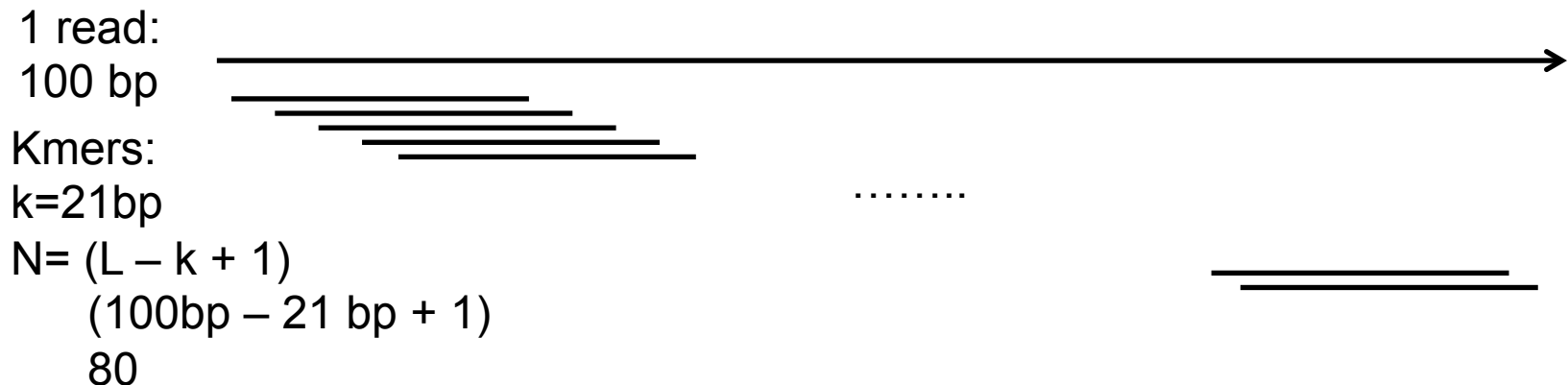


K-mer Analyses, Contamination Detection, And Mapping Based Analyses



- What is a k-mer?
 - A k-mer is a sequence of nucleotides of length k.
 - Examples of a 6-mer
 - ACGTCT
 - TGACTA
 - GATCCC
- A read of length L has L-k+1 k-mers.



K-mer Analysis

```
ACTCAGGATTATTCATACATAAATAGCCGGGG
      ATTCATACATAAATAGCCGGGG
ACTCAGGATTATTCATACATA
      ATTCATACATA
      TTCATACATAA
ACTCAGGATTIA TCATACATAAAT
CTCAGGATTIAT CATAACATAATA
TCAGGATTIATT ATACATAAATAG
CAGGATTIATTC TACATAAATAGC
AGGATTIATTCA ACATAAATAGCC
GGATTATTCAT CATAAATAGCCG
GATTATTCATA ATAATAGCCGG
ATTATTCATAC TAATAGCCGGG
TTATTCATACA AATAGCCGGGG
TATTCATACAT
      ATTCATACATA
```

- The sequence has 21 distinct 11-mers.
- The reads have 21 distinct 11-mers even though 22 11-mers are generated.
- The frequency of a distinct k-mer increases where reads overlap

K-mer Analysis

```
ACTCAGGATTATTCATACATAAATAGCCGGGG
      ATTCATACATAAATAGCCGGGG
ACTCAGGATTATTCATACATA
      GGATTATTCATACATAAATAGC
      ATTCATACATA
      TTCATACATAA
ACTCAGGATTA TCATACATAAT
CTCAGGATTAT CATAcataATA
TCAGGATTATT ATACATAATAG
CAGGATTATTC TACATAATAGC
AGGATTATTCA ACATAAATAGCC
GGATTATTCAT CATAAATAGCCG
GATTATTCATA ATAATAGCCGG
ATTATTCATAC TAATAGCCGGG
TTATTCATACA AATAGCCGGGG
TATTCATACAT
      ATTCATACATA
GGATTATTCAT
GATTATTCATA
ATTATTCATAC
TTATTCATACA
TATTCATACAT
      ATTCATACATA
      TTCATACATAA
      TCATACATAAT
      CATAcataATA
      ATACATAATAG
      TACATAATAGC
```

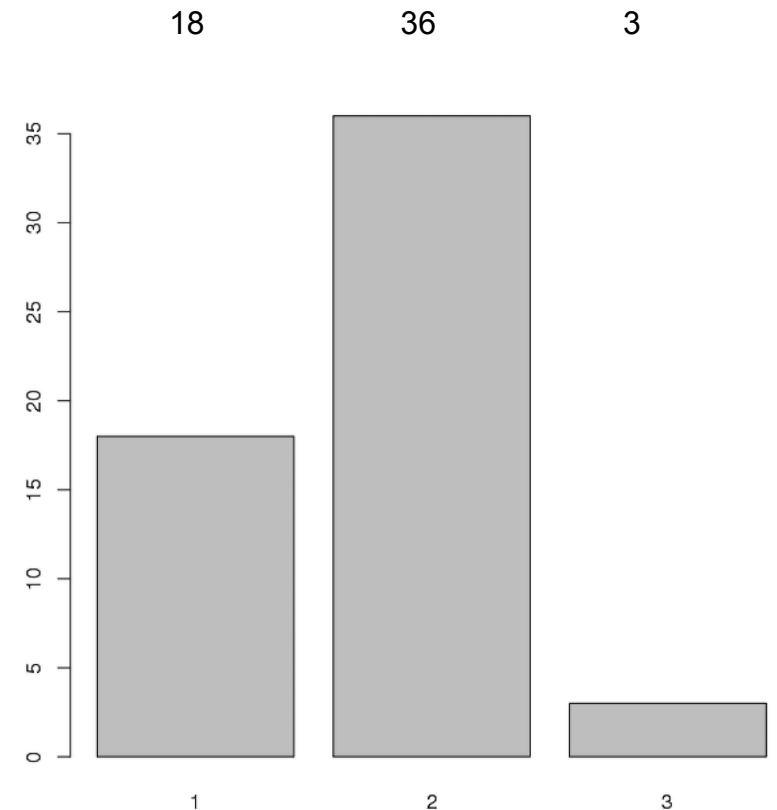
- The sequence has 21 distinct 11-mers.
- The reads have 21 distinct 11-mers even though 33 11-mers are generated.
- The frequency of a distinct k-mer increases where reads overlap

