

# Detection and inference of differentially methylated regions from bisulfite sequencing

Keegan Korthauer

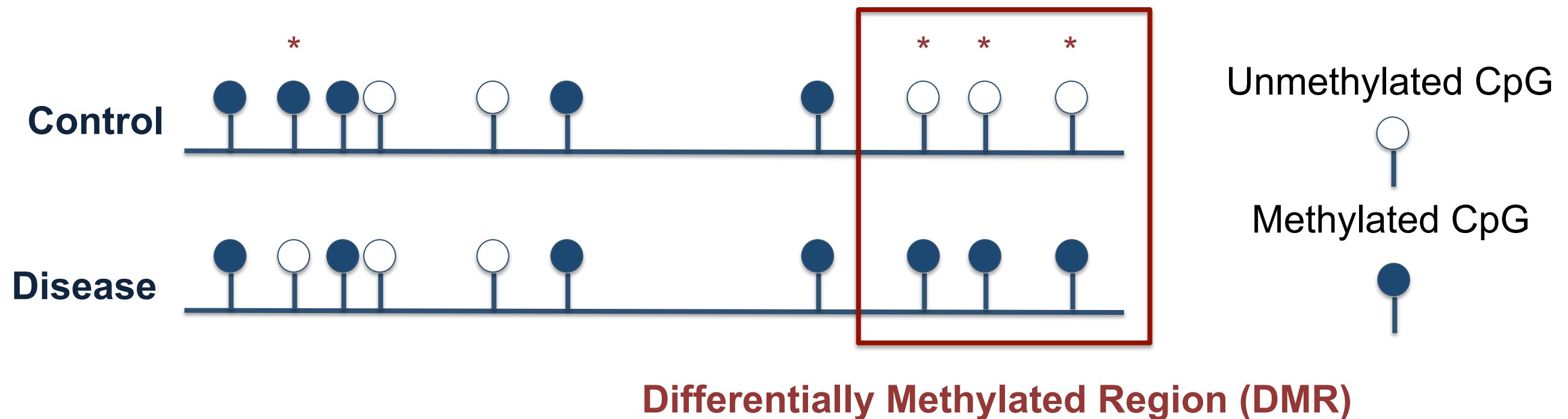
 @keegankorthauer

Bioconductor Meeting 2017

# Differential methylation

Differential methylation commonly studied in

- Cancer
- Developmental stages / Aging
- Tissue types



# Whole Genome Bisulfite Sequencing

Treat DNA with bisulfite before sequencing

- Methylated C -> **unaffected**
- Unmethylated C -> **appear as T**

<p><b>GAGCGATGGATAGCG</b></p> <p><b>CGAGTGATGGATAGC</b></p> <p><b>GTTACGAGCGACGG</b></p> <p><b>TGGTTACGAGCGATG</b></p> <p><b>. . . TCTCGGTTACGAGCGACGGATAGCG . . .</b></p>	<p>Bisulfite sequencing reads</p> <p>Reference genome</p>
--	---

