

# ChIP-Seq Data Processing and QC

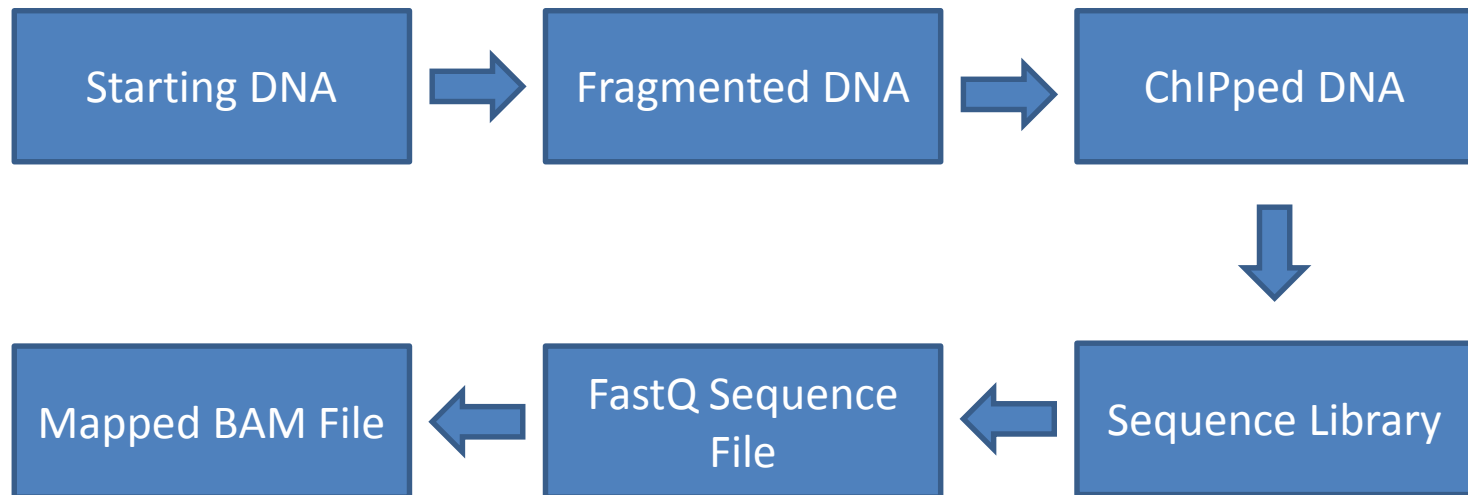
Simon Andrews

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

[@simon\\_andrews](#)

