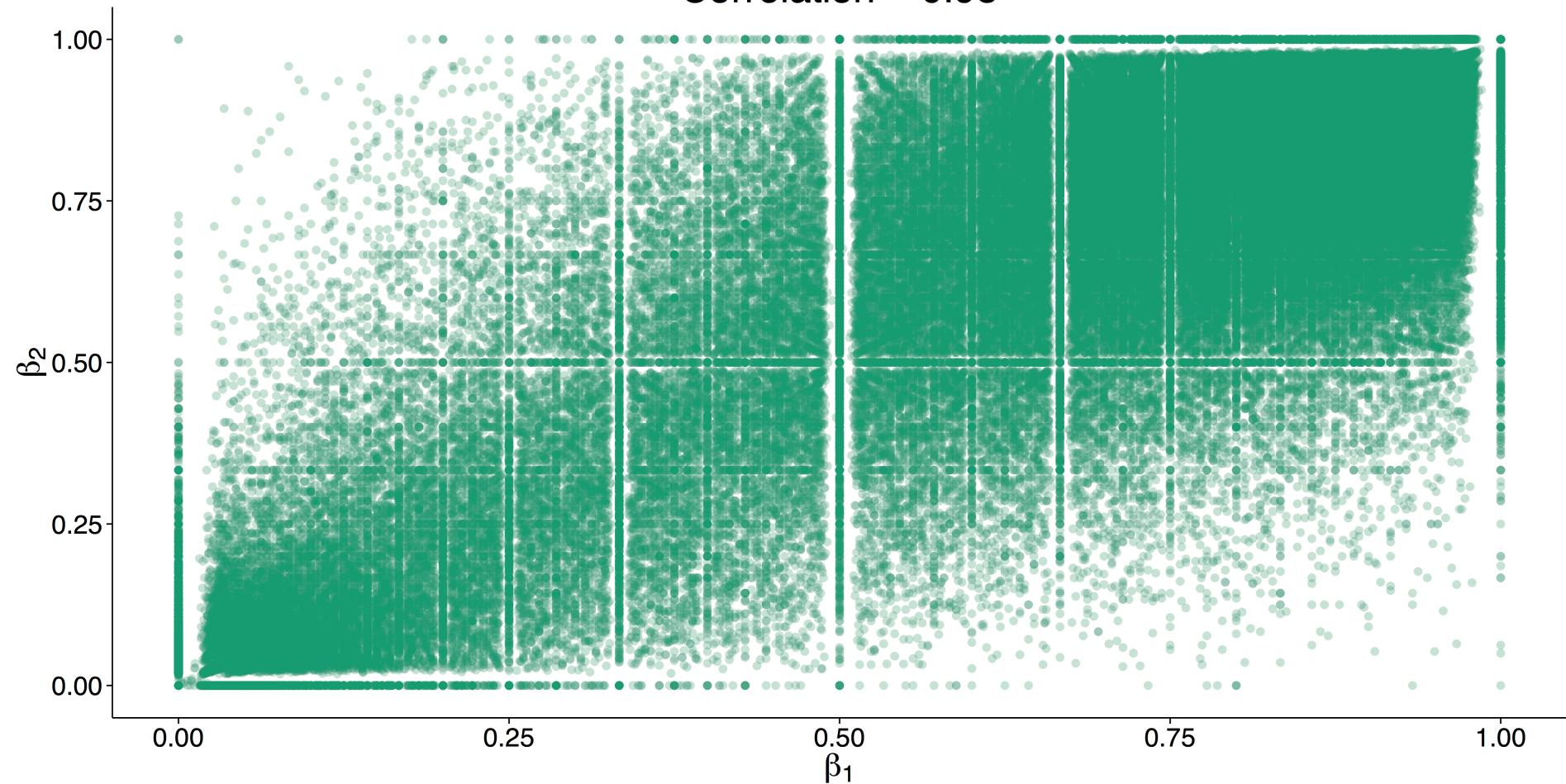


# Co-methylation = correlation

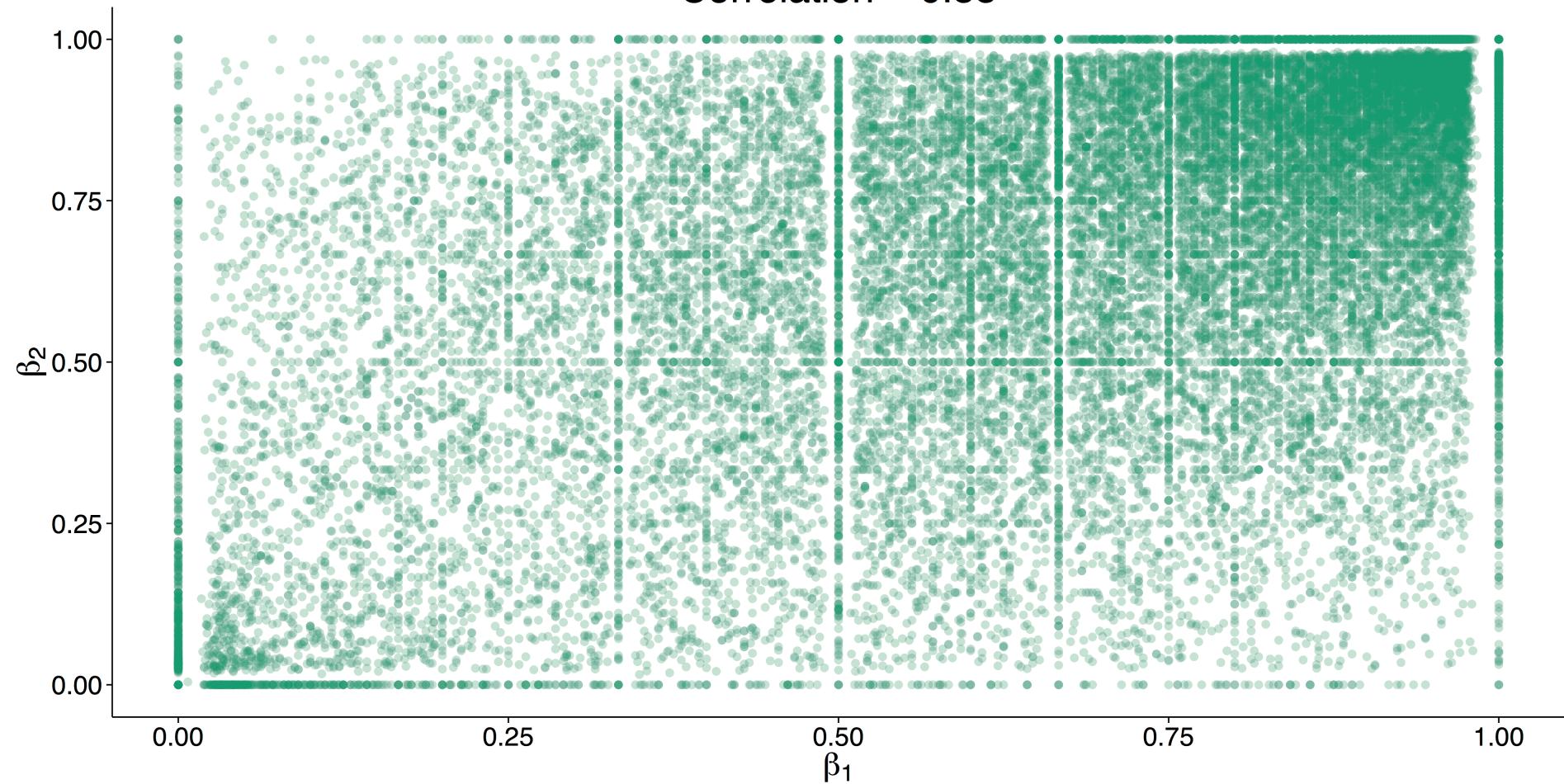
*“The relationship between the degree of methylation over distance”*

1. Within-fragment co-methylation
2. Correlation of  $\beta$ -values

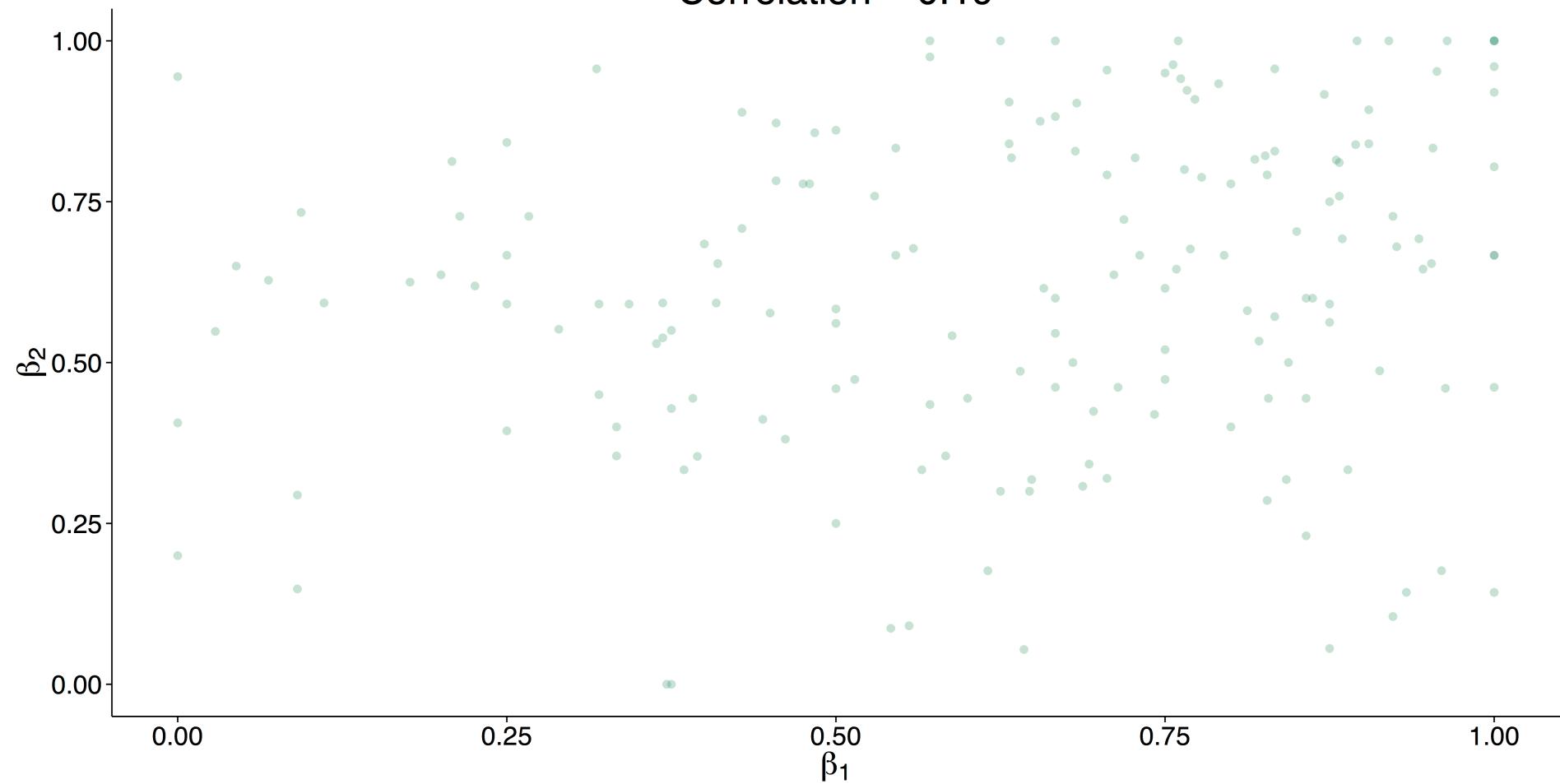
ADS: Distance = 10 bp  
Correlation = 0.95