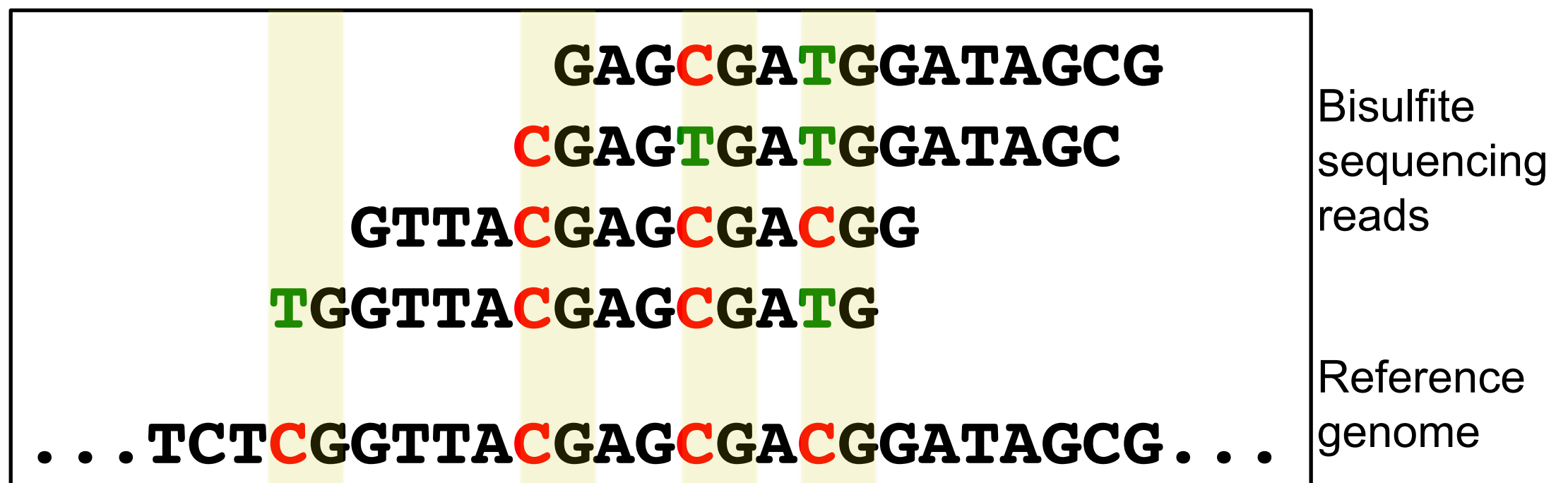
# Whole Genome Bisulfite Sequencing

Treat DNA with bisulfite before sequencing

- Methylated C -> **unaffected**
- Unmethylated C -> **appear as T**

Methylation Counts and Proportions

0/1	3/3	3/4	1/4
0	1	0.75	0.25



# Methods for DMR detection

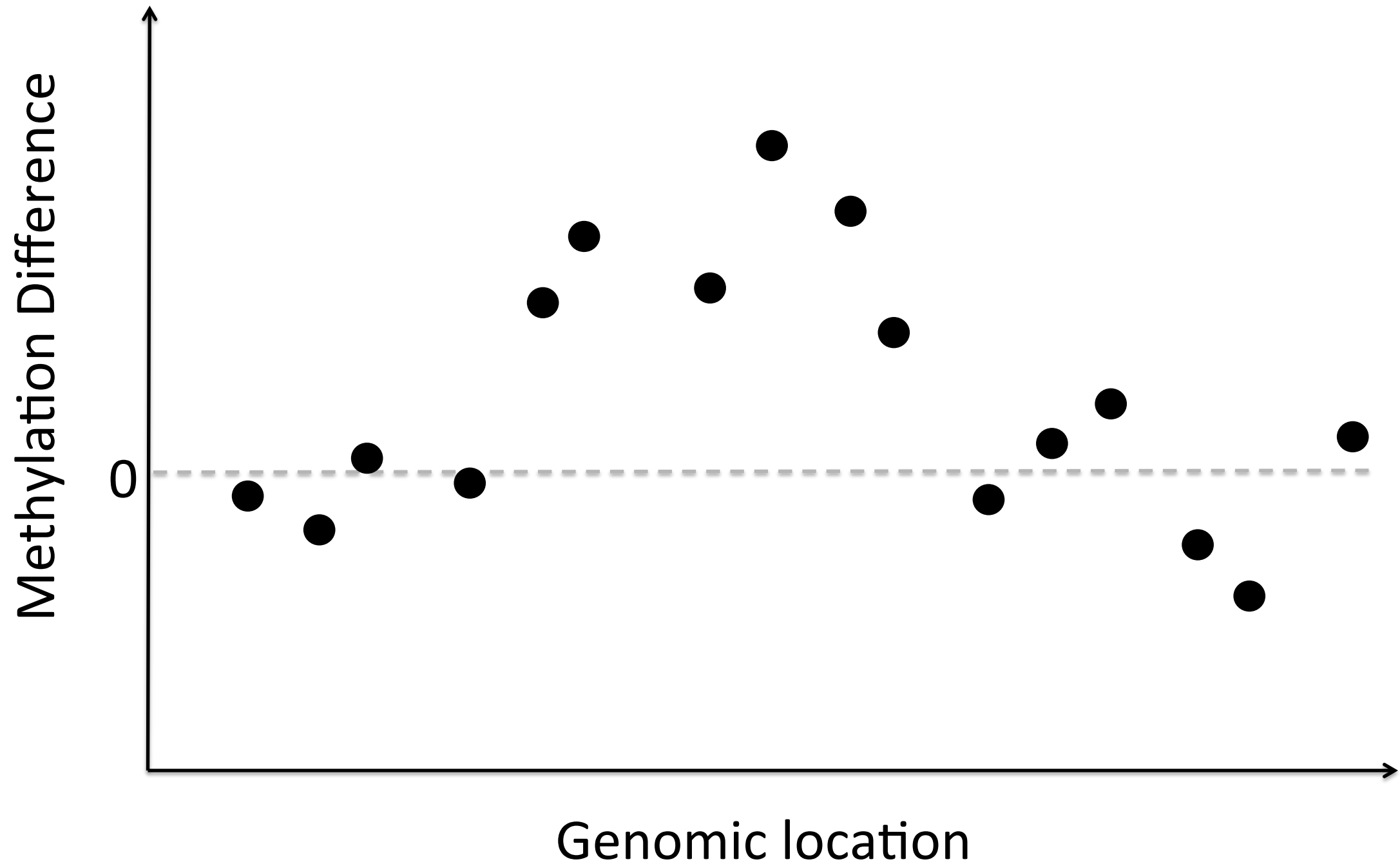
## Challenges

- Have to handle small sample sizes !
- Accommodate known sources of variability
- Detect region boundaries

## Two main strategies

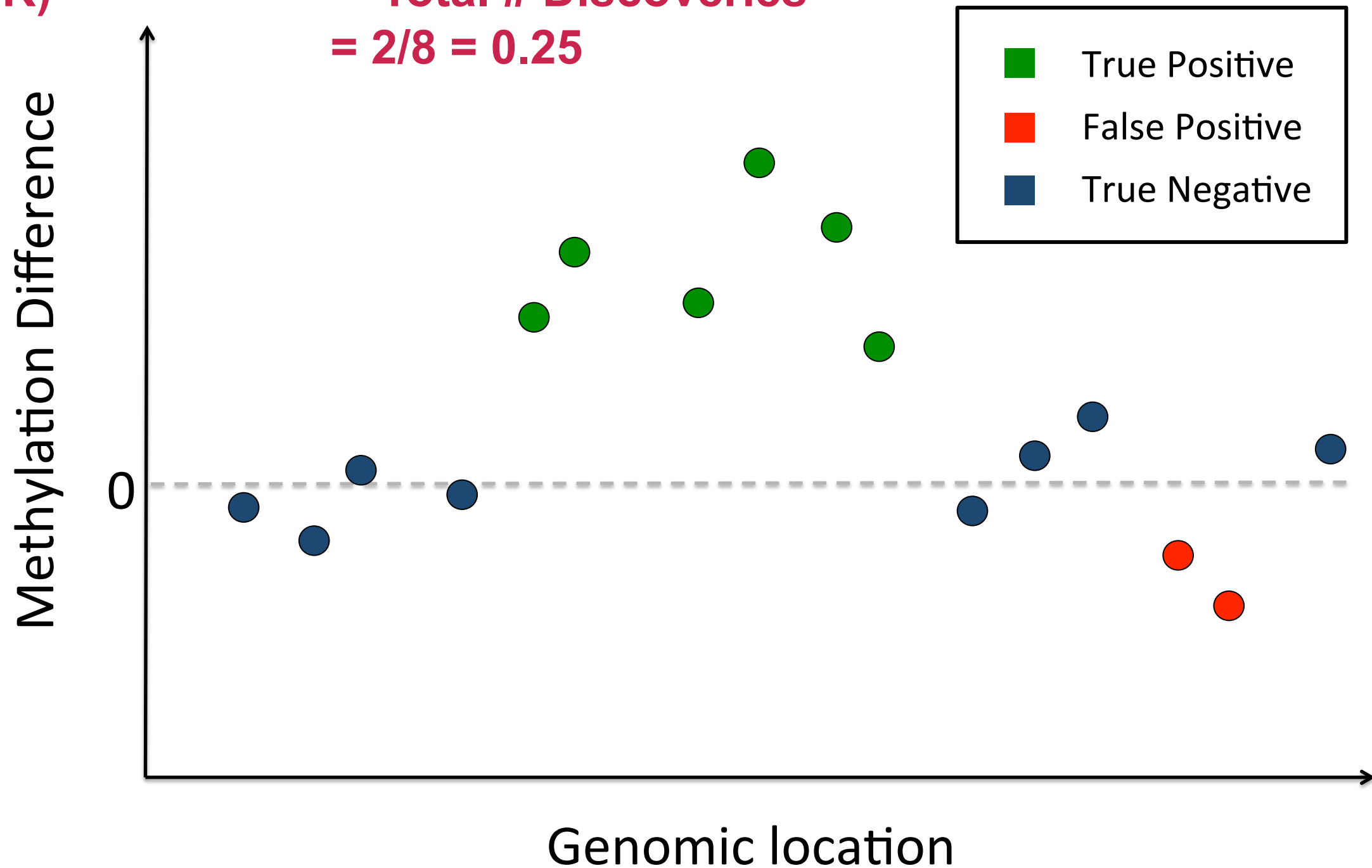
- Detect significant CpGs, then group together neighbors to form regions (DSS, BSmooth)
- Targeted regions / sliding windows specified in advance (BiSeq, MOABS)