v2018-01

# Data Creation and Processing

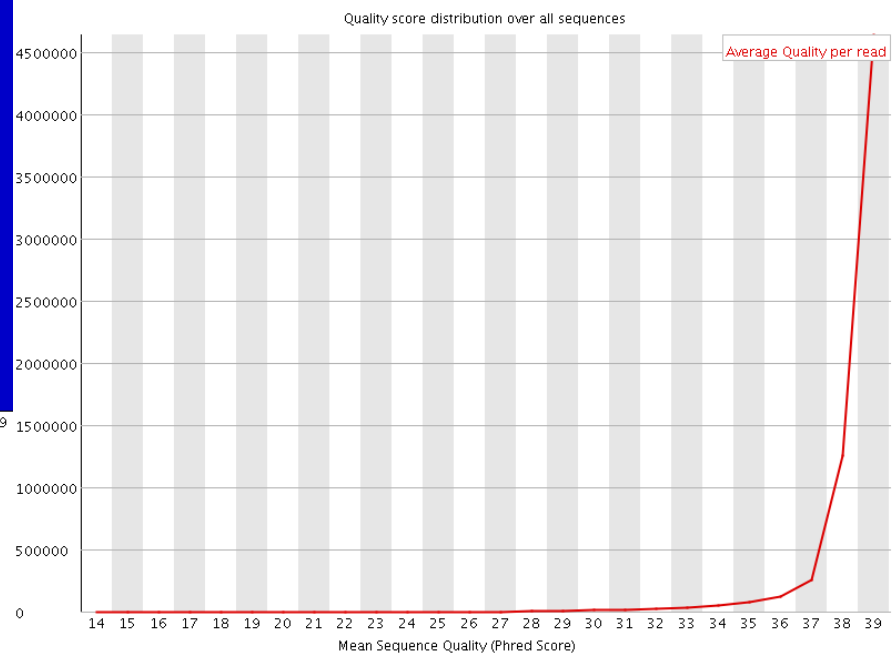
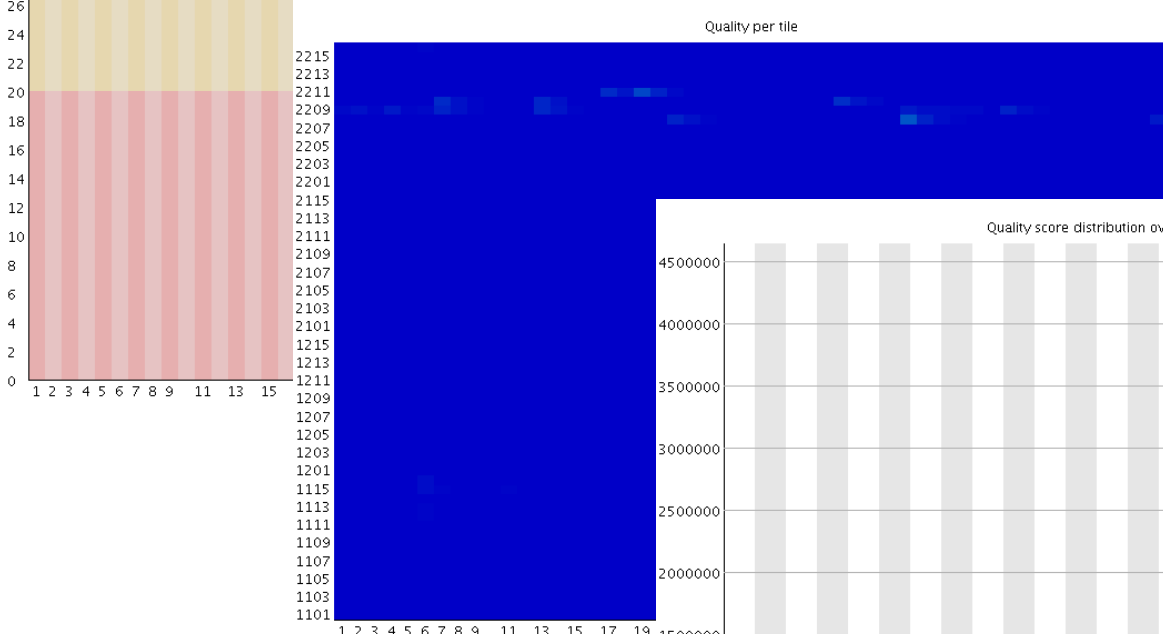
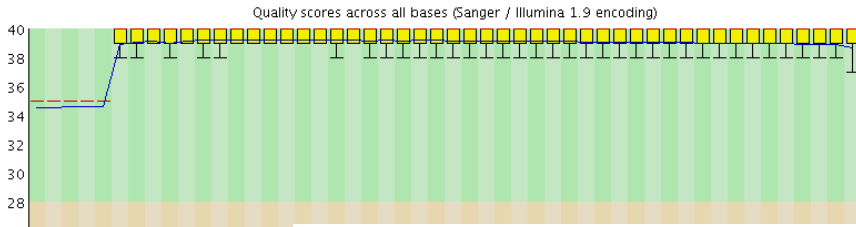


# A typical ChIP Library



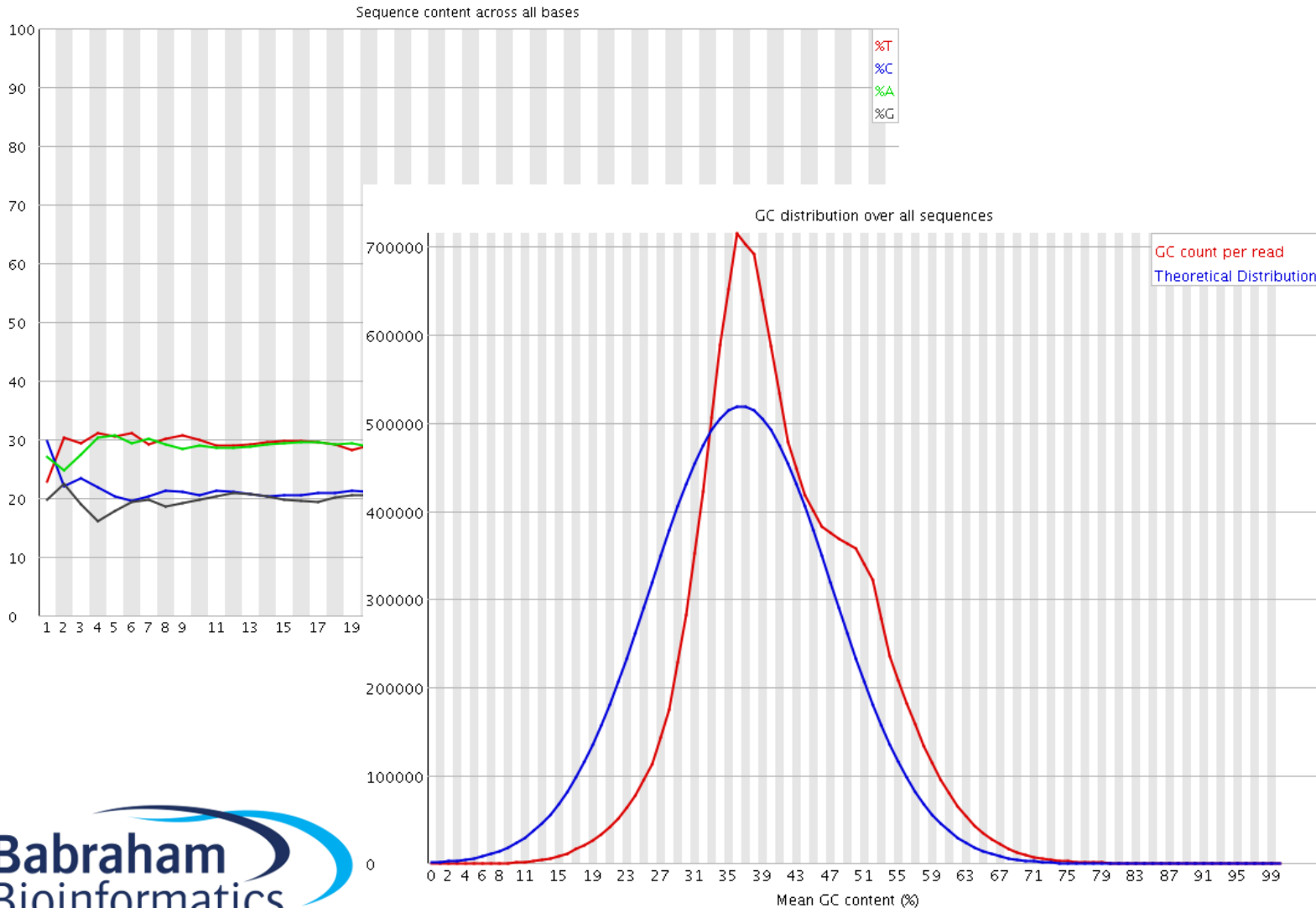
- Potential technical problems
  - Adapter contamination
  - PCR Duplication
- Potential biological problems
  - Lack of enrichment
  - Other selection biases