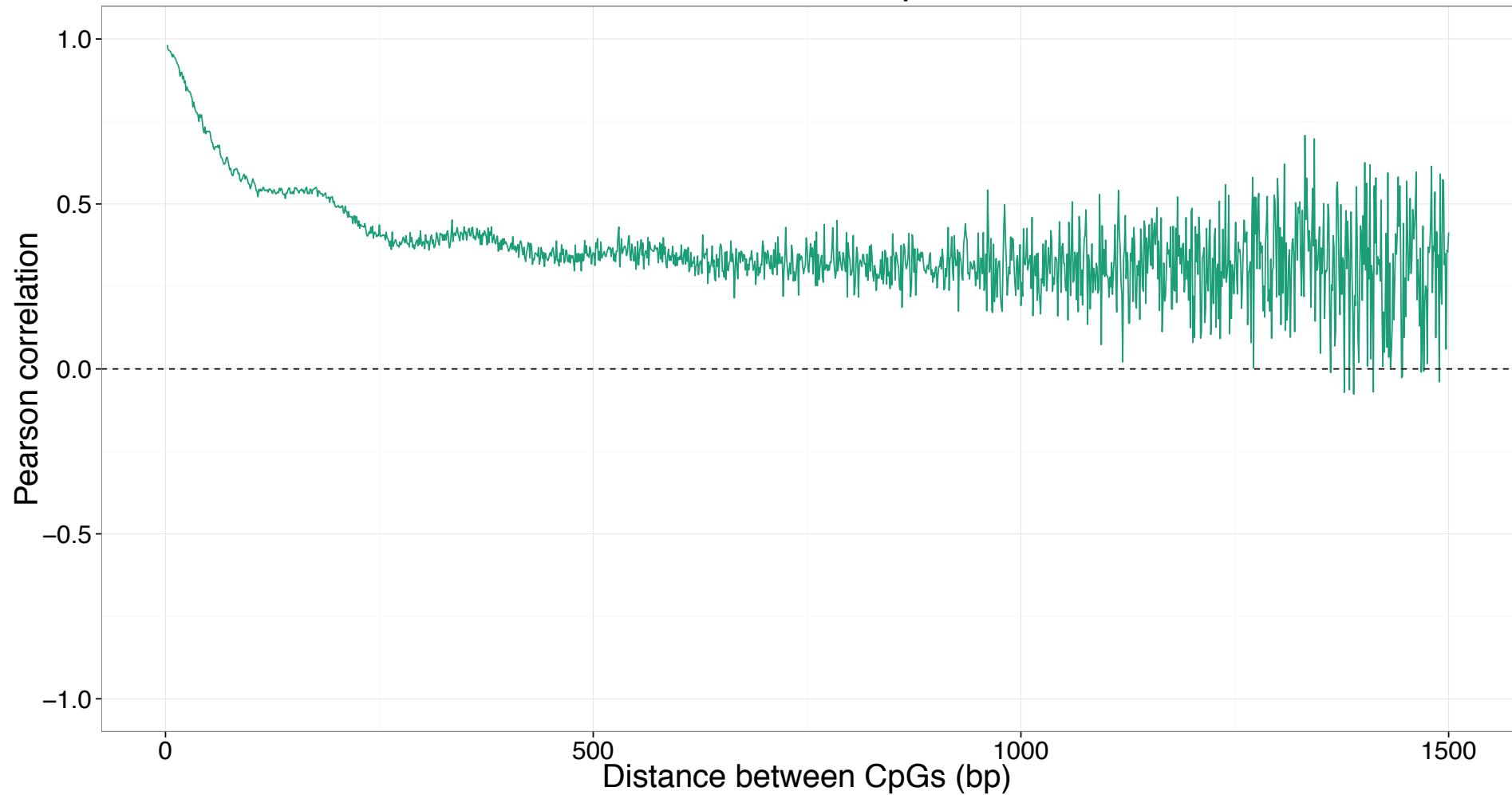
ADS: Distance = 100 bp  
Correlation = 0.55

