

Exploring and Understanding ChIP-Seq data

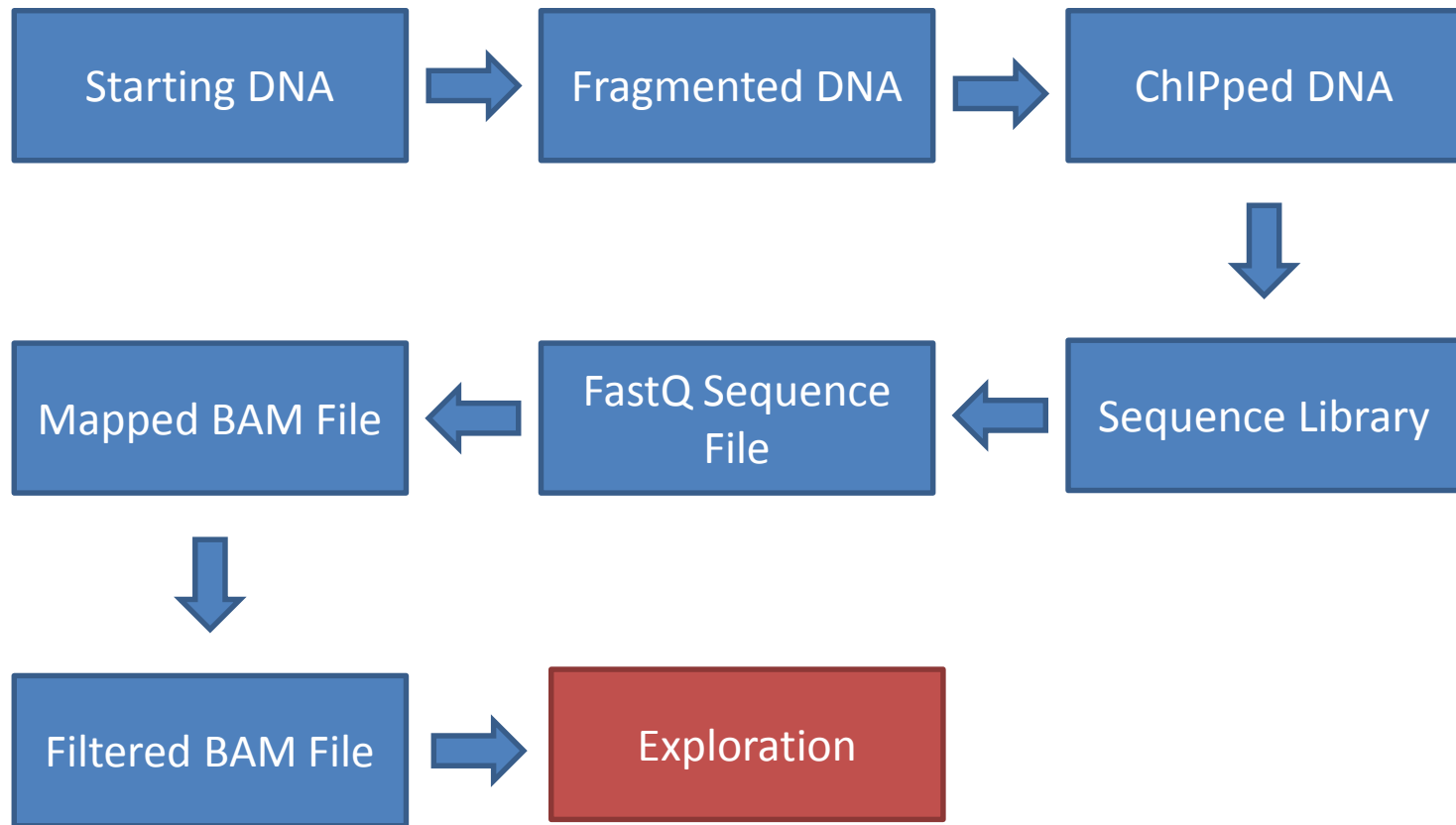
Simon Andrews

simon.andrews@babraham.ac.uk

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

v2018-03-02

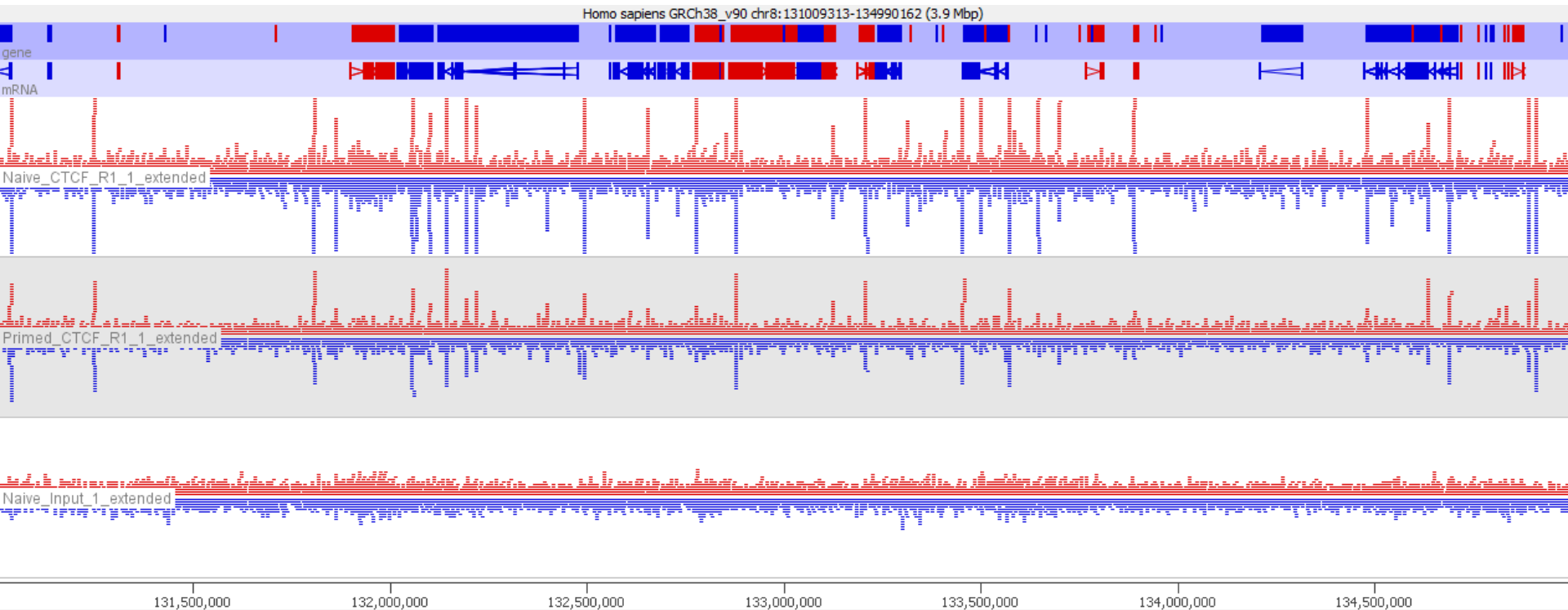
Data Creation and Processing



Some Basic Questions

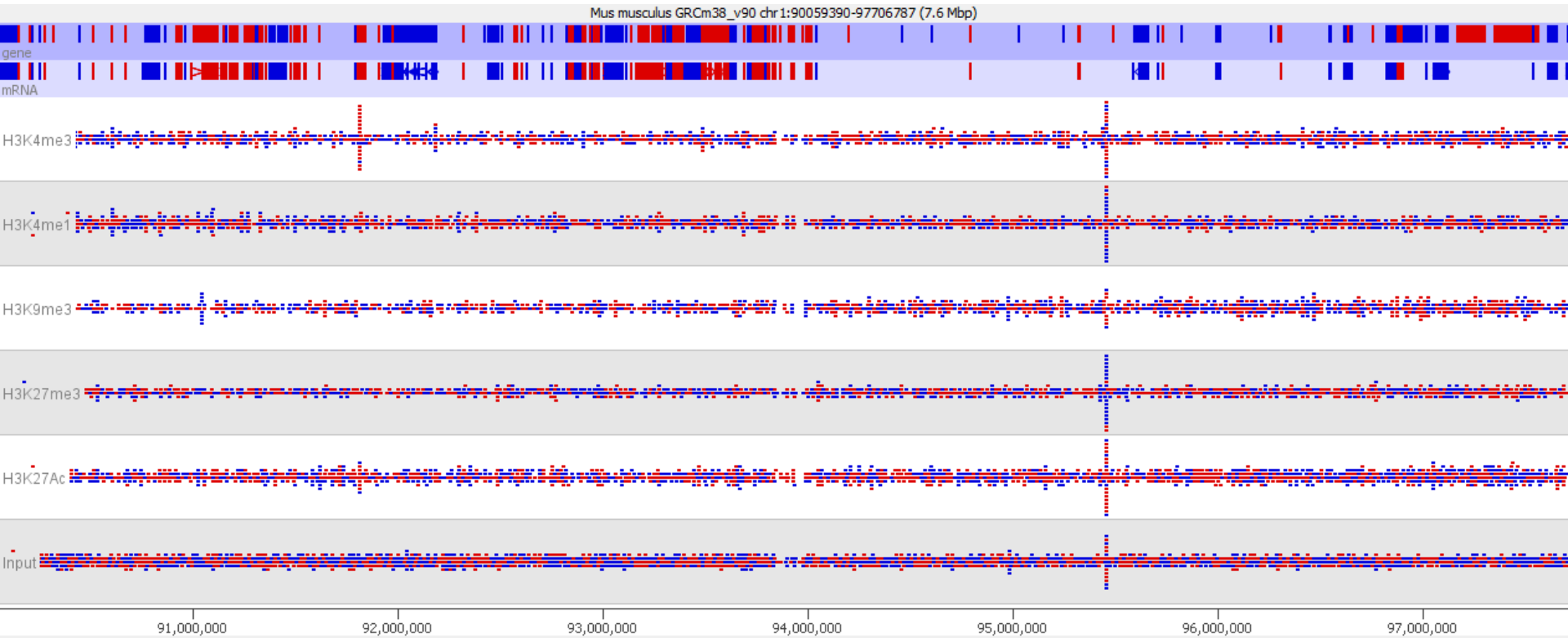
- Is there any enrichment?