# QC of raw sequence Base Call Quality



# QC of raw sequence

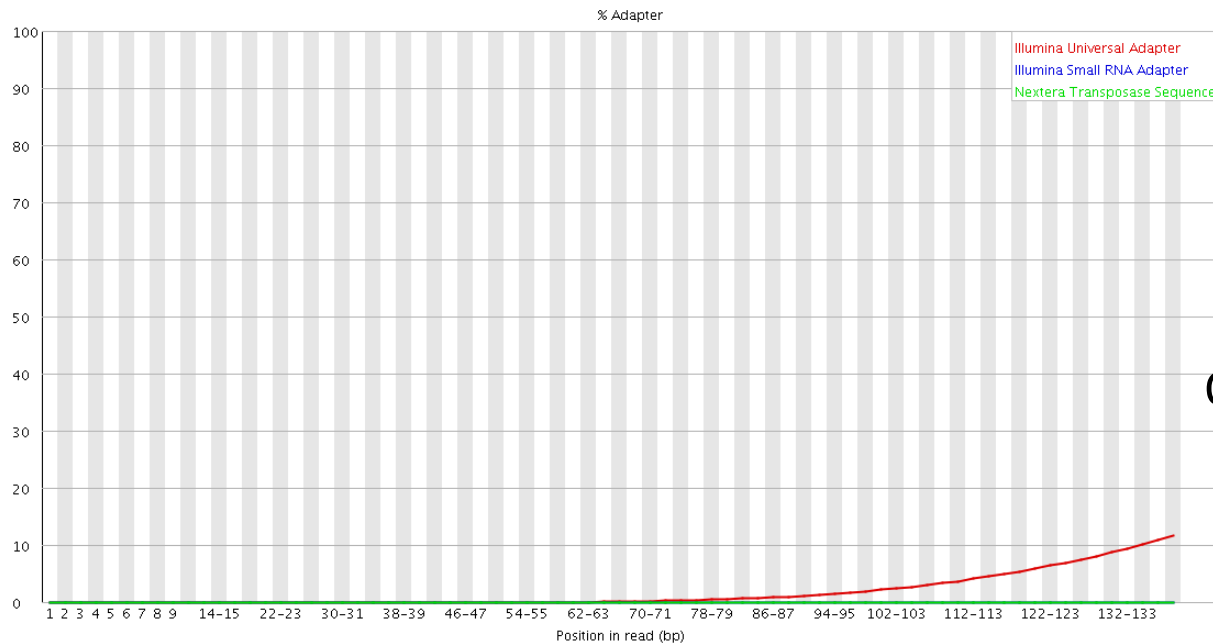
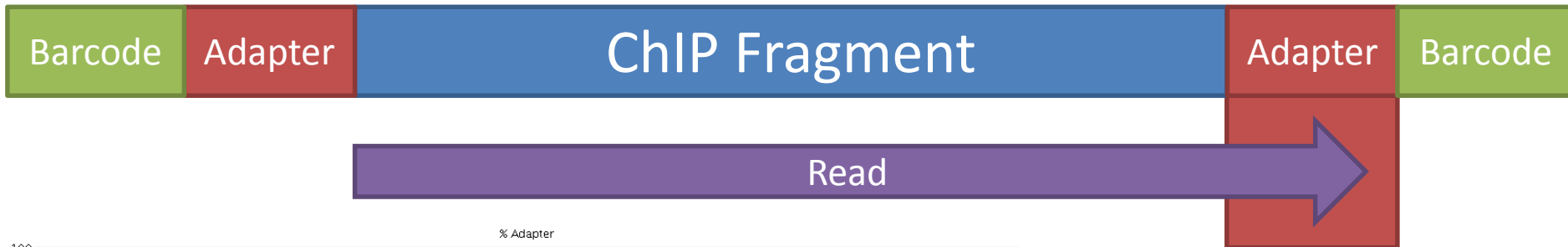
## Sequence Composition