# Danger of grouping significant CpGs



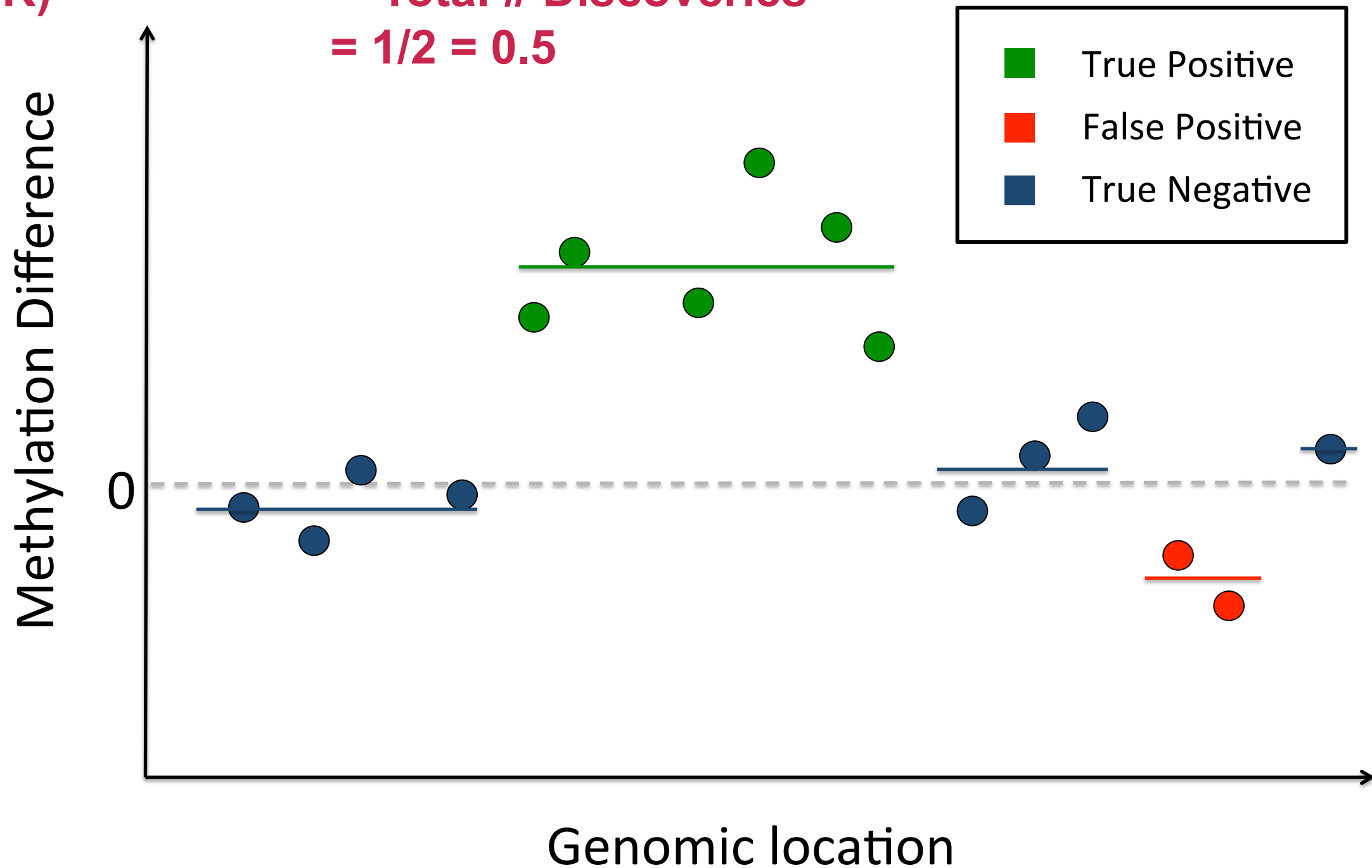
# Danger of grouping significant CpGs

**False Discovery Rate =  $\frac{\# \text{ False Discoveries}}{\text{Total \# Discoveries}}$**   
**(FDR)**  
**=  $\frac{2}{8} = 0.25$**



# Danger of grouping significant CpGs

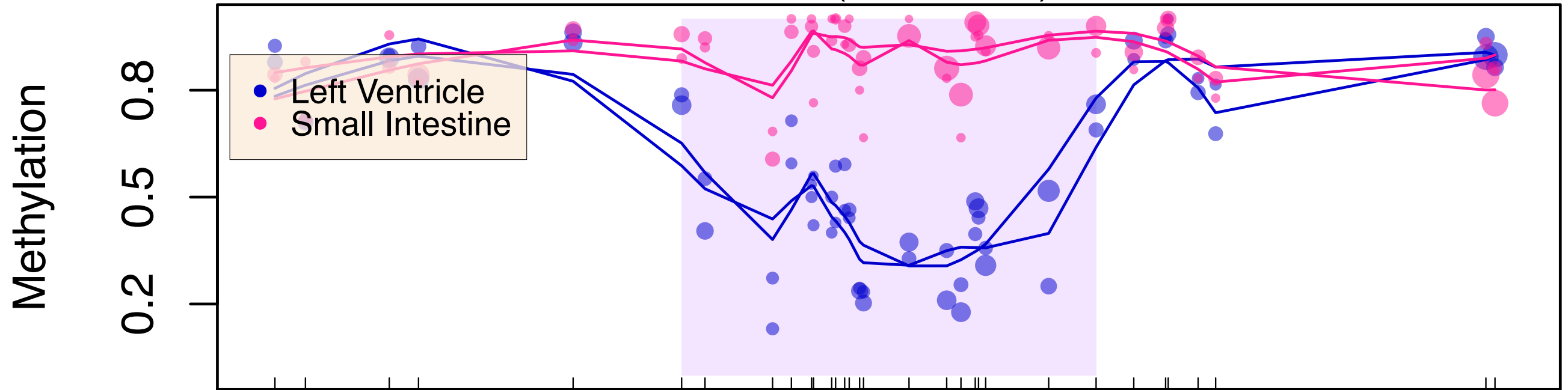
False Discovery Rate =  $\frac{\# \text{ False Discoveries}}{\text{Total \# Discoveries}}$   
(FDR)  
 $= 1/2 = 0.5$



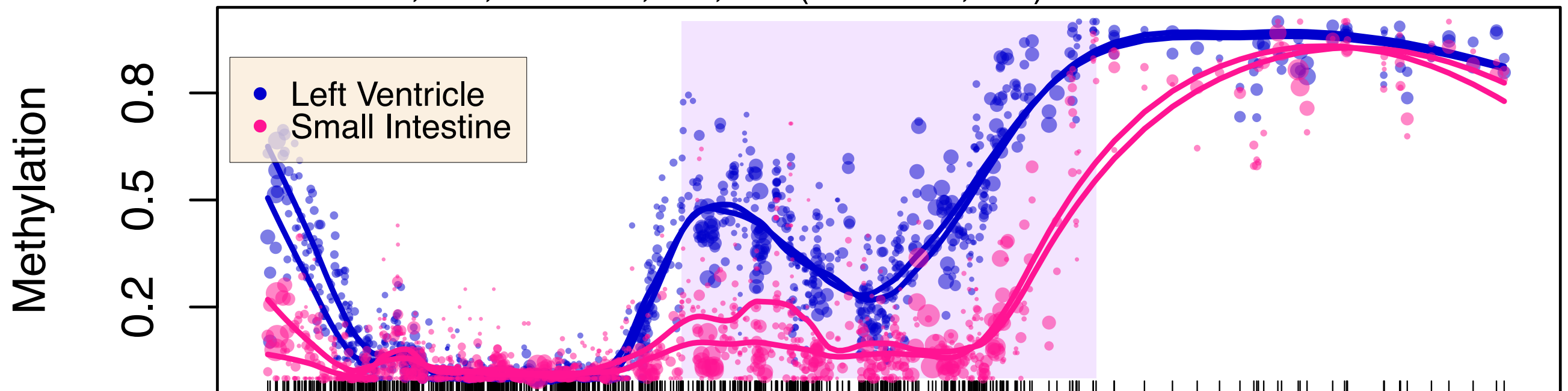


# DMRs come in all shapes and sizes

chr5: 75,654,594 – 75,655,232 (width = 639)



