

# How Sequencing Experiments Fail

v1.0

Simon Andrews

[simon.andrews@babraham.ac.uk](mailto:simon.andrews@babraham.ac.uk)

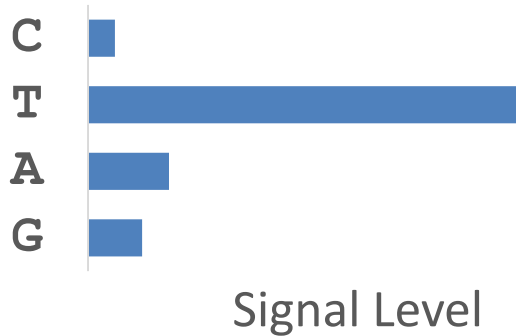
# Classes of Failure

Technical	Something went wrong with a machine
Tracking	Samples aren't what they're supposed to be
Library	Problems during sequencing library preparation
Contamination	Unexpected material in your libraries
Biological	Samples didn't behave the way you expected
Interpretation	Drawing the wrong conclusion from the data

# Technical

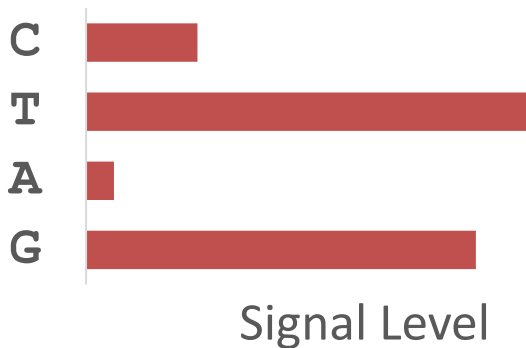
Technical

# Technical Failures



Call = T

Confidence = High



Call = T

Confidence = Low



Call = T

Confidence = Low

Technical

# Phred Scores

$$\text{Phred} = -10 \log_{10} p$$

$p$  = Probability call is incorrect

10% error

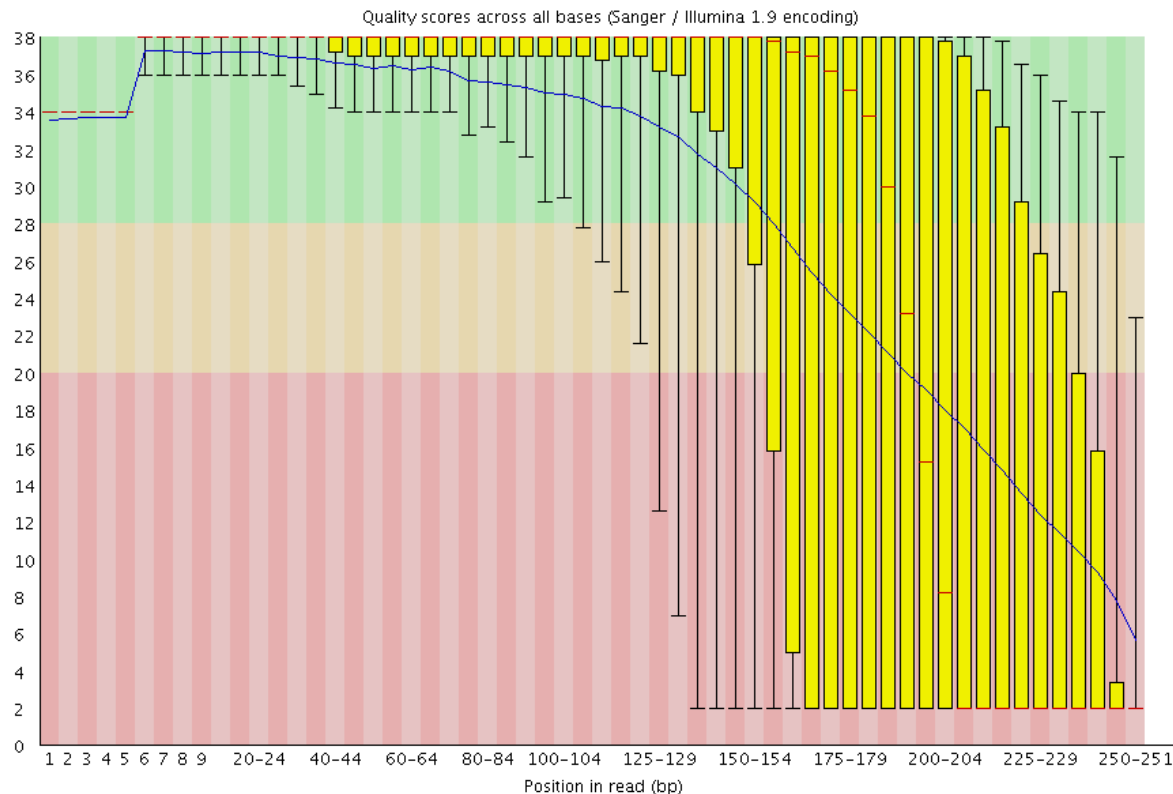
1% error

0.1% error

Phred10

Phred20

Phred30



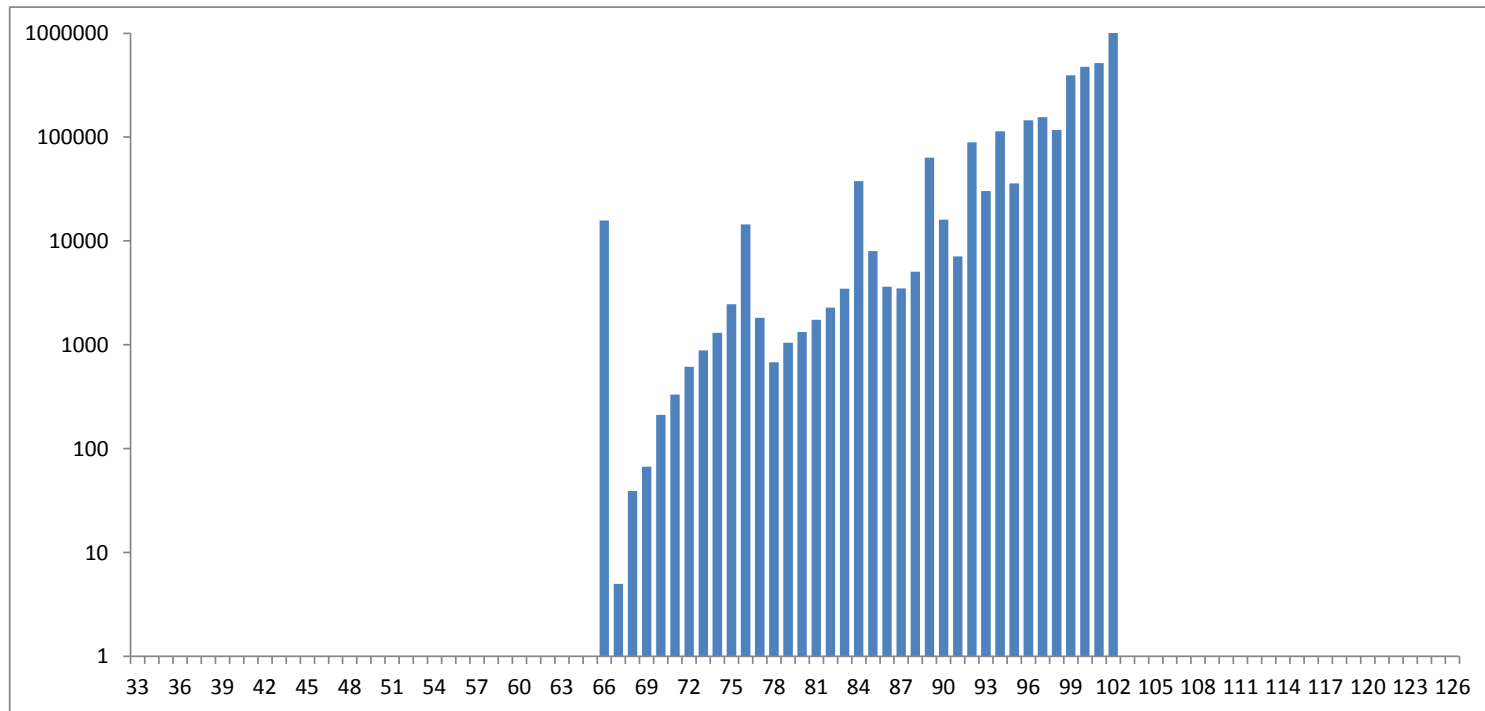
Technical

# Incorrect Encoding

Phred64

!"#\$%&'()\*+,-./0123456789:;<=>?@ABCDEFGHIJKLMN O PQRSTUVWXYZ[\]^\_`abcdefgh

Phred33

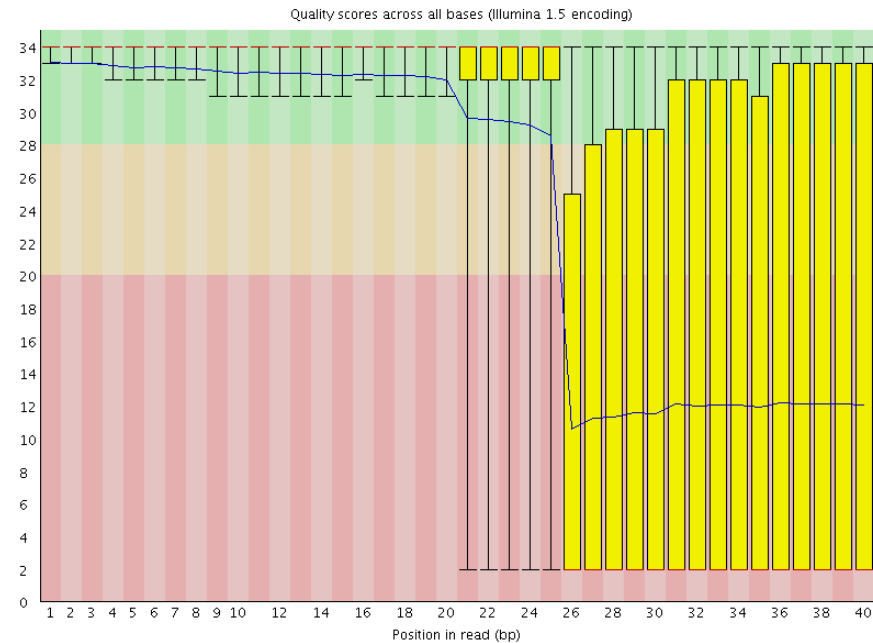
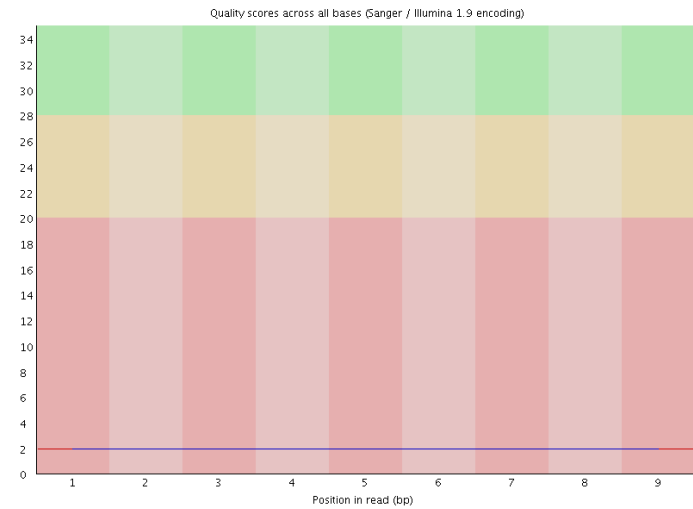
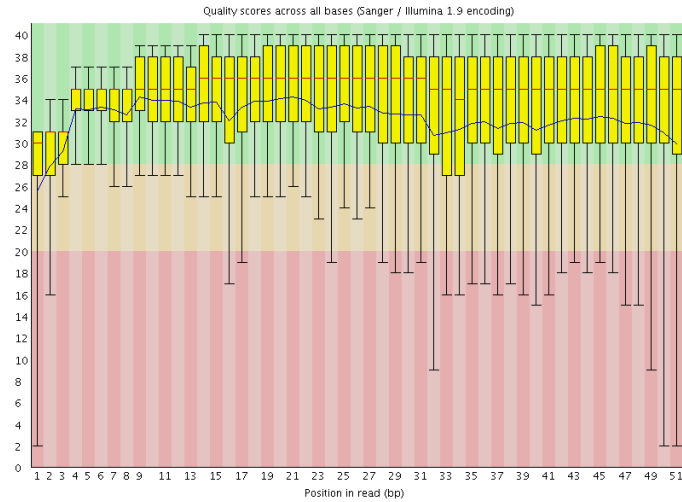


Phred64 (Illumina)

Phred33 (Sanger)

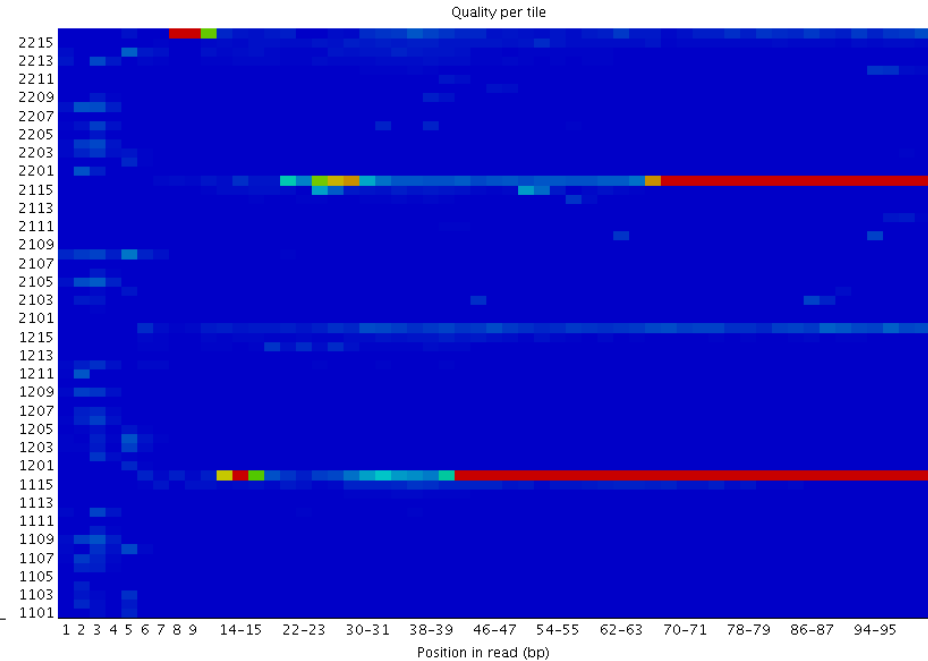
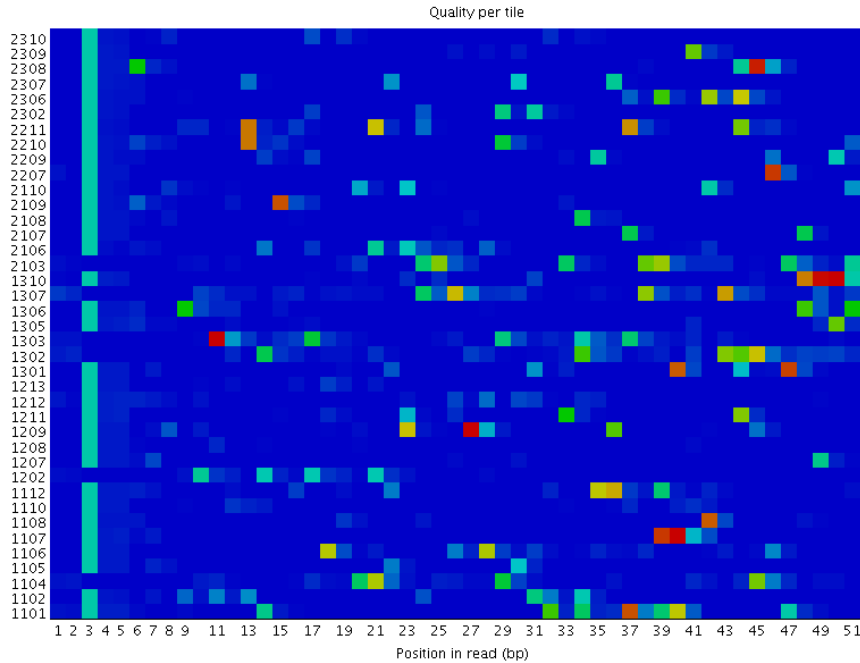
Technical