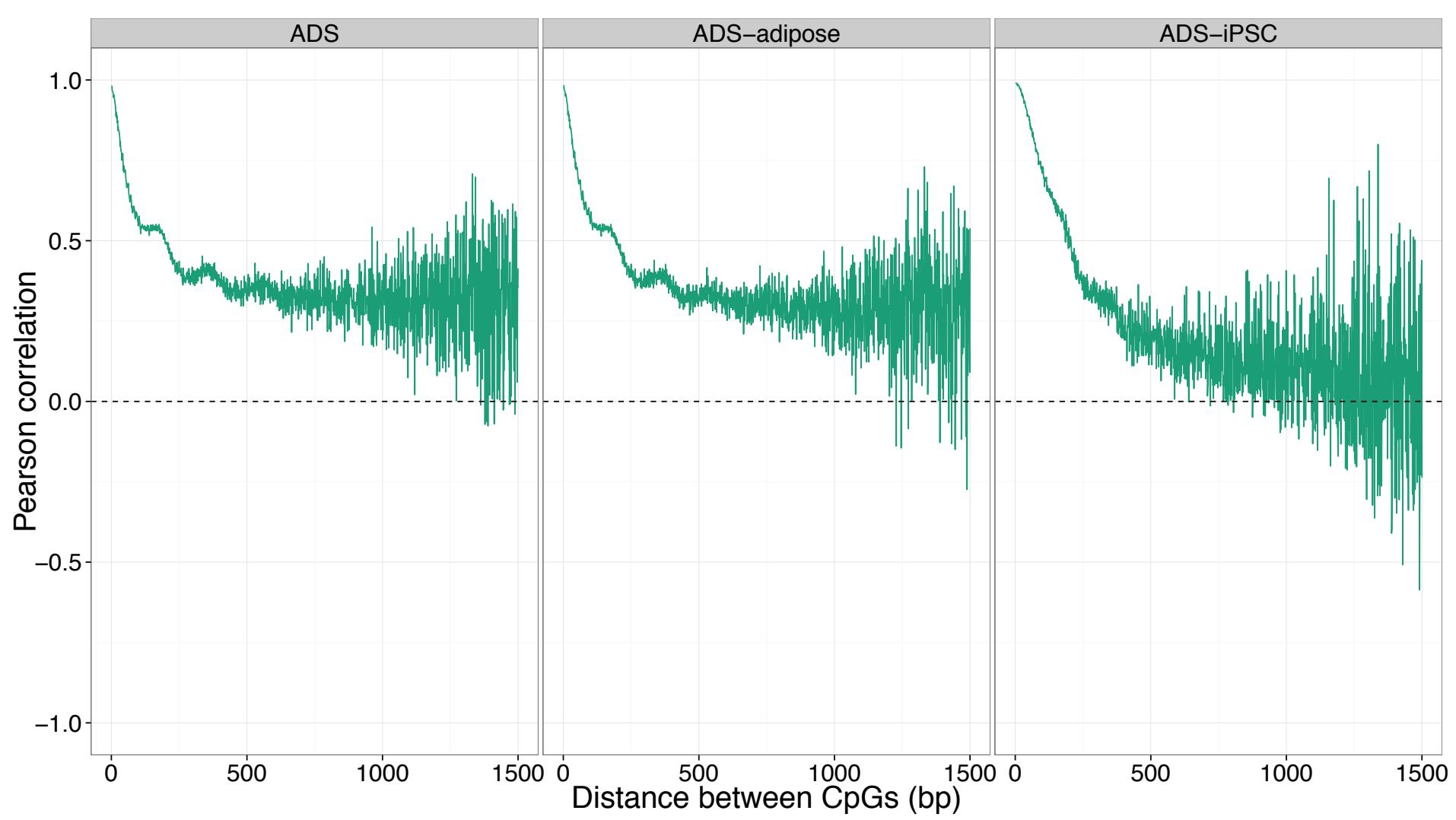


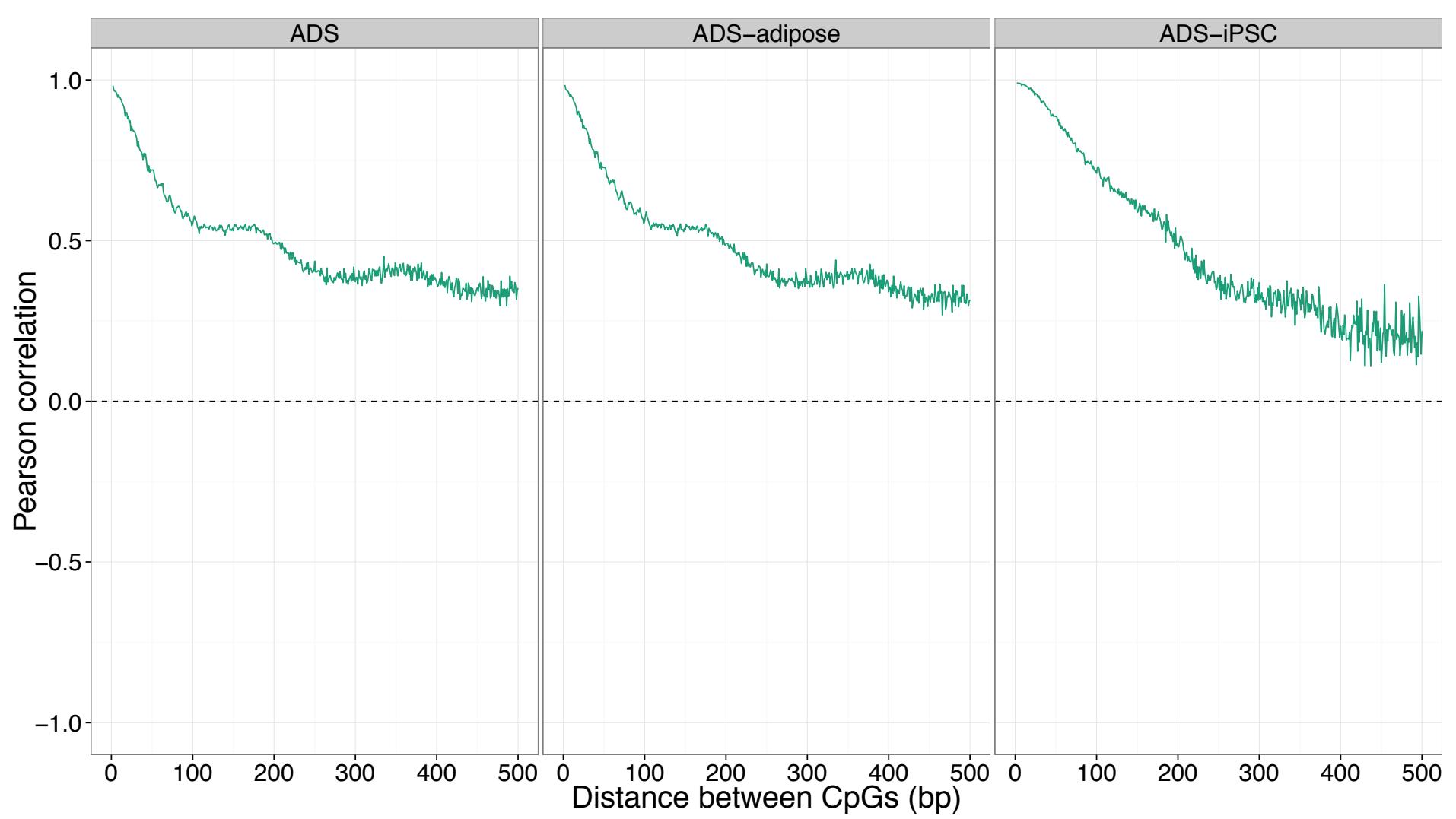
ADS: Distance = 1000 bp  
Correlation = 0.19