- What is the size / patterning of enrichment?
- How well are my controls behaving?
- What is the best way to quantitate this data?
- Are there any technical artefacts?

Start with a visual inspection



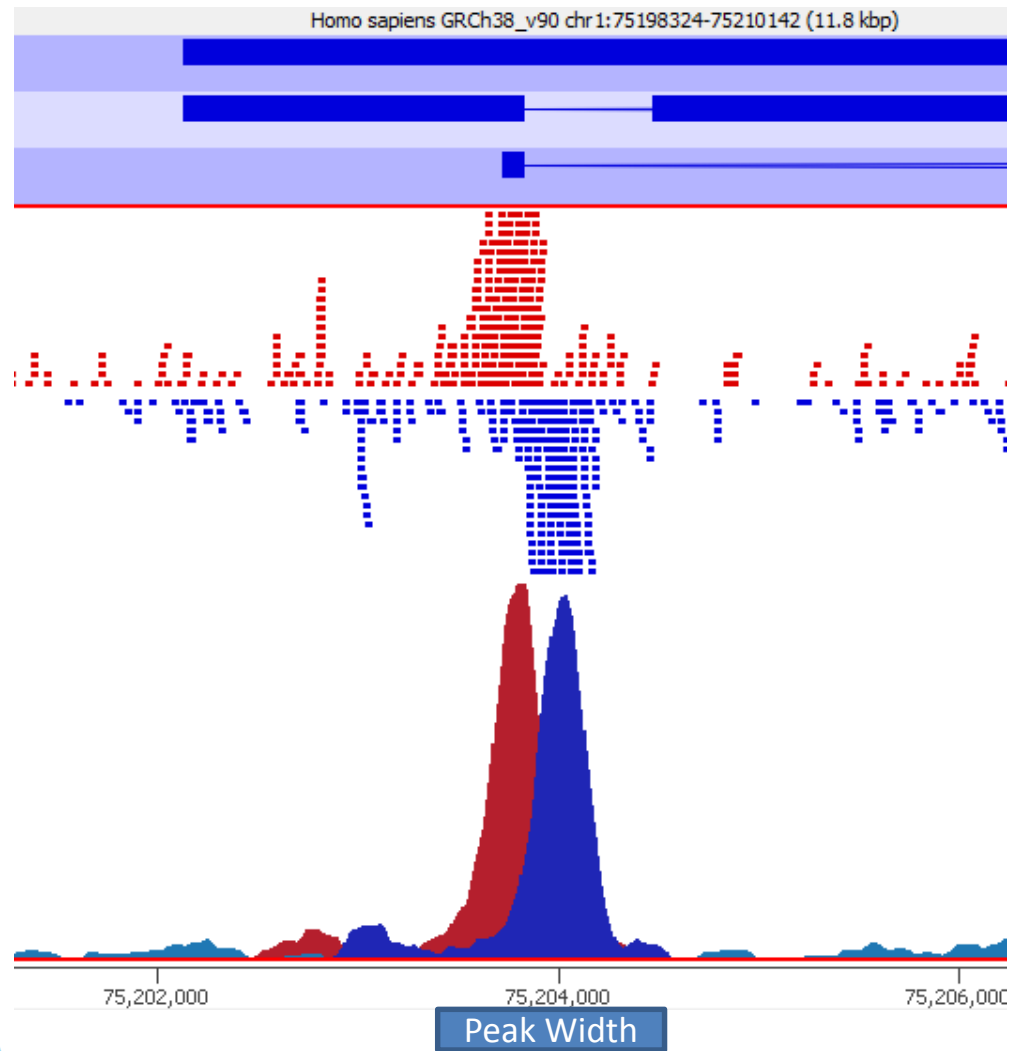
- Is there any enrichment?
- What is the size / patterning of enrichment?
- How well are my controls behaving?