# QC of raw sequence

## Adapter Contamination



**Trim Galore!**  
Quality and Adapter Trimming

# Mapping ChIP Data

- All regions should be linear genomic stretches
- Standard genomic aligners are fine
  - Bowtie2 <http://bowtie-bio.sourceforge.net/bowtie2/>
  - BWA <http://bio-bwa.sourceforge.net/>



# Example Bowtie2 Mapping

- Create Genome Index (once - slow!)

```
bowtie2-build yeast_genome.fa yeast_index
```

- Map a single FastQ file

```
bowtie \  
-x yeast_index \  
-U data.fastq.gz \  
| samtools view \  
-bS \  
-o data.bam
```

# Post Alignment QC Mapping Statistics

```
sample1: 2264052 (14.66%) aligned exactly 1 time
sample2: 2698005 (18.79%) aligned exactly 1 time
sample3: 13434392 (67.08%) aligned exactly 1 time
sample4: 1108477 (6.70%) aligned exactly 1 time
sample5: 2143911 (17.58%) aligned exactly 1 time
sample6: 2980154 (13.98%) aligned exactly 1 time
```

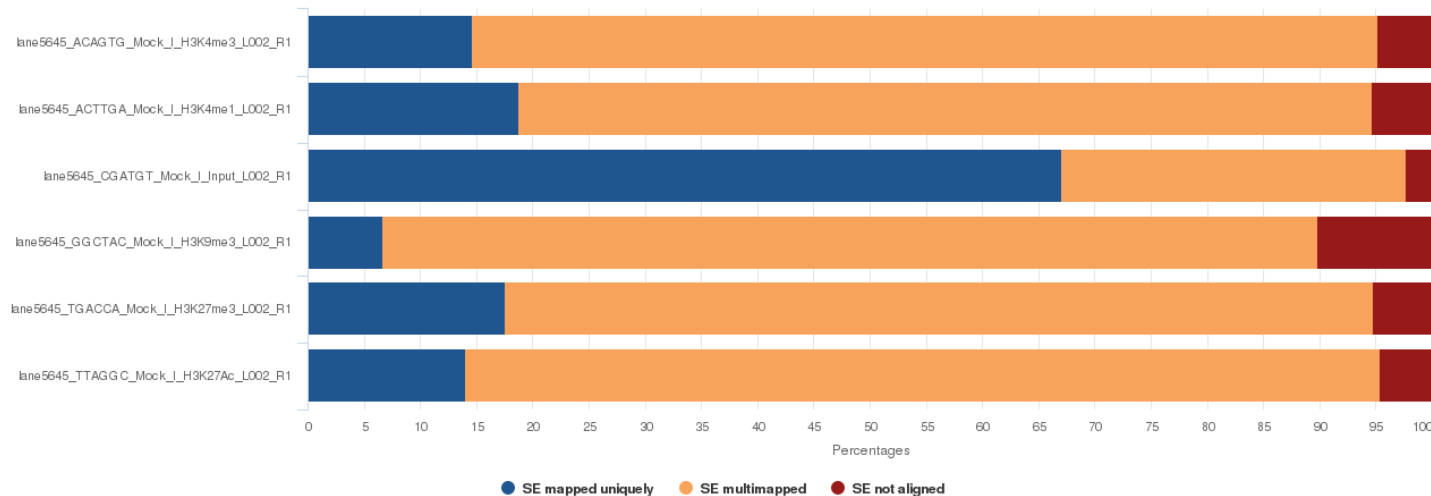
## Bowtie 2

Bowtie 2 is an ultrafast and memory-efficient tool for aligning sequencing reads to long reference sequences.