# Phred Scores



Technical

# Positional Phred Scores

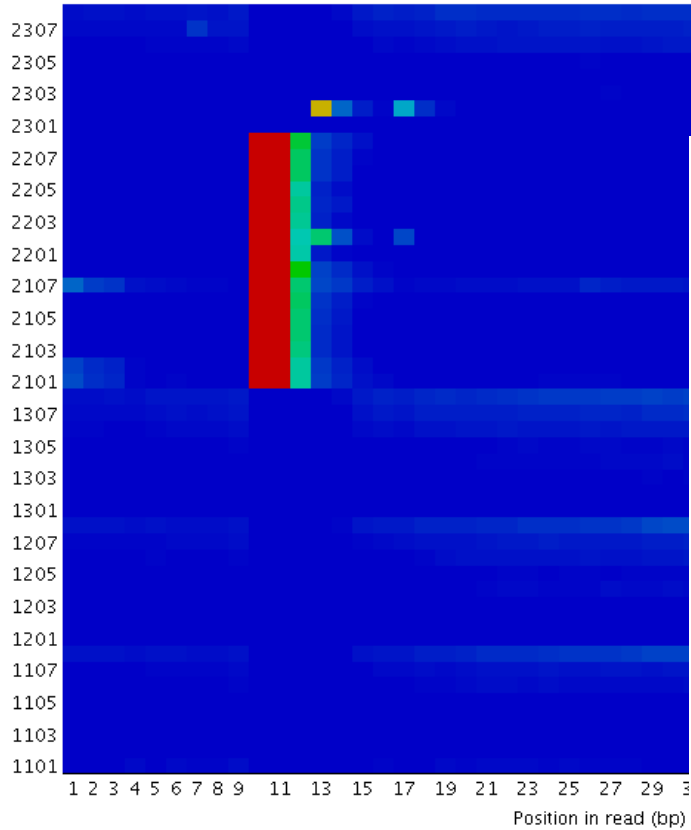




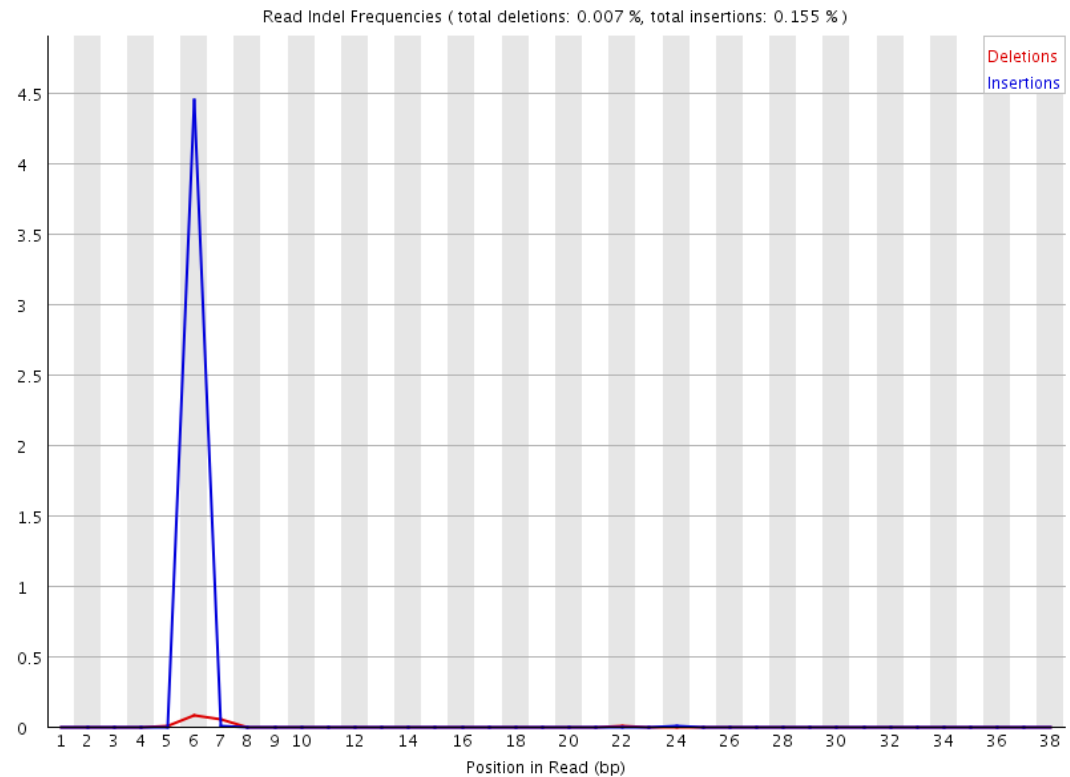
Technical

# Positional Phred Scores

Quality per tile



## Indel Frequencies



Technical

# Biased Phred Scores

HindIII sites



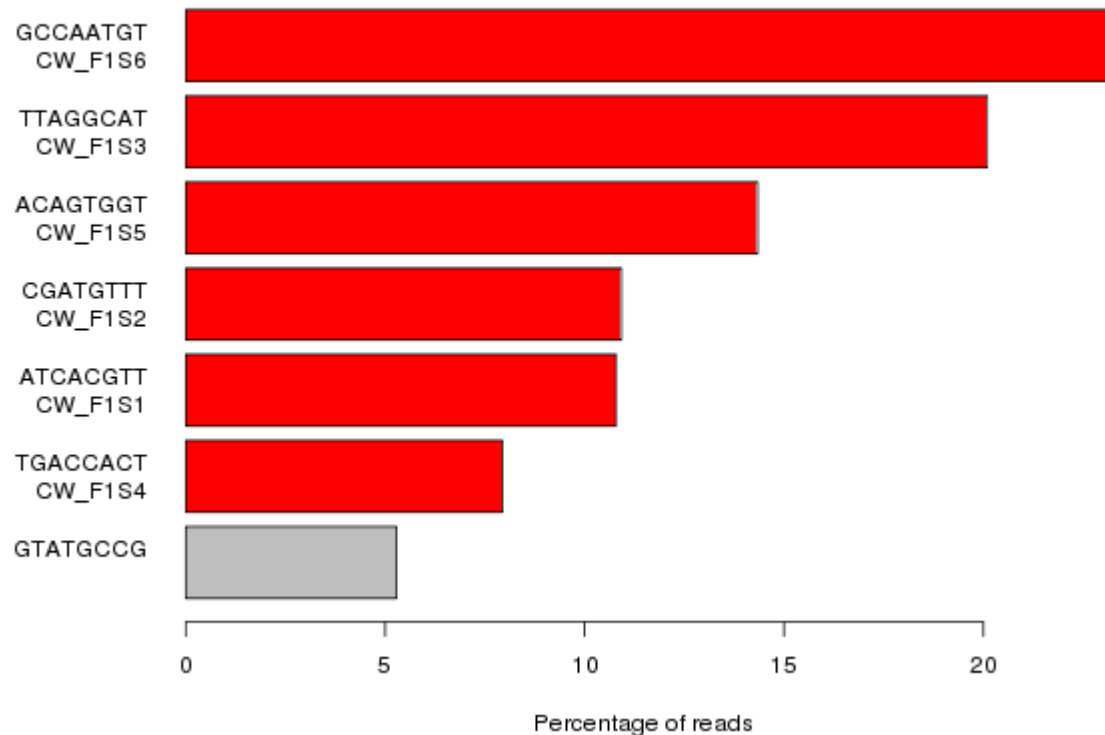
GGATCCGTATGCGATGCTAGCGT  
GGATCATATATATGCTAGCGTAT  
GGATCTATATTGCGCGATACTGG  
GGATCCCGTAGCTGCGATGCTGA  
GGATCAAGGATAGCGCGTCTAGA  
GGATCTATATAGTTGCCGTATCG  
GGATCGGAGCGGGATCGGATGCG



# Tracking

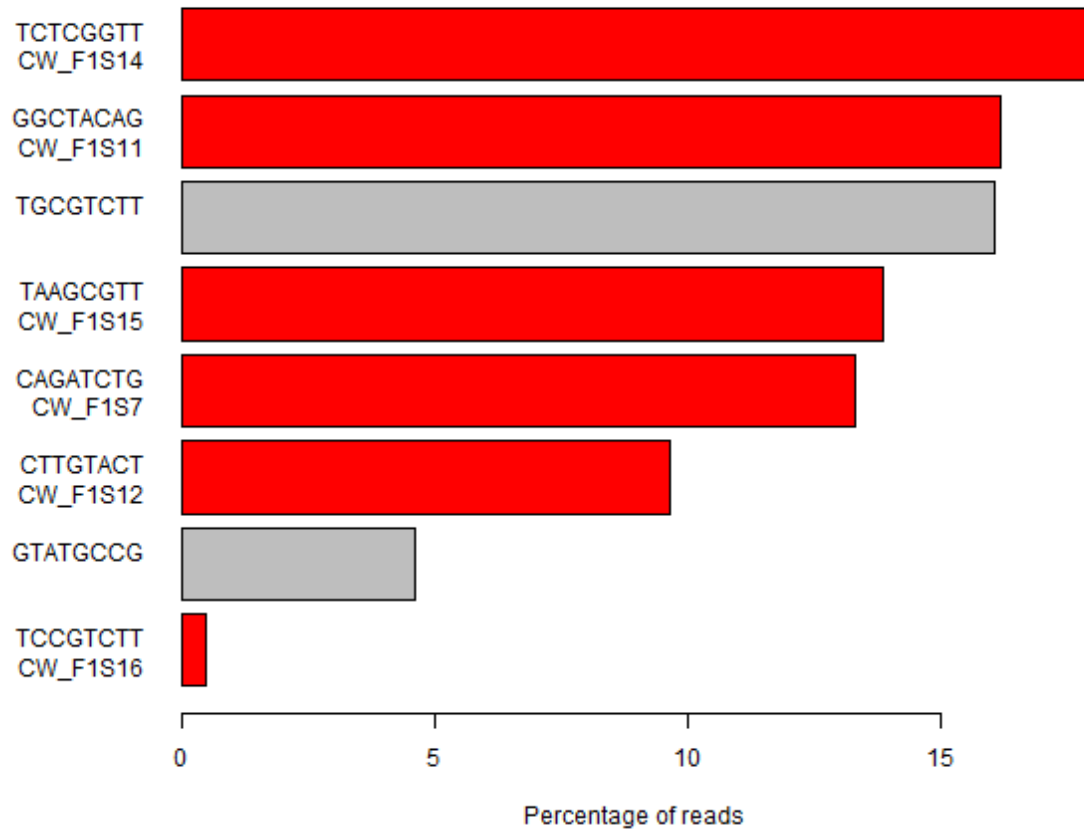
# Tracking - Barcodes

Barcodes shown explain 92% of the data



# Tracking - Barcodes

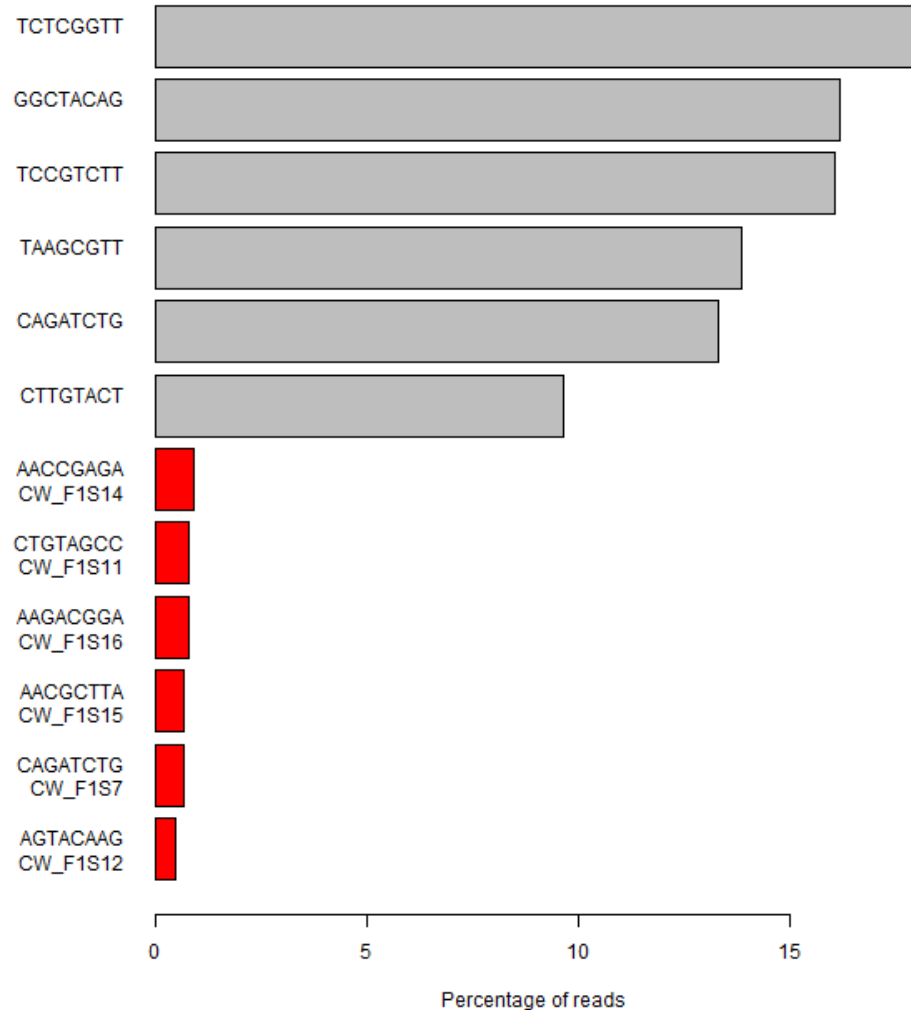
Barcodes shown explain 92% of the data



## Tracking

# Tracking - Barcodes

Barcodes shown explain 91% of the data



# Tracking Exercise

You have some barcode statistics from a set of runs from the same group.