K-mer Analysis

ACTCAGGATTATTCATACATAAATAGCCGGGGTACATTTTCATTAAGCGGCGATTACGATTACATACGT

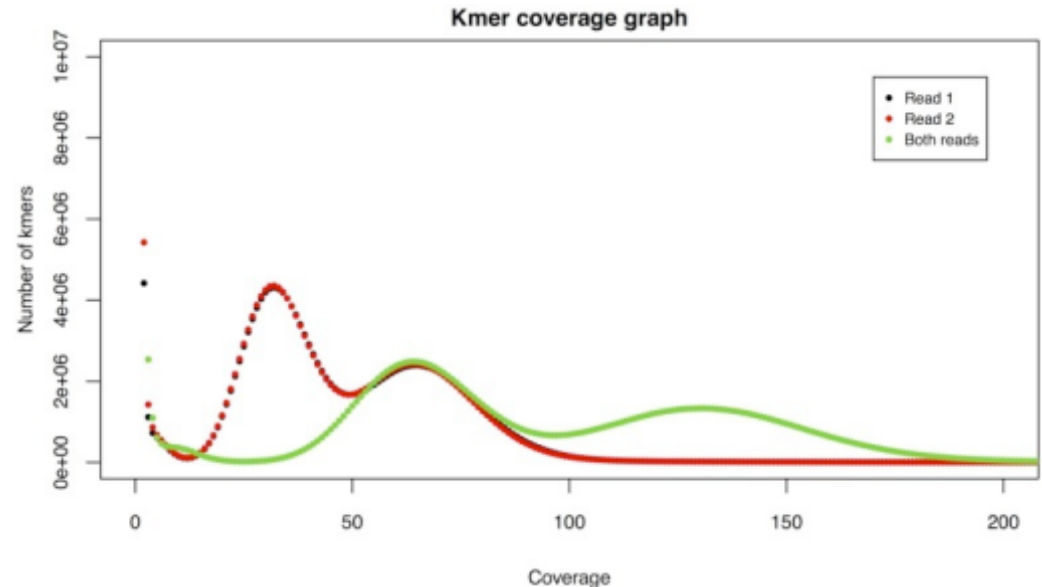
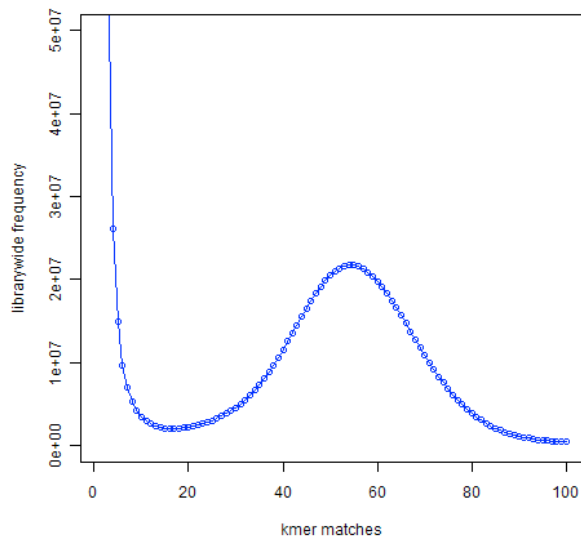
ATTCATACATAAATAGCCGGGG CATTTCATTAAGCGGCGATTA
 ACTCAGGATTATTCATACATA GCCGGGGTACATTTTCATTAAG GGCGATTACGATTACATACGT
 GGATTATTCATACATAAATAGC GGGGTACATTTTCATTAAGCGG
 ACATAATAGCCGGGGTACATT AGCGGCGATTACGATTACATA

ATTCATACATA
 TTCATACATAA
 ACTCAGGATTA TCATACATAAT
 CTCAGGATTAT CATAcataAATA GGGGTACATTT GGCGATTACGA
 TCAGGATTATT ATACATAATAG GGGTACATTT GCGATTACGAT
 CAGGATTATTC TACATAATAGC GGTACATTTCA CGATTACGATT
 AGGATTATTCA ACATAATAGCC GTACATTTTCAT GATTACGATTA
 GGATTATTCAT CATAATAGCCG TACATTTTCATT ATTACGATTAC
 GATTATTCATA ATAATAGCCGG ACATTTTCATTA TTACGATTACA
 ATTATTCATAC TAATAGCCGGG CATTTCATTA TACGATTACAT
 TTATTCATACA AATAGCCGGGG ATTTTCATTAAG ACGATTACATA
 TATTCATACAT TTTCATTAAGC CGATTACATAC
 ATTCATACATA GCCGGGGTACA TTCATTAAGCG GATTACATACG
 GGATTATTCAT CCGGGGTACAT TCATTAAGCGG ATTACATACGT
 GATTATTCATA CGGGGTACATT
 ATTATTCATAC GGGGTACATTT AGCGGCGATTA
 TTATTCATACA GGGTACATTT CCGGCGATTAC
 TATTCATACAT GGTACATTTCA CCGGCGATTACG
 ATTCATACATA GTACATTTTCAT GGCGATTACGA
 TTCATACATAA TACATTTTCATT GCGATTACGAT
 TCATACATAAT ACATTTTCATTA CGATTACGATT
 CATAcataAATA CATTTCATTA GATTACGATTA
 ATACATAATAG ATTTTCATTAAG ATTACGATTAC
 TACATAATAGC TTACGATTACA
 TACGATTACAT
 ACGATTACATA
 ACATAATAGCC CATTTCATTA
 CATAATAGCCG ATTTTCATTAAG
 ATAATAGCCGG TTTCATTAAGC
 TAATAGCCGGG TTCATTAAGCG
 AATAGCCGGGG TCATTAAGCGG
 ATAGCCGGGGT CATTAGCGGC
 TAGCCGGGGTA ATTAAGCGGCG
 AGCCGGGGTAC TTAAGCGGCGA
 GCCGGGGTACA TAAGCGGCGAT
 CCGGGGTACAT AAGCGGCGATT
 CGGGGTACATT AGCGGCGATTA



What can k-mers tell us?

- Plotting the frequencies of k-mers tells us how complicated our genome is.
 - There are a proportion of distinct k-mers that appear only once in the genome.
 - Distinct k-mers with count f (x-axis) vs Frequency of distinct k-mers with count f (y-axis).