Start with a visual inspection



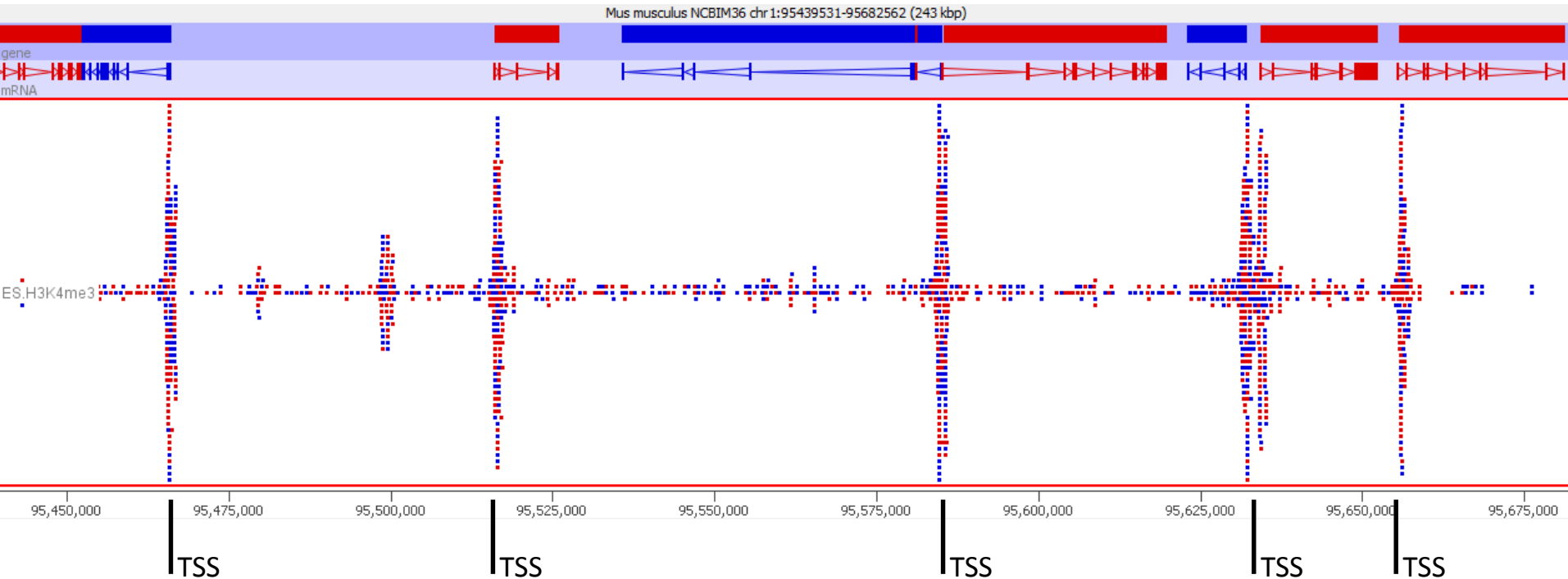
- Is there any enrichment?
- What is the size / patterning of enrichment?
- How well are my controls behaving?

Extending reads if necessary



Look for peaks

Associate with features

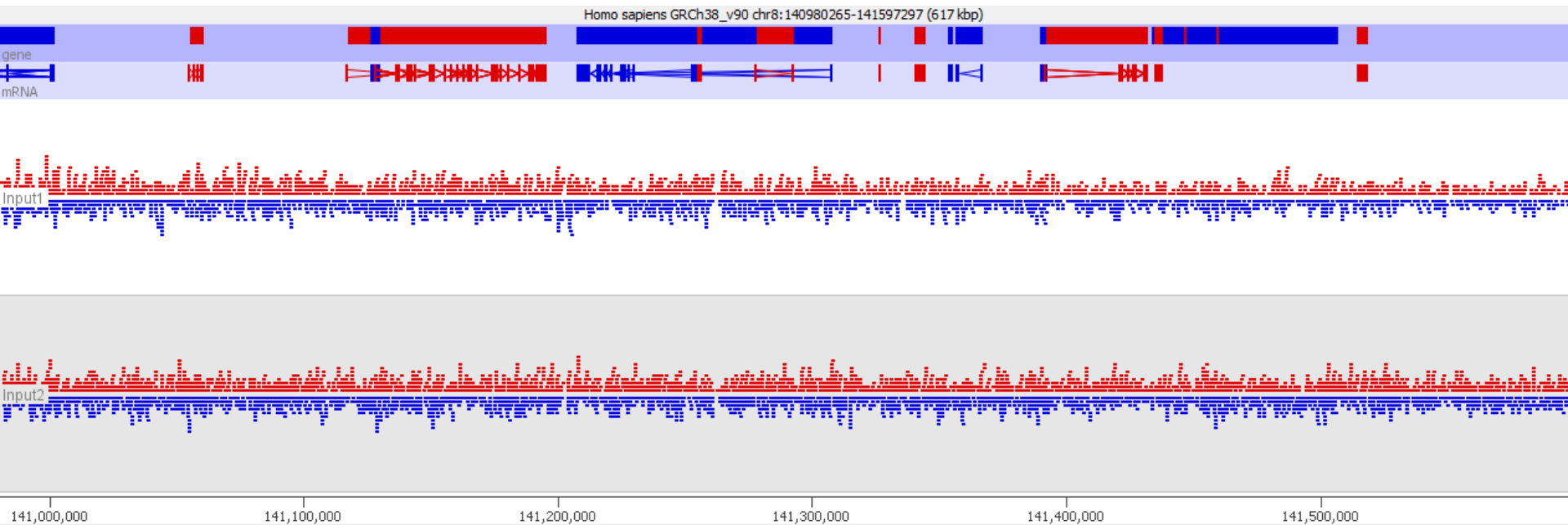


- Are my peaks narrow or broad
- Do peak positions obviously correspond to existing features?

Examine Controls

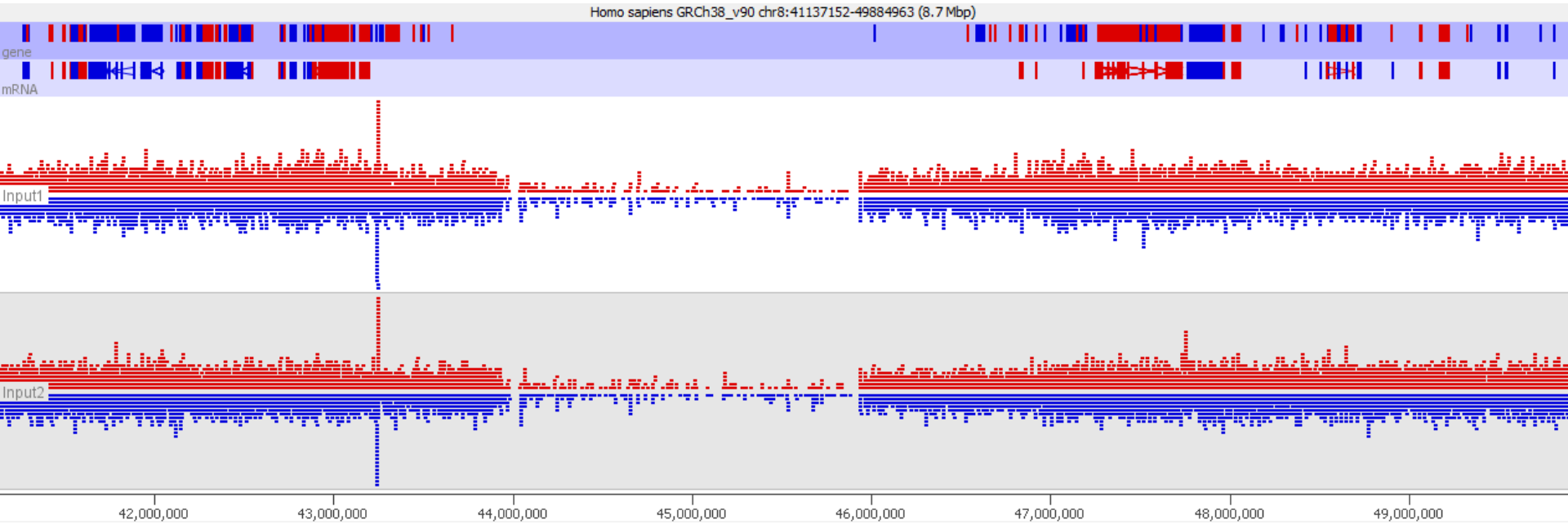
- IgG or other Mock IP
 - Good result is no material at all
 - Not worth sequencing. Reads are only informative if the ChIP hasn't worked.
- Input material (sonicated / Mnase etc)
 - Genomic library - everywhere equally
 - Technical issues can cause variation

Examine Controls



- Does the coverage look even
- If there are multiple inputs to do they look similar

Examine Controls



Why do controls misbehave?