Number of Reads Percentages

Bowtie 2 SE Alignment Scores

Export Plot



Created with MultiQC

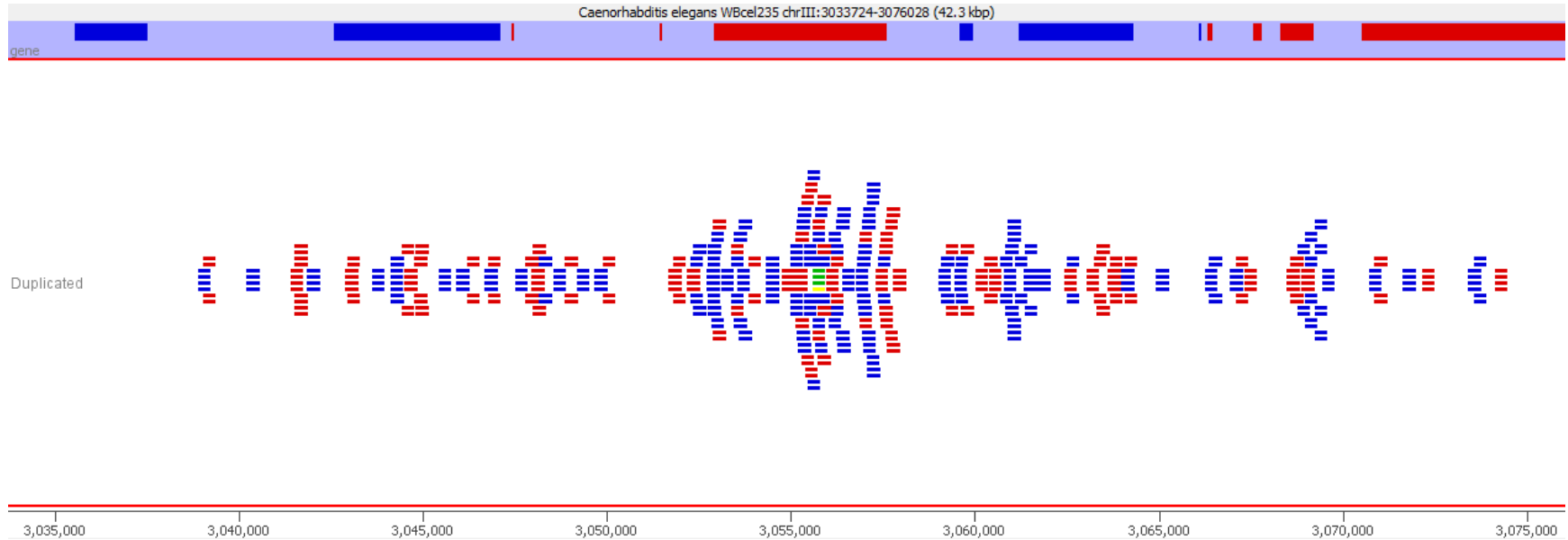
# Post Alignment Processing

## MAPQ Filtering

- ChIP-Seq relates sequences to positions in a reference genome
- You need to be confident that the reported position is correct
- Filtering on MAPQ value (likelihood of reported position being incorrect) is an easy way to do this
- MAPQ filtering should be performed in most cases

```
samtools view -q 20 -b -o filtered.bam data.bam
```

# Post Alignment Processing Deduplication



```
java -jar picard.jar SortSam \  
  INPUT=filtered.bam \  
  OUTPUT=sorted.bam \  
  SORT_ORDER=coordinate
```

```
java -jar picard.jar MarkDuplicates \  
  INPUT=sorted.bam \  
  OUTPUT=dedup.bam \  
  METRICS_FILE=metrics.txt
```