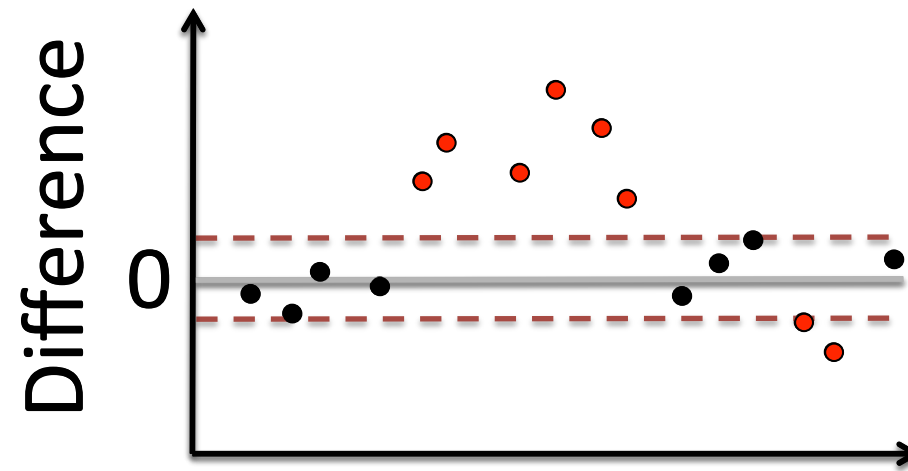
chr16: 54,288,213 – 54,292,817 (width = 4,605)



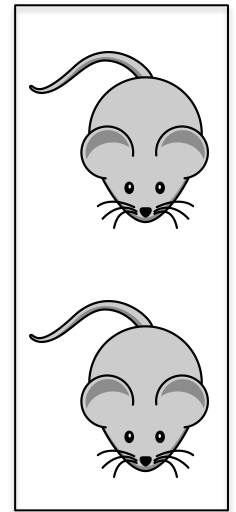
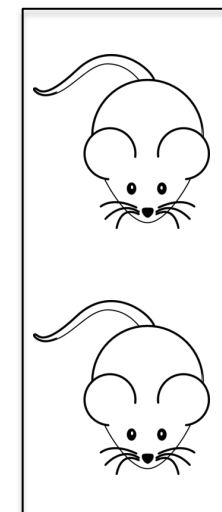
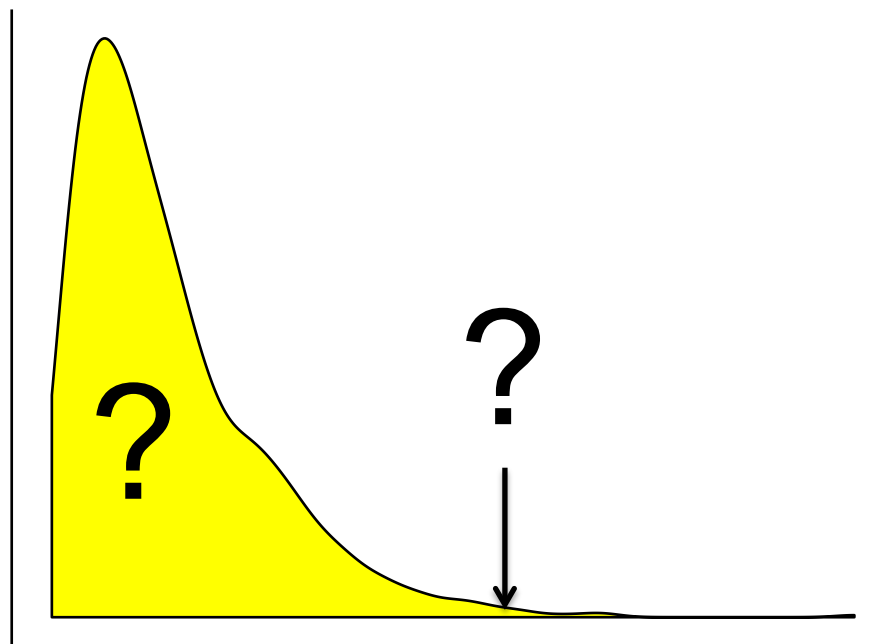
-> Detect region boundaries from the data