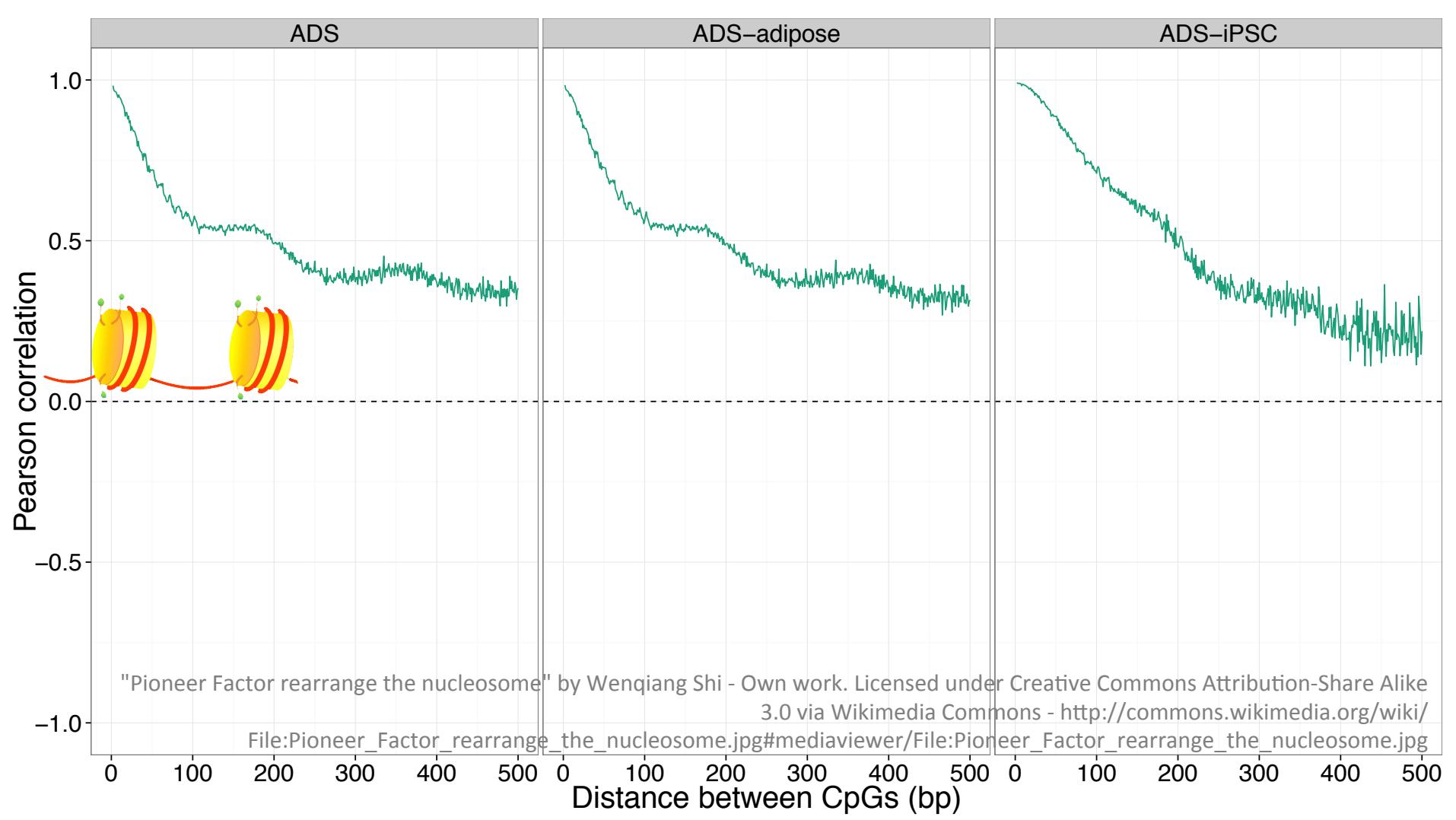


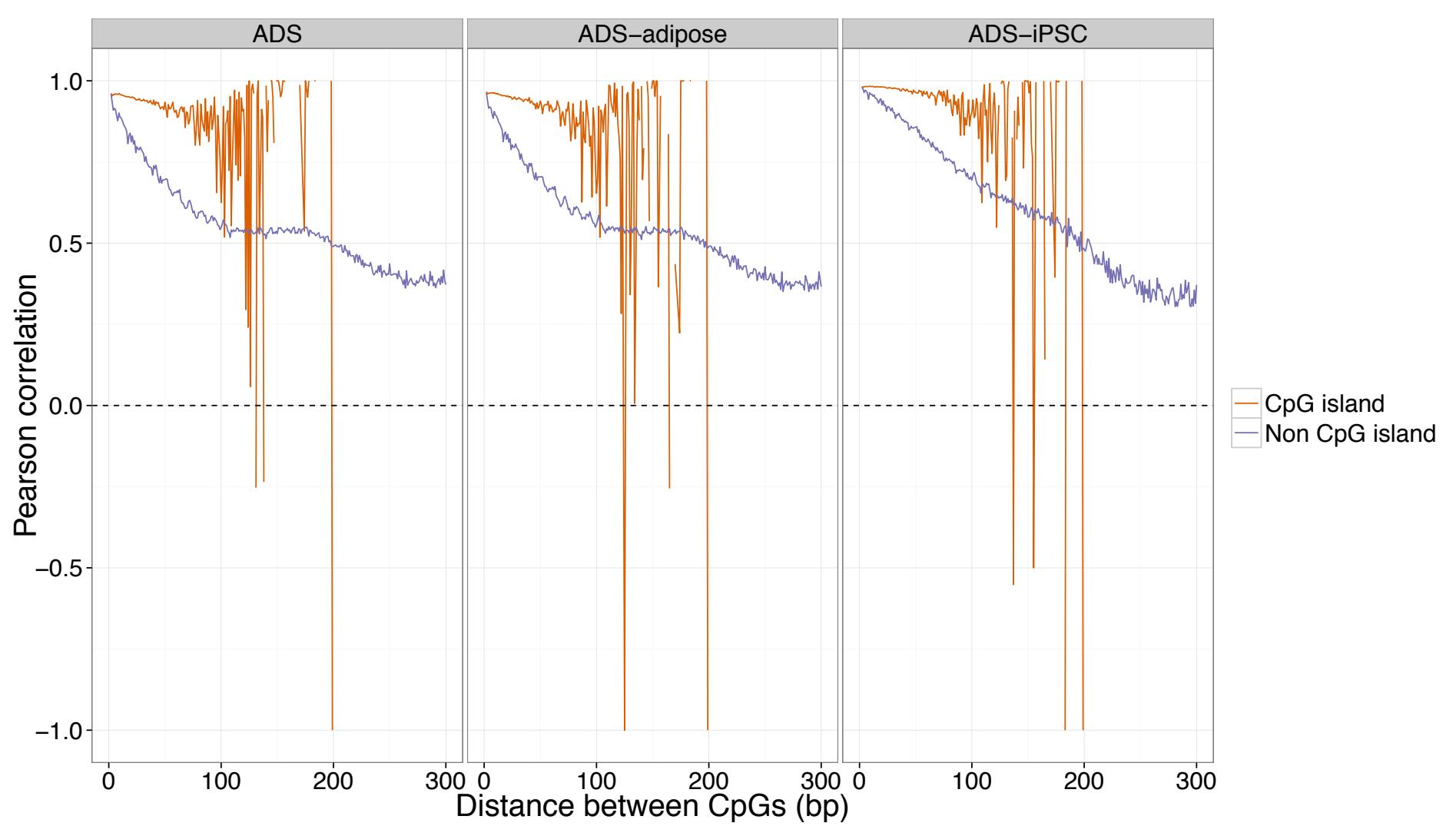
# ADS: Correlation of $\beta$ -values





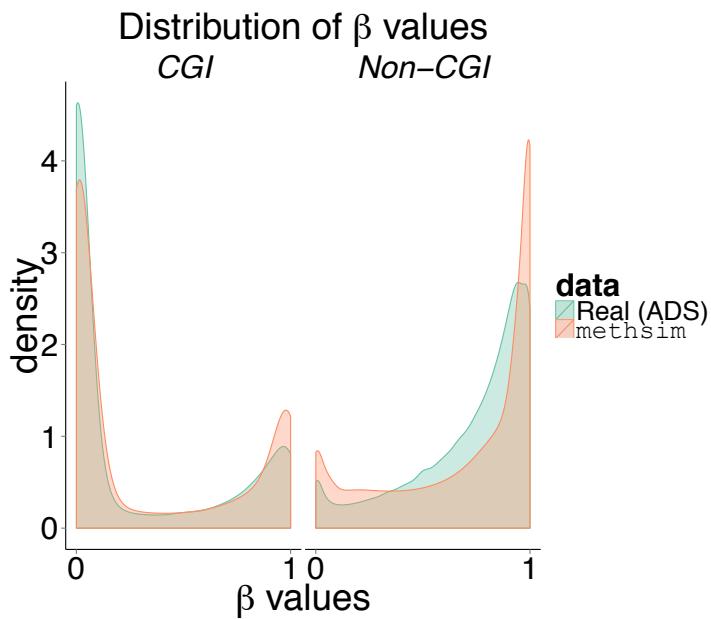




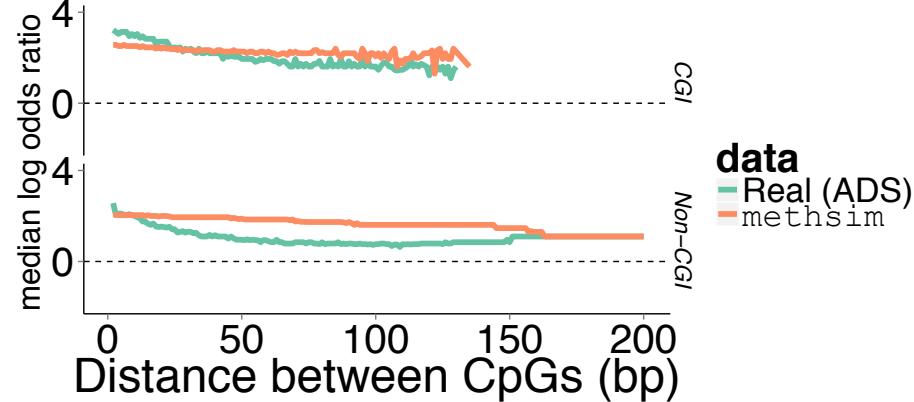


# methsim

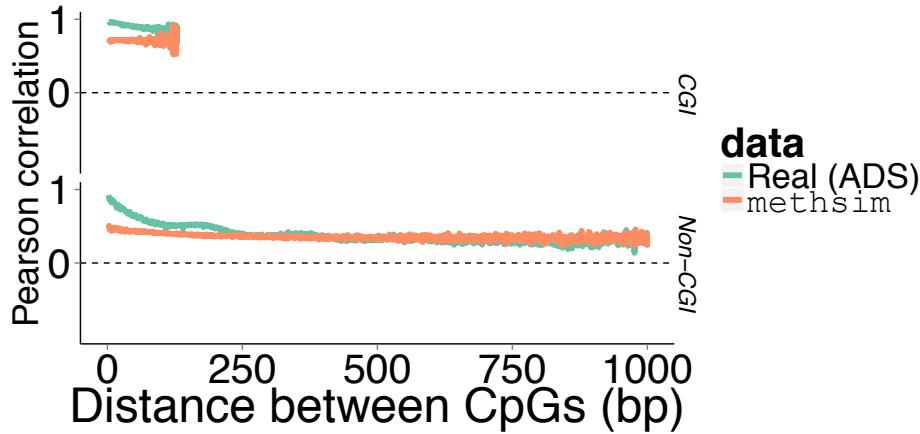
[www.github.com/PeteHaitch/methsim](https://www.github.com/PeteHaitch/methsim)



## Within haplotype co-methylation at neighbouring CpGs

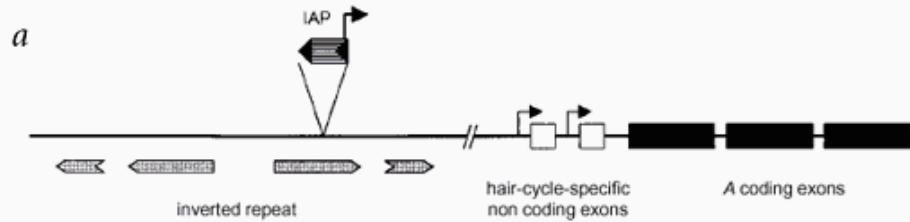


## Correlations of pairs of $\beta$ values

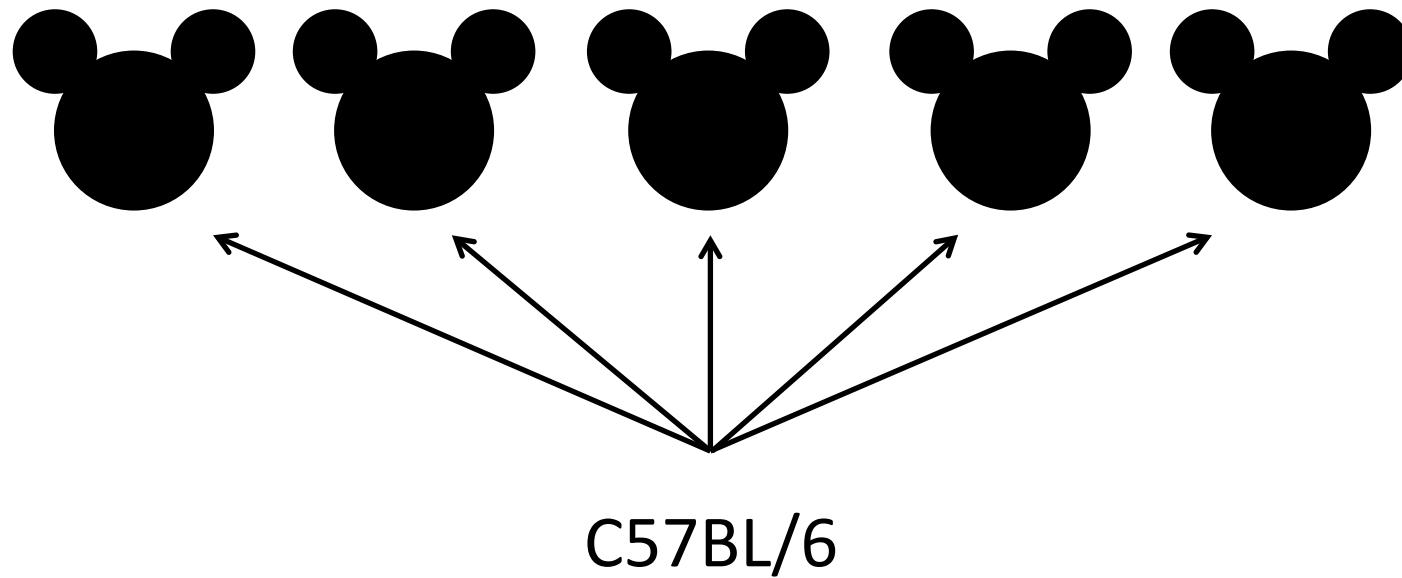


Some of these mice are not like the others (we hope...)

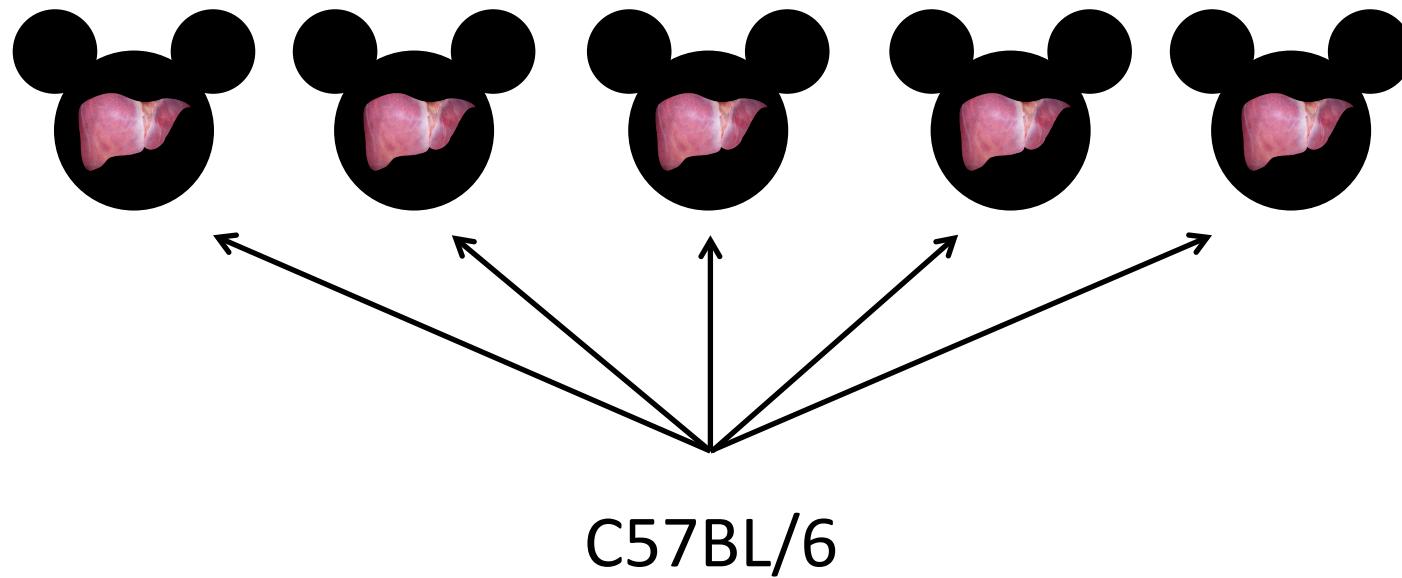
# **Methylome of the agouti viable yellow mouse (A<sup>vy</sup>)**