**DO NOT DEDUPLICATE AS A MATTER OF COURSE! THINK FIRST!**

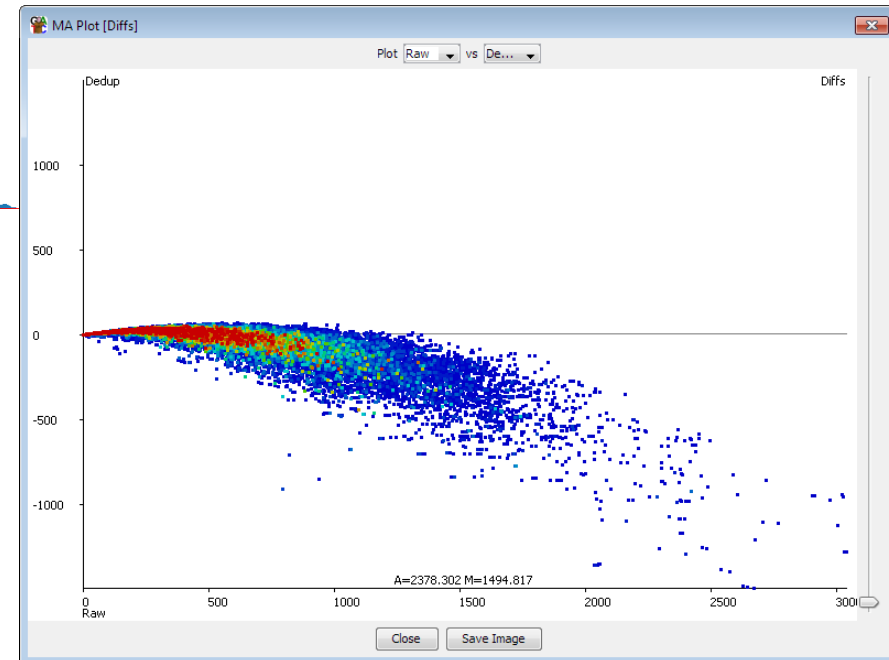
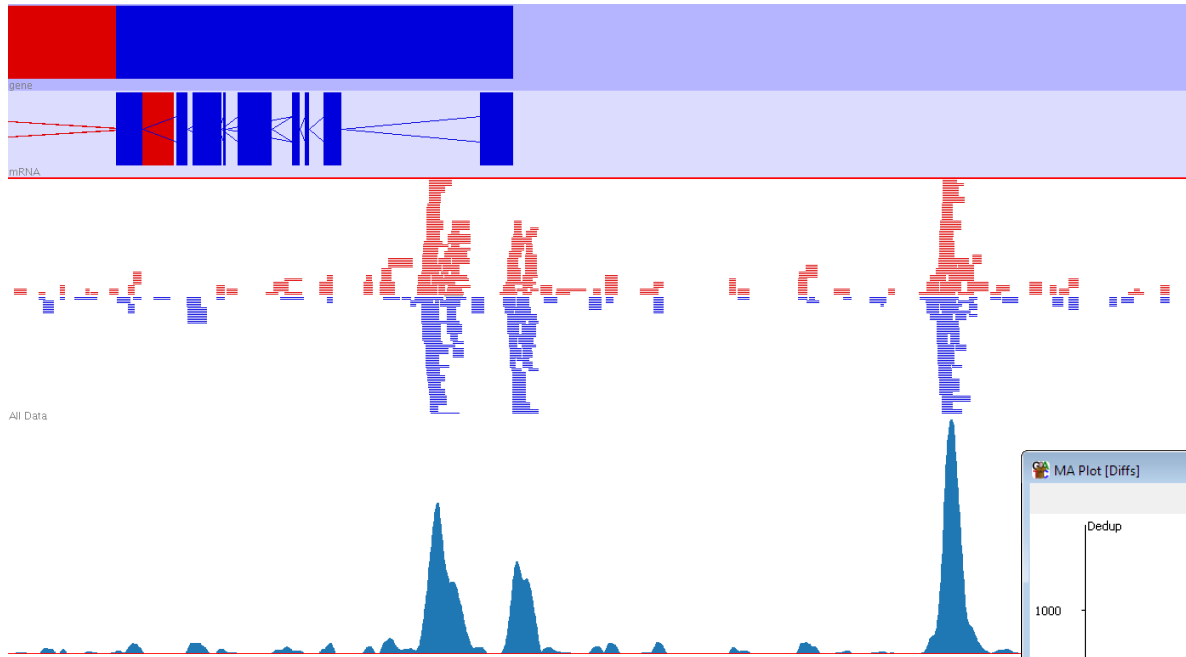
# To Deduplicate or Not?

- Deduplication can make enrichment visually clearer and help to spot truly enriched regions
- Why not just deduplicate everything?
  - Quantitation compression
  - Repeat enrichment

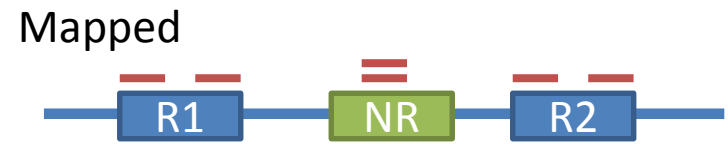
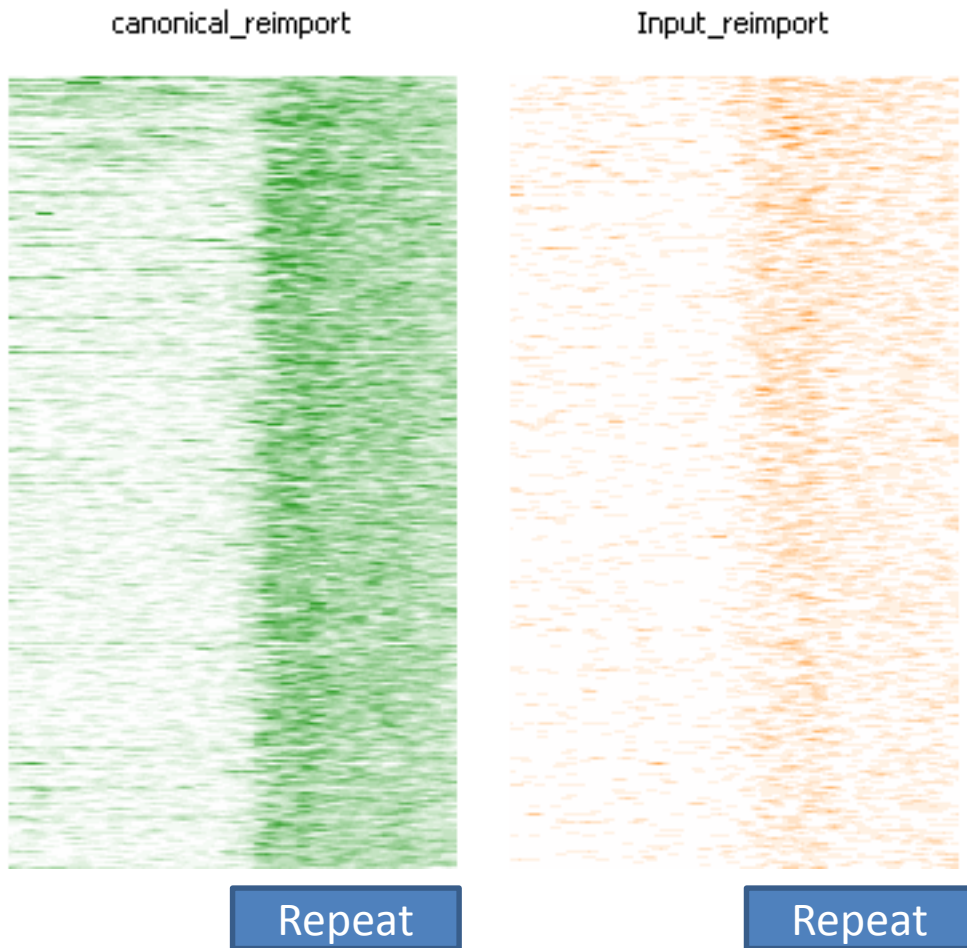
# Deduplication