K-mer Analysis

ACTCAGGATTATTCATACATAAATAGCCGGGG
ACTCCGGATTATTCATACATAAATAGCCGGGG

ACTCAGGATTATTCATACATA
ACTCCGGATTATTCATACATA
 ATTCATACATAAATAGCCGGGG
 ATTCATACATAAATAGCCGGGG

ATTCATACATA
TTCATACATAA

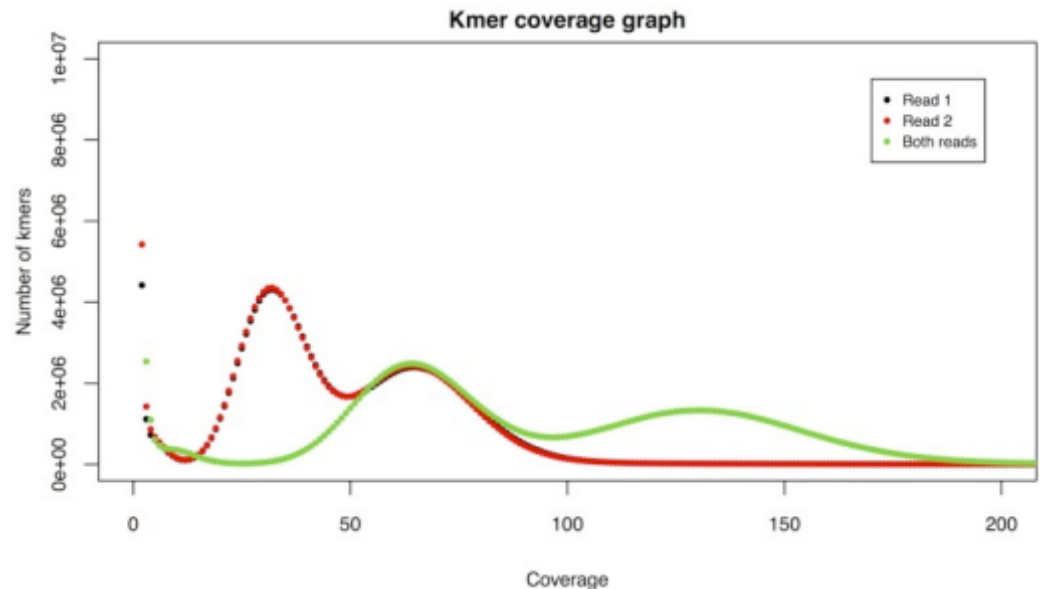
ACTCAGGATTA TCATACATAAT
CTCAGGATTAT CATAcataAATA
TCAGGATTATT ATACATAATAG
CAGGATTATTC TACATAATAGC
AGGATTATTCA ACATAATAGCC
GGATTATTCAT CATAAATAGCCG
GATTATTCATA ATAATAGCCGG
ATTATTCATAC TAATAGCCGGG
TTATTCATACA AATAGCCGGGG
TATTCATACAT

ATTCATACATA
ATTCATACATA

TTCATACATAA

ACTCCGGATTA TCATACATAAT
CTCCGGATTAT CATAcataAATA
TCCGGATTATT ATACATAATAG
CCGGATTATTC TACATAATAGC
CGGATTATTCA ACATAATAGCC
GGATTATTCAT CATAAATAGCCG
GATTATTCATA ATAATAGCCGG
ATTATTCATAC TAATAGCCGGG
TTATTCATACA AATAGCCGGGG
TATTCATACAT
ATTCATACATA

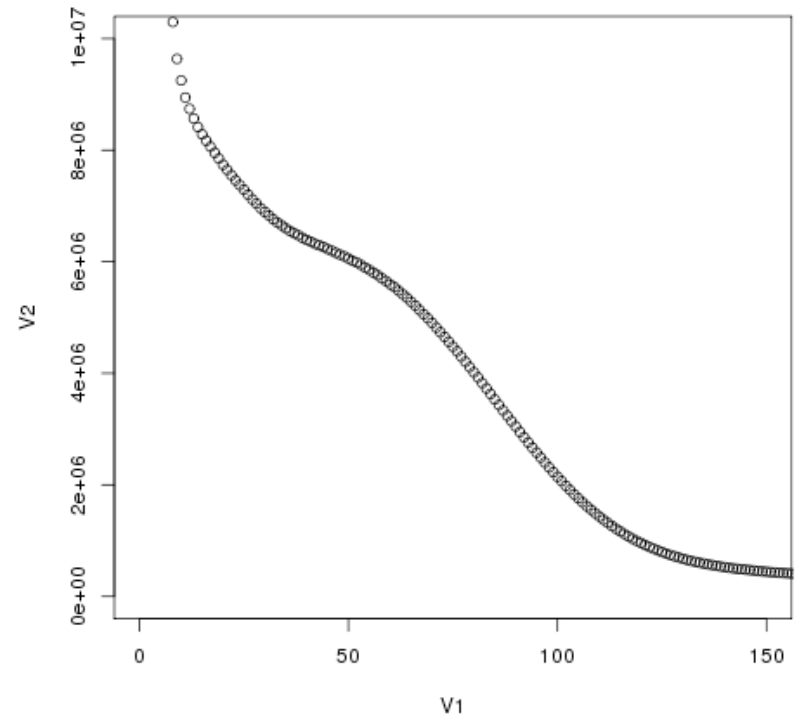
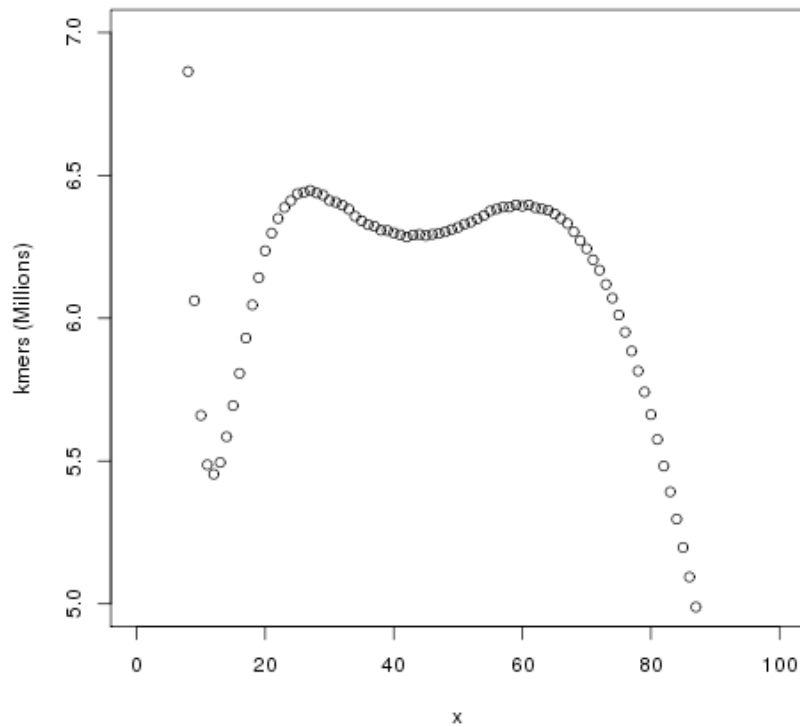
- Variation generates additional distinct k-mers
- Variants are at half the frequency of the rest of the distinct genome



