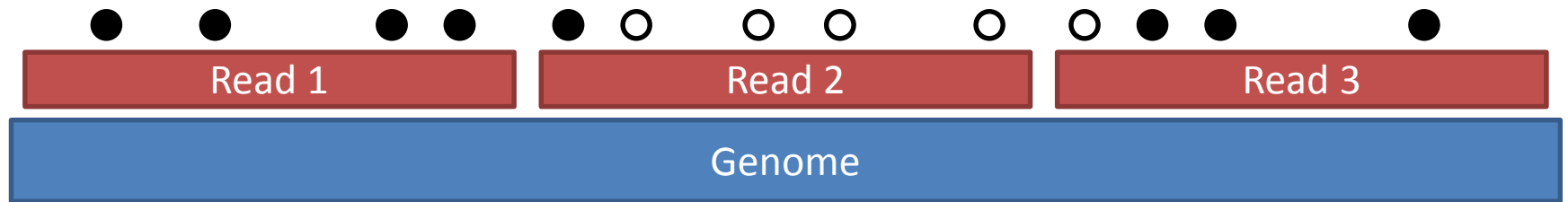


Visualising and Exploring BS-Seq Data

Simon Andrews
simon.andrews@babraham.ac.uk
@simon_Andrews

V2017-01

Starting Data



`L001_bismark_bt2_pe.deduplicated.bam`

`CHG_OB_L001_bismark_bt2_pe.deduplicated.txt.gz`

`CHG_OT_L001_bismark_bt2_pe.deduplicated.txt.gz`

`CHH_OB_L001_bismark_bt2_pe.deduplicated.txt.gz`

`CHH_OT_L001_bismark_bt2_pe.deduplicated.txt.gz`

`CpG_OB_L001_bismark_bt2_pe.deduplicated.txt.gz`

`CpG_OT_L001_bismark_bt2_pe.deduplicated.txt.gz`

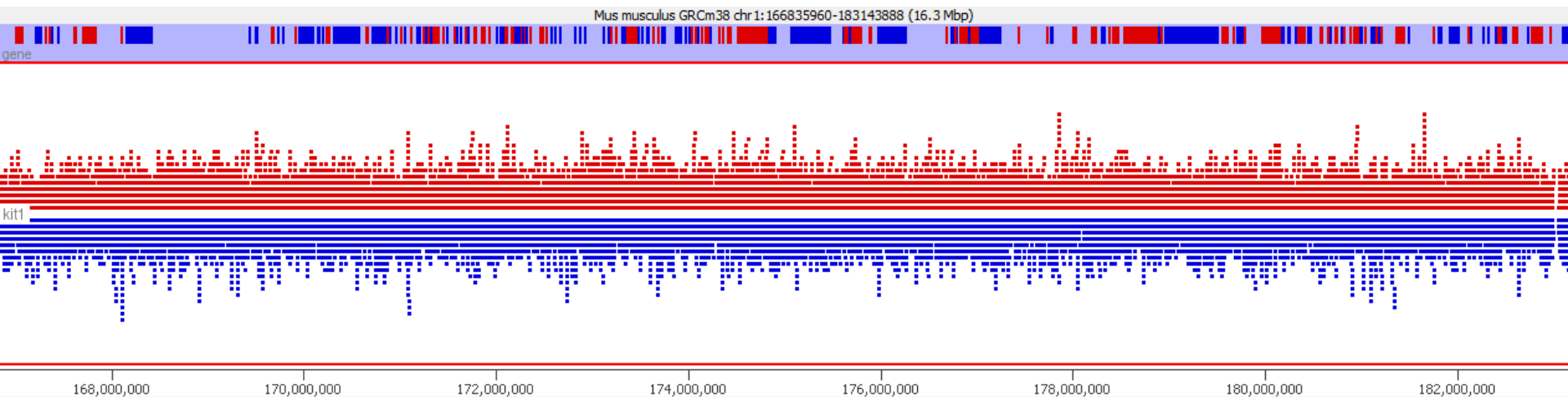
`L001_bismark_bt2_pe.deduplicated.cov.gz`

Which data to use?

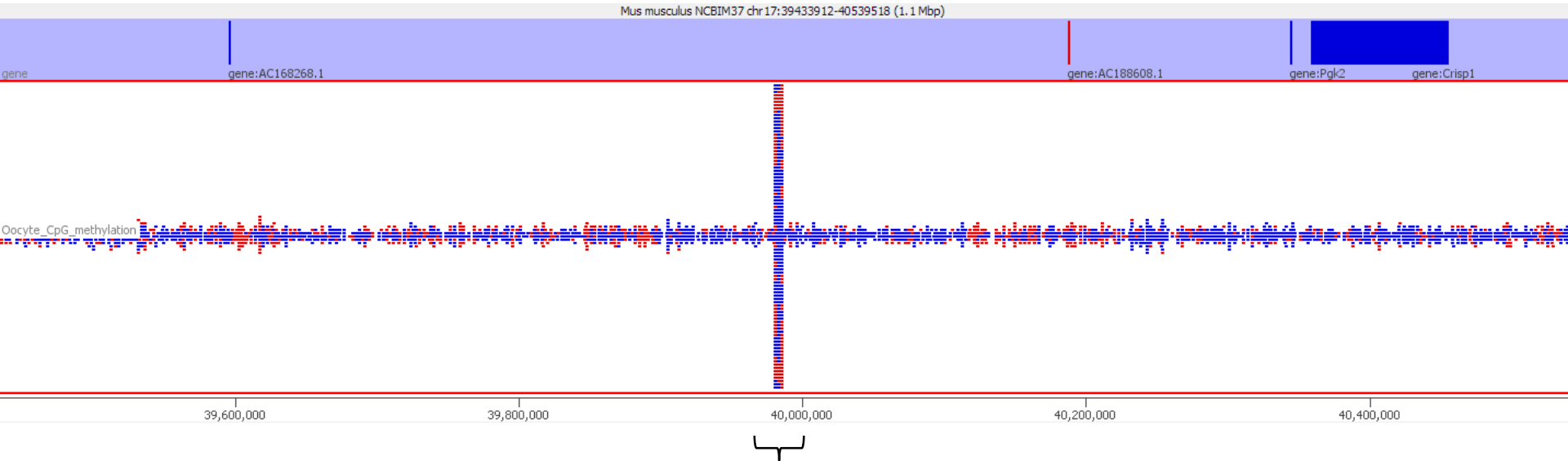
- Methylation contexts
 - CpG: Only generally relevant context for mammals
 - CHG: Only known to be relevant in plants
 - CHH: Generally unmethylated
- Methylation strands
 - CpG methylation is generally symmetric
 - Normally makes sense to merge OT / OB strands

Coverage Expectations

- Whole genome covered equally
- Both strands present equally



Coverage Outliers



Around 600x average genome density

Coverage Outliers



Coverage Outliers

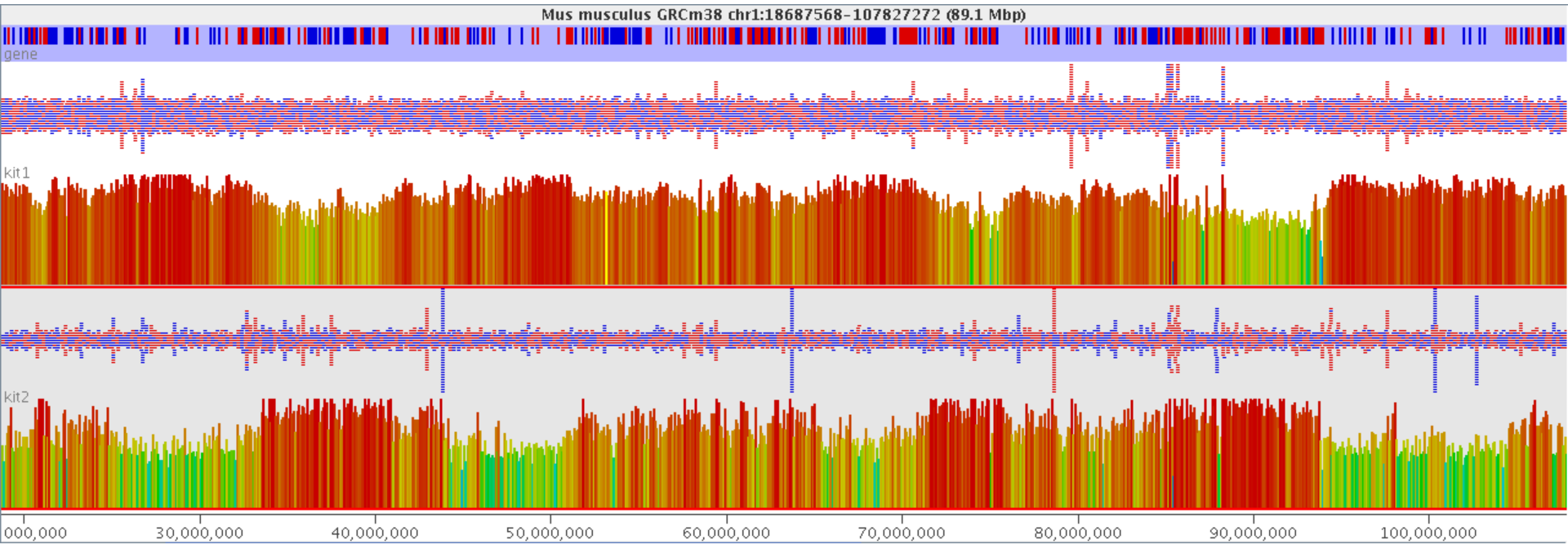
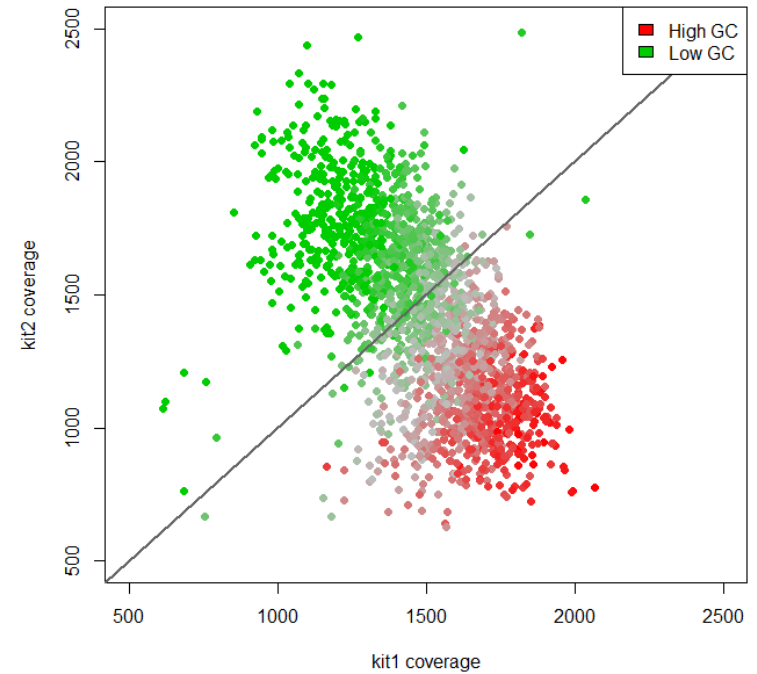
- Normally the result of mis-mapping repetitive sequences not in the genome assembly
- Centromeric / telomeric sequences are common
- Can be a significant proportion of all data
- Can throw off calculations of overall methylation
- Should be removed