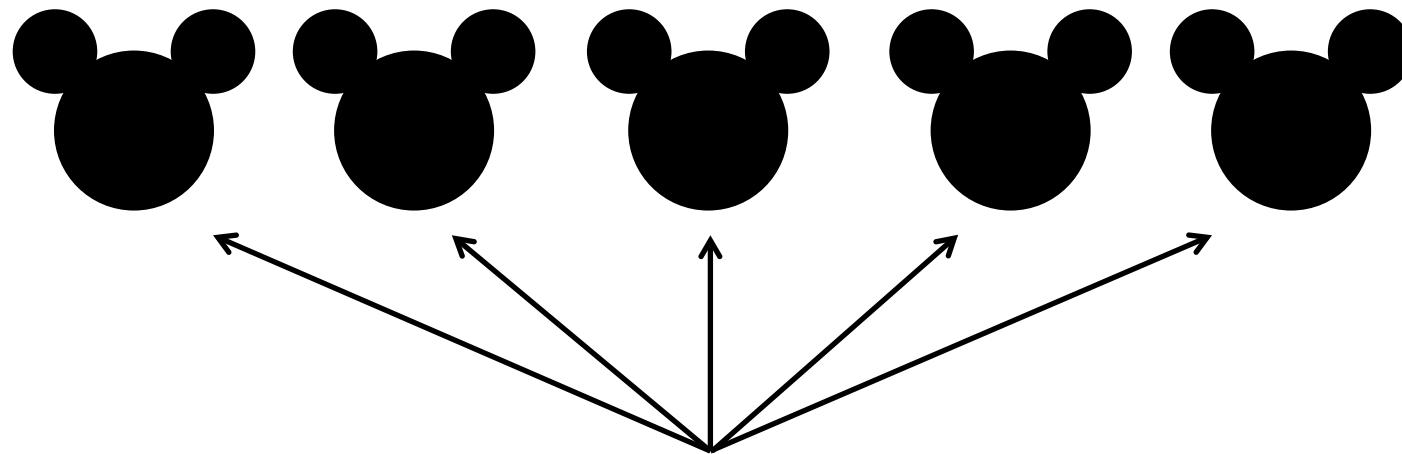
# Experimental design



# Experimental design

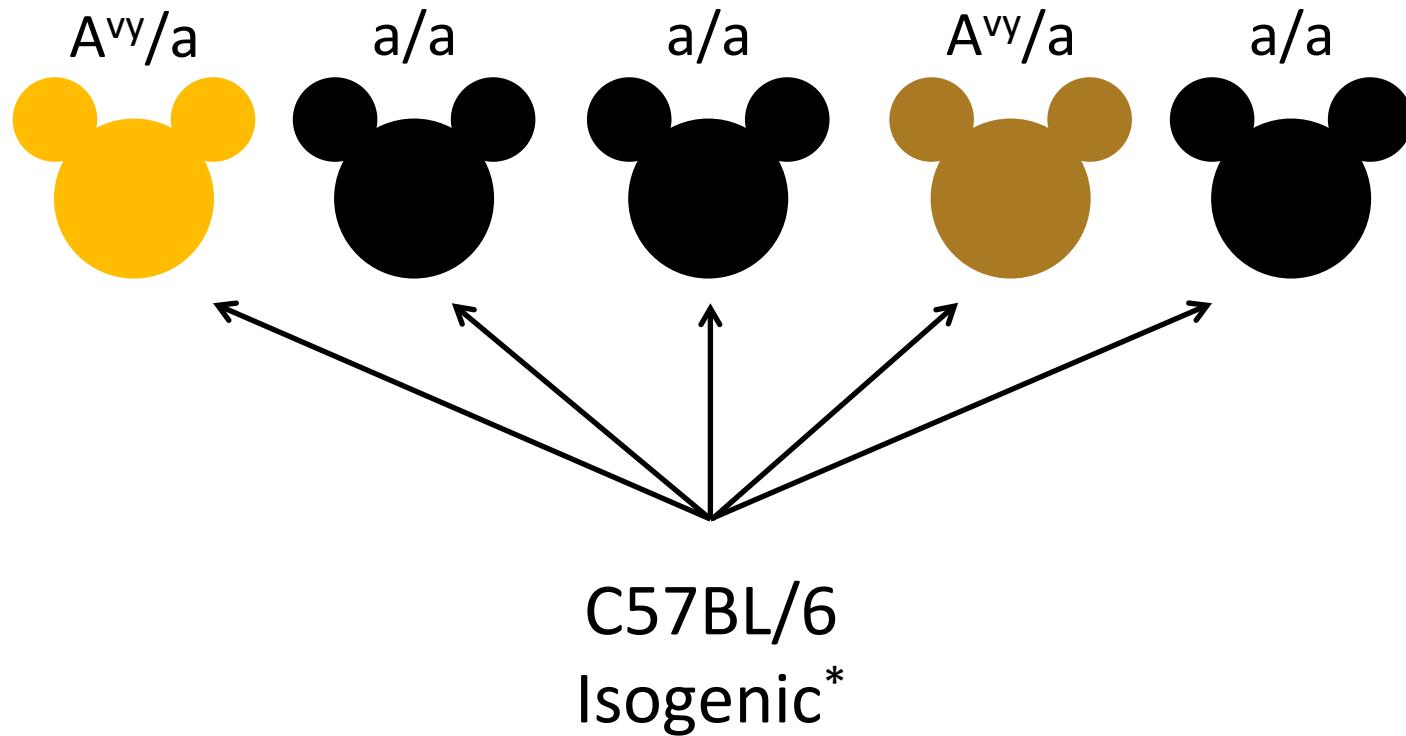


# Experimental design

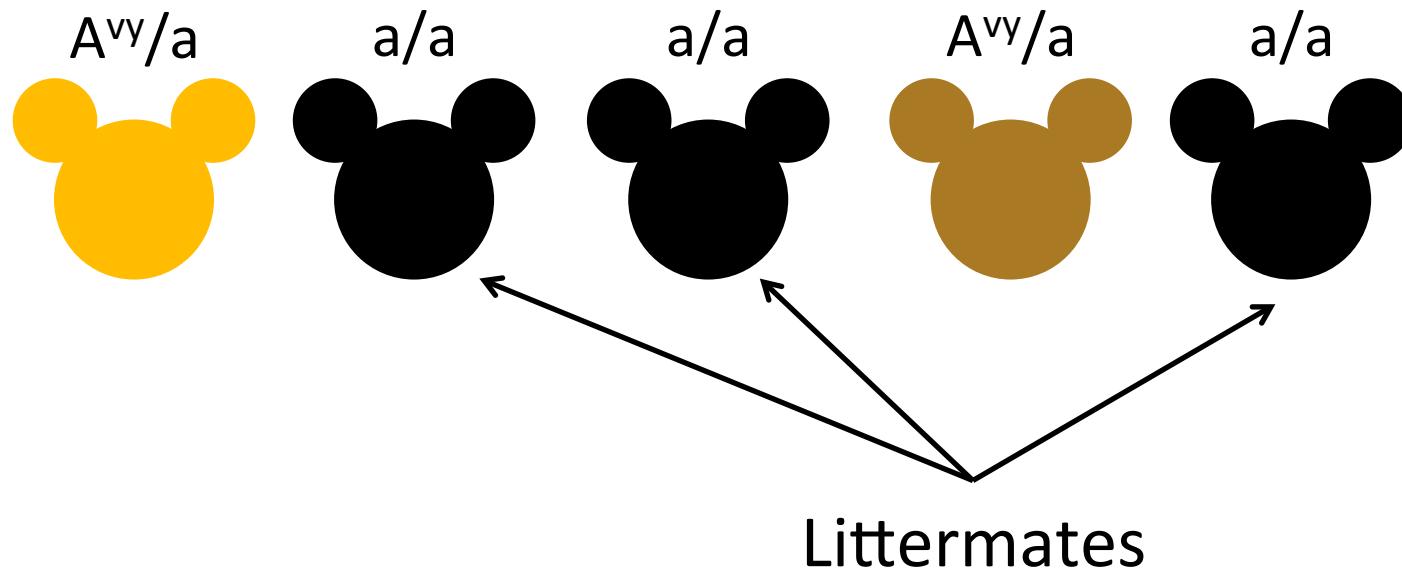


C57BL/6  
Isogenic\*

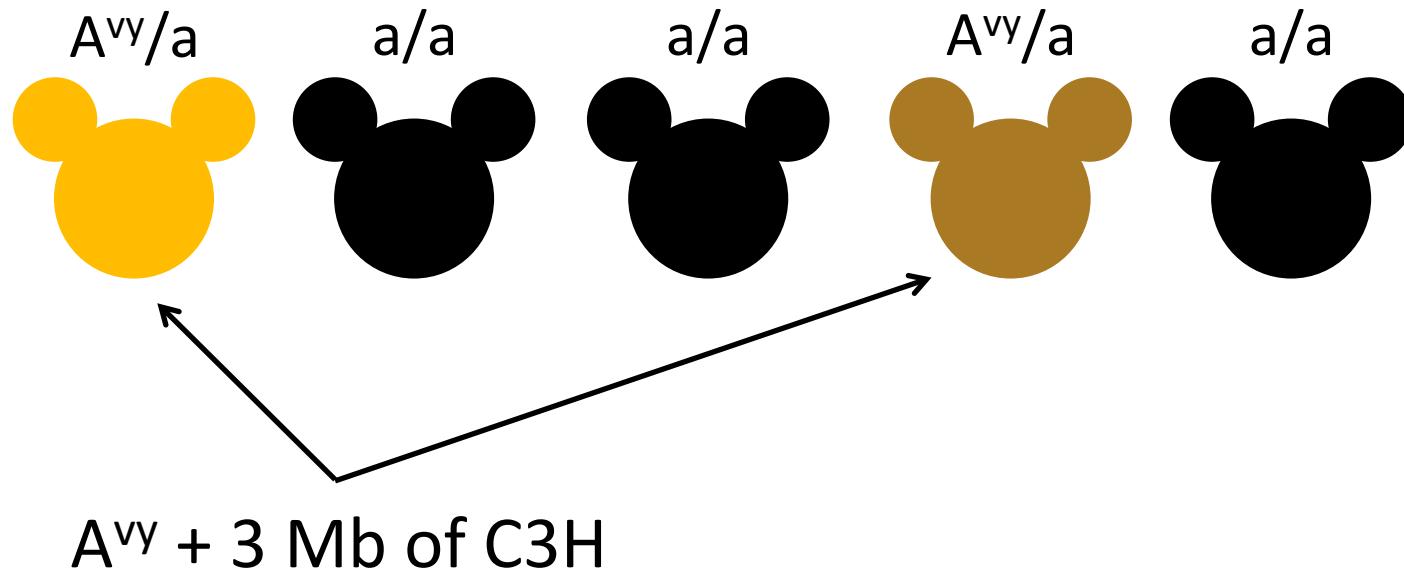
# Experimental design



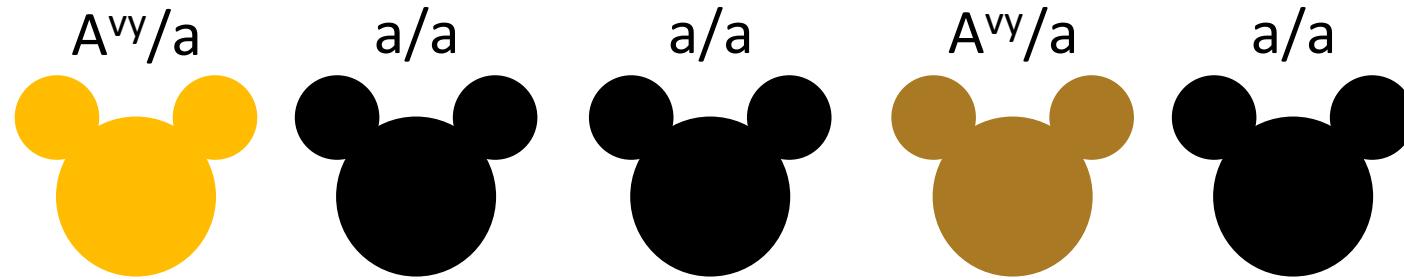
# Experimental design



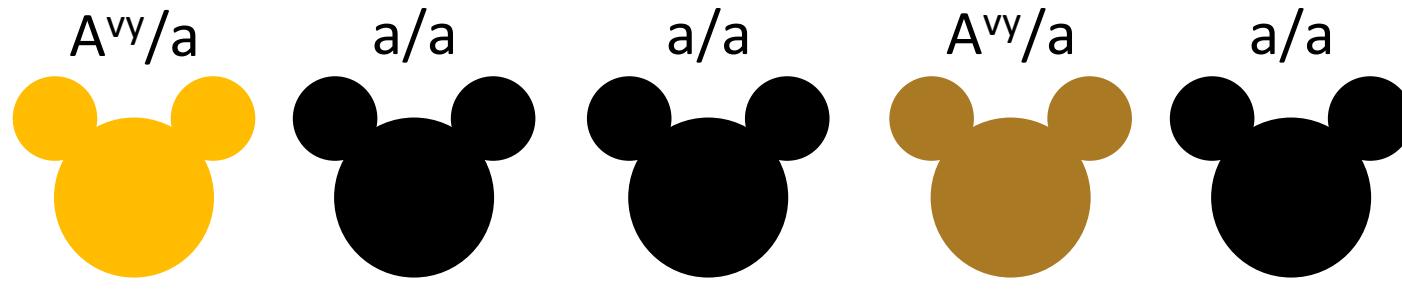
# Experimental design



# Experimental design



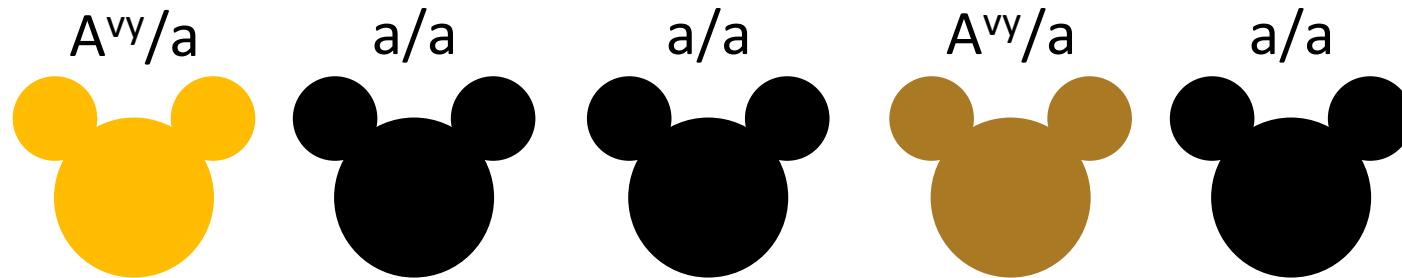
# Experimental design



$30\times$  whole-genome bisulfite-sequencing

=

# Experimental design



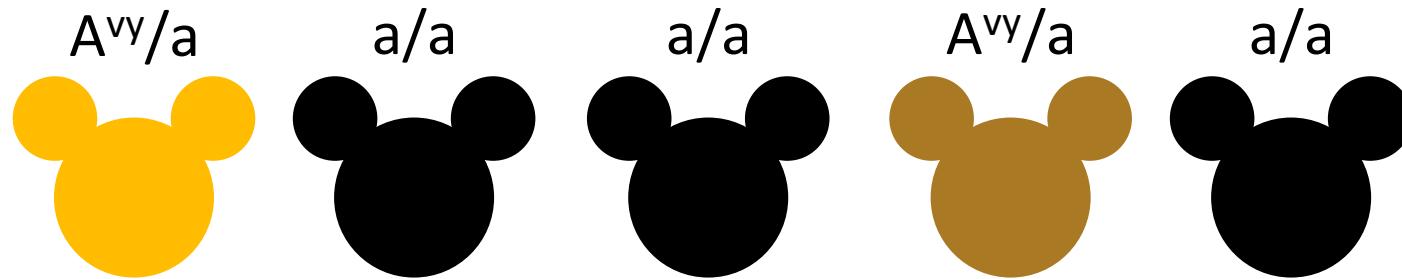
+

30× whole-genome bisulfite-sequencing

=

\$\$\$

# Experimental design



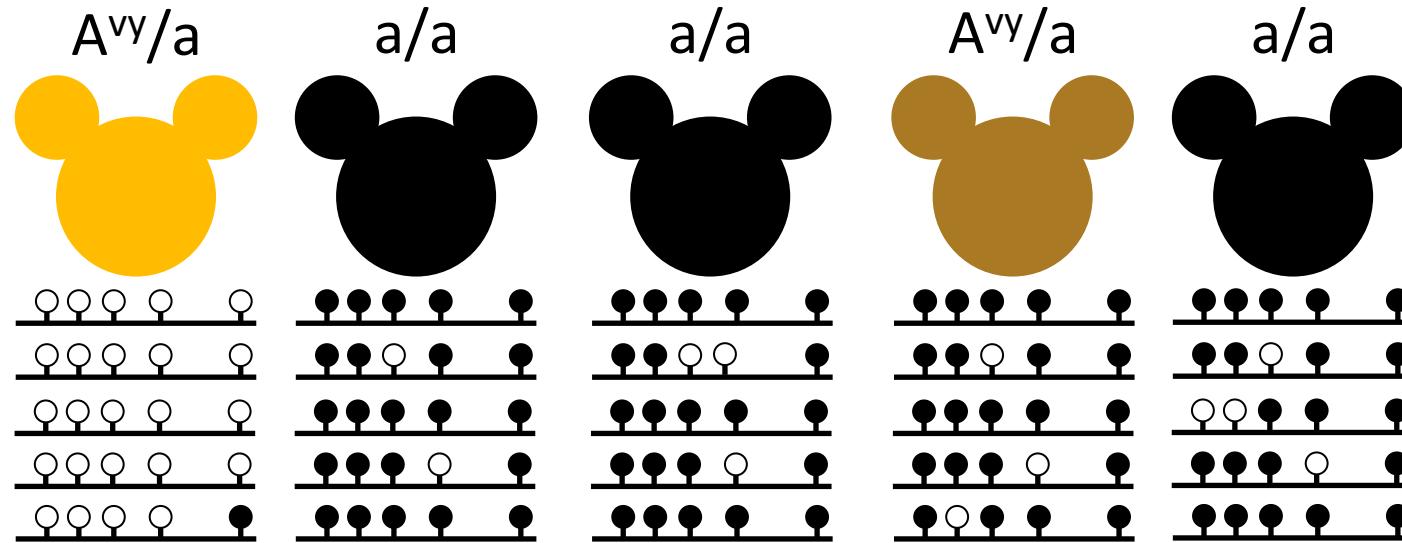