- We can use coverage to estimate genome size
 - $N = M * L / (L - K + 1)$
 - N is Depth of Read Coverage
 - M is mean k-mer coverage
 - L is read length
 - K is k-mer size
 - $G = T / N$
 - G is the genome size
 - T is the total number of bases

Estimating Genome size

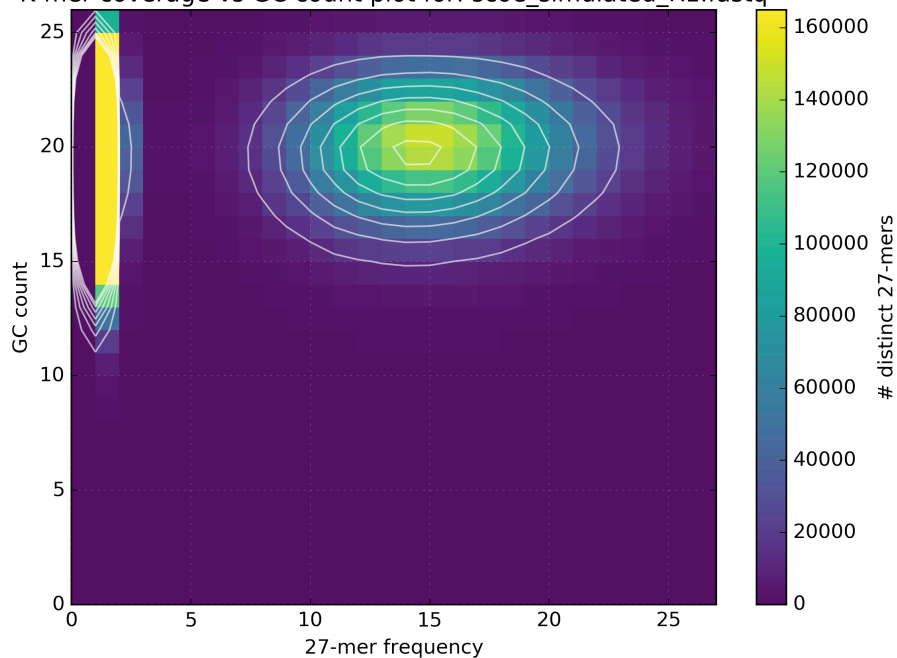
- Not so easy: estimating complexity



- Monitor GC Bias

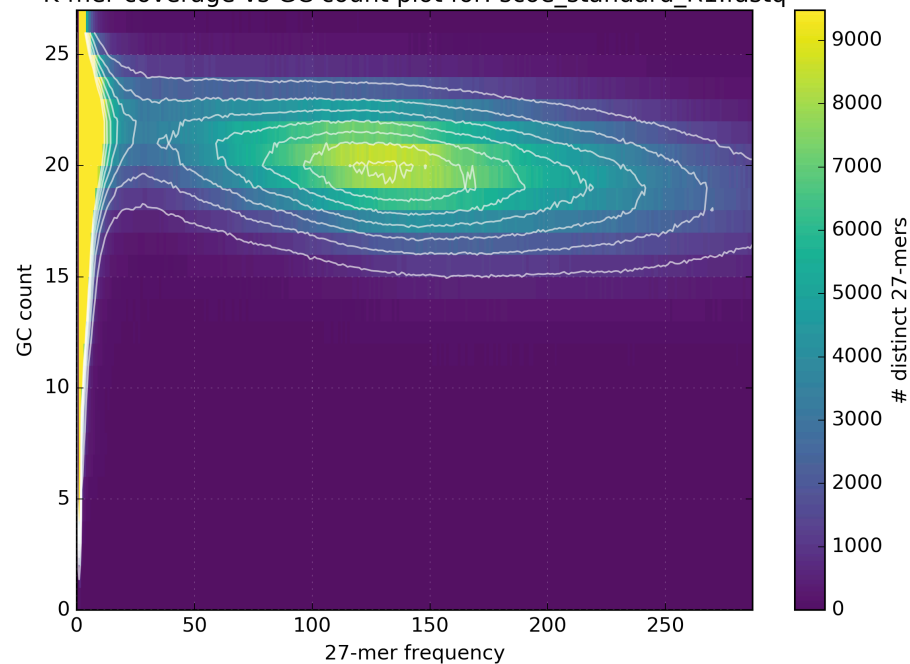
Simulated data

K-mer coverage vs GC count plot for: scoe_simulated_R1.fastq

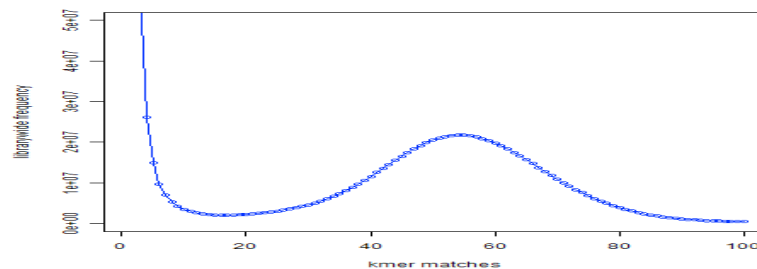