Red = Expected barcode

Grey = Unexpected barcode

Can you see if you can spot a problem within this data set,  
and say how many of the lanes it might affect?

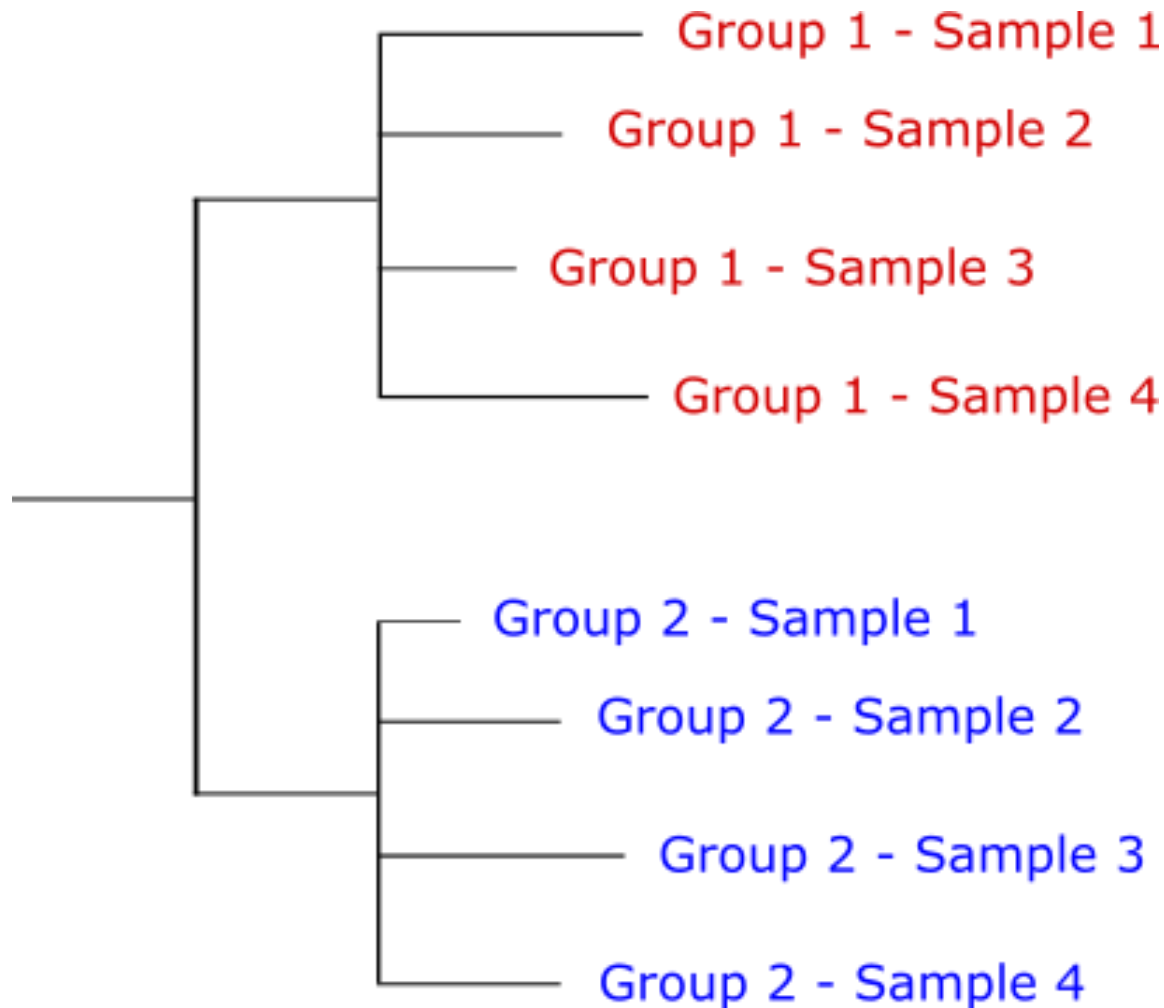
# Tracking – Swapped Samples

- Swapped between users
  - Different sample type
  - Different species
  - See later Contamination / Biology sections
- Swapped within experiment
  - Look for consistent biological signal
  - Use other knowledge of the samples to validate



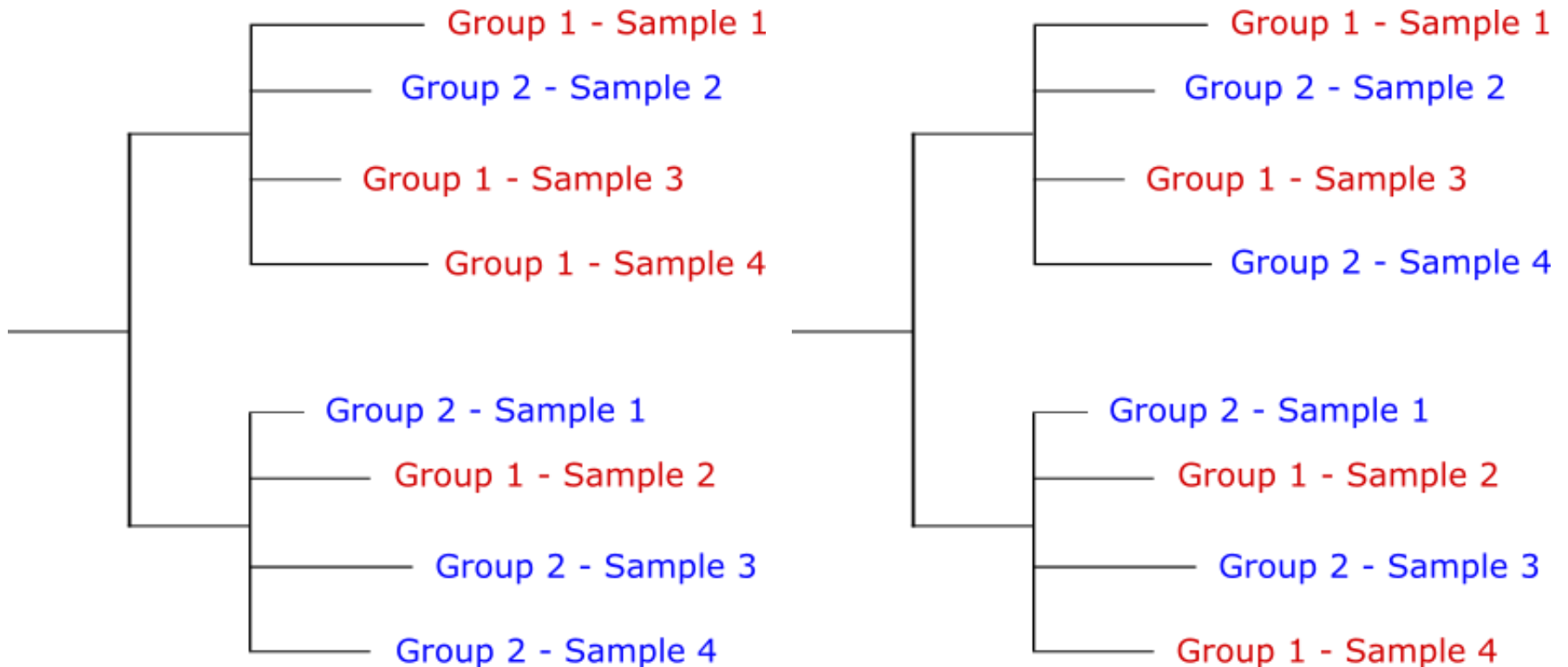
Tracking

# Tracking – Sample groups



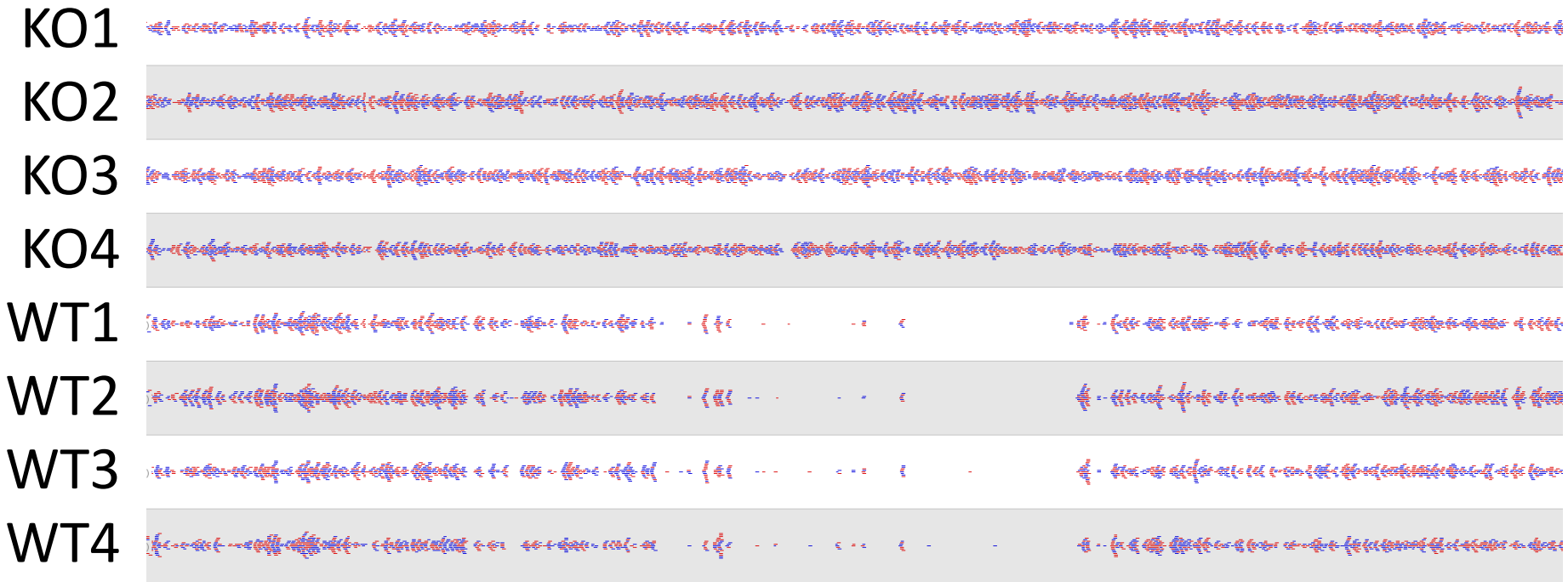
Tracking

# Tracking – Sample groups



Tracking

# Tracking – Sample swaps



# Library

# Library Problems

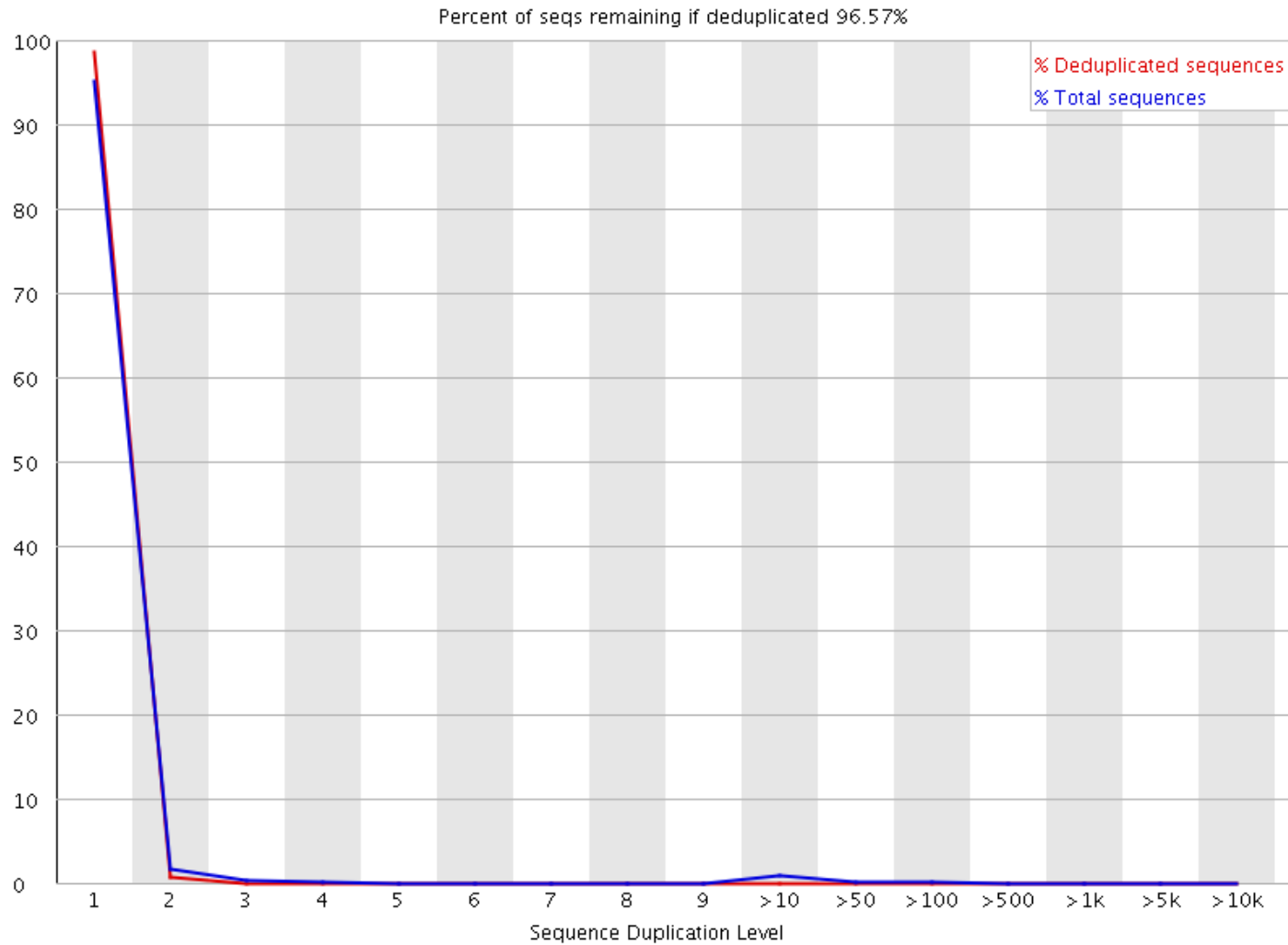
- Material Lost
  - Overamplification
    - Duplication
- Biases in selection
  - Priming bias
  - GC bias
  - Methylation bias
  - Size selection bias
- Technical contamination
  - Read through adapter
  - Adapter dimers