+

30× whole-genome bisulfite-sequencing

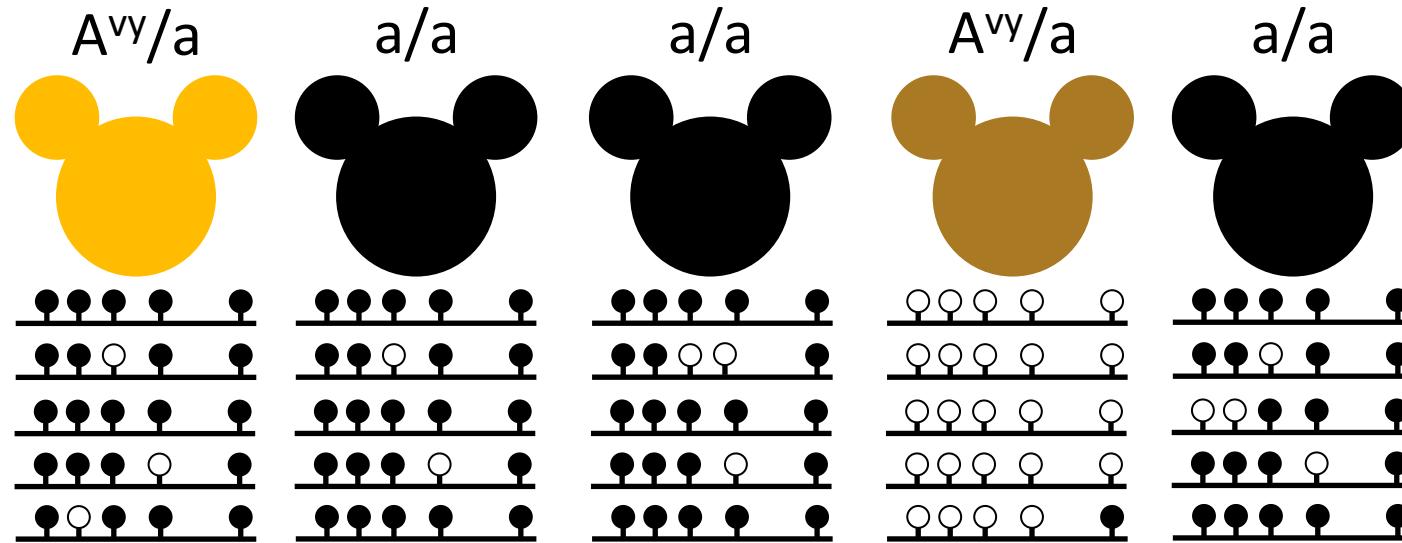
=

*epialleles*

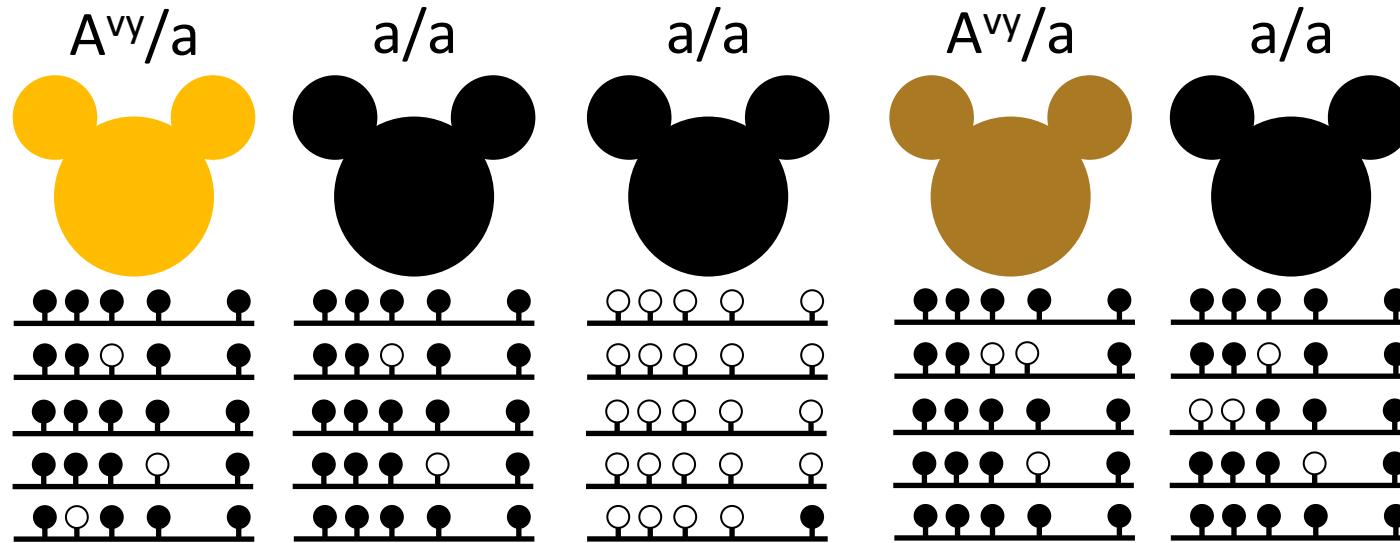
# Epialleles



# Epialleles



# Epialleles



# Method

1. Some of these mice are not like the others?
2. Is my neighbour also different?
3. Is my neighbour different in the same way as me?

# Some of these mice are not like the others?

	Mouse1	Mouse2	Mouse3	Mouse4	Mouse5
Methylated	17	31	15	23	9
Unmethylated	1	3	0	1	1

P-value = 0.76

# Some of these mice are not like the others?

	Mouse1	Mouse2	Mouse3	Mouse4	Mouse5
Methylated	38	79	59	69	44
Unmethylated	1	2	1	2	46

P-value =  $2 \times 10^{-25}$

# Some of these mice are not like the others?

	Mouse1	Mouse2	Mouse3	Mouse4	Mouse5