Coverage Biases

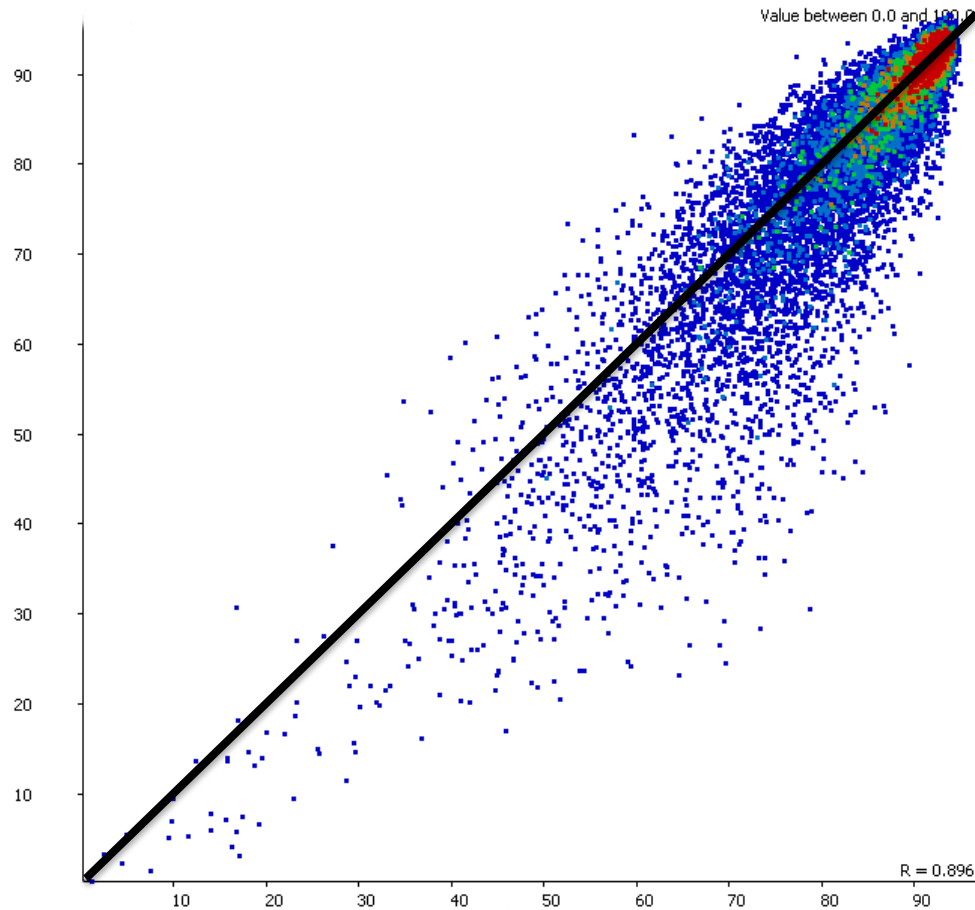
- BS-Seq read level coverage should be flat
- Different kits have different biases
- Most bias is based on composition (normally GC)
- Link between methylation and GC content can cause problems.

Coverage Bias

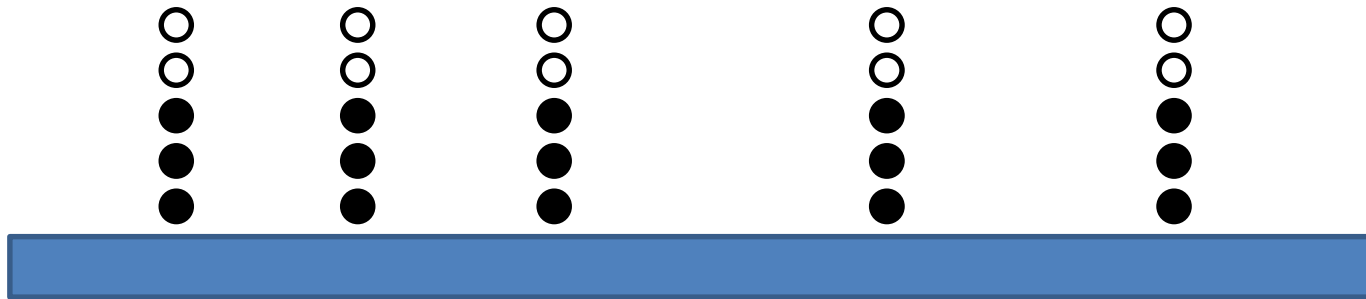
GC Content is most likely but others could exist