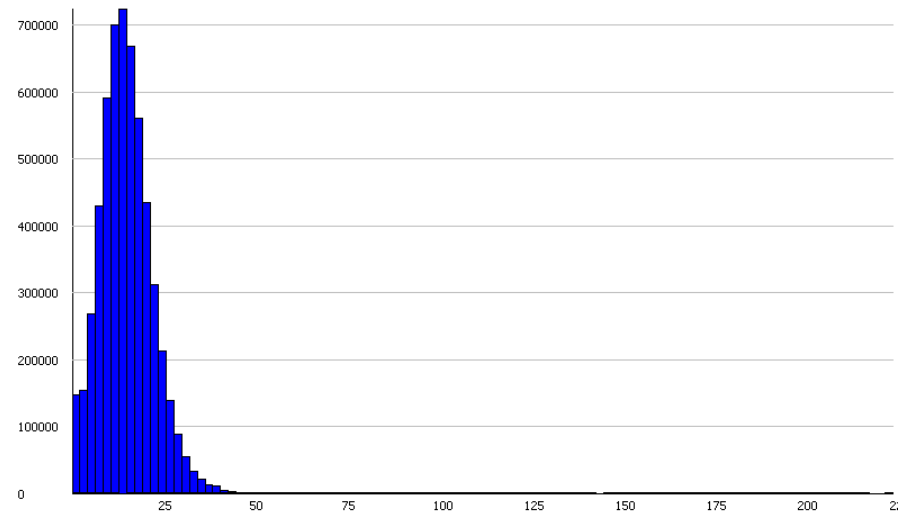
- Low coverage
 - Repetitive unmappable regions
 - Holes in the assembly
- High coverage
 - Mismapped reads from outside the assembly
- Biases
 - DNA stability
 - GC content
 - Segmental Duplication

Types of input problem

- Categorical
 - Mismapped reads
 - Indicates that the region can't be trusted
 - Blacklist and remove - don't try to correct
- Quantitative
 - Other biases (most often GC)
 - Some potential to correct, but difficult
 - Hopefully consistent, so will cancel out

Making Blacklists

- Unusual Coverage
 - Outlier detection (boxplots etc.)
 - Often only filter over-representation (maybe also zero counts)
- Pre-built lists
 - Large projects often build these (ENCODE)
 - Not for all species



Pre-Existing Blacklists

- ENCODE
- modENCODE
- UCSC



Impact of artifact removal on ChIP quality metrics in ChIP-seq and ChIP-exo data

Thomas S. Carroll^{1*}, Ziwei Liang^{2†}, Rafik Salama^{1†}, Rory Stark¹ and Ines de Santiago^{1*}

¹ Cambridge Institute CRUK, University of Cambridge, Cambridge, UK

² Lymphocyte Development, MRC Clinical Sciences Centre, Imperial College, London, UK

- Check assembly versions

Comparison of samples

Initial Quantitation

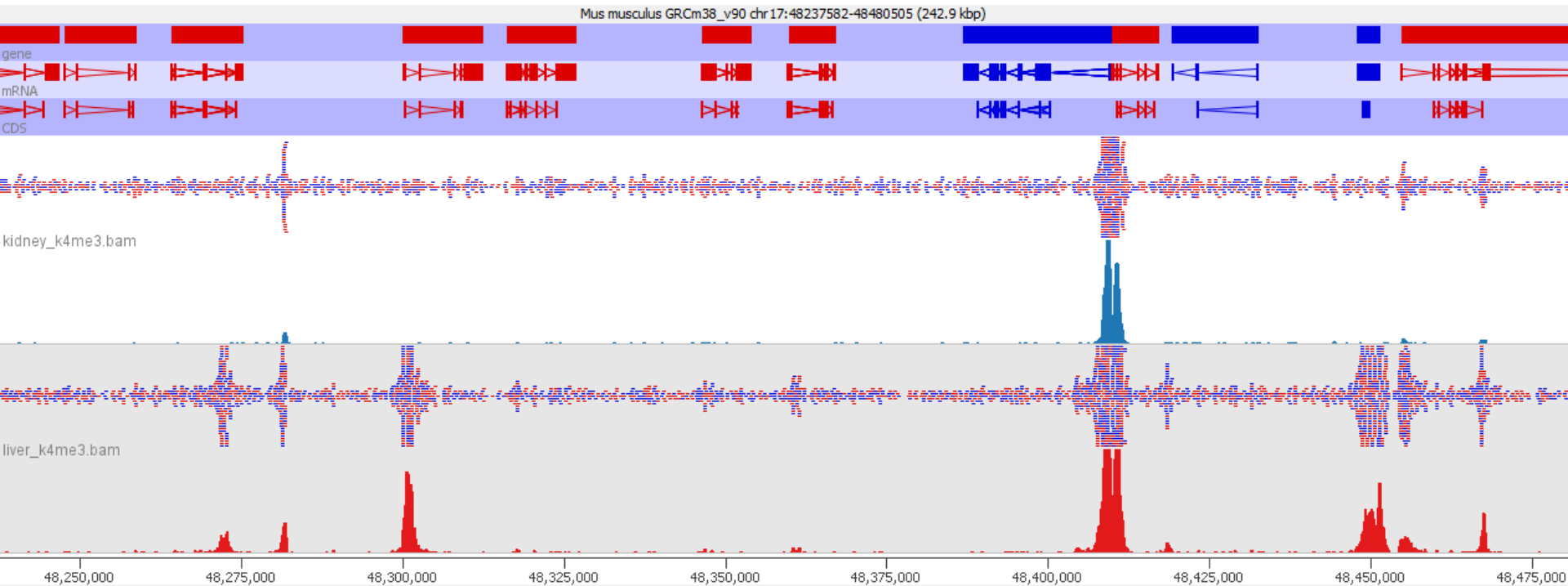
- Always start with a simple unbiased quantitation (not focussed on features/peaks)
- Tiled measures over the whole genome
 - Use approximate insert size as window size
 - Something around 500bp is normally sensible
- Linear read count quantitation corrected for total library size

Quantitation Questions

- Initially similar to what we did with the raw data, but with some values behind it
 - Can we observe enrichment in our ChIP above our controls
 - Do we see similar enrichment between replicates, or between conditions