GC bias in Standard Illumina protocol

K-mer coverage vs GC count plot for: scoe_standard_R1.fastq

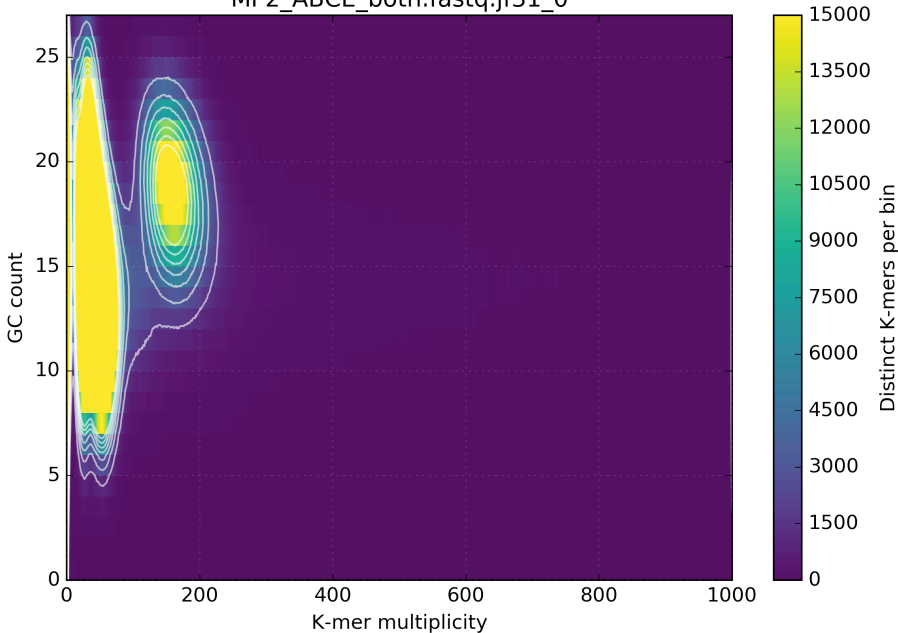


ACTCAGGATTA GC=4 Count=1
AATAGCCGGGG GC=7 Count=2

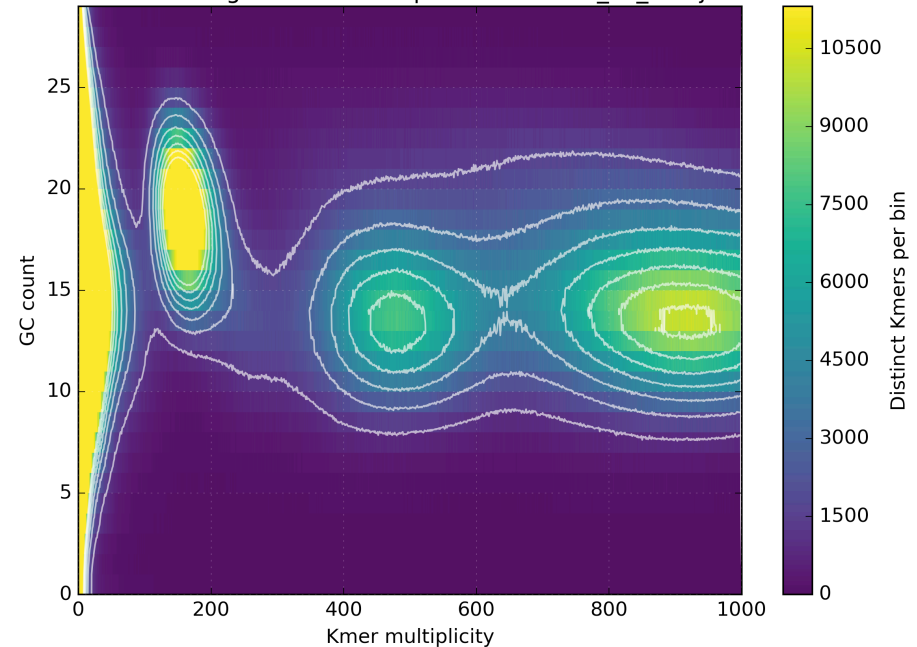


- Uncover contamination
 - Separate bacteria from eukaryote
 - Separate organelle from nuclear genome

K-mer coverage vs GC count plot for:
MP2_ABCE_both.fastq.jf31_0

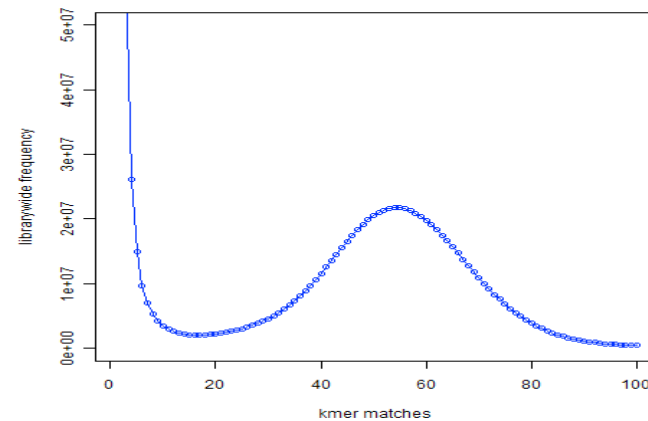
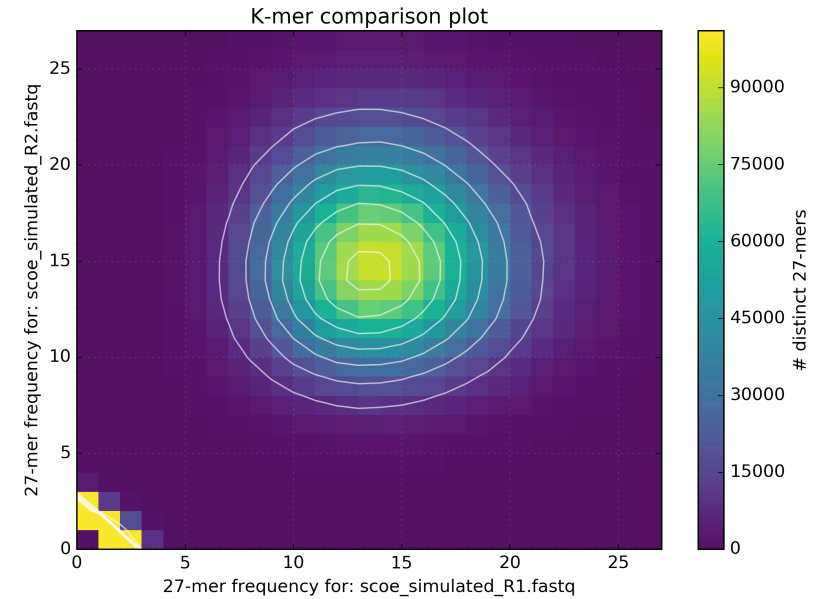
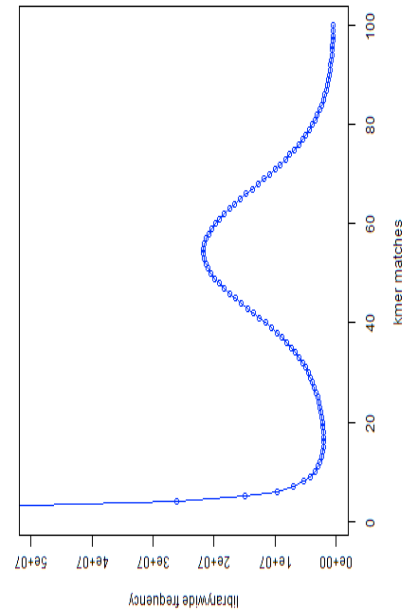