Coverage bias can lead to Methylation bias



Quantitating Methylation

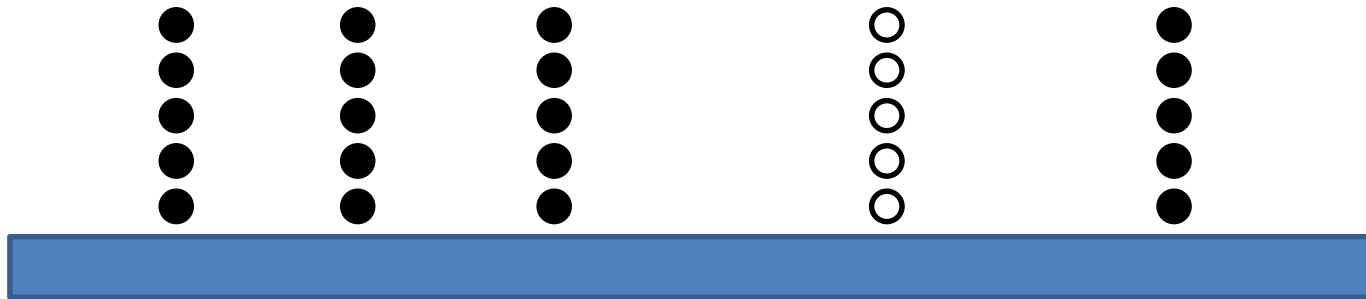


Total methylated calls = 15

Total unmethylated calls = 10

Methylation level = $(15/(15+10)) * 100 = 60\%$

Quantitating Methylation

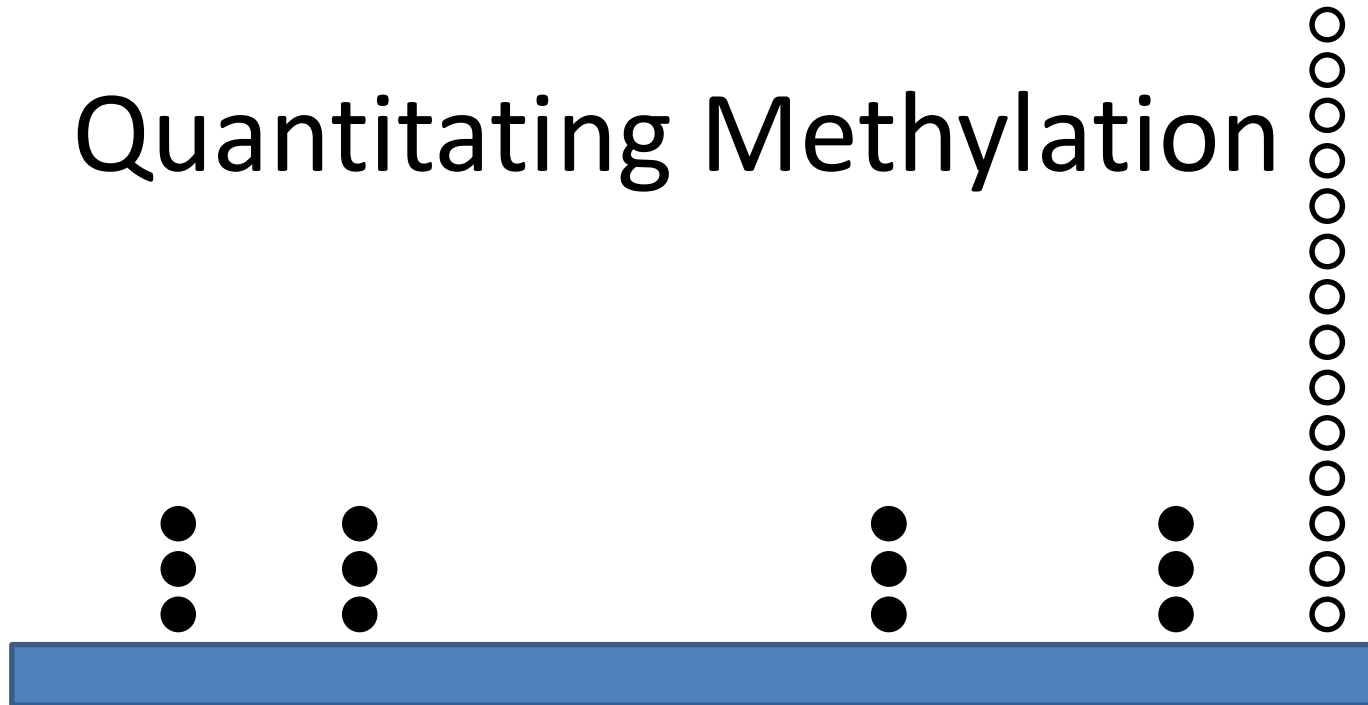


Total methylated calls = 20

Total unmethylated calls = 5

Methylation level = $(20/(20+5))*100 = 80\%$

Quantitating Methylation

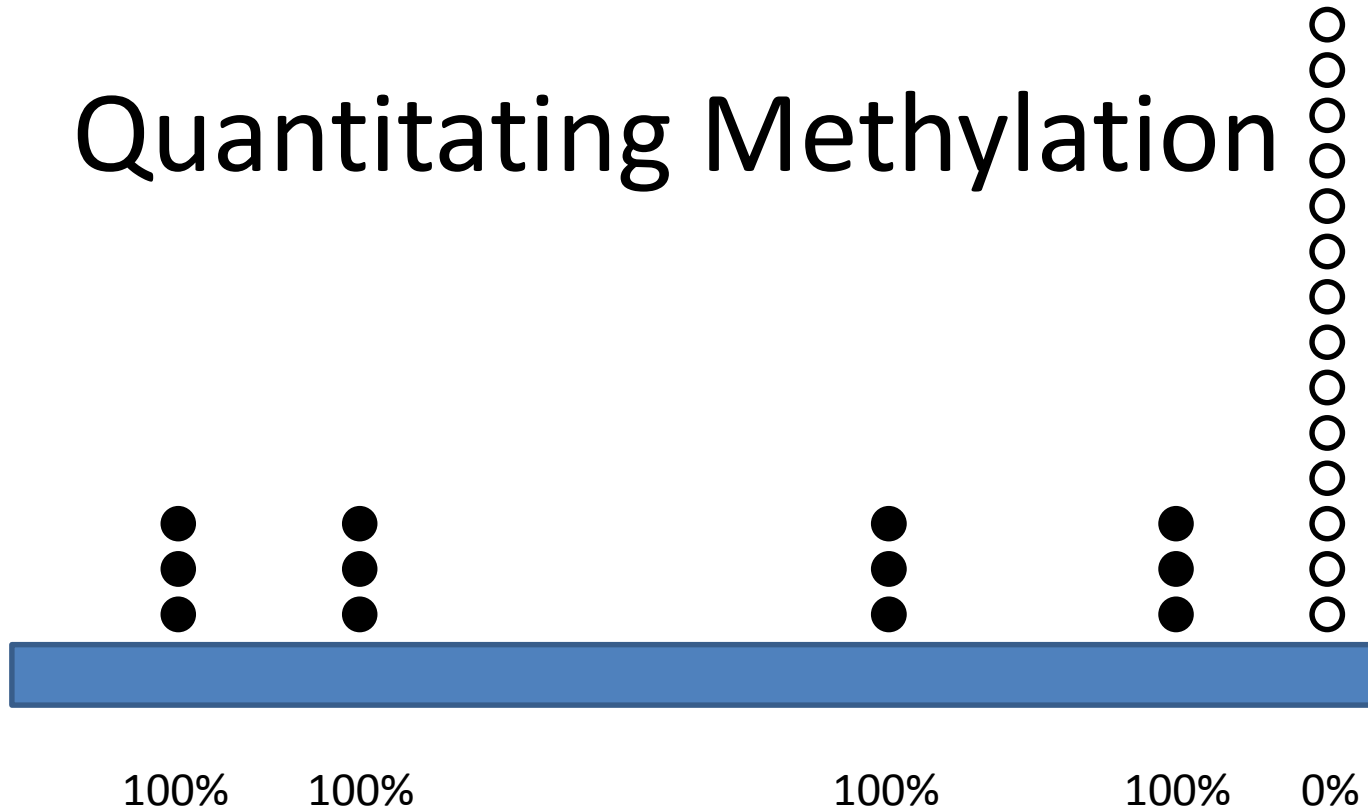


Total methylated calls = 12

Total unmethylated calls = 14

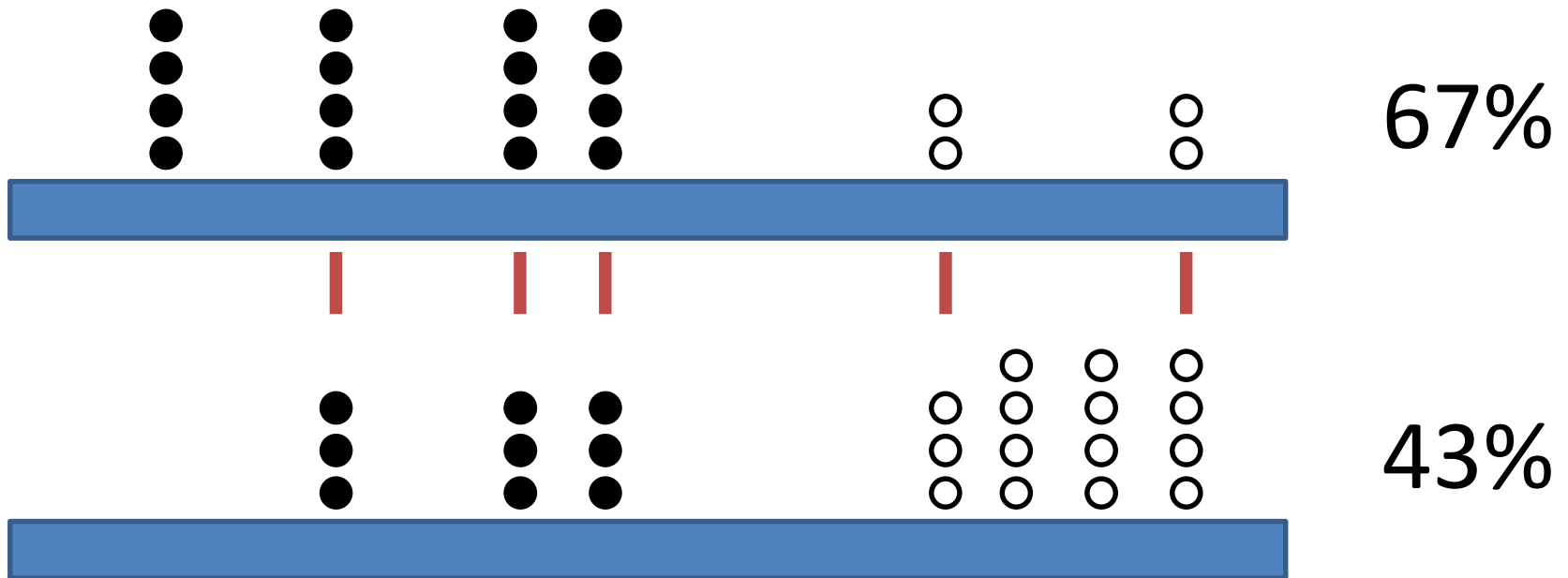
Methylation level = $(12/(12+14)) * 100 = 46\%$?

Quantitating Methylation



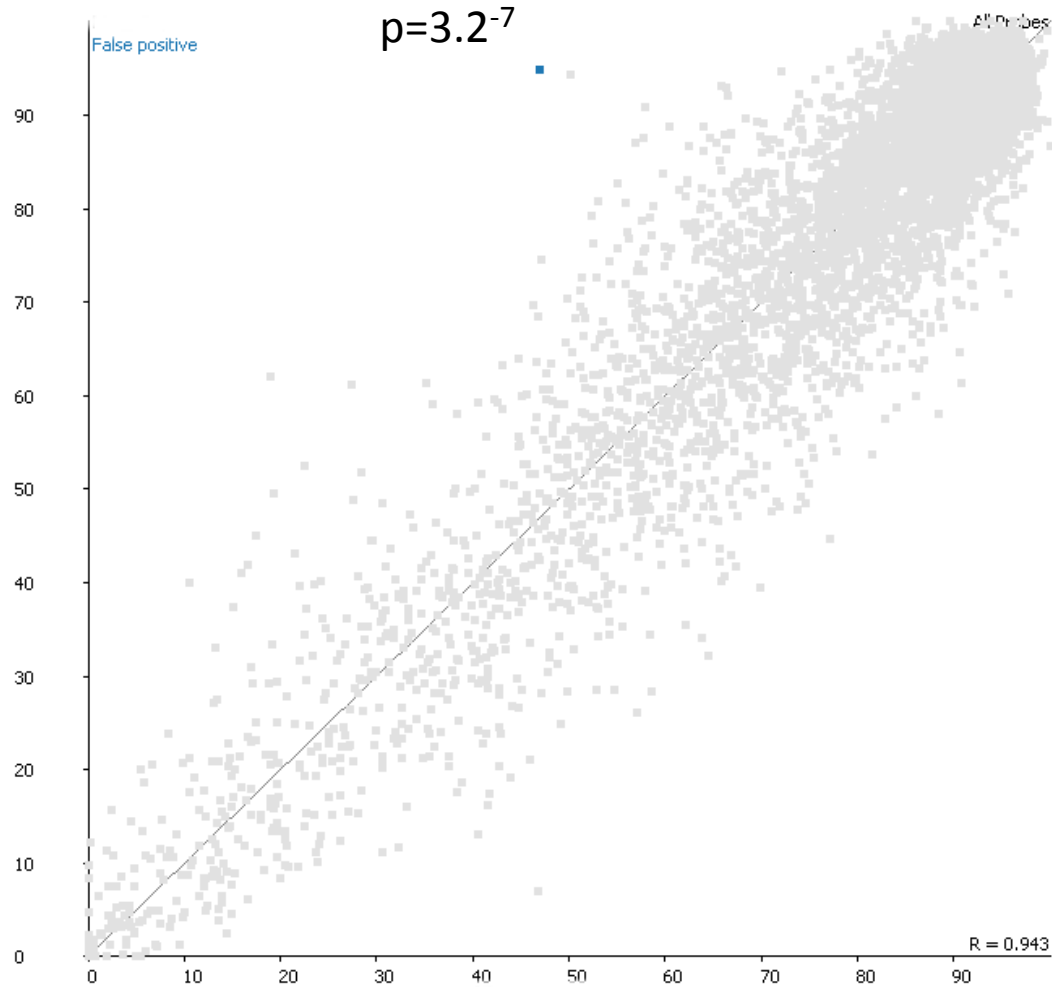
$$\text{Methylation level} = (100+100+100+100+0)/5 = 80\%$$