# Duplication

- Over-sequencing of library complexity
- Too little material or too much PCR
- Can be difficult to assess
- Why does duplication matter?
  - Potentially biased
  - Over-estimates measurement accuracy

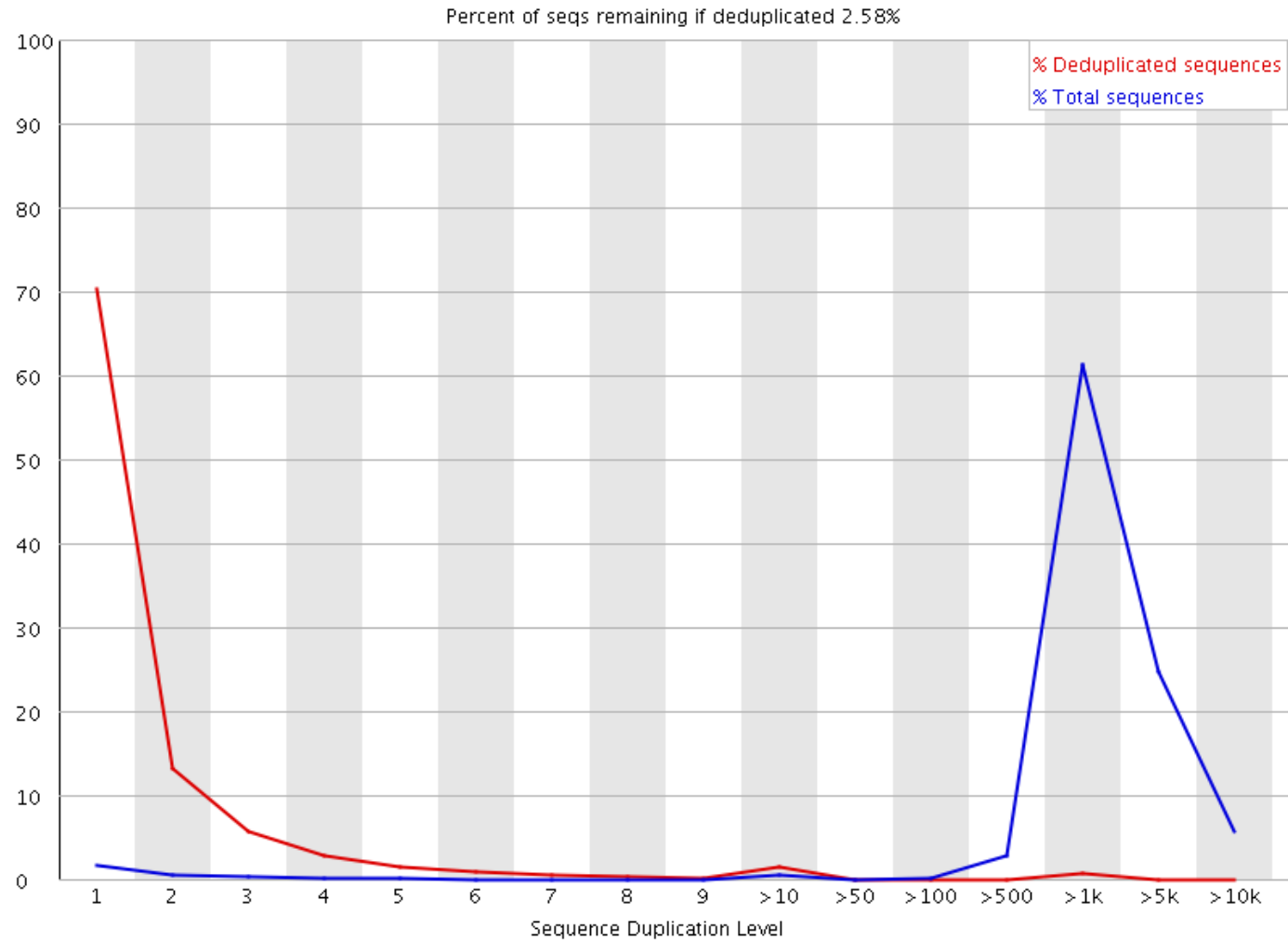
Library

# Duplication



Library

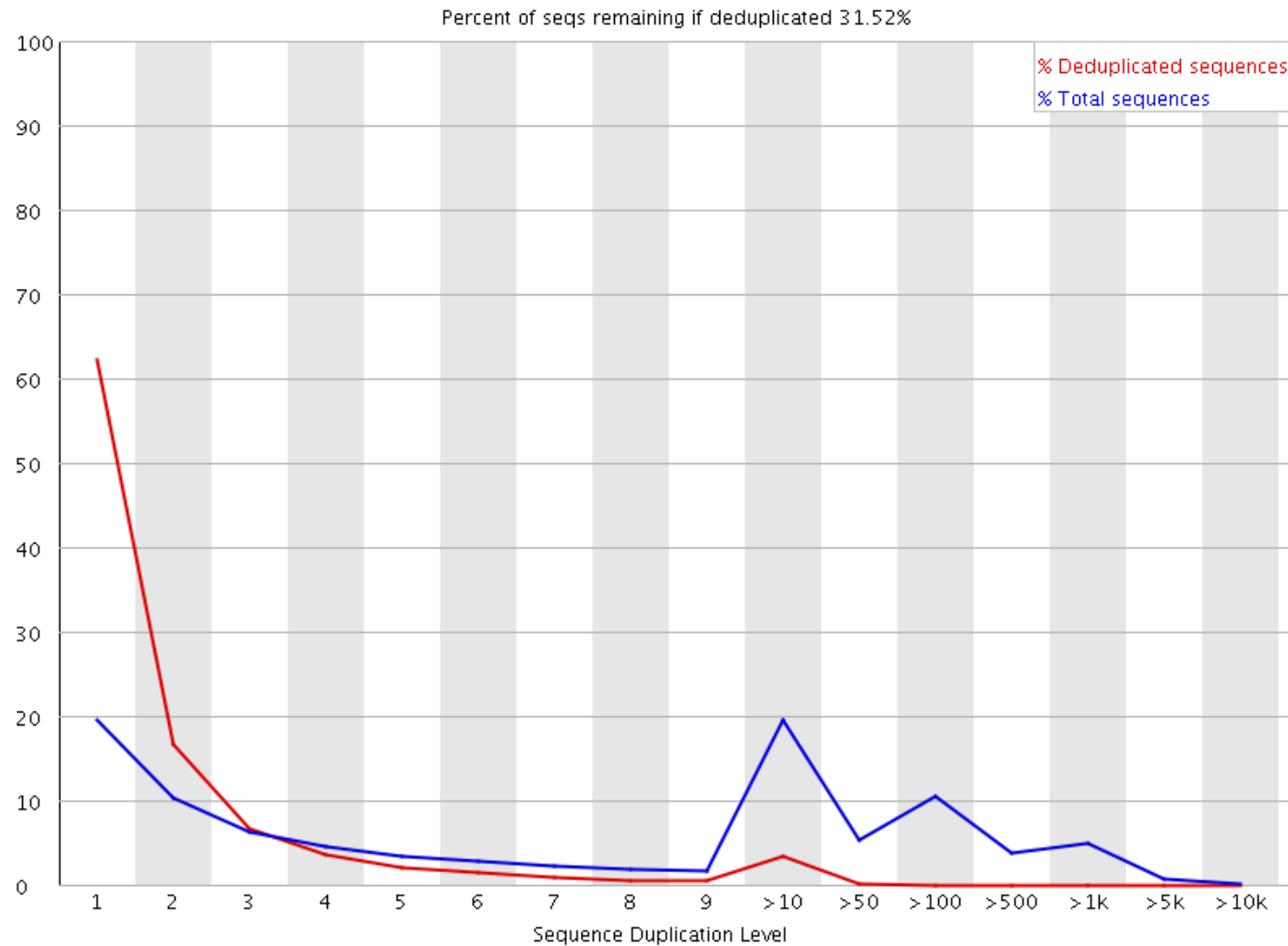
# Duplication





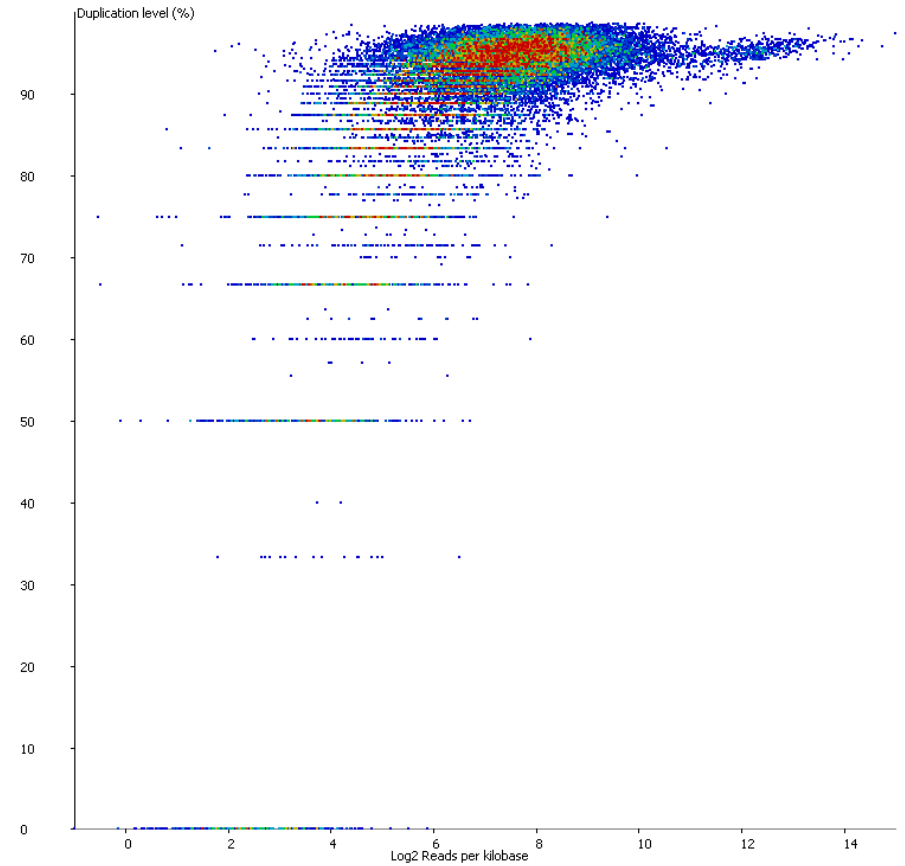
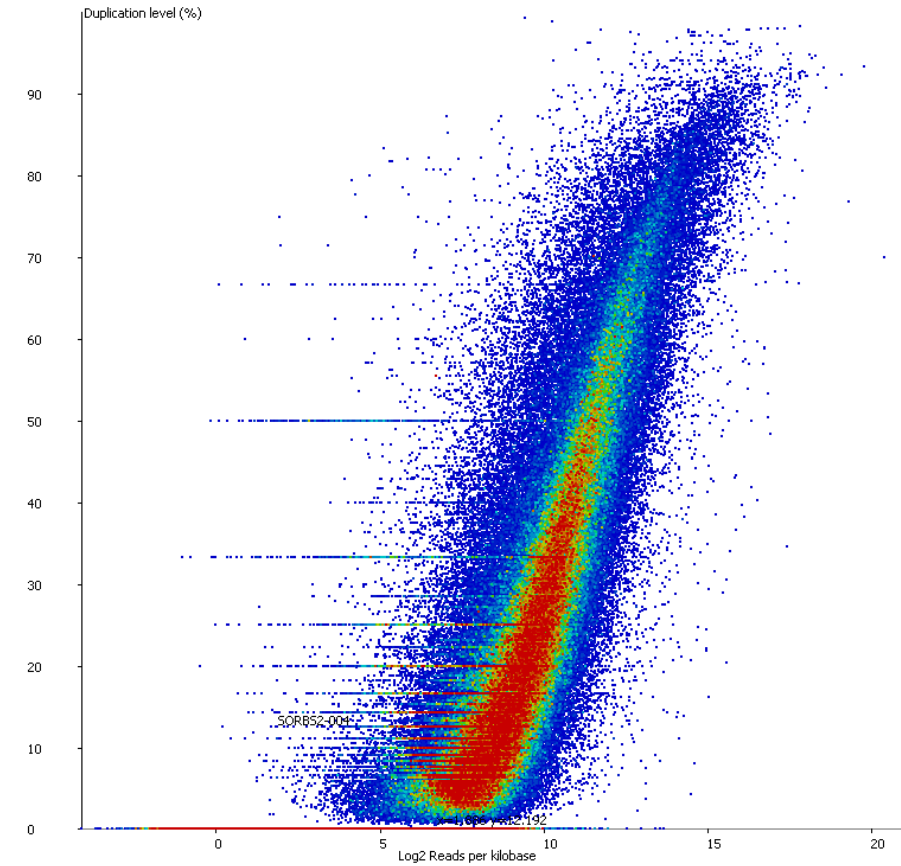
Library