# Quantitation Compression from Deduplication



# Repeat Enrichment from Deduplication

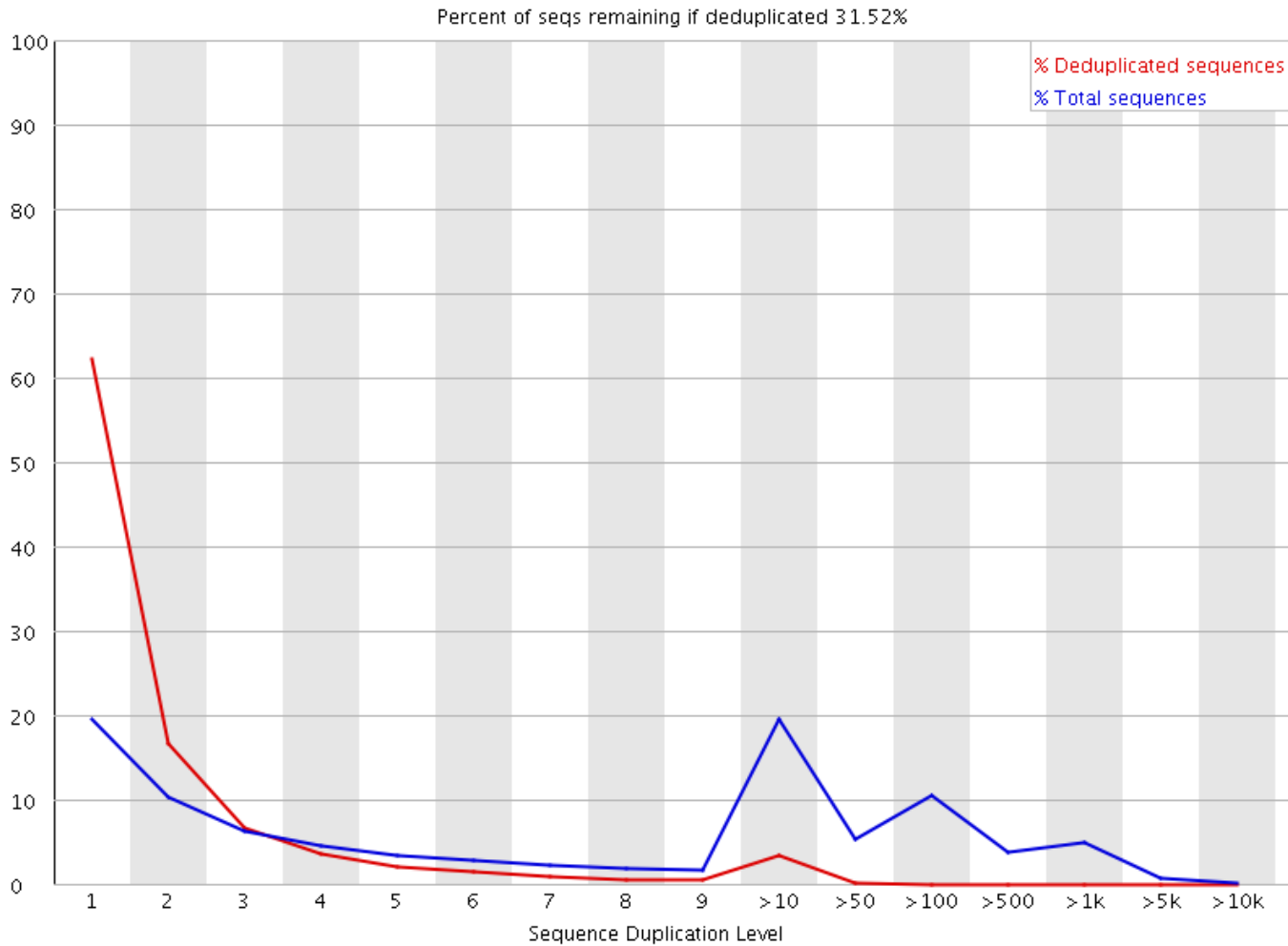


Peak callers (MACS for example) deduplicate internally, so you don't have to consciously do this.