Quantitating Methylation

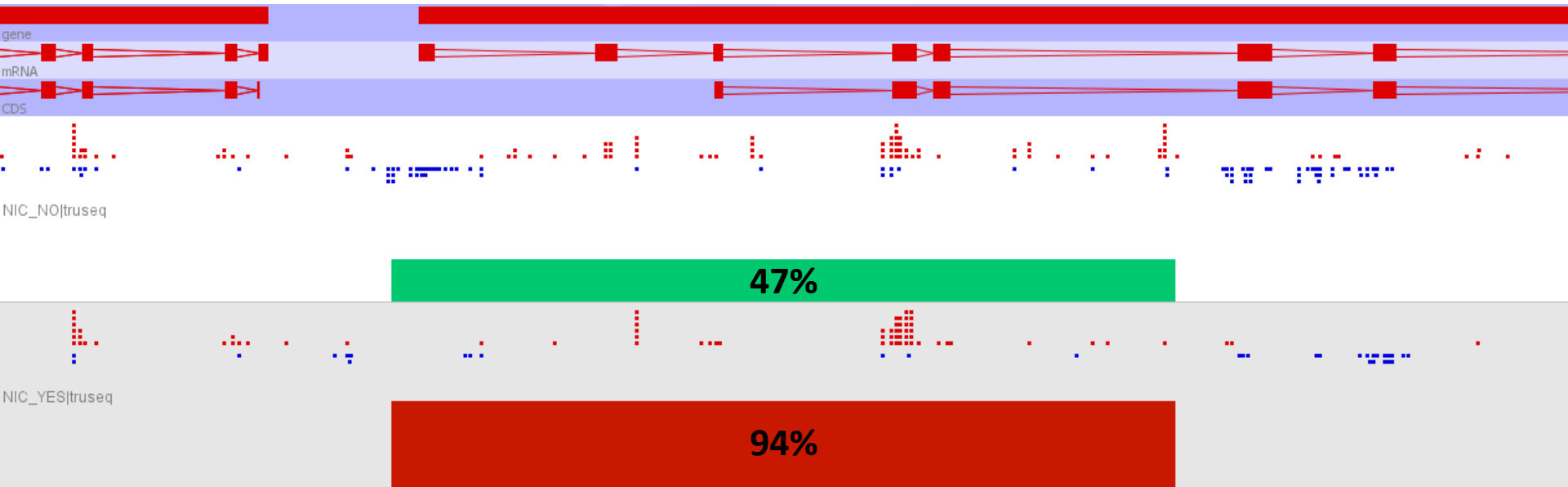


Common = 60% in both

Example



Example



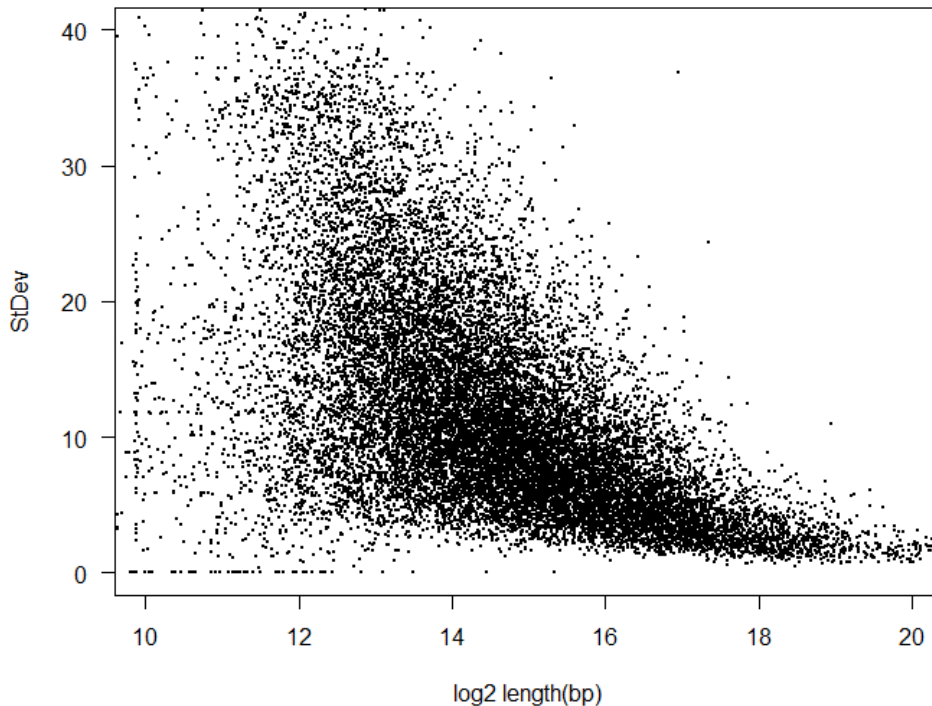
More Complex Methods

- Factors to consider during quantitation
 - Surrounding methylation levels
 - Assume that methylation doesn't change over very short distances
 - Coverage of individual bases
 - Down-weight very low or high values
 - Density of CpGs
 - Apply the level to all bases within a region

Where to make measures

- Per base
 - Very large number of measures
 - Poor accuracy for individual bases
- Unbiased windows
 - Tiled over whole genome
 - Need to decide how they will be defined
- Targeted regions
 - Which regions
 - What context

Making Comparable Measures

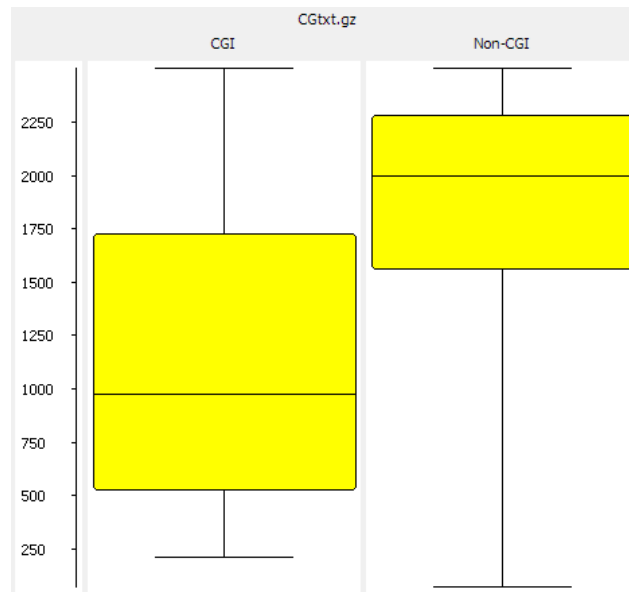
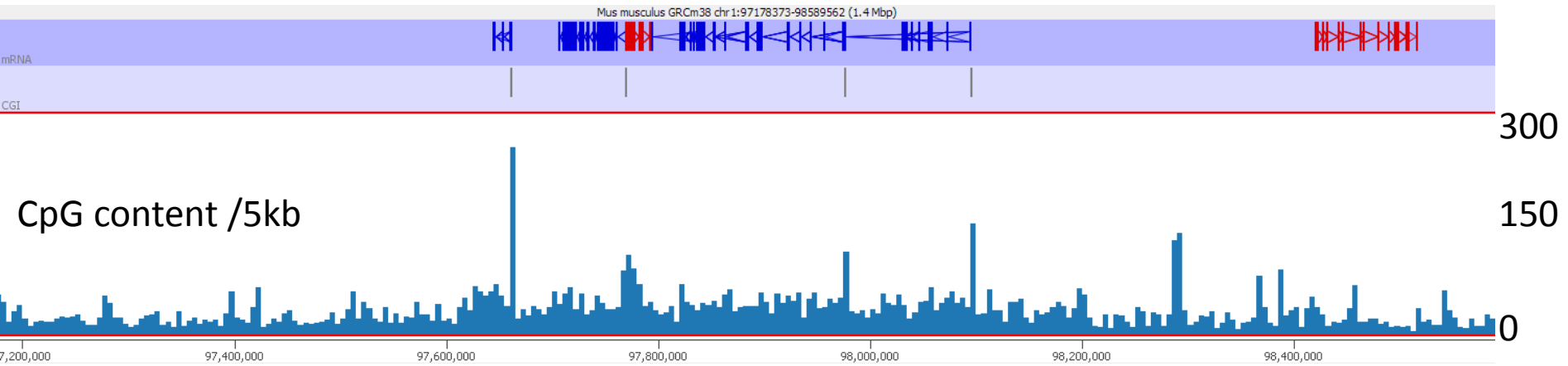


- Observation level correlates with stability.
- Want to try to have similar amounts of data in each measurement window.
- Equalises noise for visualisation and power for analysis.