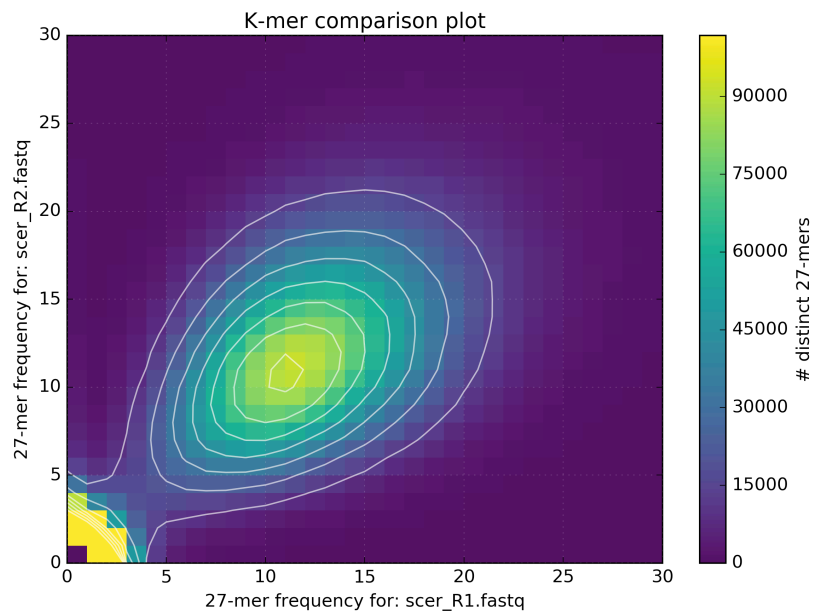
Kmer coverage vs GC count plot for: diatom_all_k31.jf31



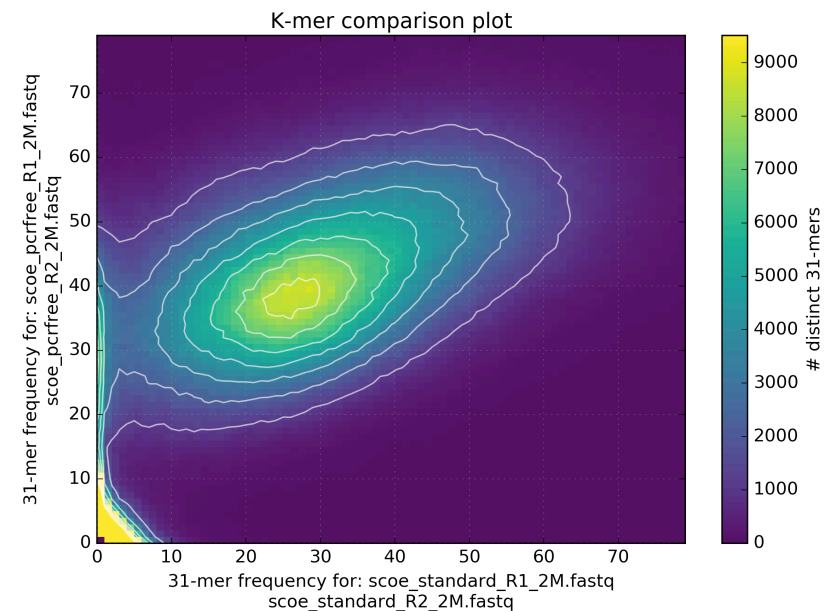
What can k-mers tell us?

- Comparing k-mer counts between data reveals biases
 - R1 vs R2
 - Lib1 vs Lib2





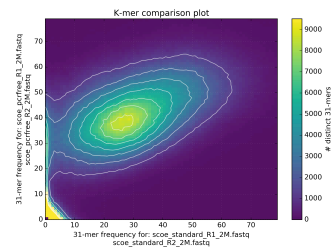
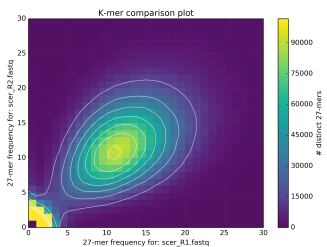
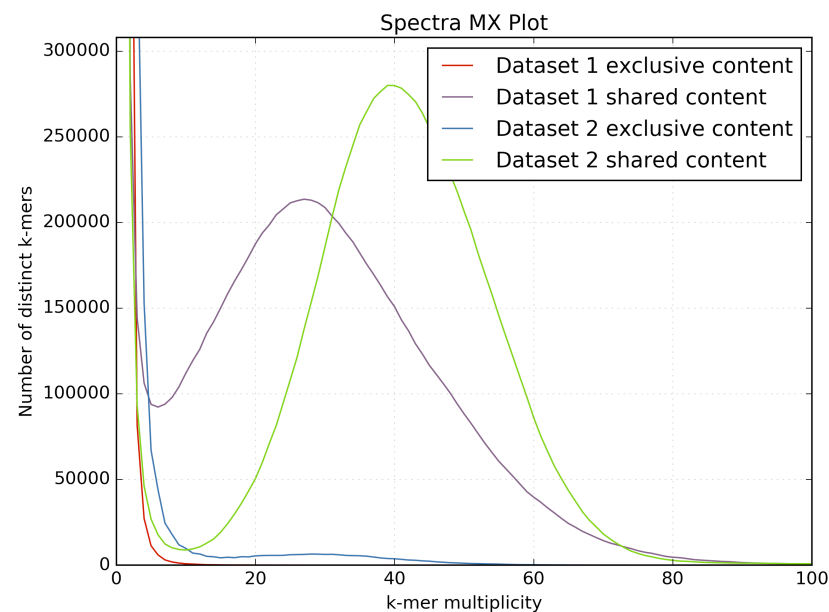
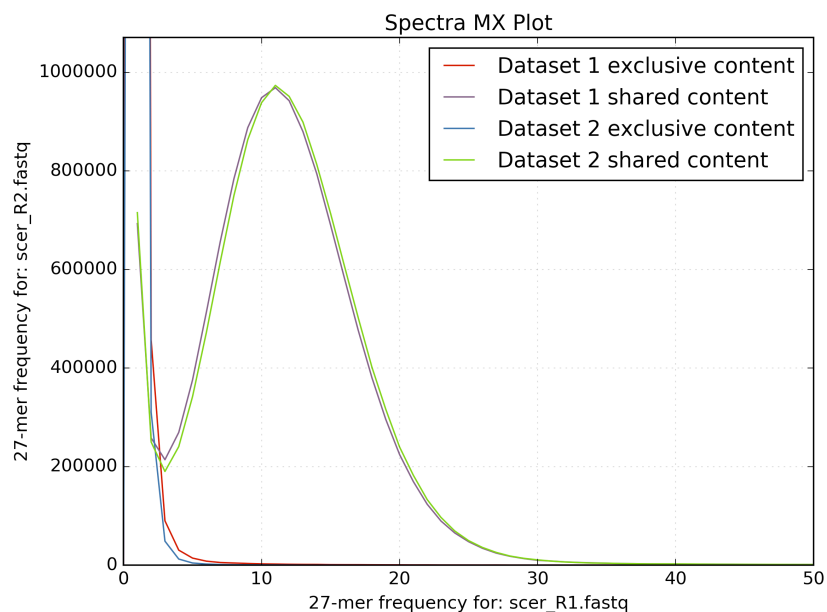
R1 vs R2