Methylated	38	79	59	69	44
Unmethylated	1	2	1	2	46

$$P\text{-value} = 2 \times 10^{-25}$$

$P\text{-value} < \text{threshold}$

→ (candidate) differentially methylated CpG (**DMC**)

Is my neighbour also different?



"Run-DMC Logo" Licensed under Public domain via Wikimedia Commons - [http://commons.wikimedia.org/wiki/File:Run-DMC\\_Logo.svg#mediaviewer/File:Run-DMC\\_Logo.svg](http://commons.wikimedia.org/wiki/File:Run-DMC_Logo.svg#mediaviewer/File:Run-DMC_Logo.svg)

# Is my neighbour also different?

## 1. Find runs of CpGs

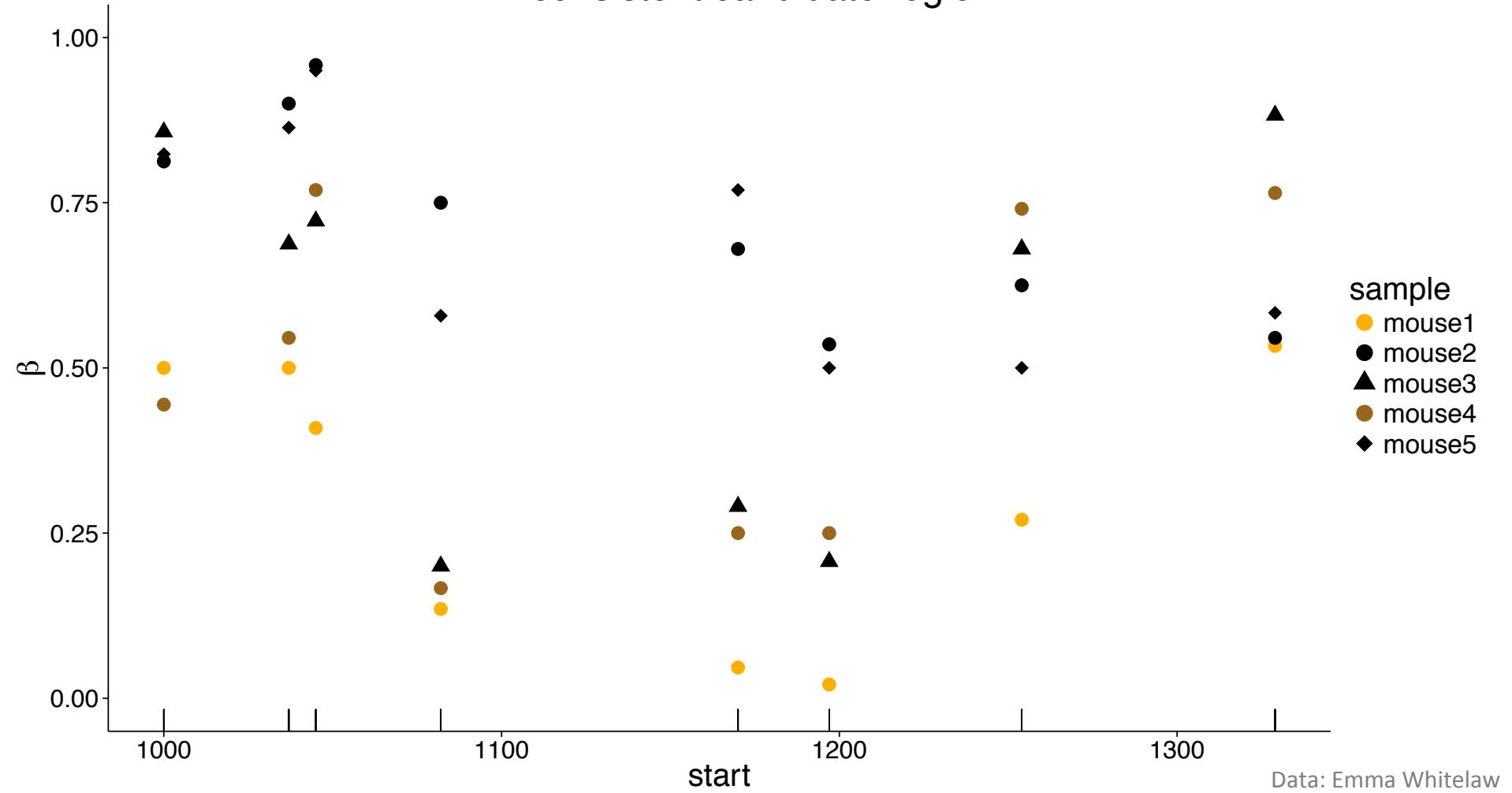
- P-value < *threshold*
- Within *distance* of next CpG
- Some allowance for missing or “insignificant” CpGs