Unbiased analysis

- What basis do you use for creating the windows?
 - Fixed size?
 - Uneven CpG distribution can be problematic
 - Fixed CpG count?
 - Different resolution
- Need to tailor the window sizes to the data density
- Can relate windows to features later

Fixed size vs Fixed CpG

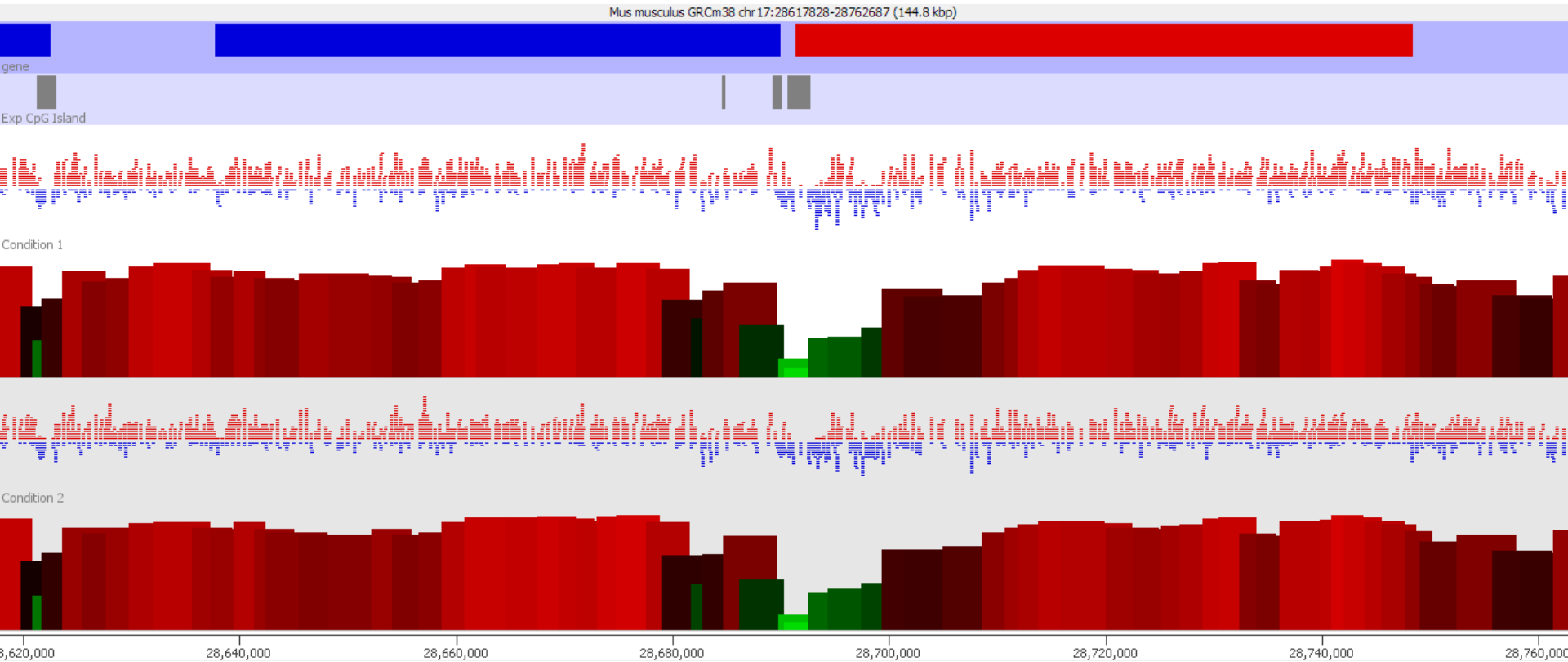


50 CpG window lengths

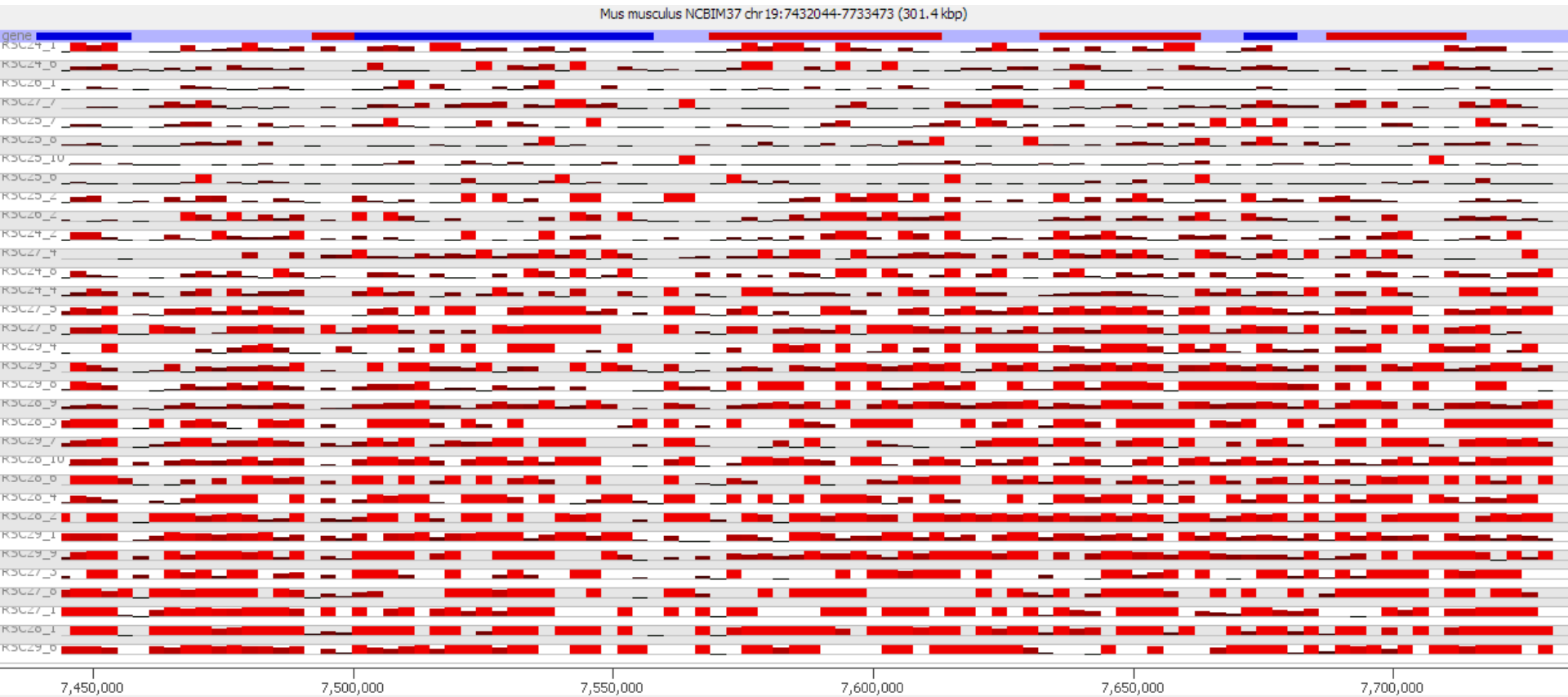
Targeted Quantitation

- Measure over features
 - CpG islands
 - Be careful where you get your locations
 - Try to fix sizes
 - Promoters
 - Should probably split into CpG island and non-CpG island
 - Try to fix sizes
 - Gene bodies
 - Filter by biotype to remove small RNA genes?