# Our approach: dmrseq

1. Detect regions by scanning the genome for candidate regions

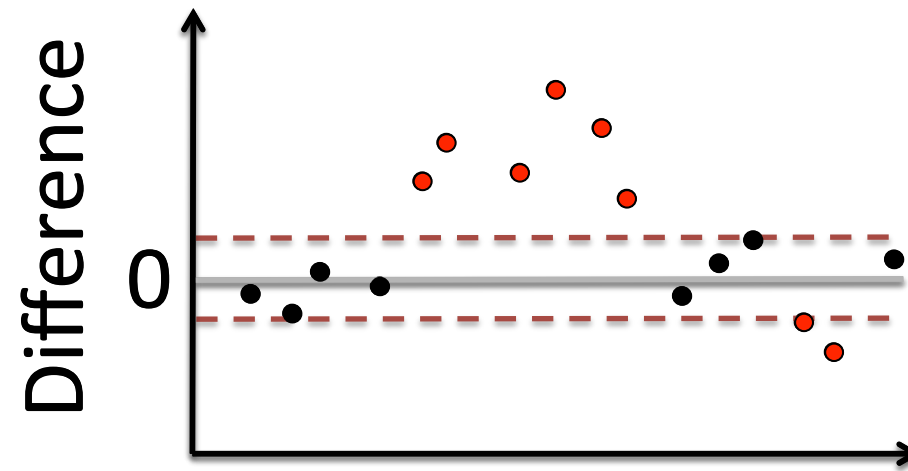


2. Region-level summary scores compared to a null to evaluate significance

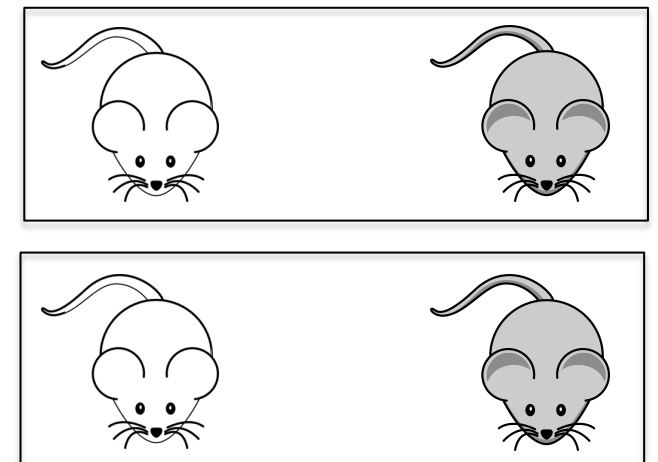
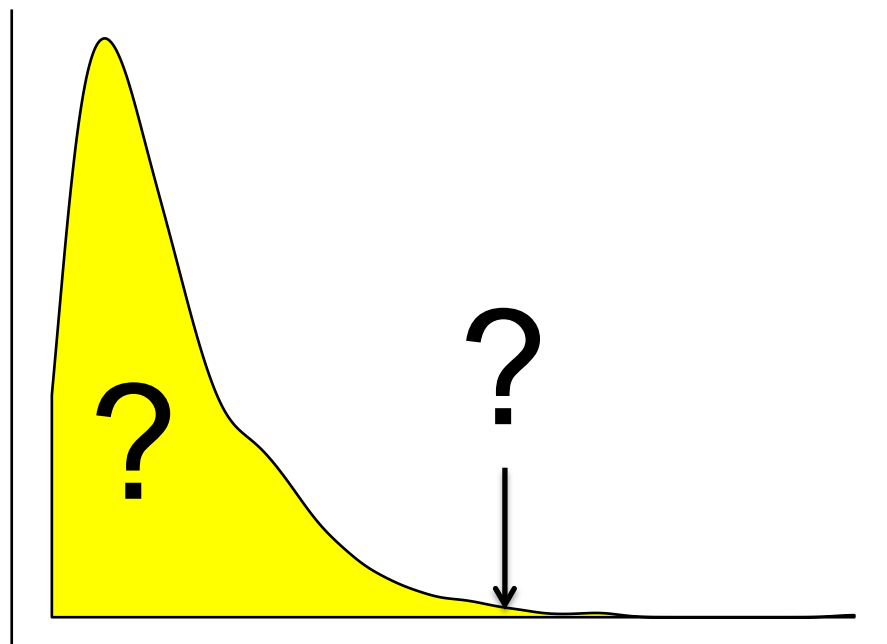


# Our approach: dmrseq

1. Detect regions by scanning the genome for candidate regions

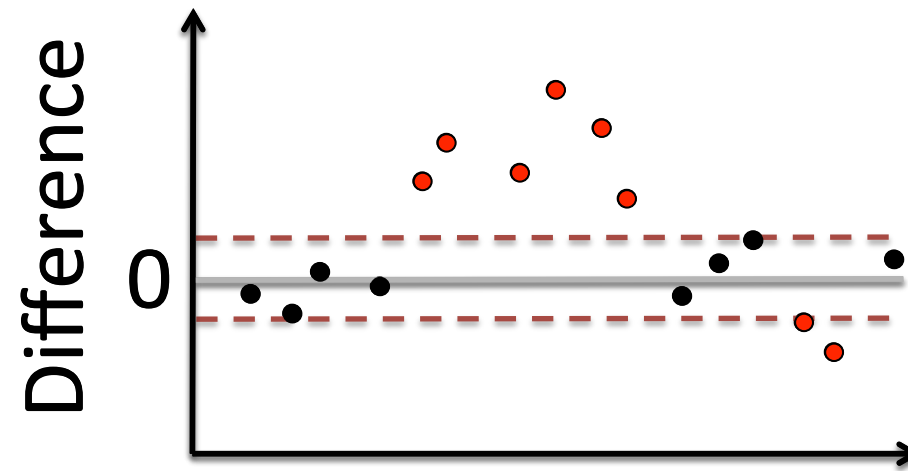


2. Region-level summary scores compared to a null to evaluate significance

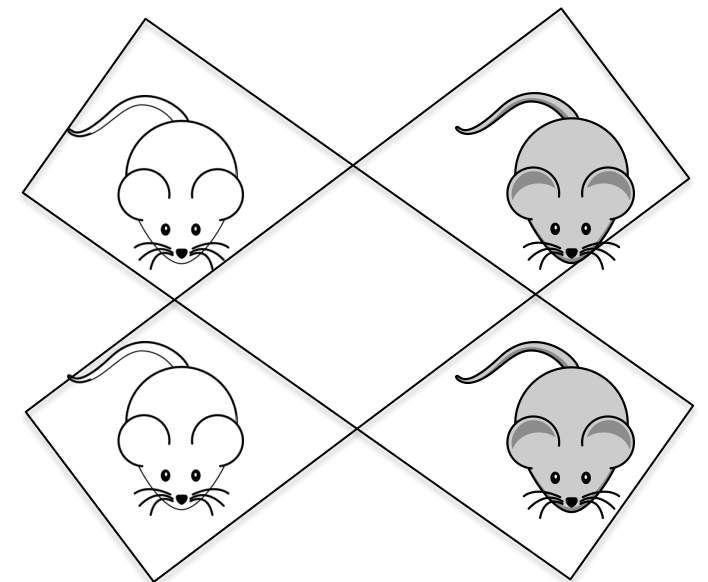
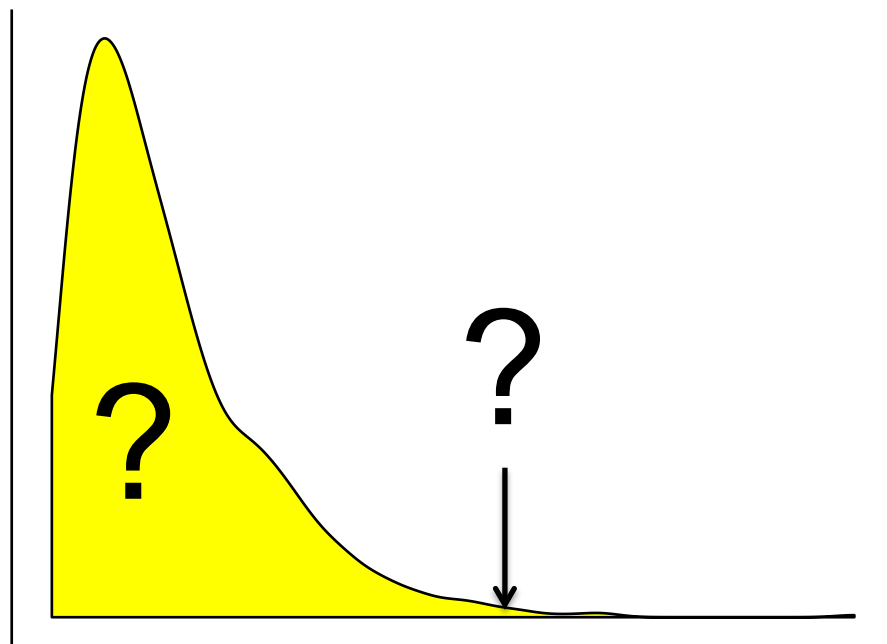


# Our approach: dmrseq

1. Detect regions by scanning the genome for candidate regions



2. Region-level summary scores compared to a null to evaluate significance



# Region level Summary Scores

Summary statistic that is approximately exchangeable across the genome so we can generate a **pooled null**

- Biological variability among samples
- Correlation among nearby CpGs
- Higher variability in CpGs with lower coverage