# Duplication



Library

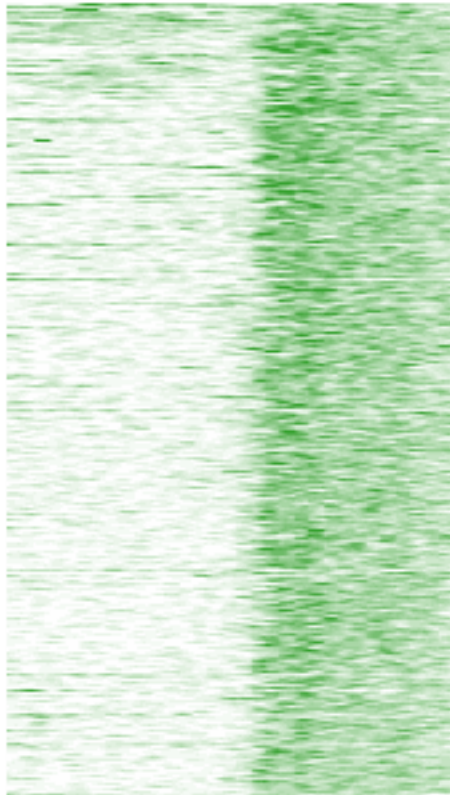
# Duplication



Library

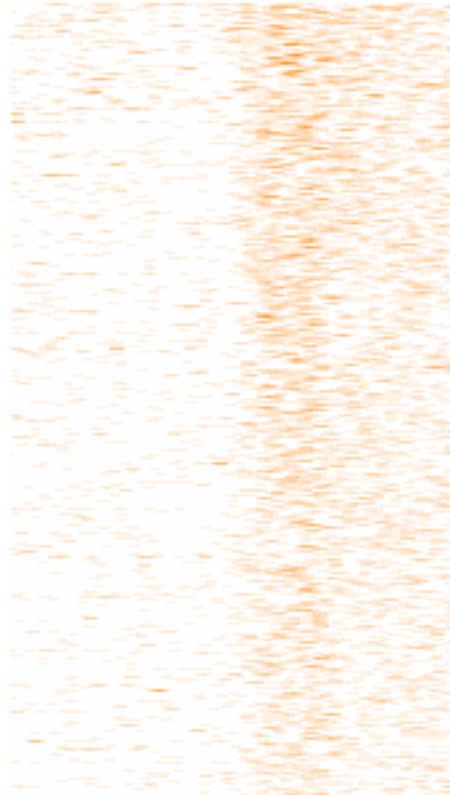
# Duplication - Repeats

canonical\_reimport



Repeat

Input\_reimport



Repeat

Real



Mapped



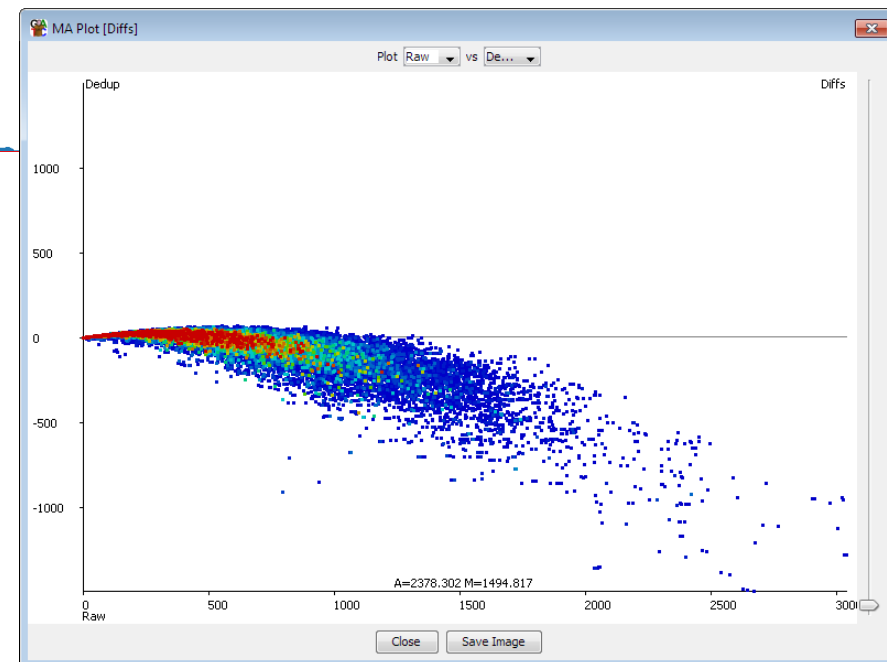
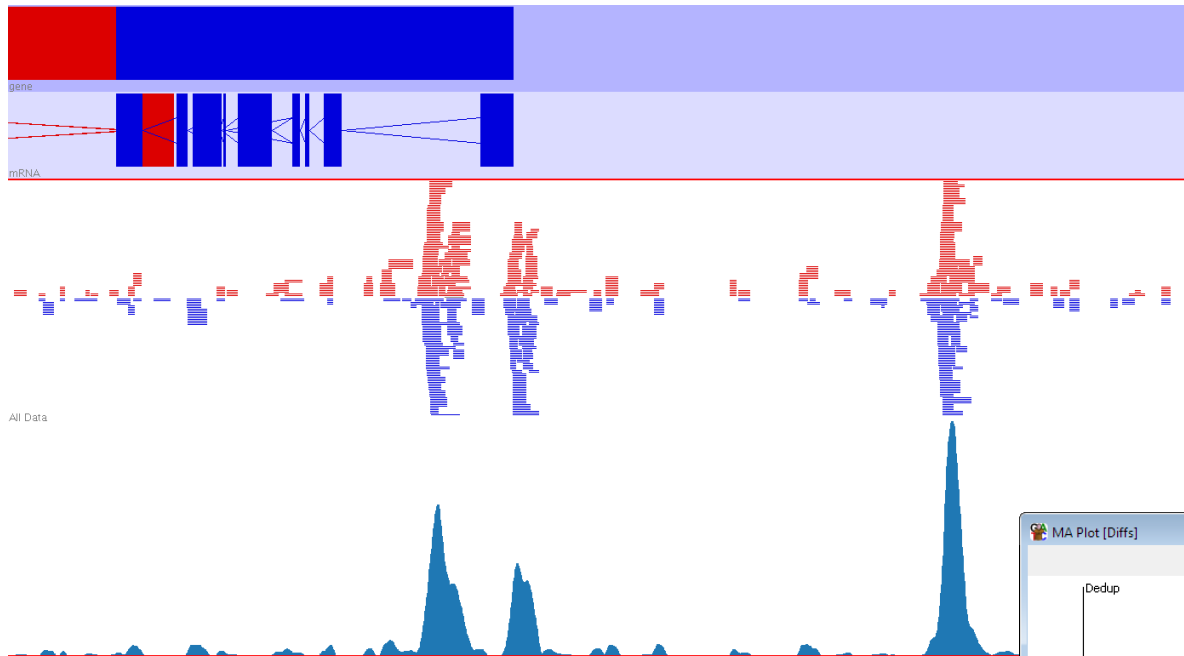
Deduplicated



Peak callers (MACS for example) deduplicate internally, so you don't have to consciously do this.

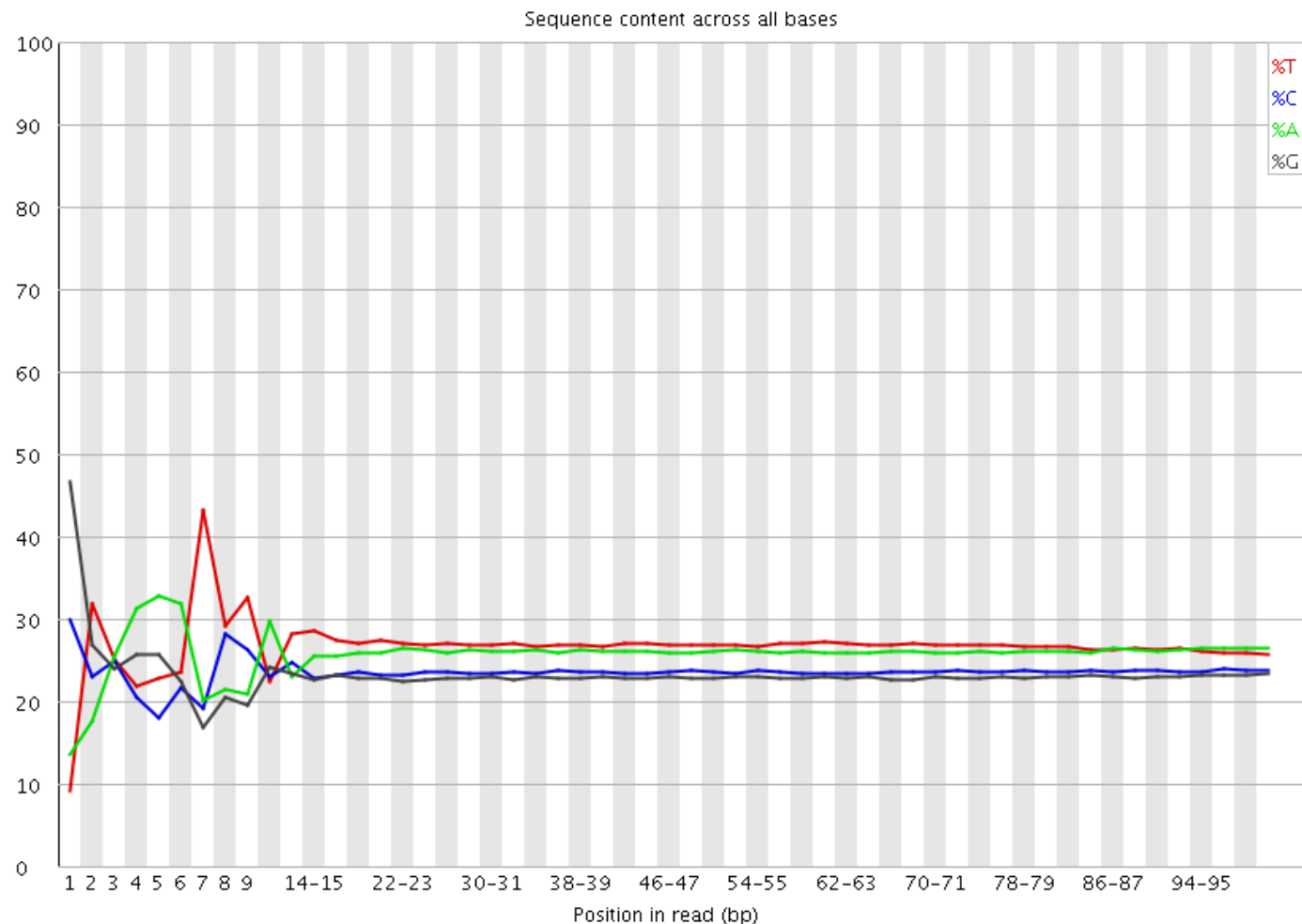
Only avoided by using uniquely mapped reads.

# Deduplication



Library

# Priming Bias



# Contamination

# Contamination

- Different kinds
  - Technical contamination
    - Adapter dimers
  - Contamination with a species you might expect
    - *E.coli* in a mouse sample
  - Contamination with something unexpected
  - Contamination with the wrong material
    - DNA in an RNA-Prep
  - Mixed samples

# Mapping Efficiency

- Know what to expect
  - Data type (genomic / transcriptomic)
  - How good / complete is the genome
- Distinguish unique / multi-mapped reads
  - Understand the mapping process

Reads:

```
Input:      5725730
Mapped:     4703342 (82.1% of input)
of these:   471516 (10.0%) have multiple alignments
                                   (471516 have >1)
```

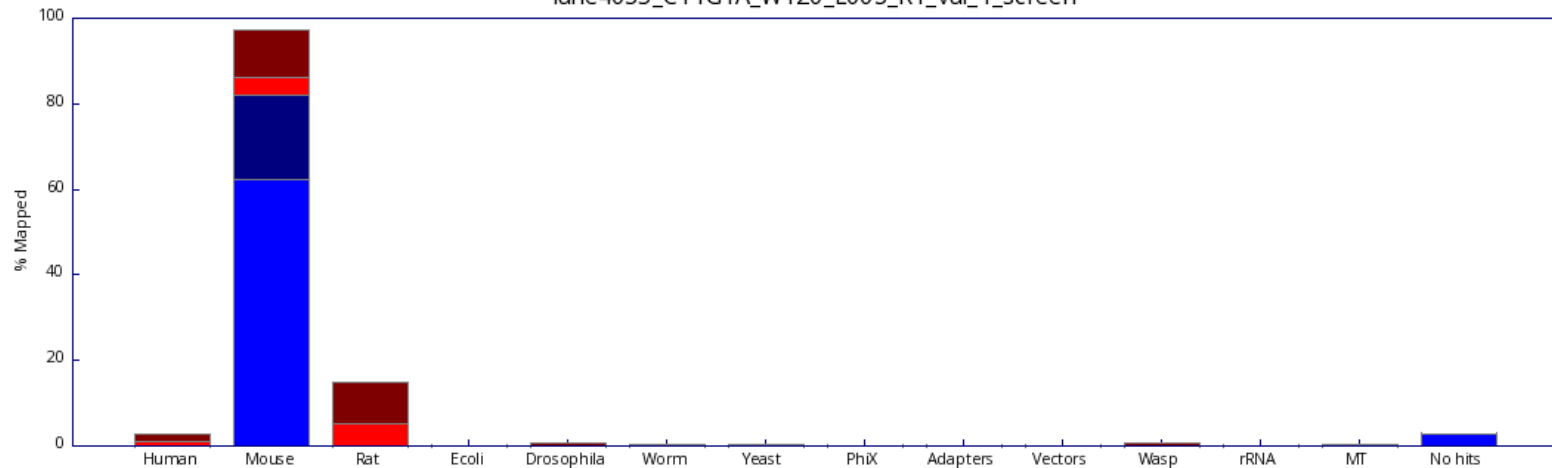
82.1% overall read alignment rate.