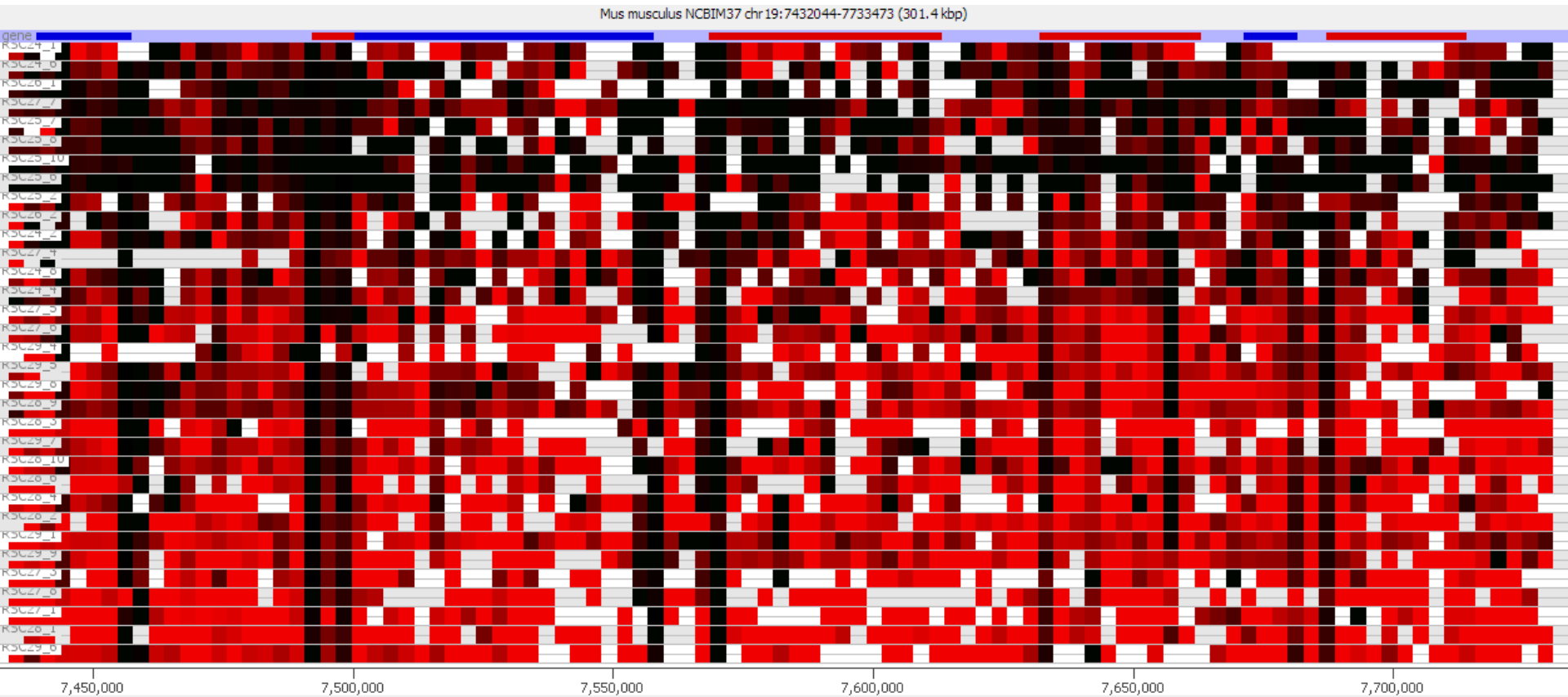
Viewing quantitated methylation



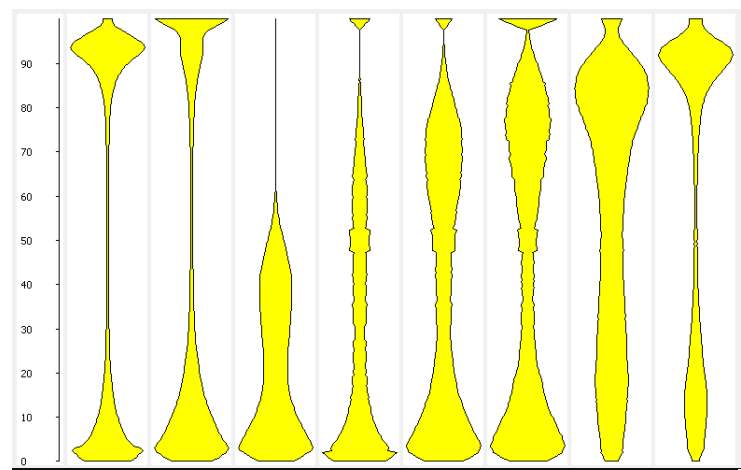
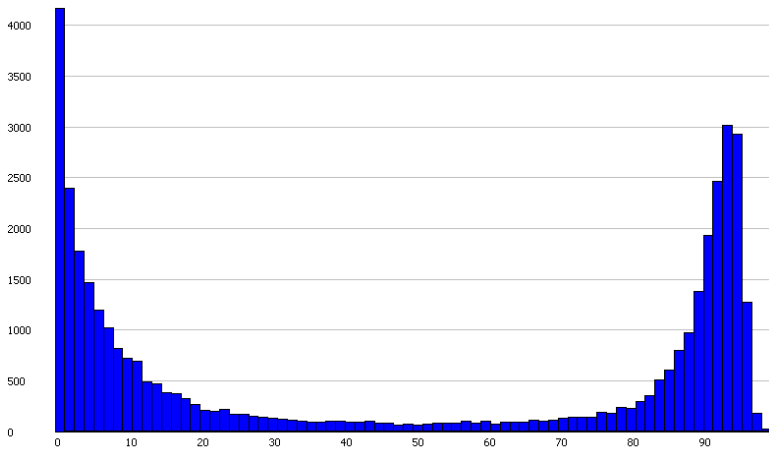
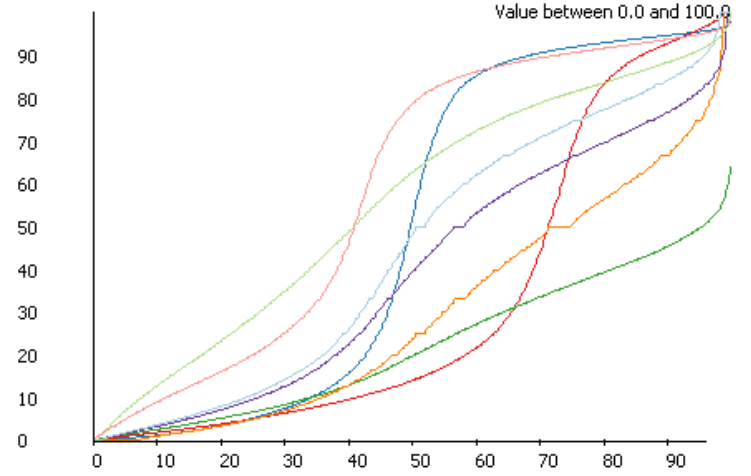
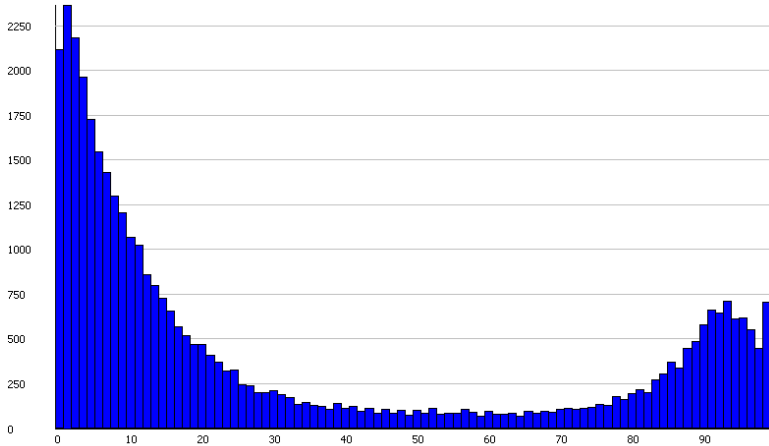
Large sample numbers



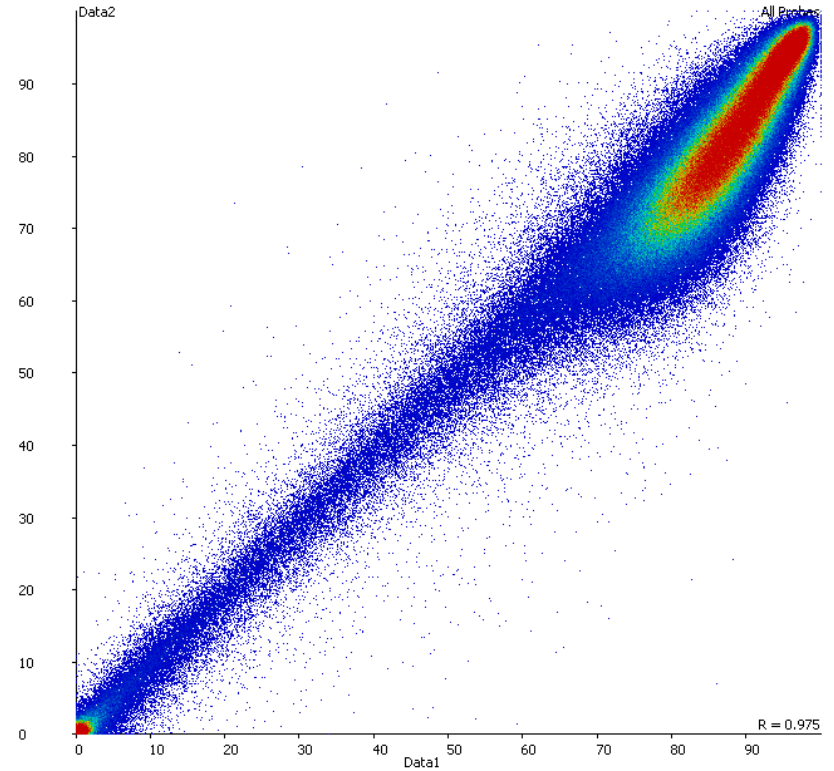
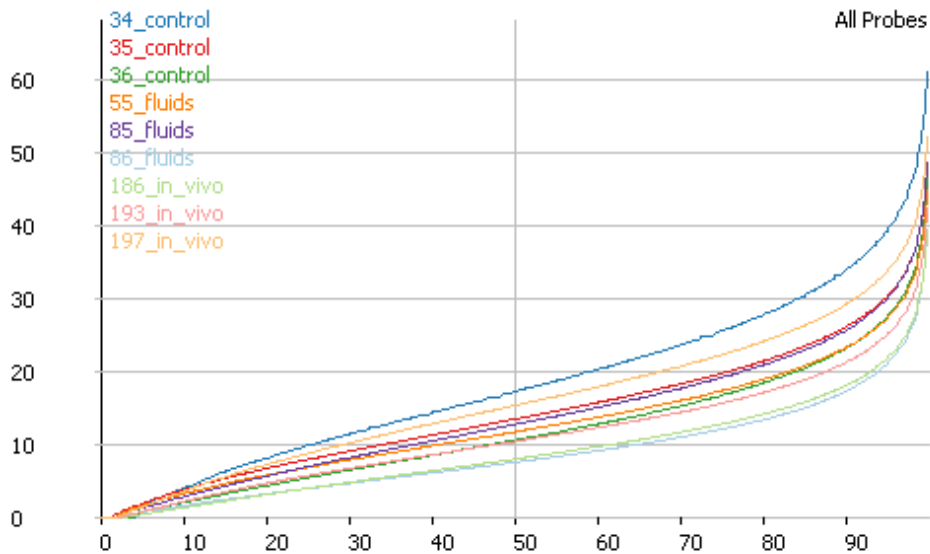
Large sample numbers