Compare samples

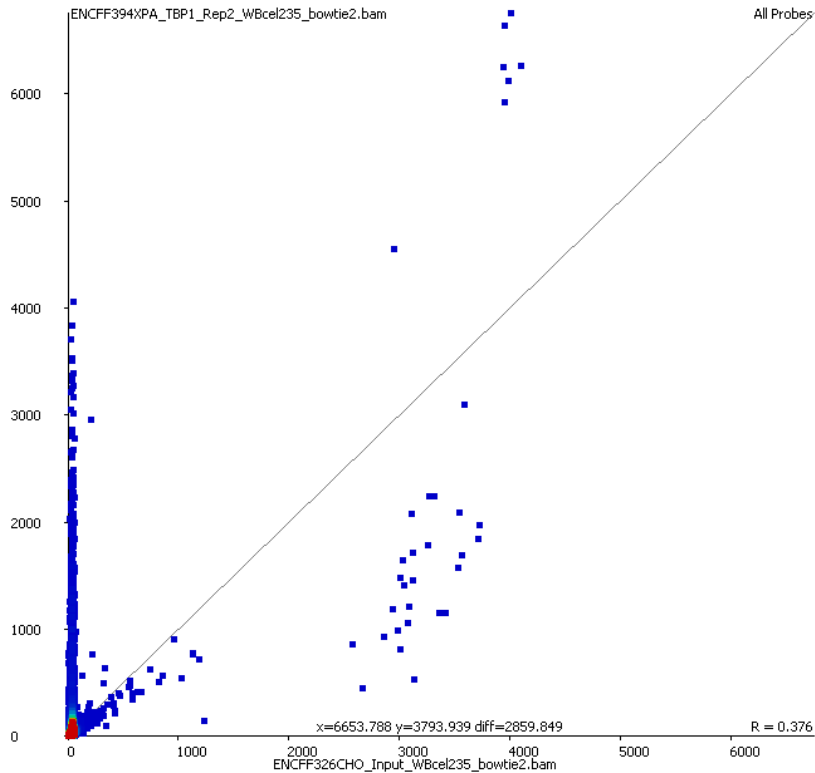
Visual comparison against raw data



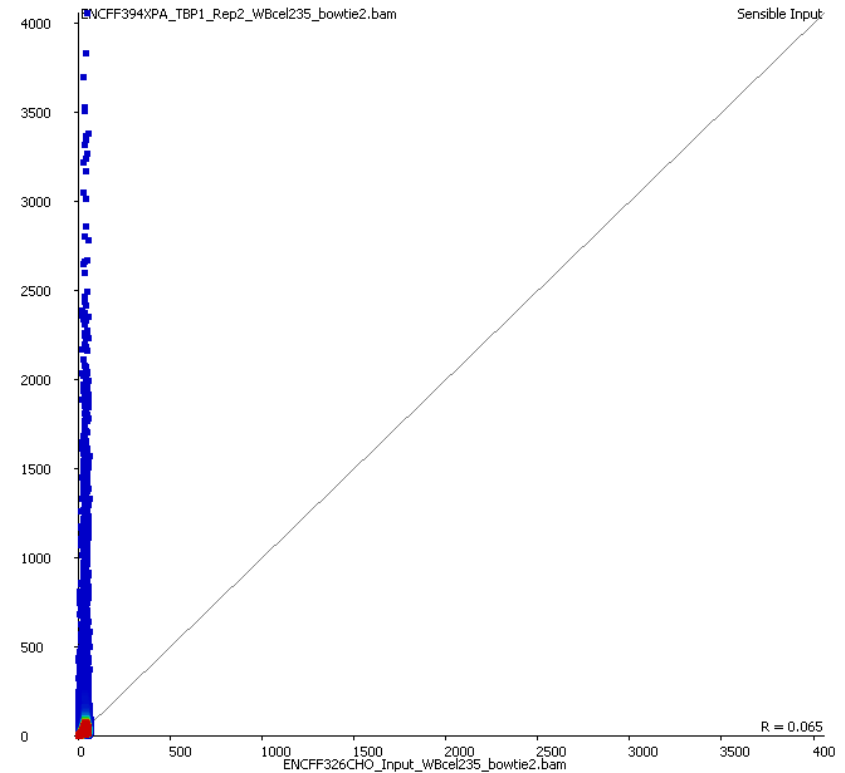
- Similar apparent overall enrichment
- Any obvious differences

Compare samples

Scatterplot input vs ChIP



Raw

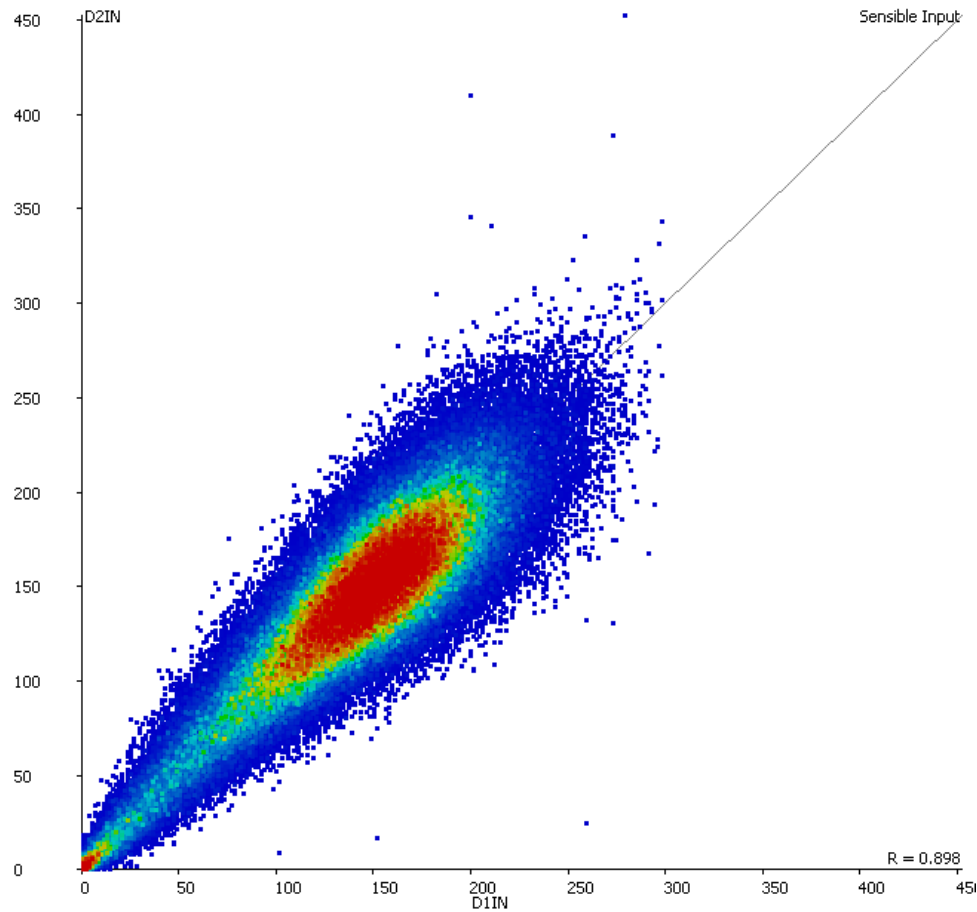


Filtered

- Any relationship between input and ChIP
- What proportion enriched
- How enriched over input