Standard vs PCR free

- PCR free captures data missing in standard protocol

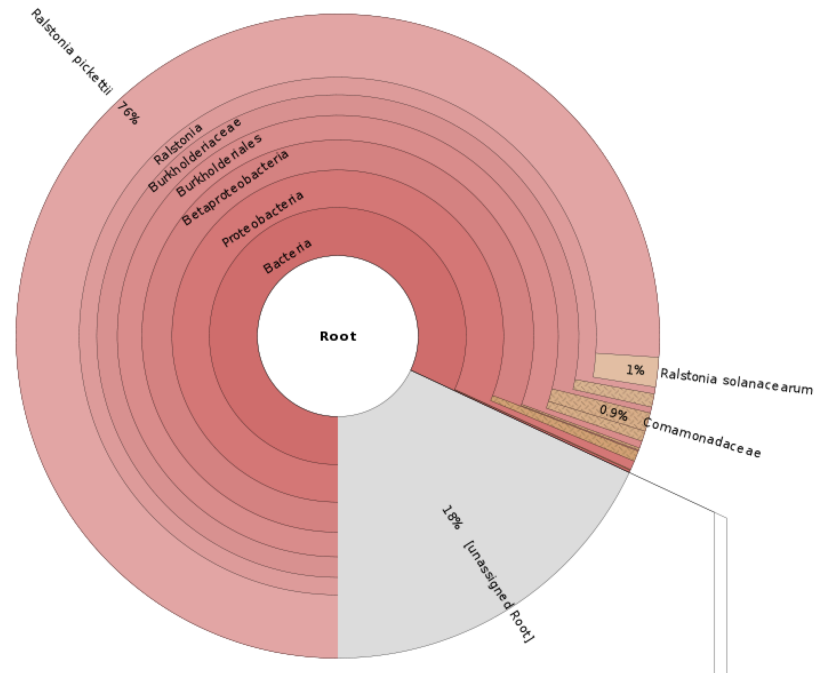
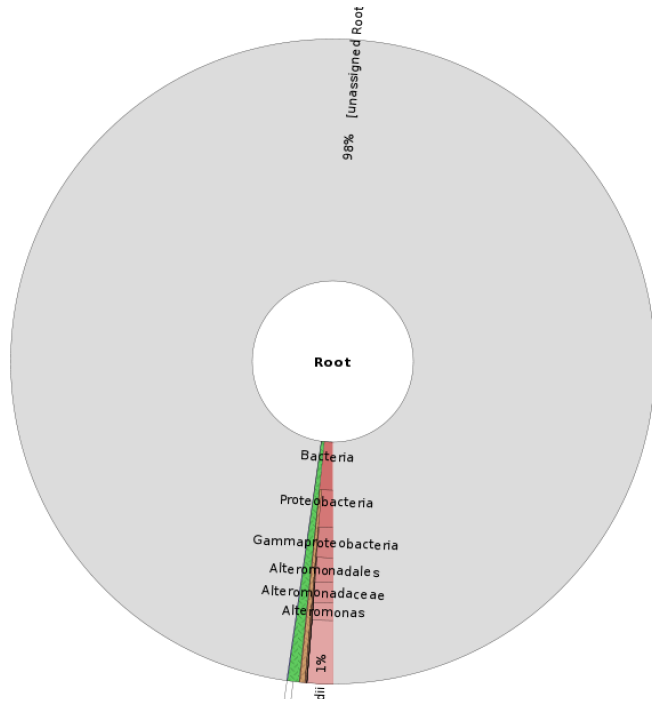
Data comparison