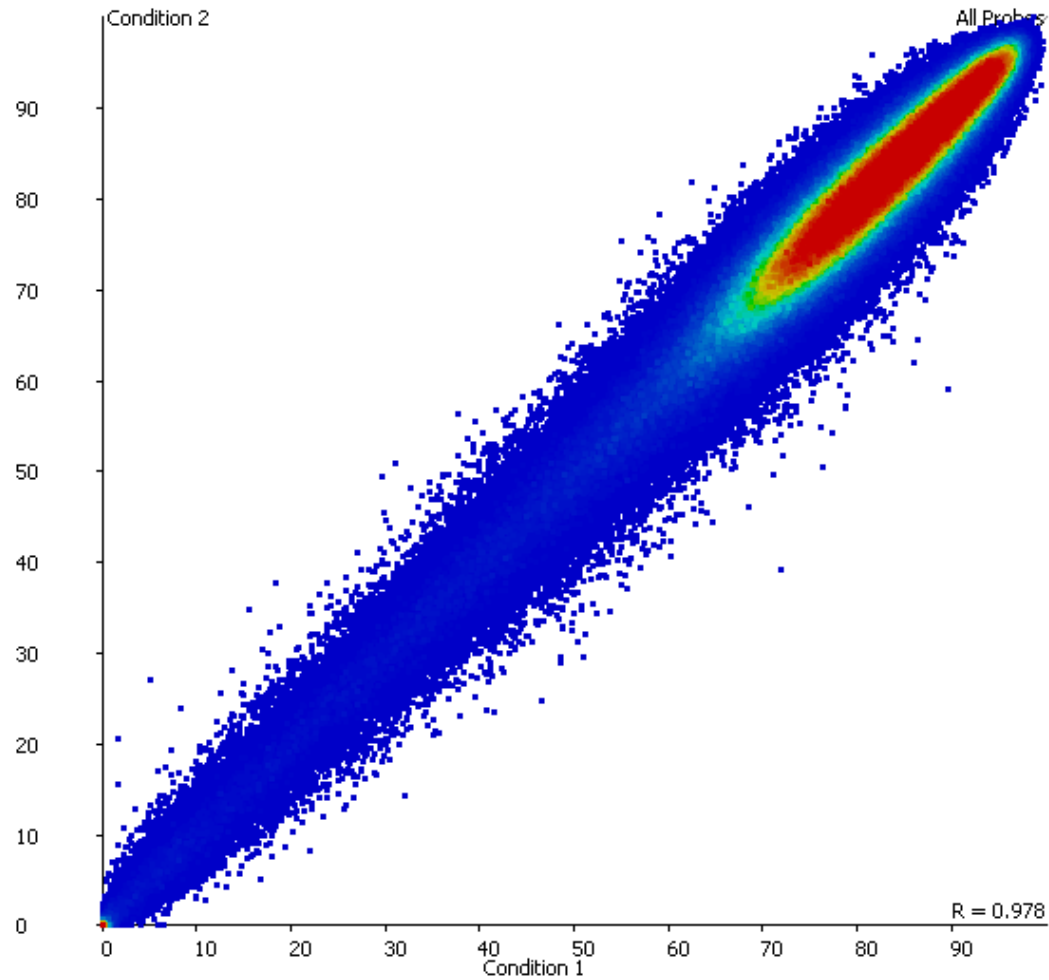
Viewing and Comparing Distributions



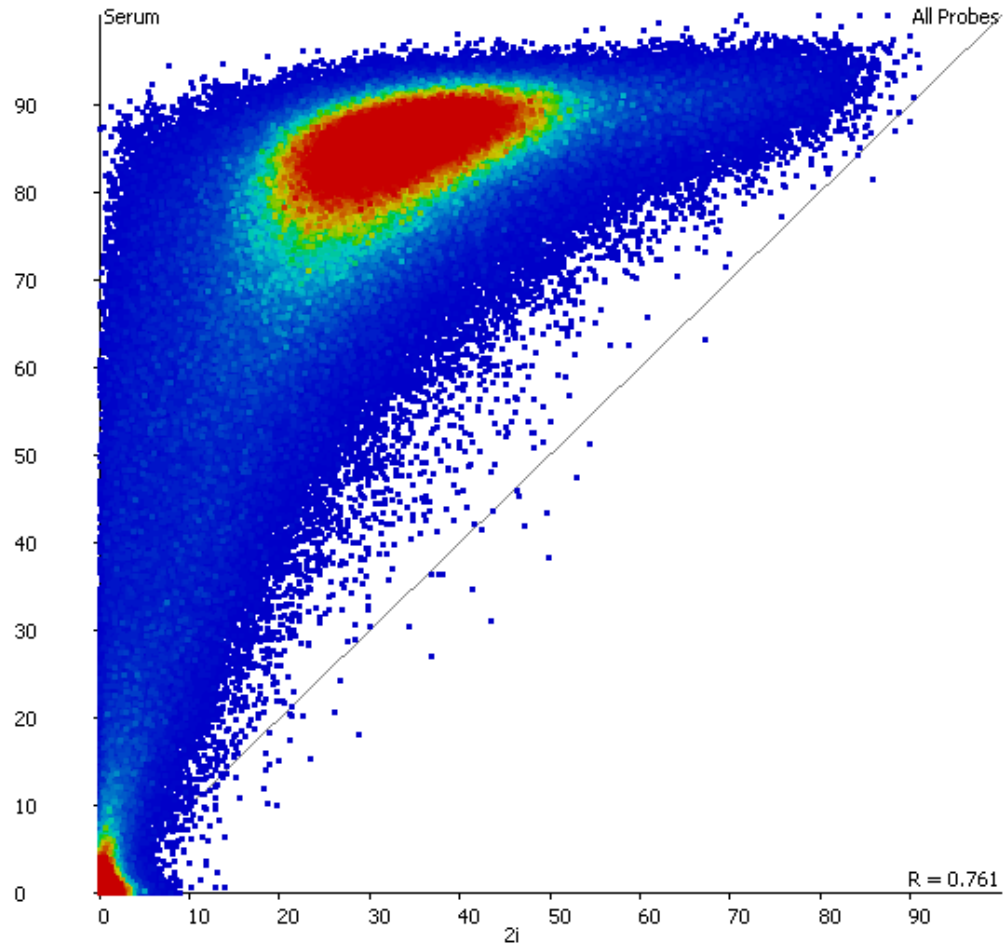
Normalising Methylation Values



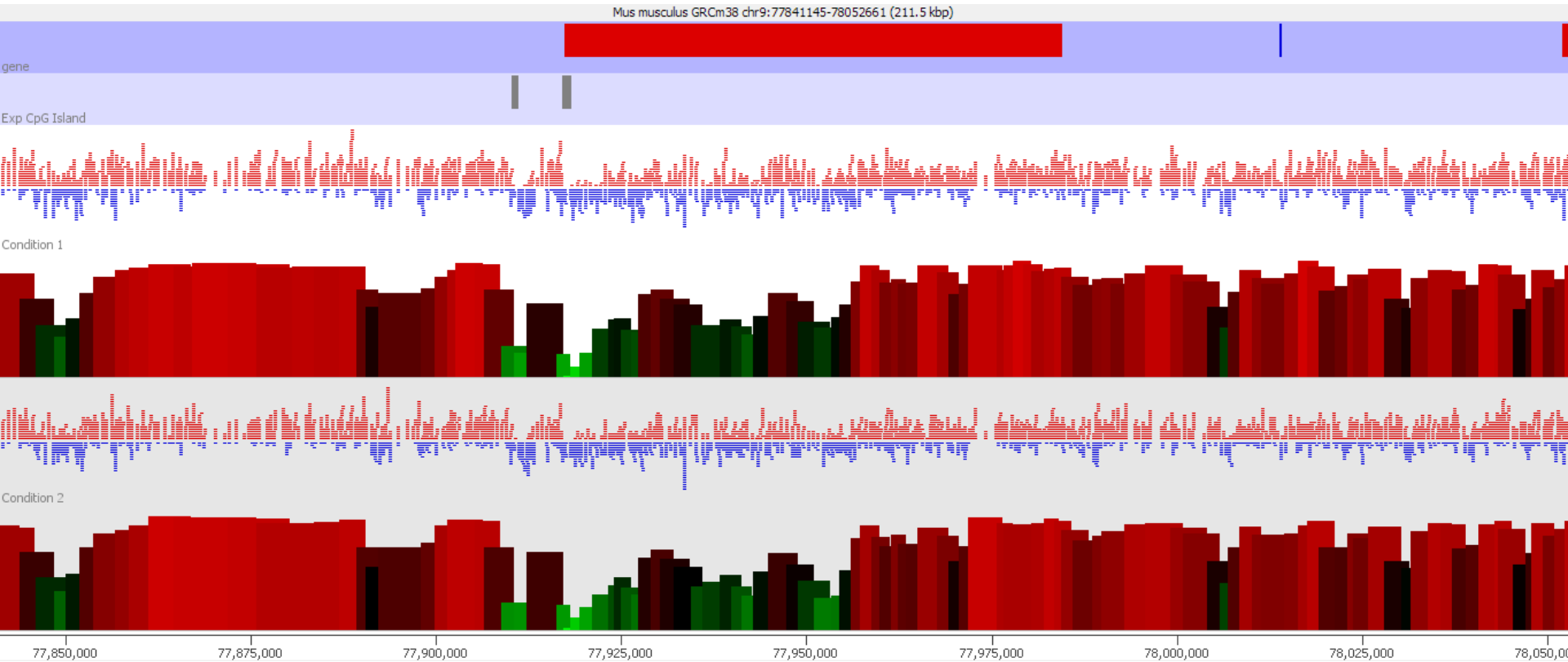
Viewing comparisons



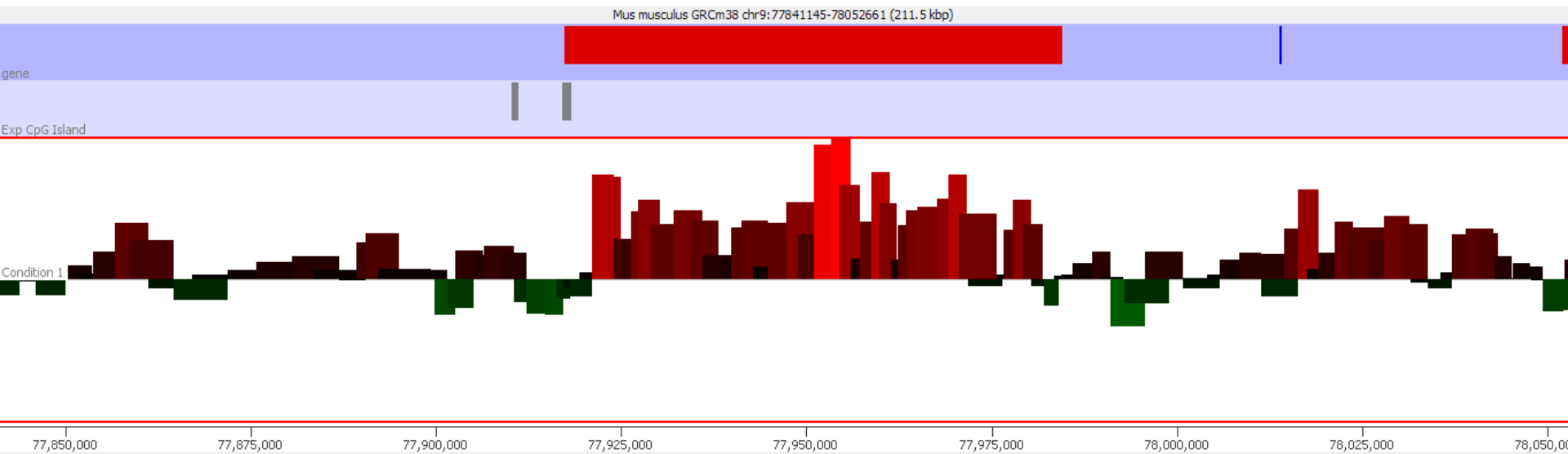
Viewing comparisons



Viewing differences



Viewing differences



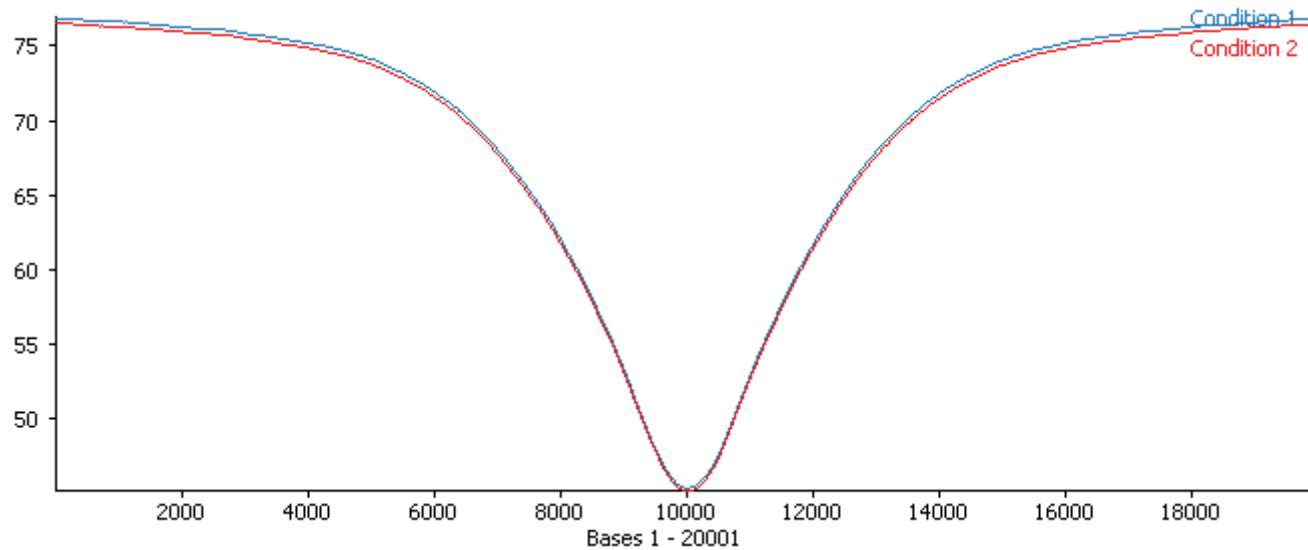
Trends

- Effects at individual loci can be subtle
- Want to find more generalised effect
- Collate information across whole genome
- Look for general trend

Considerations for trend plots

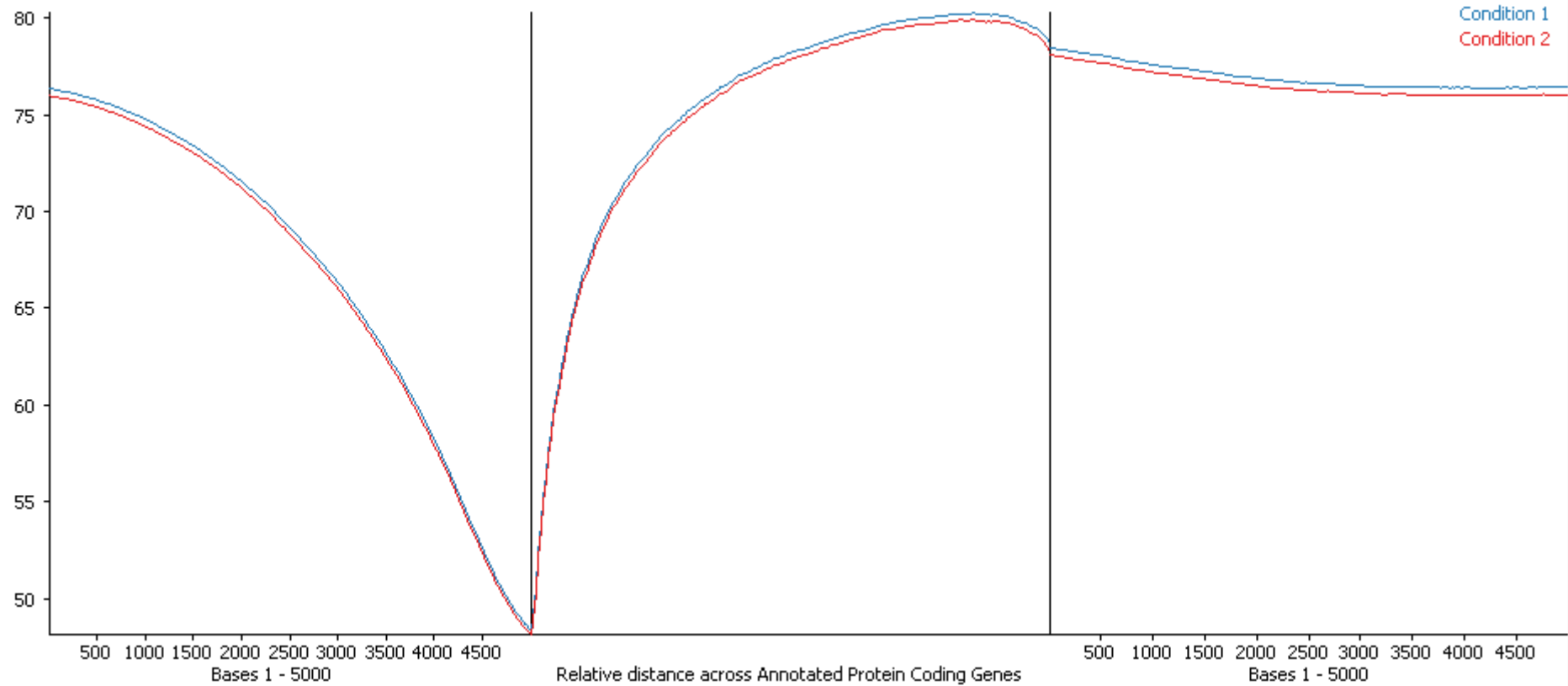
- What features to use
 - Fixed vs relative scale