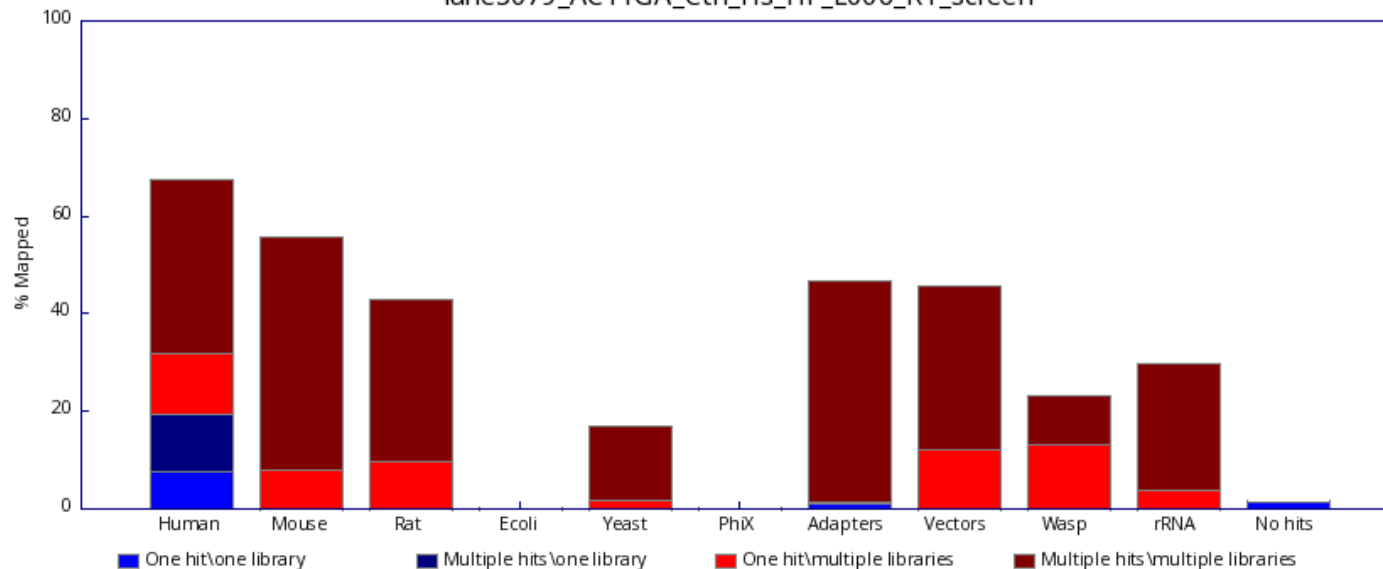
## Contamination

# Species Screen

lane4035\_CTTGTA\_WT20\_L003\_R1\_val\_1\_screen

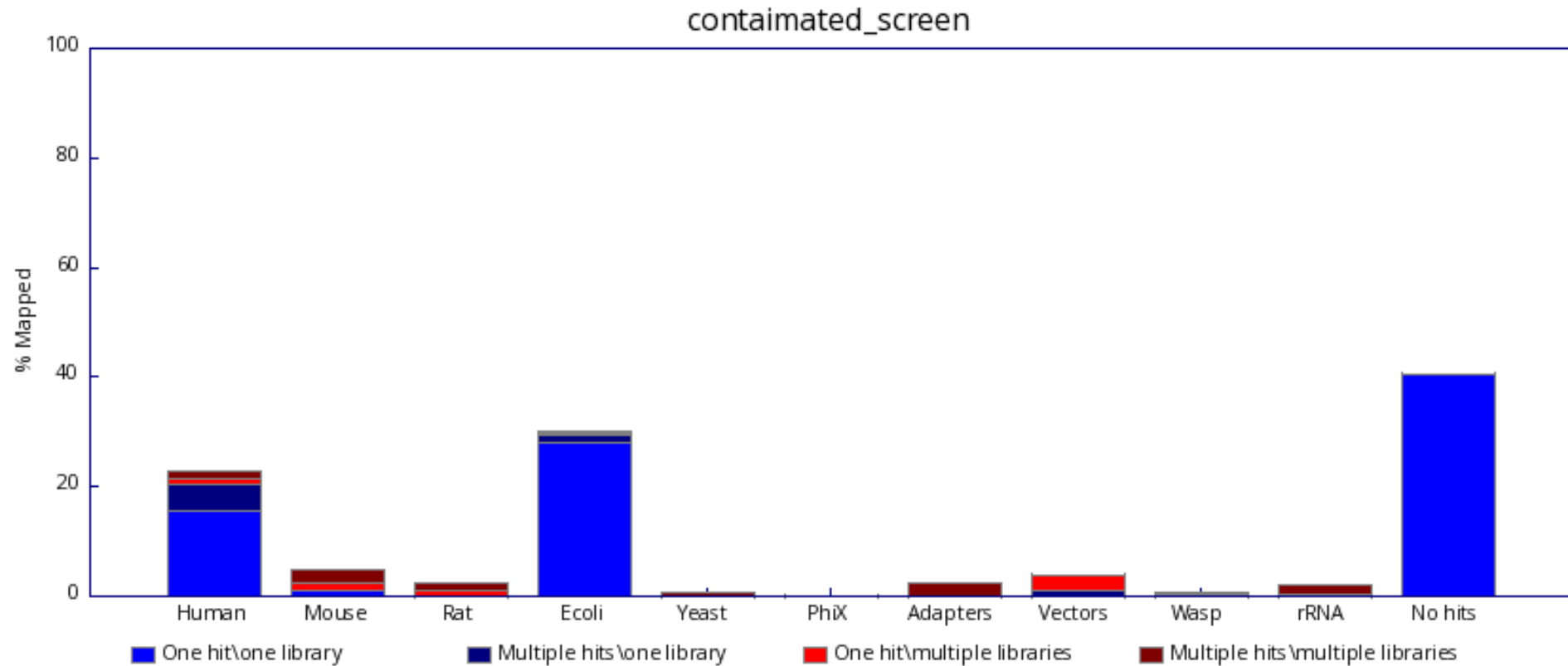


lane3079\_ACTTGA\_Ctrl\_Hs\_HF\_L006\_R1\_screen



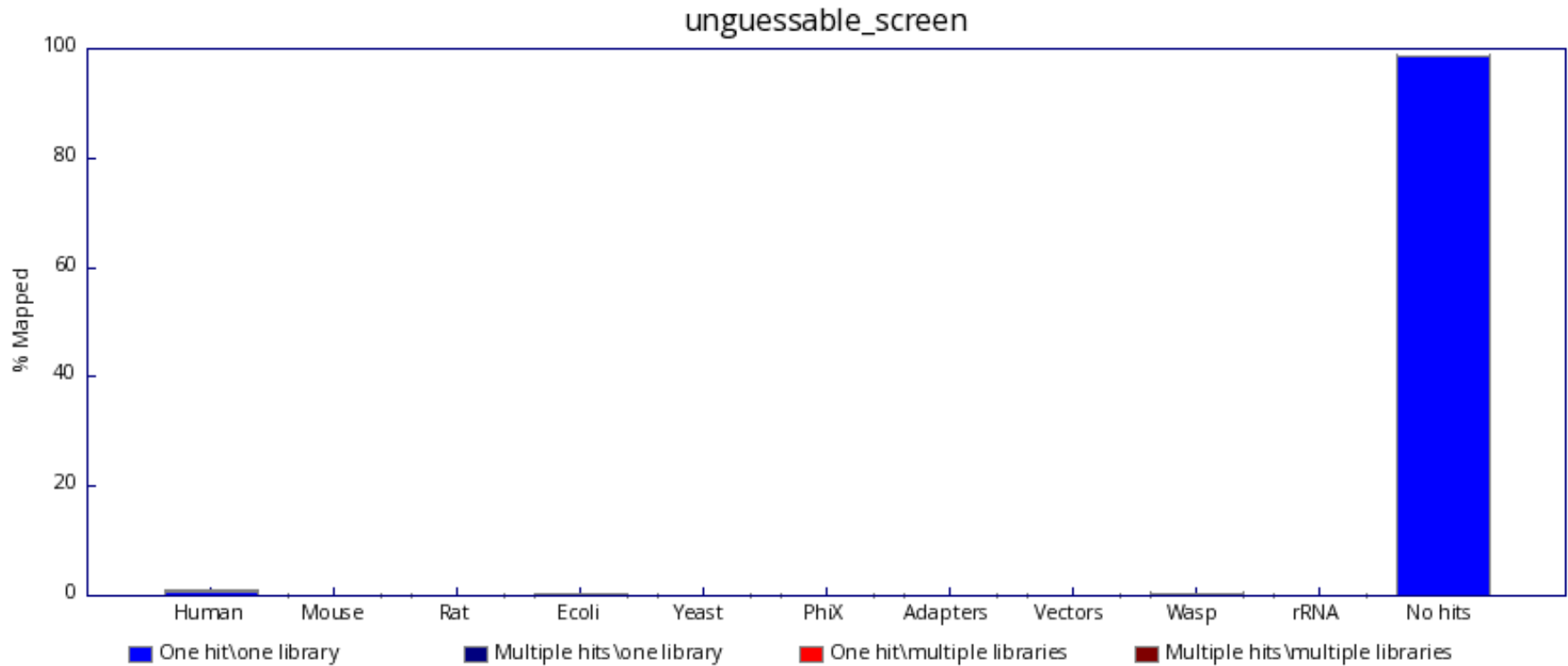
Contamination

# Species Screen



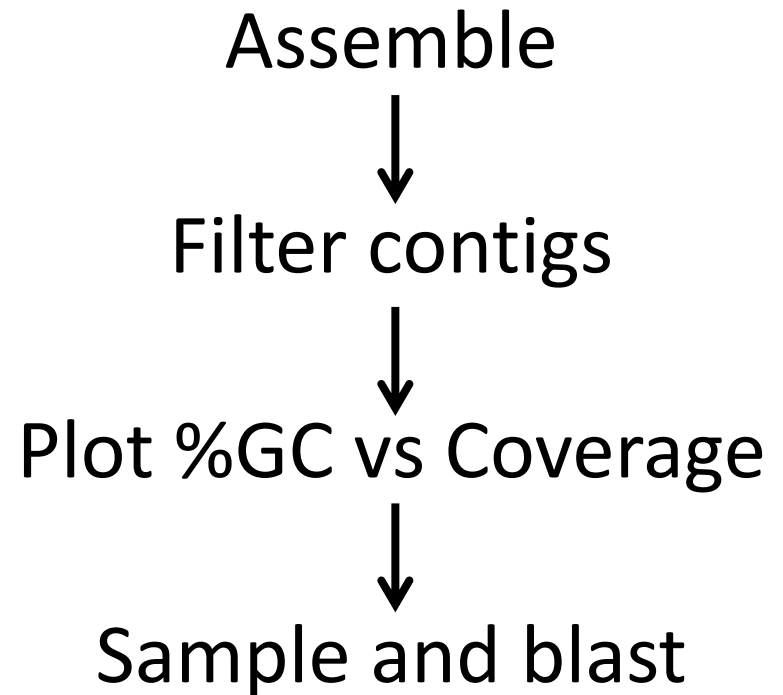
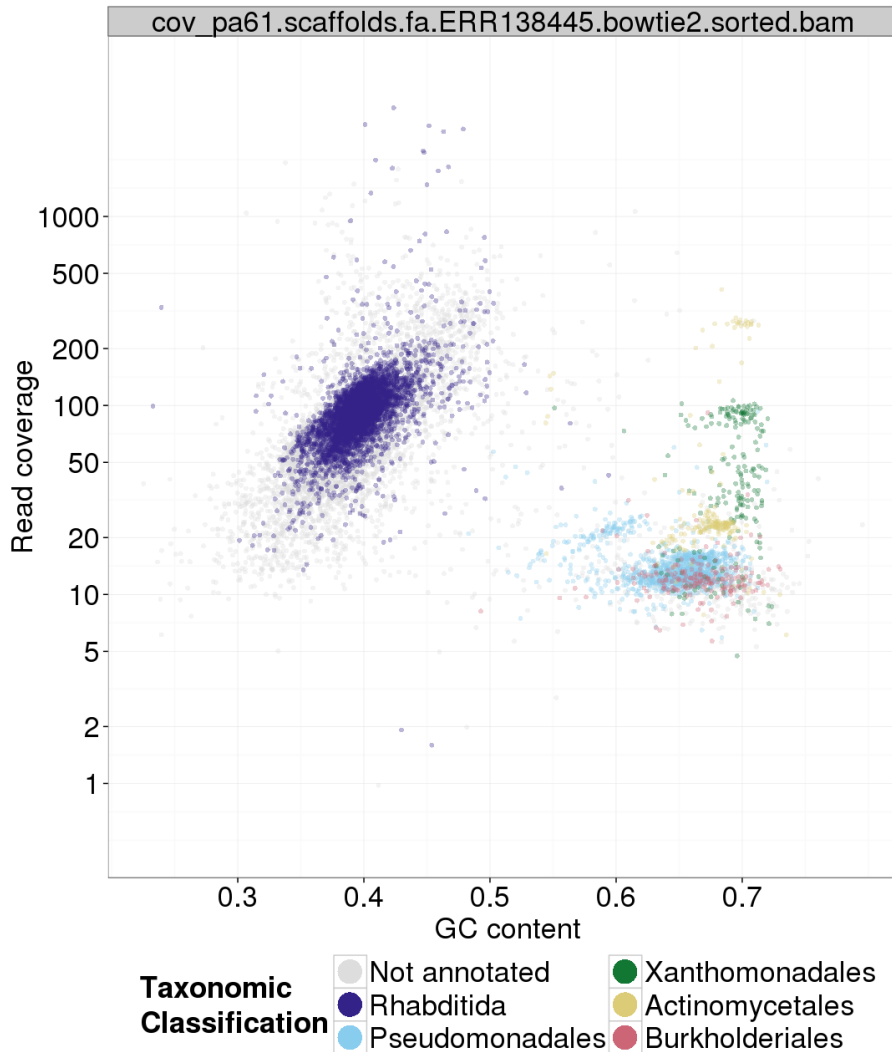
Contamination

# Species Screen



## Contamination

# TAGC Plots



Salter et al. *BMC Biology* 2014, **12**:87  
<http://www.biomedcentral.com/1741-7007/12/87>