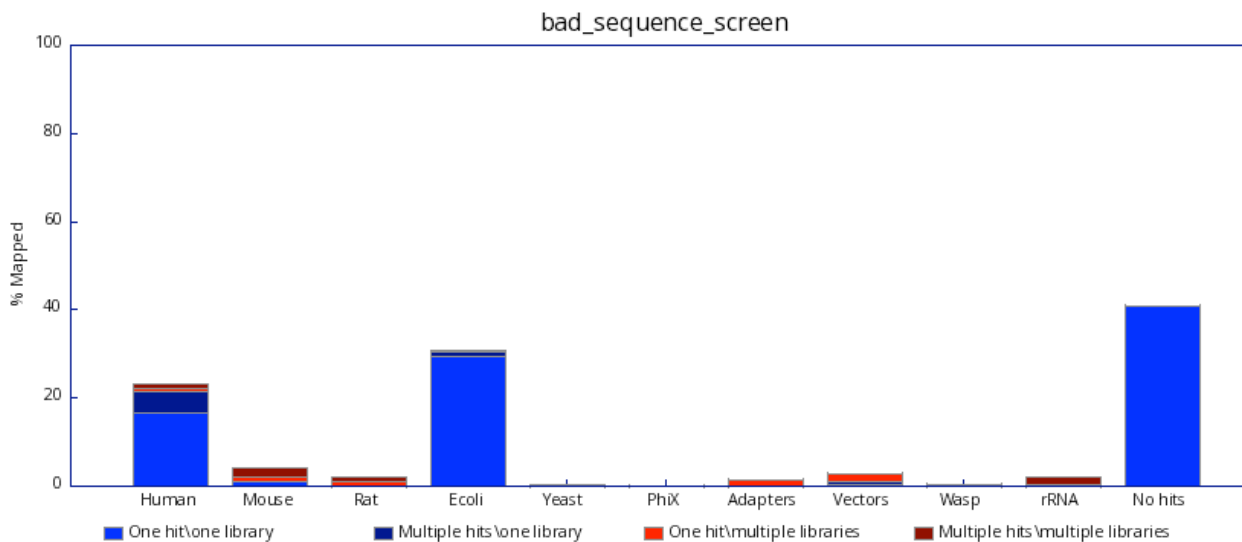
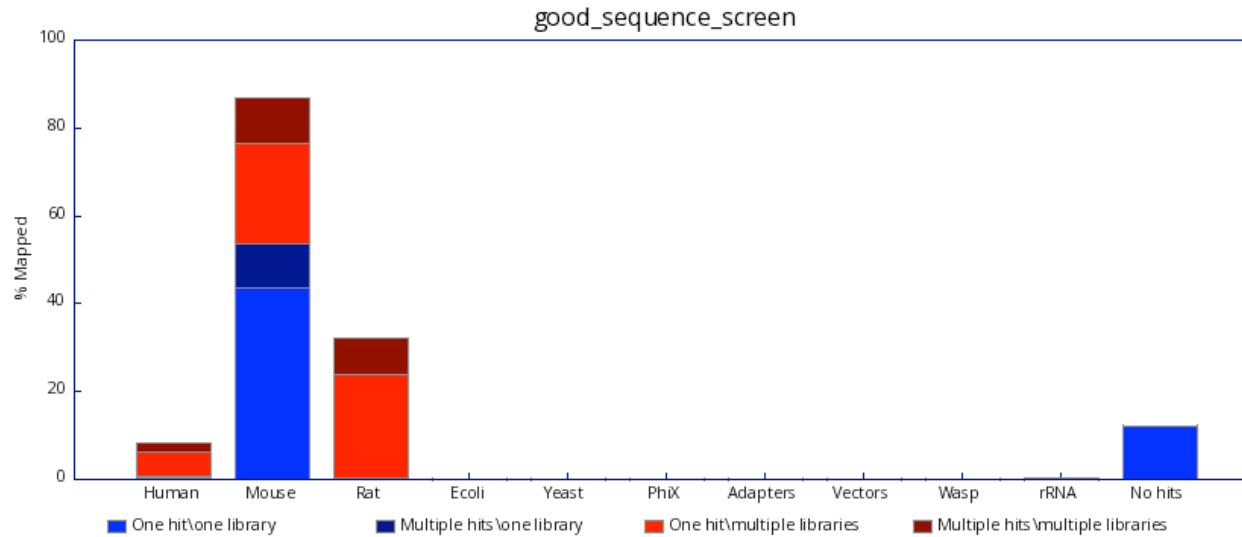


- Methods:
 - Digital Normalisation
 - Error Correction
 - Cut-off based
- The result is difficult to predict
 - Removes low coverage data
 - Overcorrects repeats
 - Removal of similar alleles
 - Large computation time for potentially little gain

- Read based contamination analyses are tricky
 - Entirely dependent on your reference database
 - Short string matching increases alignment to multiple targets
 - Unrelated organisms can contain similar strings of nucleotides

Kraken taxonomic classification

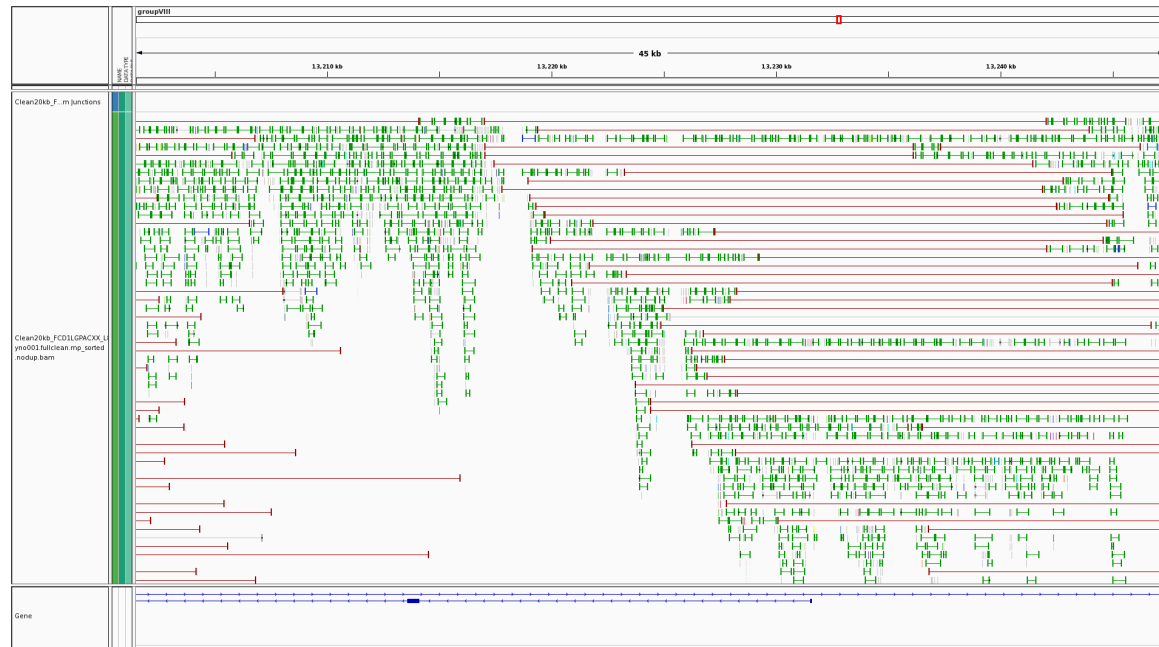




- Alignment to a known close relative.
 - Check % alignment.

```
5351663 + 0 in total (QC-passed reads + QC-failed reads)
0 + 0 secondary
19003 + 0 supplementary
0 + 0 duplicates
543253 + 0 mapped (10.15% : N/A)
5332660 + 0 paired in sequencing
2666330 + 0 read1
2666330 + 0 read2
0 + 0 properly paired (0.00% : N/A)
524250 + 0 with itself and mate mapped
0 + 0 singletons (0.00% : N/A)
12048 + 0 with mate mapped to a different chr
160 + 0 with mate mapped to a different chr (mapQ>=5)
```

- Other QC issues
 - Incorrect fragment size
 - Alignment -> Picard Metrics / Genome Browser



- Incorrect Demultiplexing
 - Check data proportions

- Check data integrity
 - Check basic statistics
 - Check for contaminants, adapter, duplication
 - Check data sets share the same information content
 - Check platform specific diagnostics
 - Trash in -> trash out.
-
- If something looks wrong, talk to your sequencing provider
 - (but get a second opinion on the answer they give you)