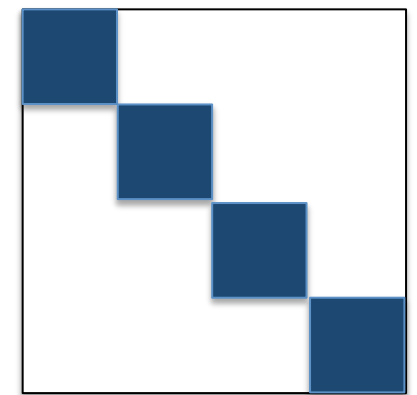


Generalized Least Squares (GLS) regression

$$\arcsin(2\pi_{ijr} - 1) = \sum_l \alpha_{lr} 1_{[i=l]} + \beta_r X_j + \varepsilon_{ijr}$$

- Nested autoregressive correlated error structure

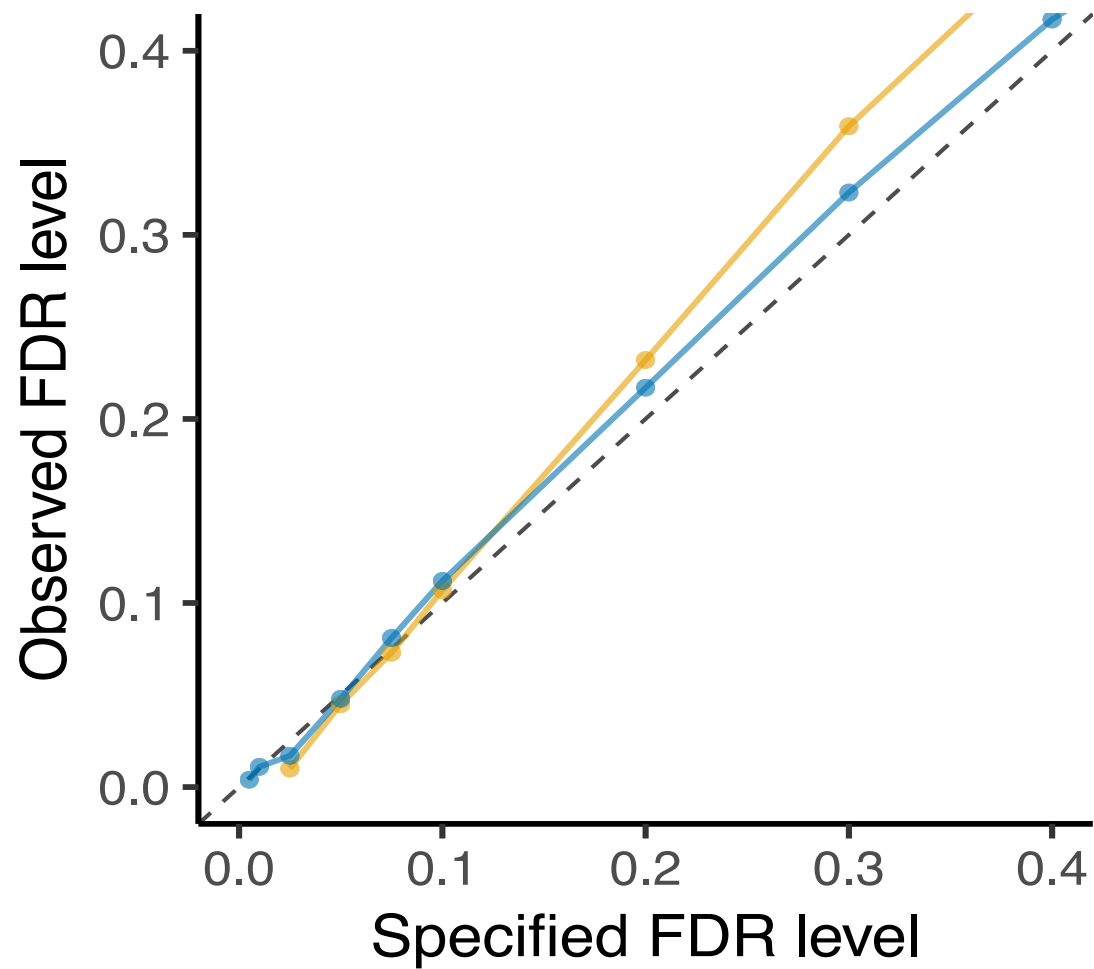
$$\text{Var}(\varepsilon_{ijr}) =$$



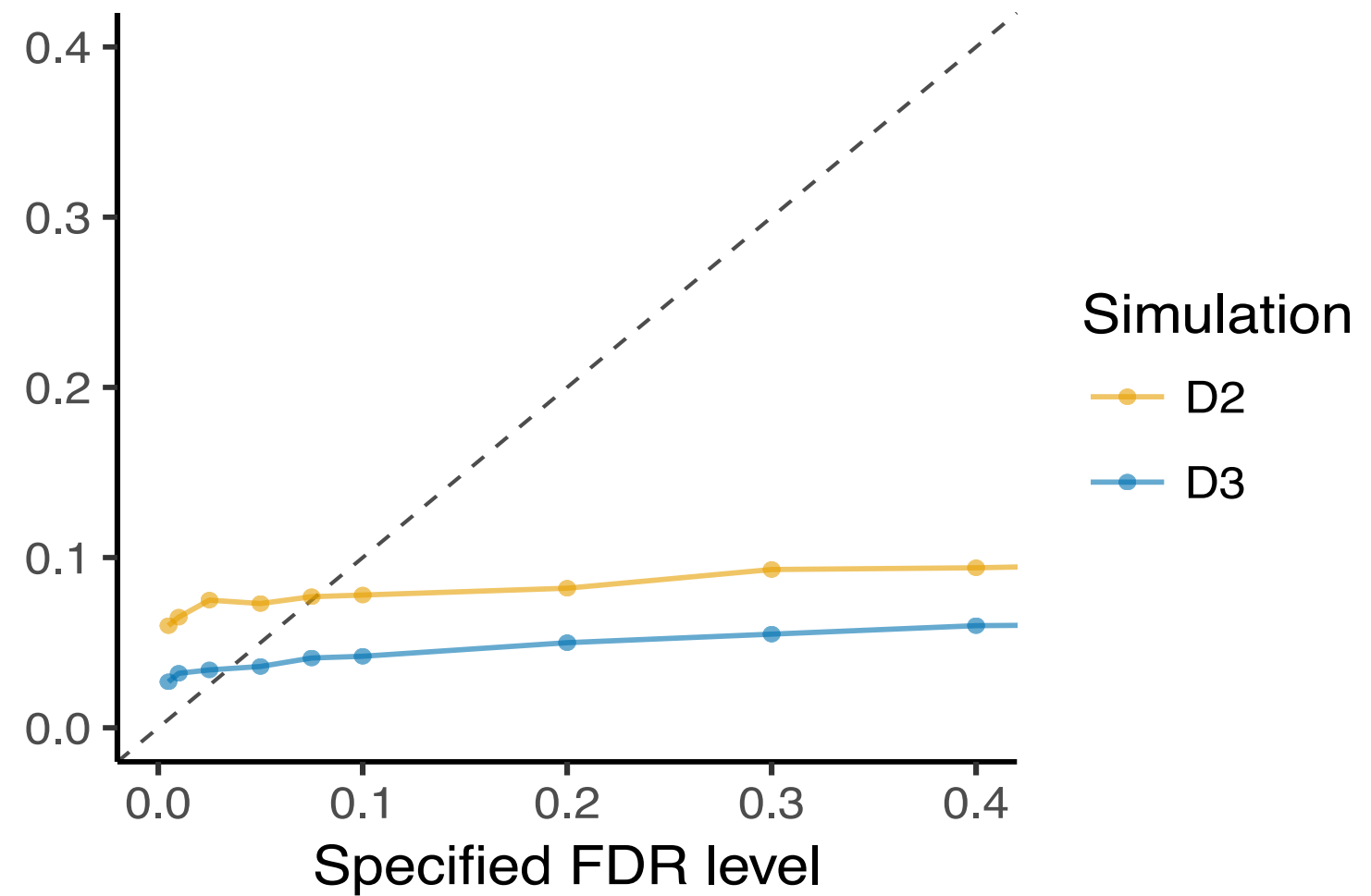
- Coverage-weighted variance

# dmrseq accurately controls FDR

**(A) dmrseq**

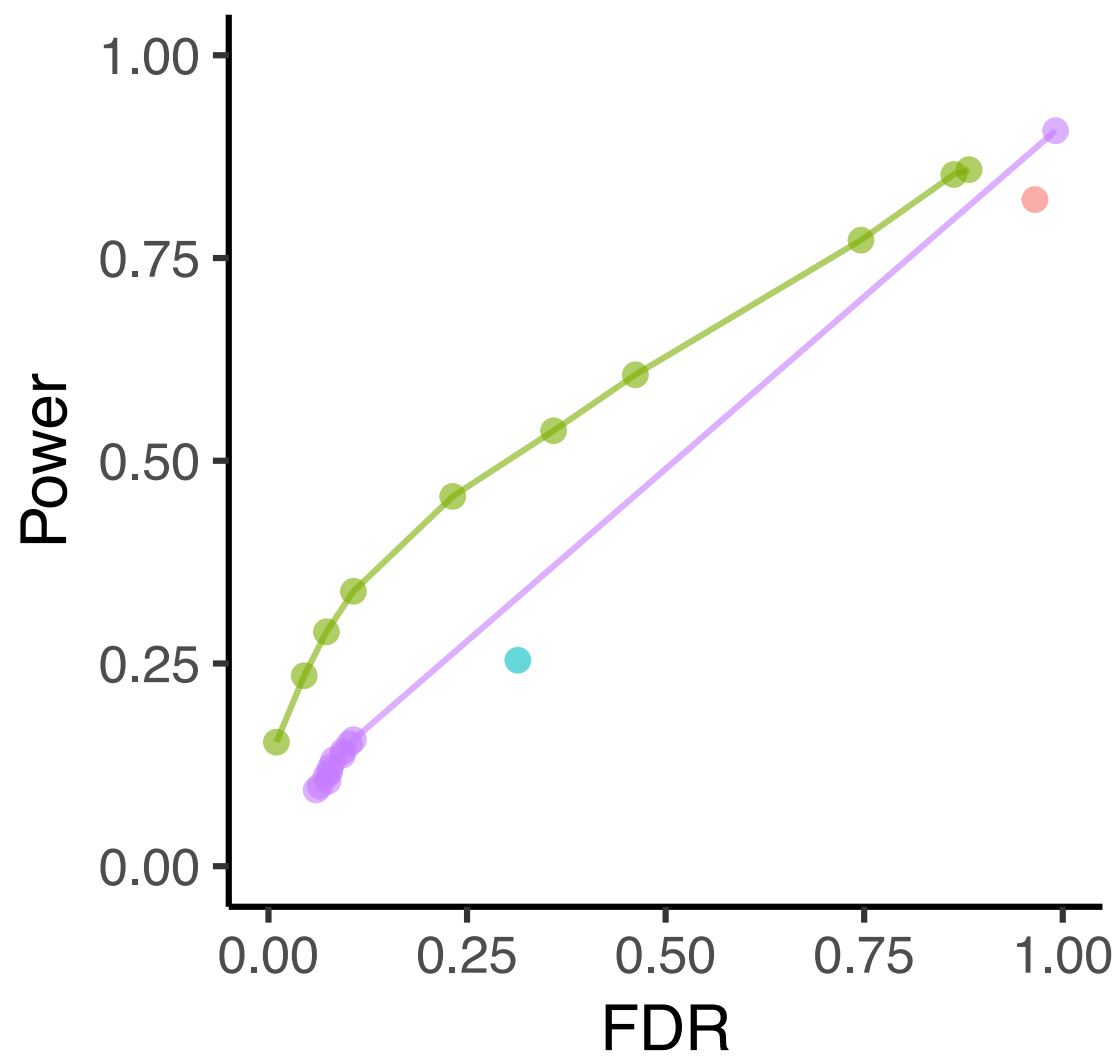


**(B) metilene**

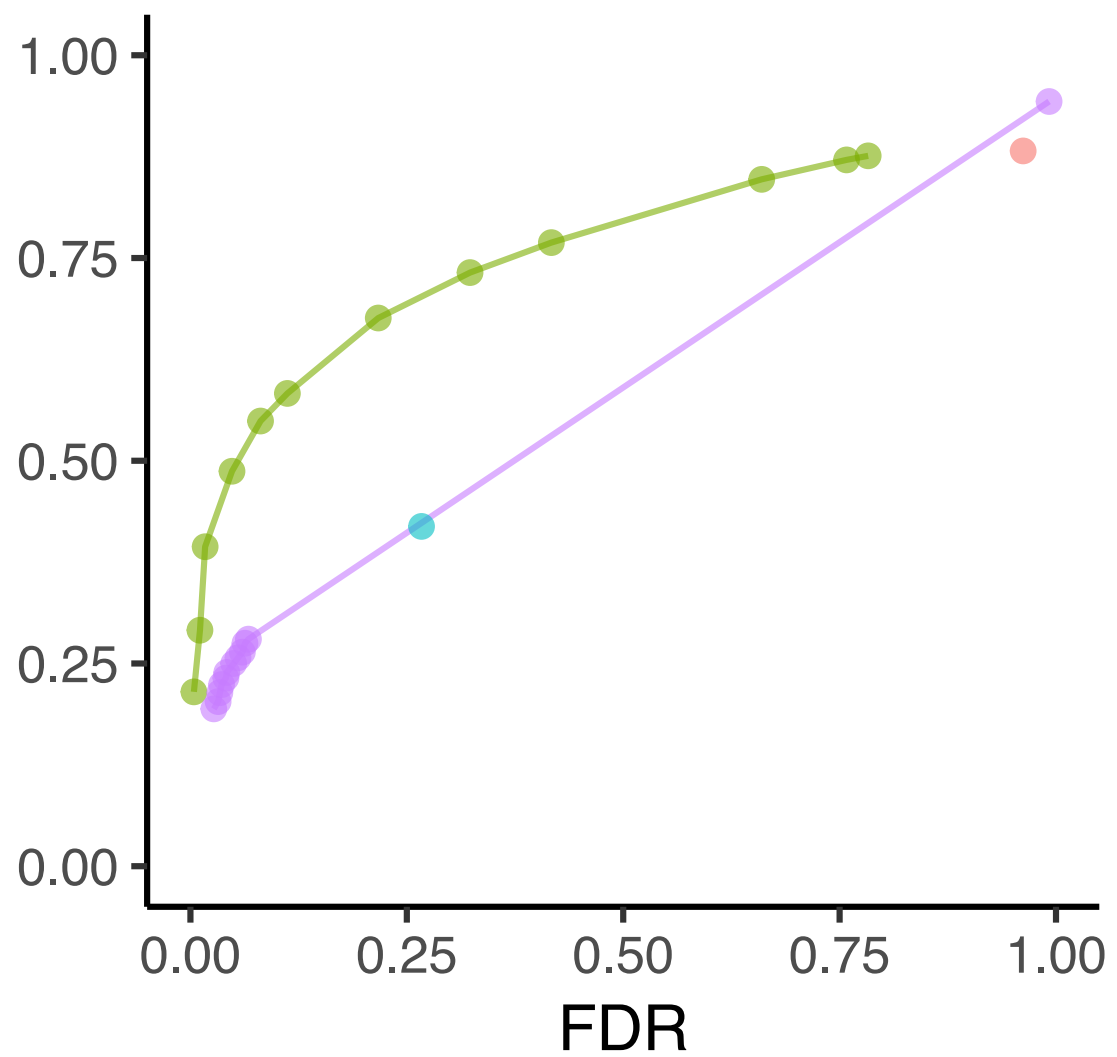