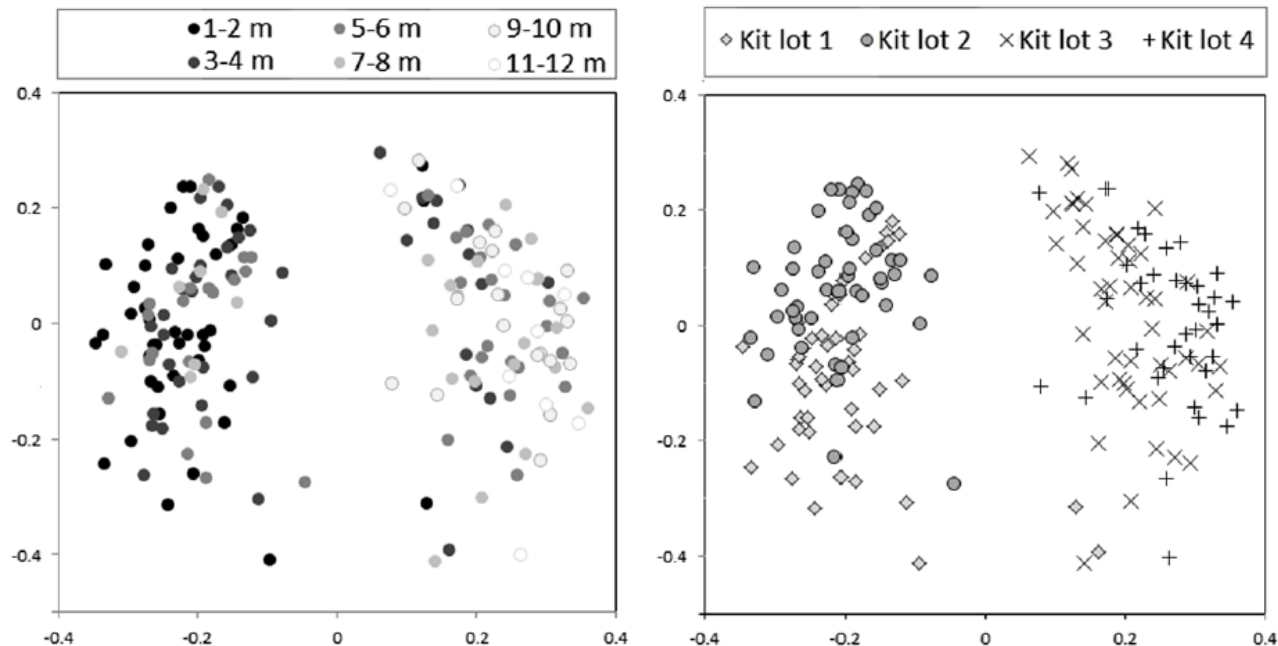


## RESEARCH ARTICLE

## Open Access

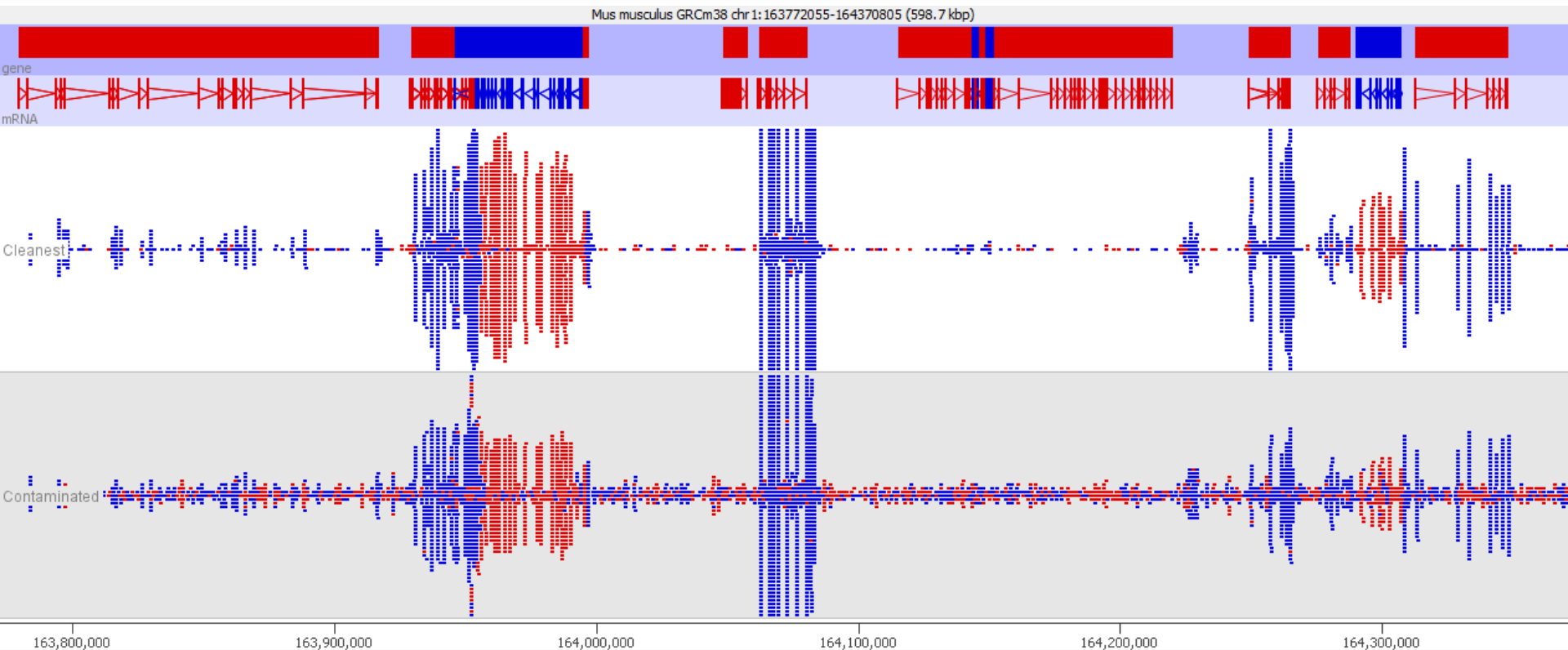
## Reagent and laboratory contamination can critically impact sequence-based microbiome analyses

Susannah J Salter<sup>1\*</sup>, Michael J Cox<sup>2</sup>, Elena M Turek<sup>2</sup>, Szymon T Calus<sup>3</sup>, William O Cookson<sup>2</sup>, Miriam F Moffatt<sup>2</sup>, Paul Turner<sup>4,5</sup>, Julian Parkhill<sup>1</sup>, Nicholas J Loman<sup>3</sup> and Alan W Walker<sup>1,6\*</sup>



Contamination

# Internal Contamination



# Biological

# Samples Don't Behave

- All samples come with a set of expectations
  - Biological effect
  - Sample source
  - Rough biological behaviour
- If these aren't met
  - Samples may not be what you expect
  - Statistical analyses may be invalid
  - Larger biological picture may be missed



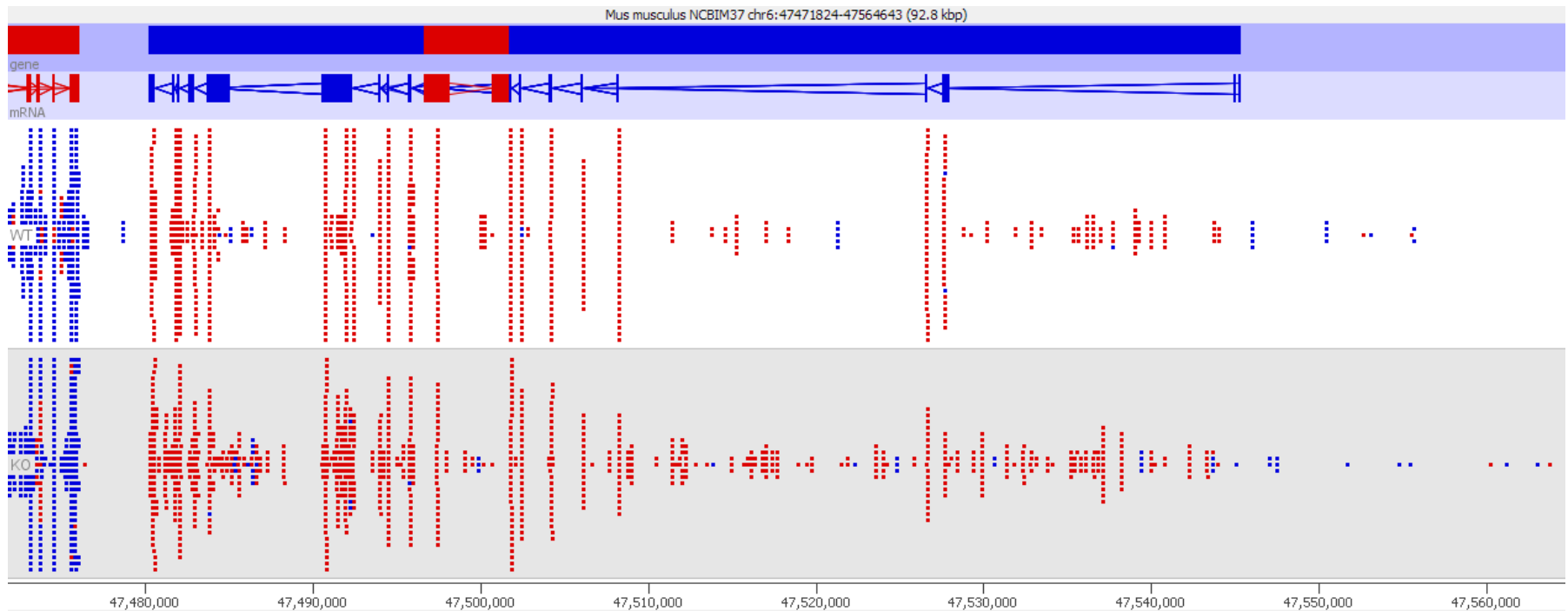
# Exercise

You have been given a set of QC and visualisation results for a knockout in male black6 mice (same genotype as the reference) of a single gene.

Have a look through the plots and see if there is anything which would cause you concern regarding the behaviour of the samples.