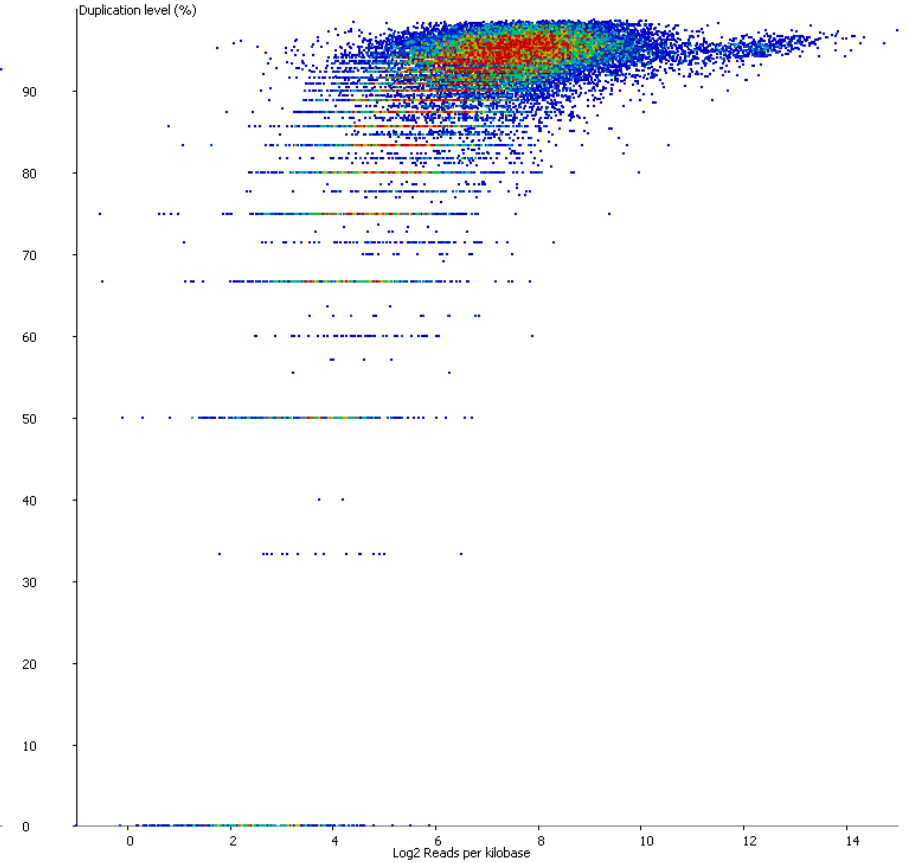
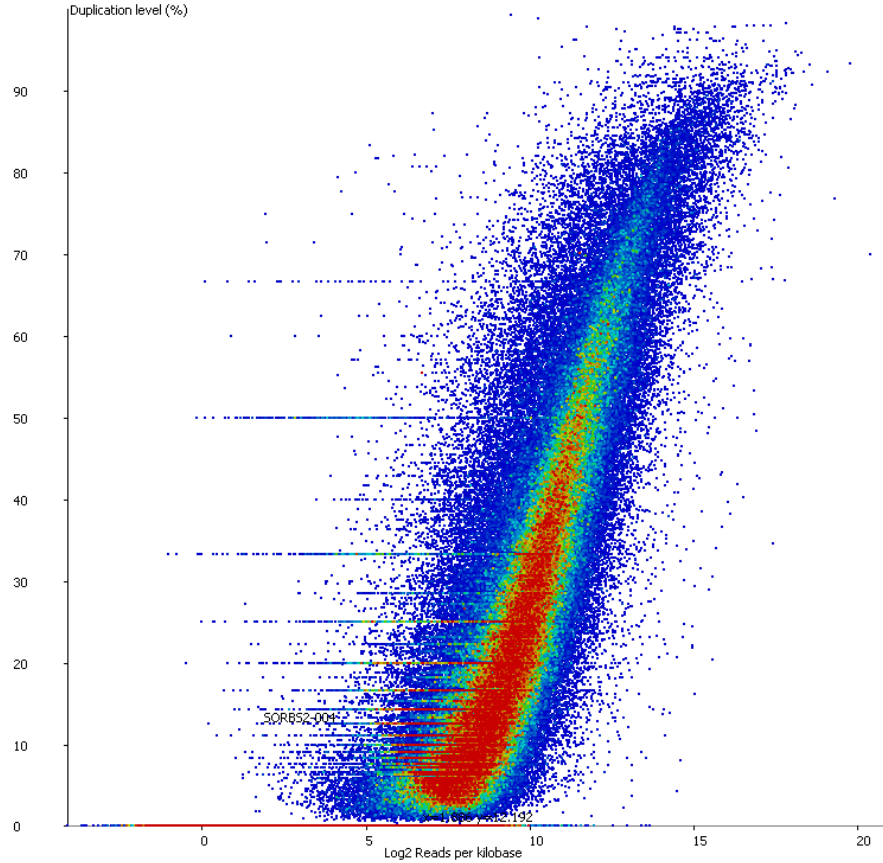
Only avoided by using uniquely mapped reads.



# Assessing Duplication



# Assessing Duplication



# Standard Processing Workflow