Compare samples

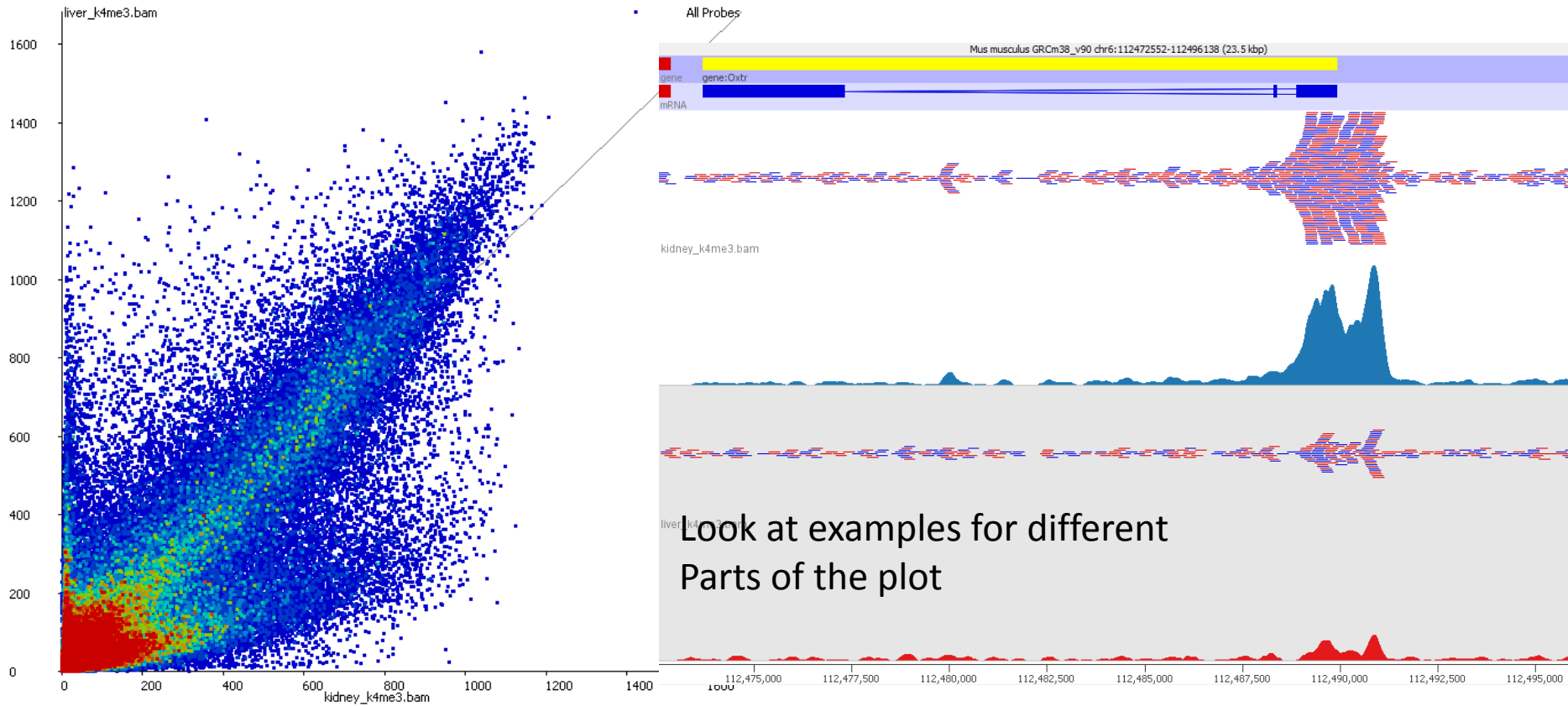
Scatterplot input vs input



- Any suggestion of differential biases in inputs
- Can we merge them to use as a common input

Compare samples

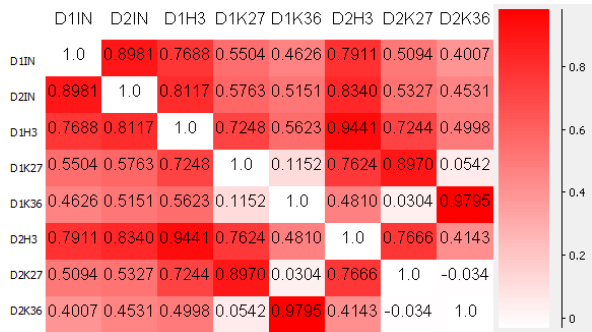
Scatterplot CHIP vs CHIP



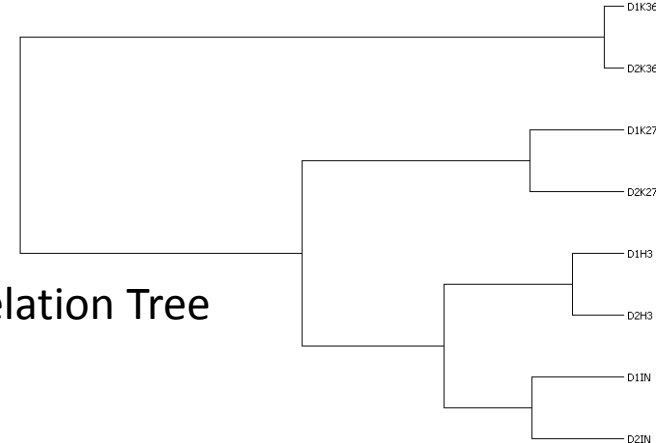
- Look for outgroups (differentially enriched)
- Compare level of enrichment (compare to diagonal)