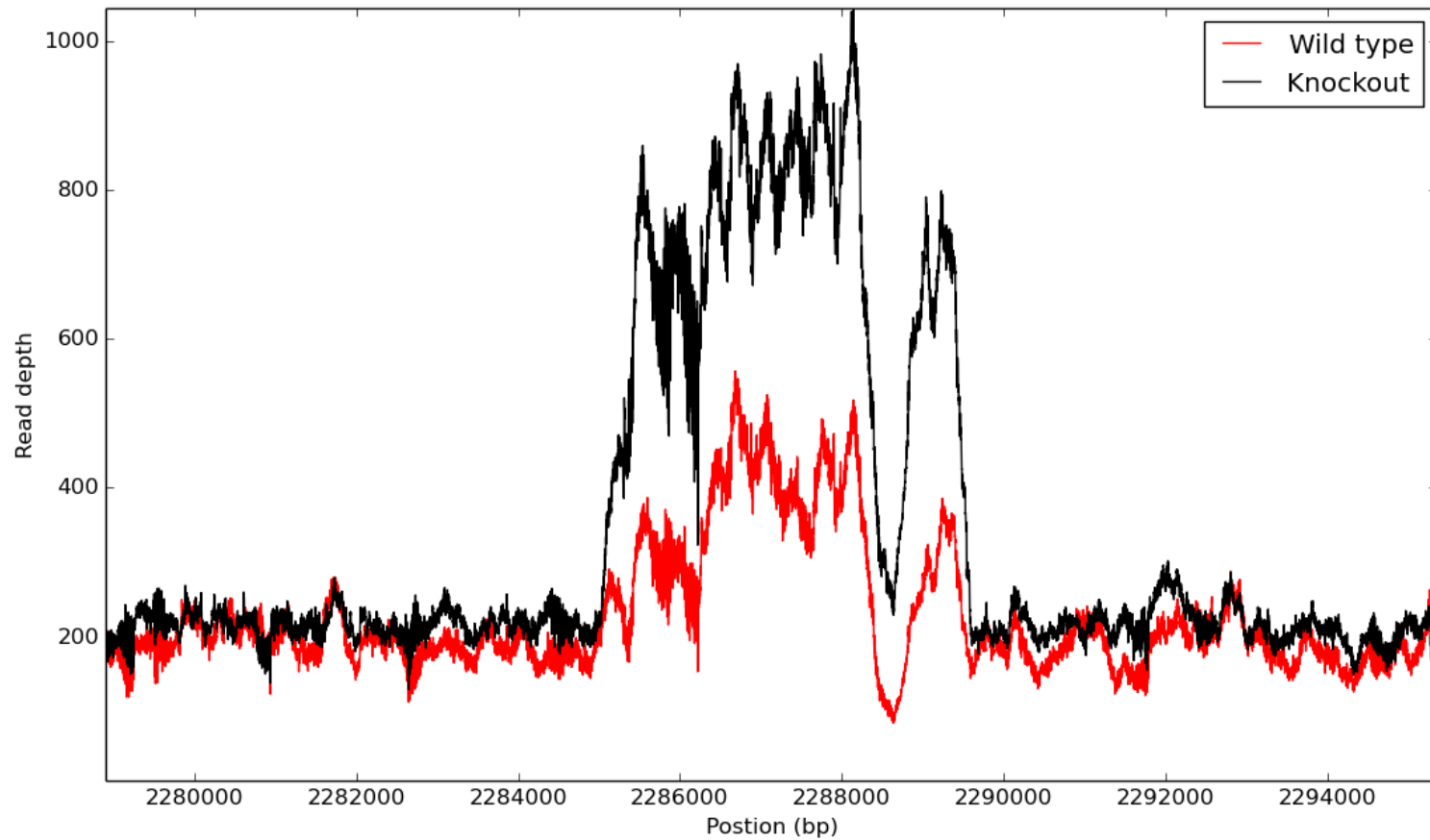
Biological

# Expected effects missing



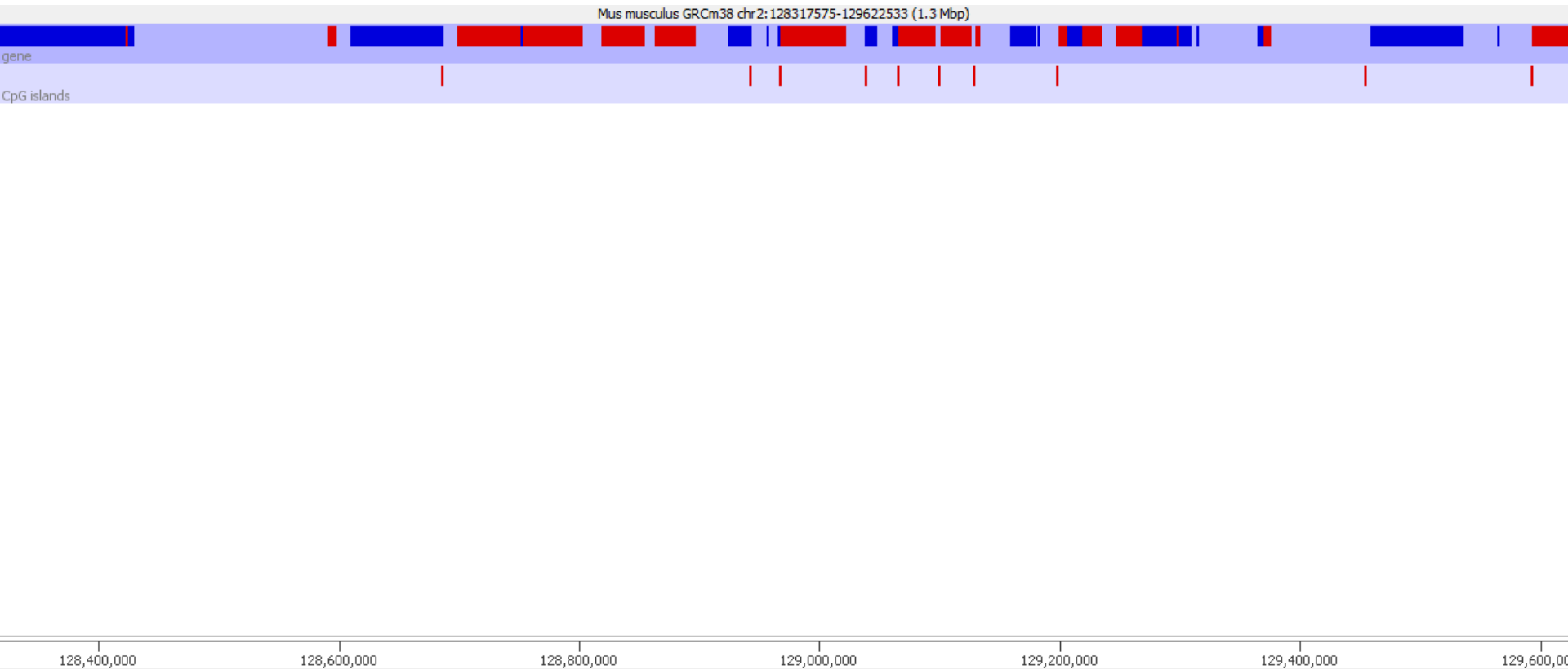
Biological

# Confounded effects



Biological

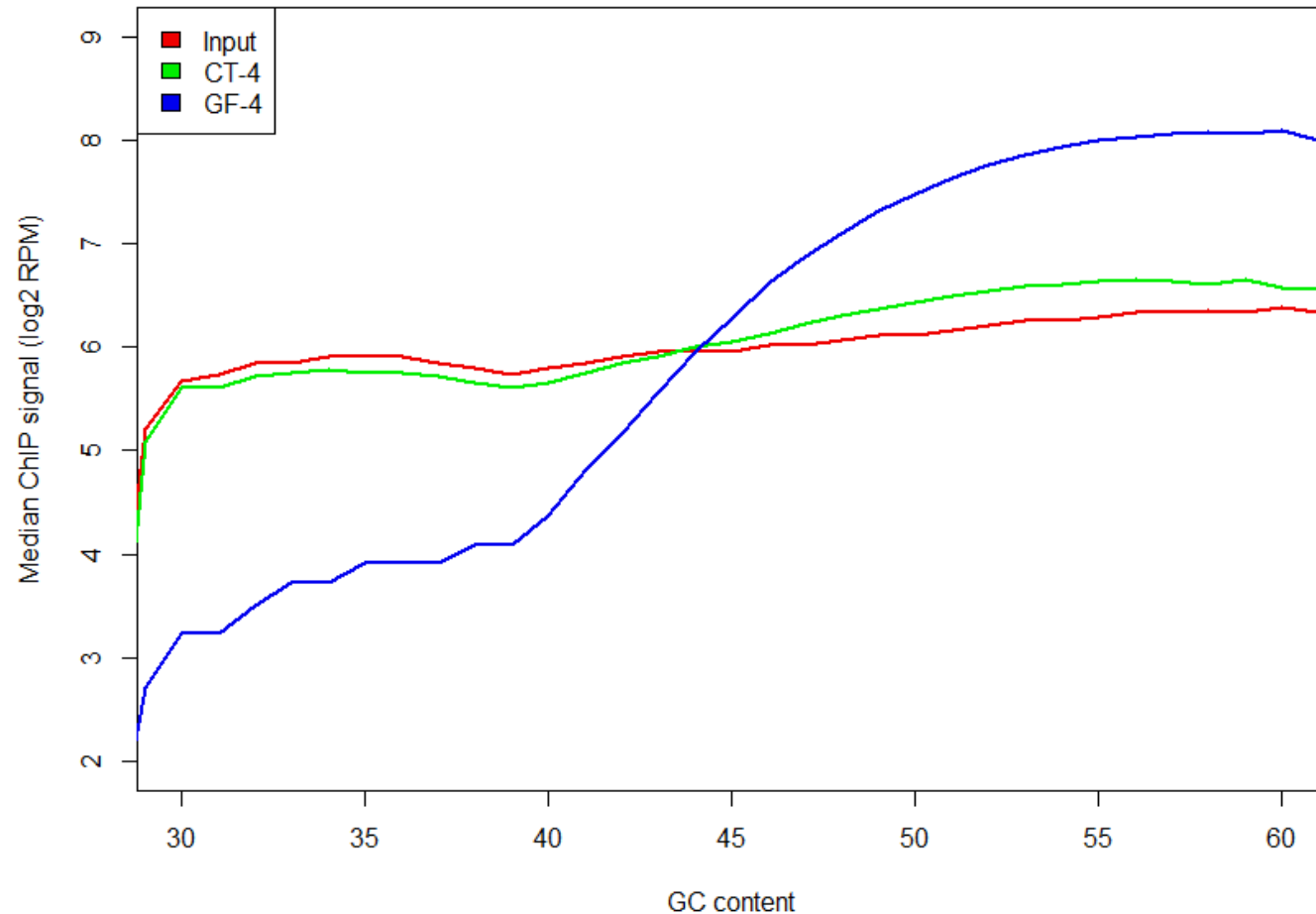
# ChIP doesn't behave



Biological

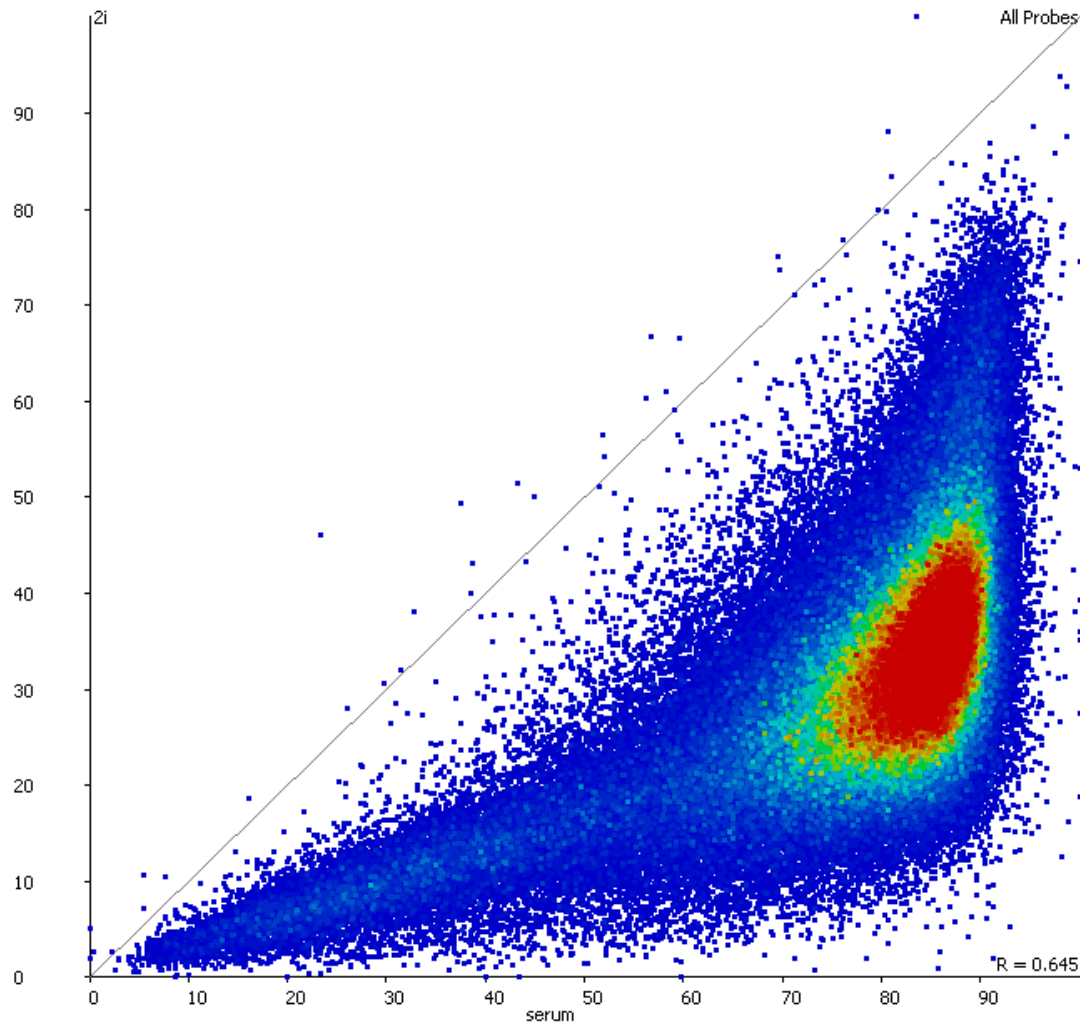
# ChIP doesn't behave

ChIP vs GC Content



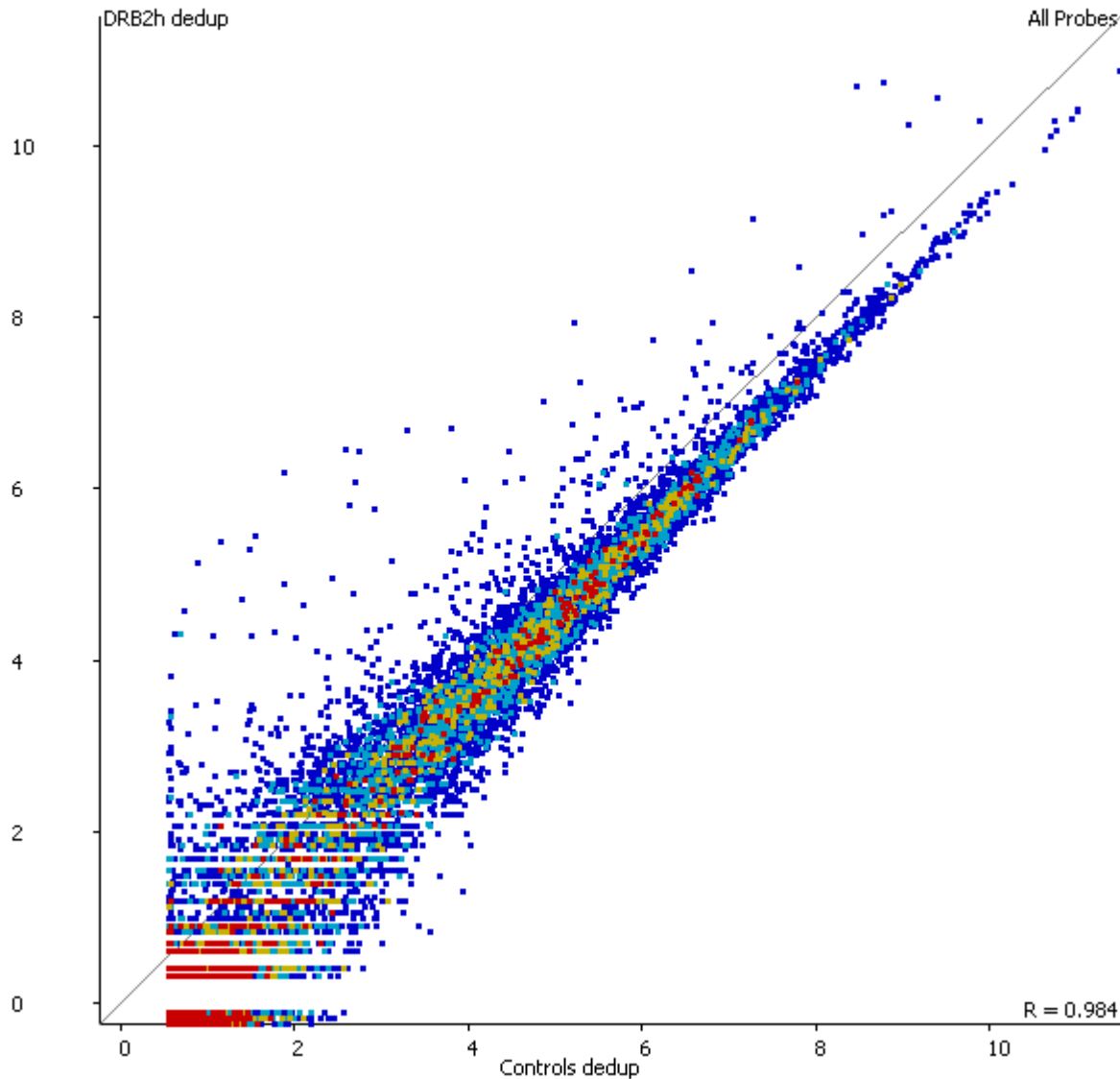
Biological

# Methylation doesn't behave



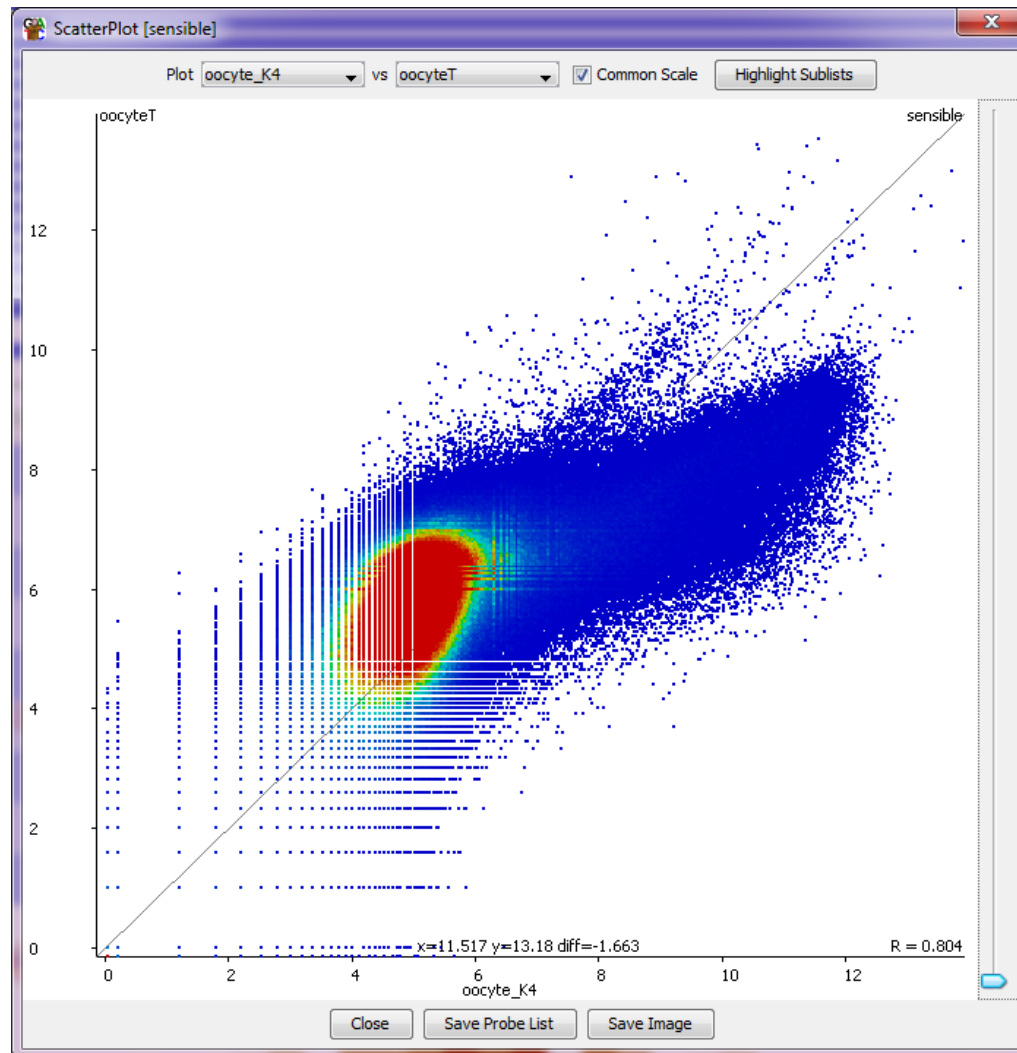
Biological

# RNA-Seq doesn't behave



Biological

# Multiple subgroups





# Interpretation