- How much context
 - Variable scales
- How to calculate base measures
 - What window size
 - Aligned vs unaligned windows
- Missing values
- Scale

Simple Example



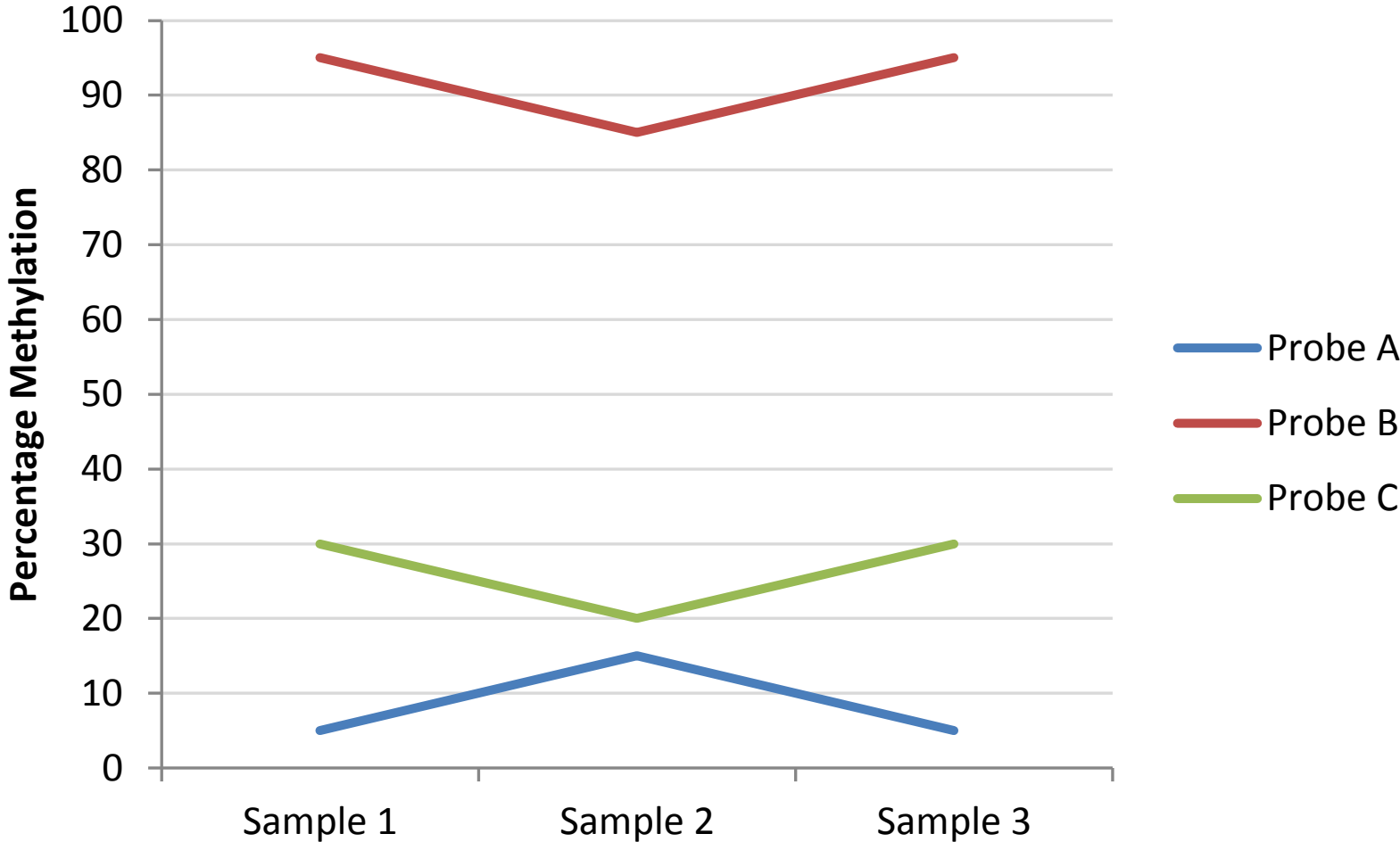
Methylation profile centred on CpG islands +/- 10kb

More complex example



Methylation profile over genes +/- 5kb

Clustering



Clustering

- Correlation Clustering
 - Focusses on the differences between conditions
 - Absolute values not important
 - Look for similar trends
 - Show median normalised values
- Euclidean Clustering
 - Focusses on absolute differences between conditions
 - Look for similar levels
 - Show raw values

Clustering