Compare samples

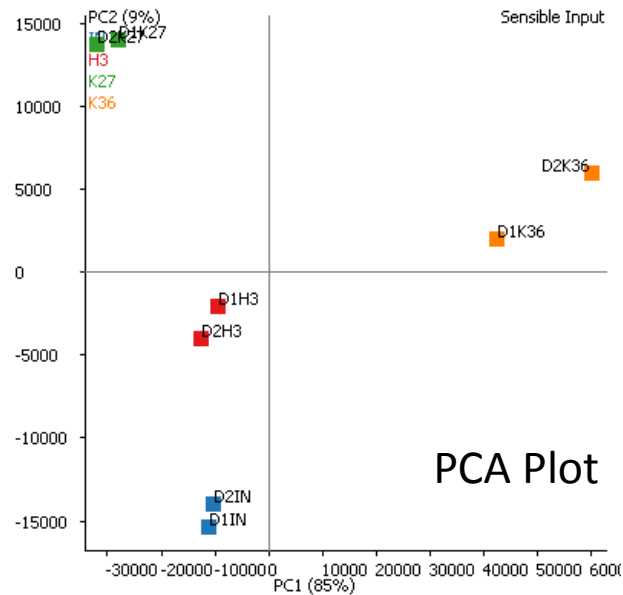
Higher level clustering



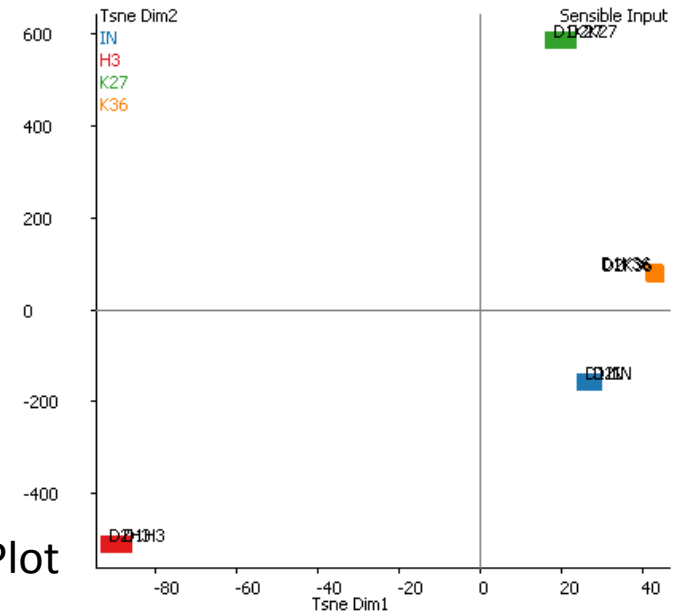
Correlation Matrix



Correlation Tree



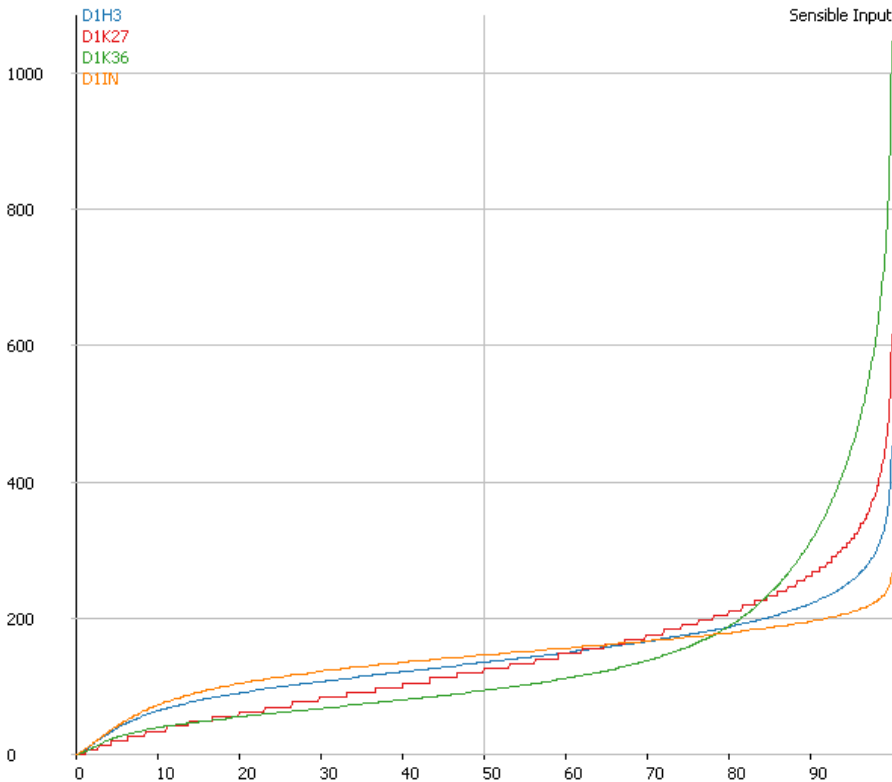
PCA Plot



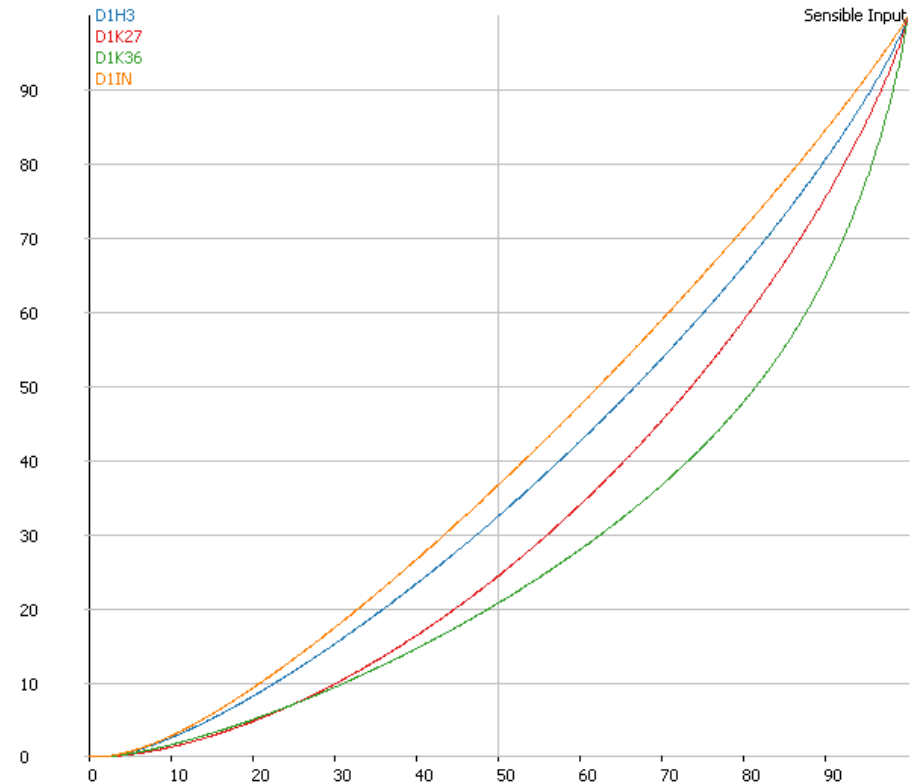
tSNE Plot

Compare samples

Summarise distributions



Cumulative Distribution Plot



QQ Plot

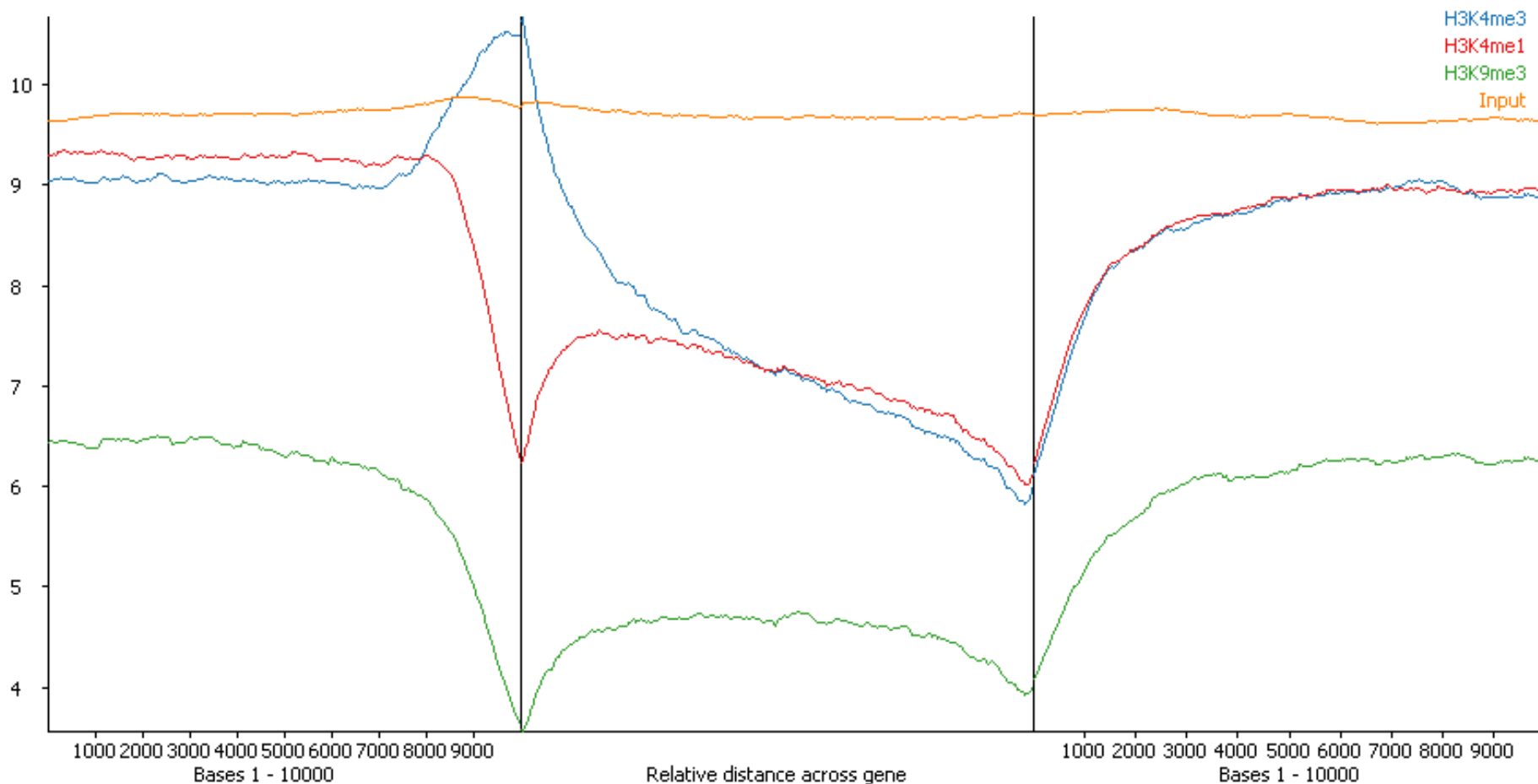
- Flatness of input
- Separation of ChIPs
- Consistency of ChIPs