# dmrseq has high sensitivity

**(A) Simulation D2**



**(B) Simulation D3**



Method

BSmooth

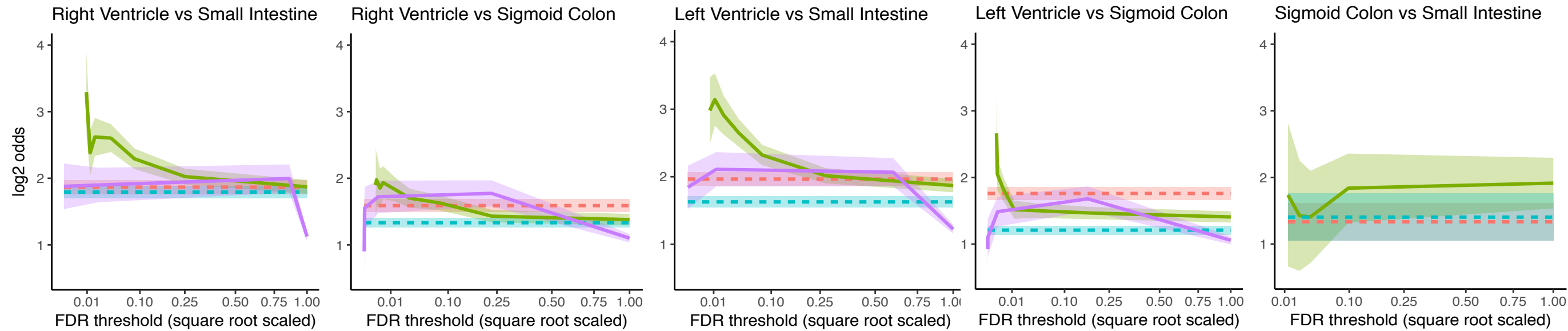
dmrseq

DSS

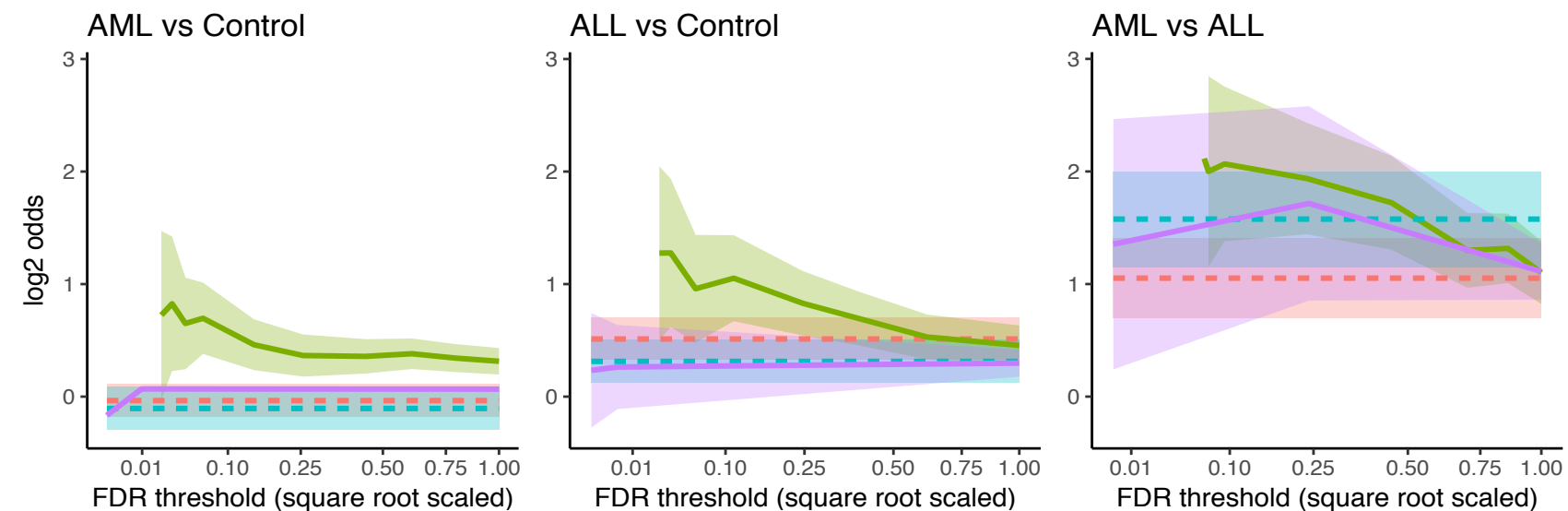
metilene

# DMRs enriched for associations with differential expression

## (A) Roadmap Tissue Comparisons



## (B) Murine Leukemia Models



Method — BSmooth — dmrseq — DSS — metilene



# dmrseq Summary

- dmrseq **identifies and prioritizes DMRs** from bisulfite sequencing experiments
- Computes region summary statistics that account for known sources of variability across the genome
- Achieves **accurate False Discovery Rate control** by generating a null distribution that pools information across the genome

# Learn More

## Slides

[goo.gl/MwQz5f](https://goo.gl/MwQz5f)

## dmrseq available on GitHub

<https://github.com/kdkorthauer/dmrseq>



# Contact



[keegan@jimmy.harvard.edu](mailto:keegan@jimmy.harvard.edu)



[@keegsdur](https://twitter.com/keegsdur)

# Acknowledgements



## Rafa Lab

Rafael Irizarry

Chinmay Shukla

## Collaborators

Sutirtha Chakraborty

Yuval Benjamini