## 2. Filter candidate runs

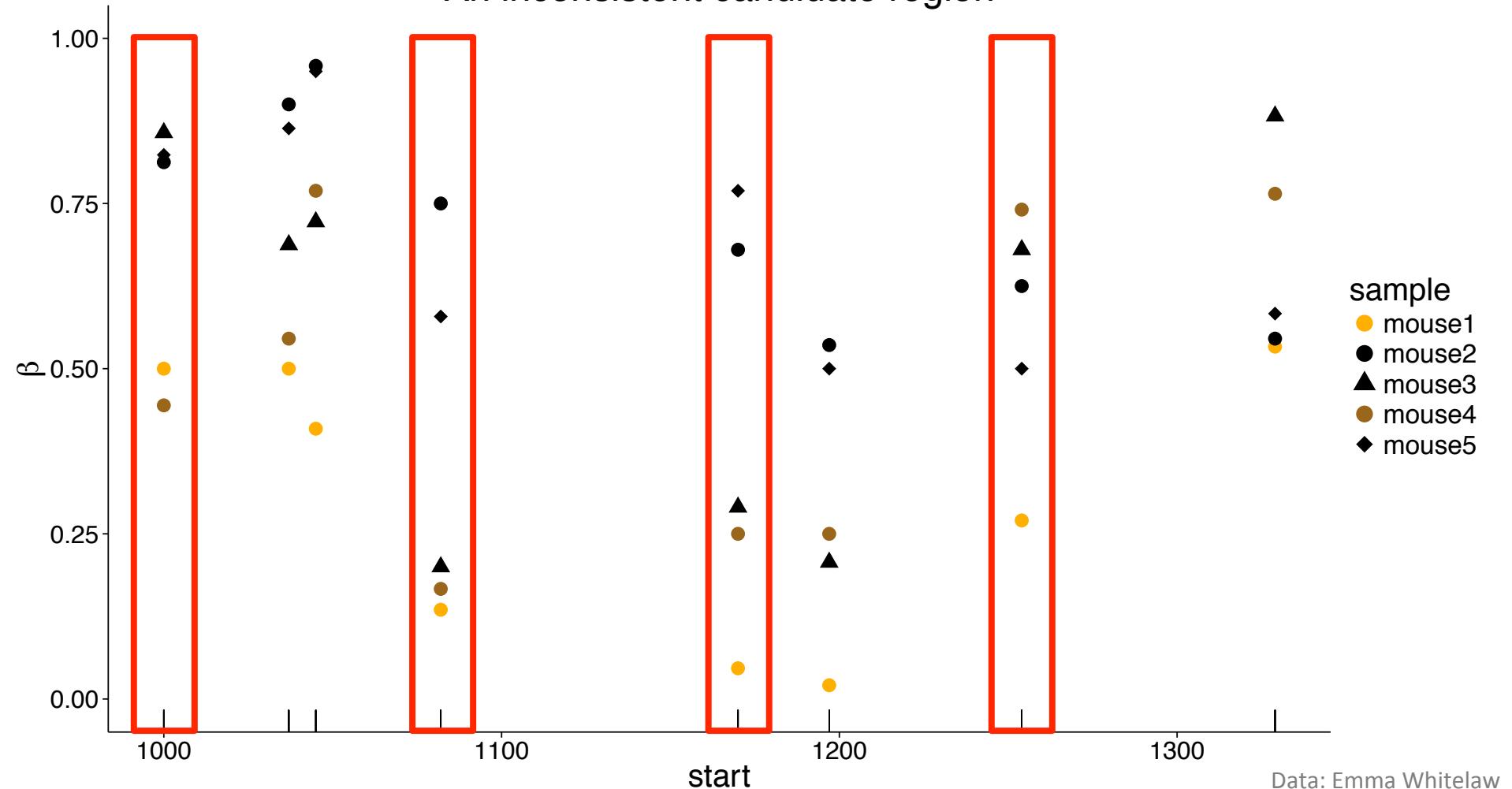
- Run contains enough CpGs

Is my neighbour different in the same  
way as me?

# An inconsistent candidate region

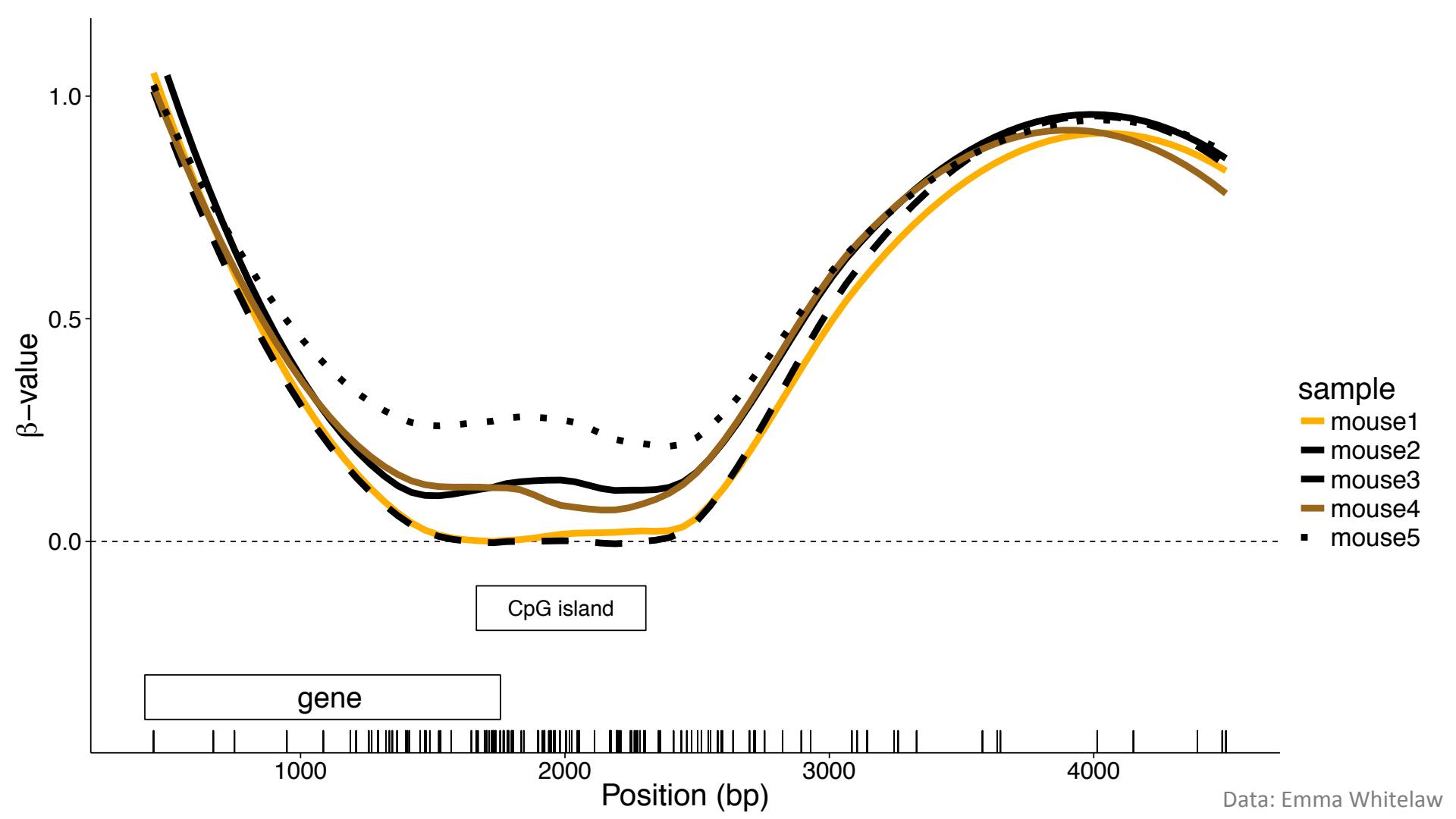


# An inconsistent candidate region



# Is my neighbour different in the same way as me?

- Flag regions with 3-way interaction between sample  $\times$  methylation level  $\times$  position
  - Not quite what we want
  - So plot, plot, plot



How I work, what I found, and what I'm proud of

# Summary

*“Doesn’t the gardener lavish more care on  
the thorns than on the flowers”*



- *Hartman in Metamorphosis* by S.Y. Agnon  
Via @erichlya

"Agnon" of Unknown - The David B. Keidan Collection of Digital Images from the Central Zionist Archives (via Harvard University Library). Licensed under the Public Domain via Wikimedia Commons - <http://commons.wikimedia.org/wiki/File:Agnon.jpg#mediaviewer/File:Agnon.jpg>

*“You can observe a lot by watching”*



- Yogi Berra

"Yogi Berra 1956" by unknown - Baseball Digest, front cover, September 1956 issue. [1]. Licensed under Public domain via Wikimedia Commons  
[http://commons.wikimedia.org/wiki/File:Yogi\\_Berra\\_1956.png#mediaviewer/File:Yogi\\_Berra\\_1956.png](http://commons.wikimedia.org/wiki/File:Yogi_Berra_1956.png#mediaviewer/File:Yogi_Berra_1956.png)



**COLLABORATE  
AND LISTEN**

*- Robert Van Winkle*



*- Robert Van Winkle  
a.k.a. Vanilla Ice*

# What I found

- Estimated strong spatial dependence of DNA methylation

# What I found

- Estimated strong spatial dependence of DNA methylation
- Cell-type differences in dependence structure

# What I found

- Estimated strong spatial dependence of DNA methylation
- Cell-type differences in dependence structure
- Evidence of higher order chromatin structure in spatial dependence data

# What I'm proud of

[www.github.com/PeteHaitch](https://www.github.com/PeteHaitch)

# Acknowledgements



Terry Speed



Peter Hall

# Acknowledgements

## Data

- Ryan Lister et al. (UWA, Salk Institute)
- Sue Clark, Aaron Statham (Garvan Institute)
- Emma Whitelaw, Harry Oey (La Trobe)
- Kasper Hansen, Rafael Irizarry (Johns Hopkins, Harvard)
- **Everyone who makes their data publicly available**

# Acknowledgements

## Data

- Ryan Lister et al. (UWA, Salk Institute)
- Sue Clark, Aaron Statham (Garvan Institute)
- Emma Whitelaw, Harry Oey (La Trobe)
- Kasper Hansen, Rafael Irizarry (Johns Hopkins, Harvard)
- **Everyone who makes their data publicly available**

## Methodology & technology

- Kasper Hansen, Rafael Irizarry (Johns Hopkins, Harvard)
- Felix Krueger (Babraham Institute)
- Toby Sargeant (WEHI)
- Keith Satterley (WEHI)
- Bioconductor developers
- WEHI Bioinformatics
- **Everyone who makes their software open source**

# Acknowledgements

## Data

- Ryan Lister et al. (UWA, Salk Institute)
- Sue Clark, Aaron Statham (Garvan Institute)
- Emma Whitelaw, Harry Oey (La Trobe)
- Kasper Hansen, Rafael Irizarry (Johns Hopkins, Harvard)
- **Everyone who makes their data publicly available**

## Methodology & technology

- Kasper Hansen, Rafael Irizarry (Johns Hopkins, Harvard)
- Felix Krueger (Babraham Institute)
- Toby Sargeant (WEHI)
- Keith Satterley (WEHI)
- Bioconductor developers
- WEHI Bioinformatics
- **Everyone who makes their software open source**

**Funding:** APA and VLSCI

# Acknowledgements

## Data

- Ryan Lister et al. (UWA, Salk Institute)
- Sue Clark, Aaron Statham (Garvan Institute)
- Emma Whitelaw, Harry Oey (La Trobe)
- Kasper Hansen, Rafael Irizarry (Johns Hopkins, Harvard)
- **Everyone who makes their data publicly available**

## Methodology & technology

- Kasper Hansen, Rafael Irizarry (Johns Hopkins, Harvard)
- Felix Krueger (Babraham Institute)
- Toby Sargeant (WEHI)
- Keith Satterley (WEHI)
- Bioconductor developers
- WEHI Bioinformatics
- **Everyone who makes their software open source**

**Funding:** APA and VLSCI

**Sanity:** Family and friends