Associate enrichment with features

Trend Plots

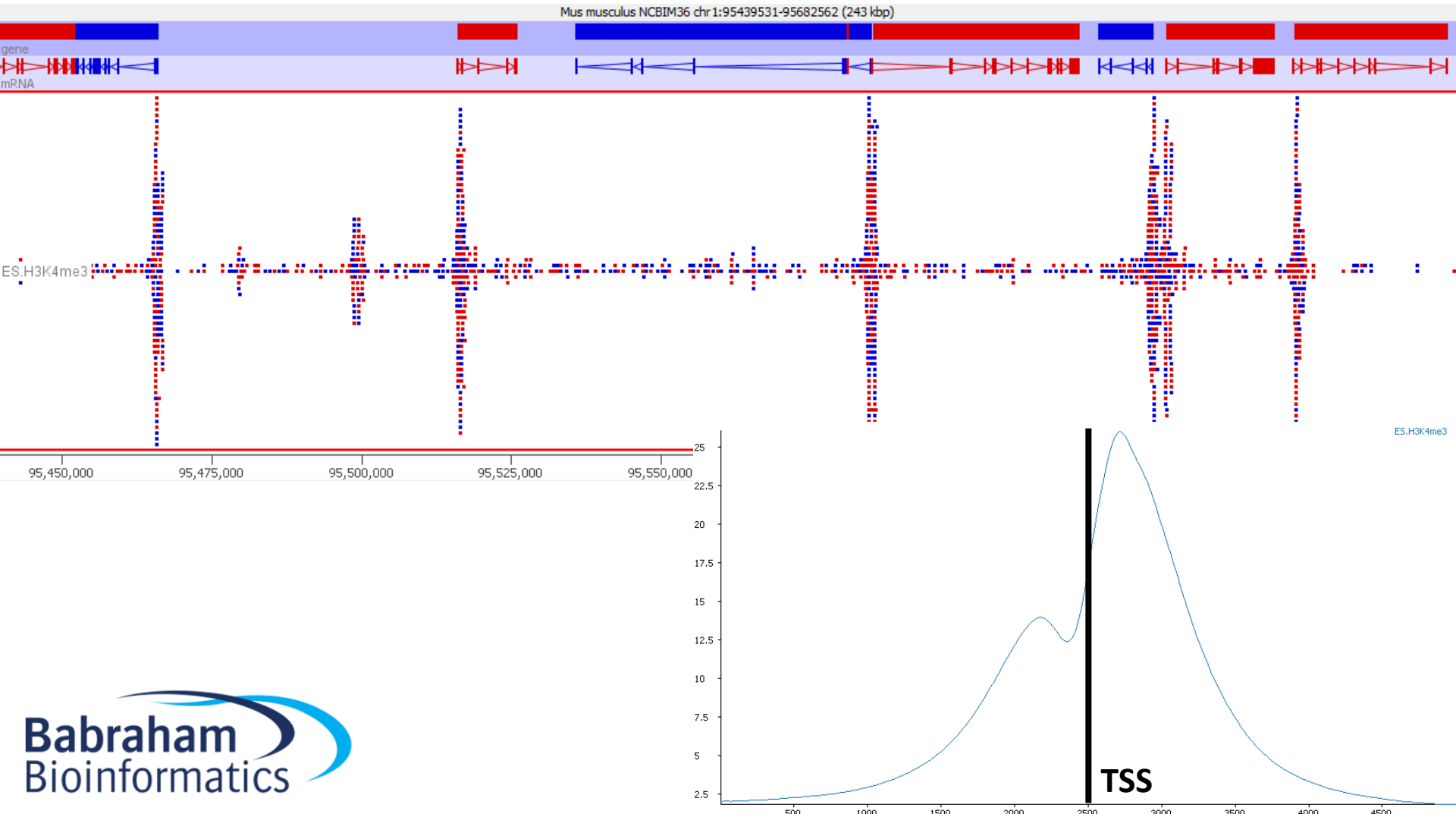
- Graphical way to look at overall enrichment relative to positions in features
 - Gene bodies
 - Promoters
 - CpG islands
- May influence how we later quantitate and analyse the data
 - Analyse per feature
 - Look for exceptions to the general rule

Trend Plot Example

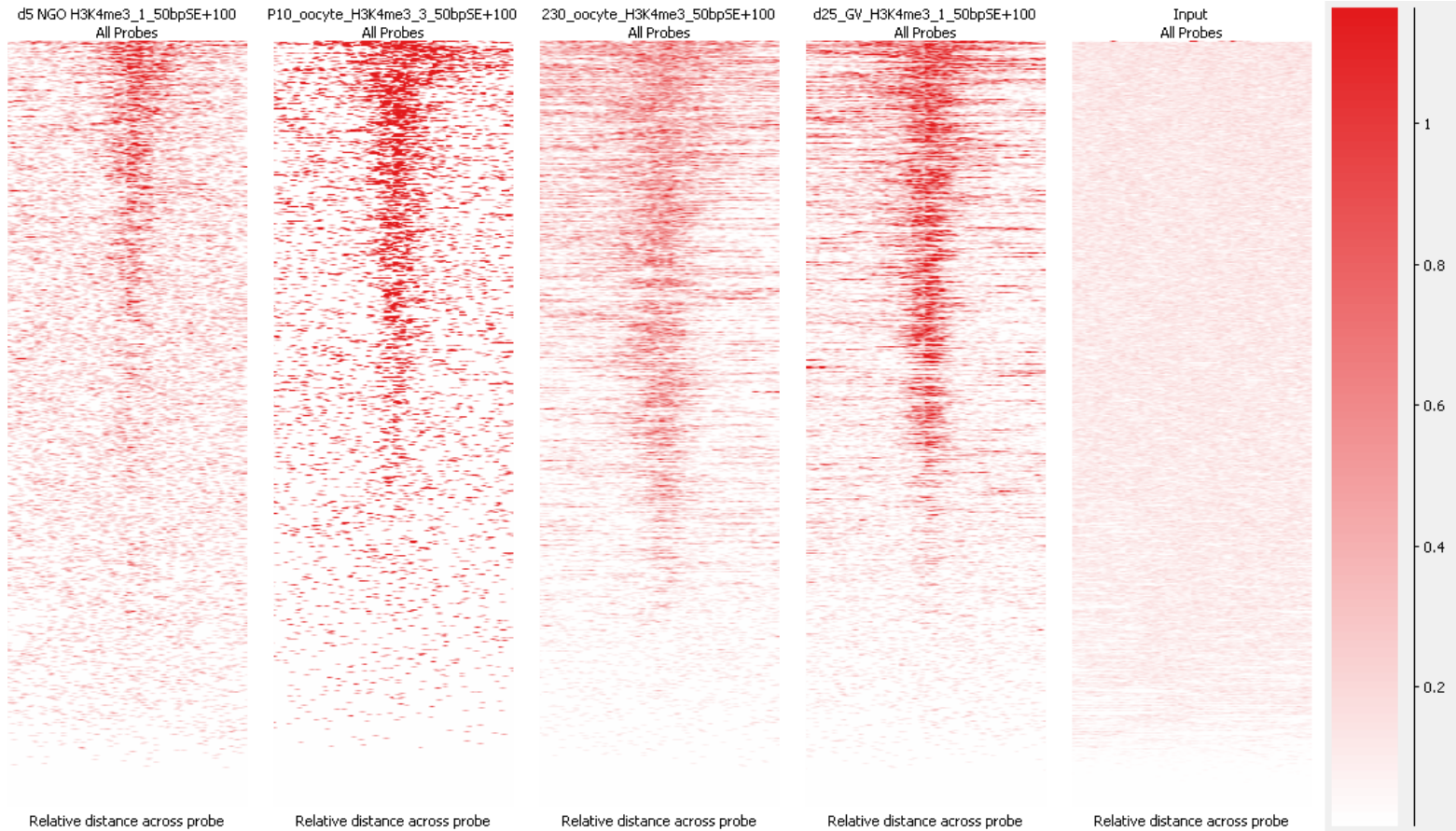


- Overall average
- Says nothing about proportion of features affected

Trend plots should match the data



Aligned Probes Plots give more detail



- Information per feature instance
- Comparison of equivalent features in different marks/samples

After exploration you should..

- Know whether your ChIP is really enriched
- Know the nature / shape of the enrichment
- Know whether your controls behave well
- Know whether you're likely to have differential enrichment
- Know if you will need additional normalisation
- Know the best strategy to measure your data