

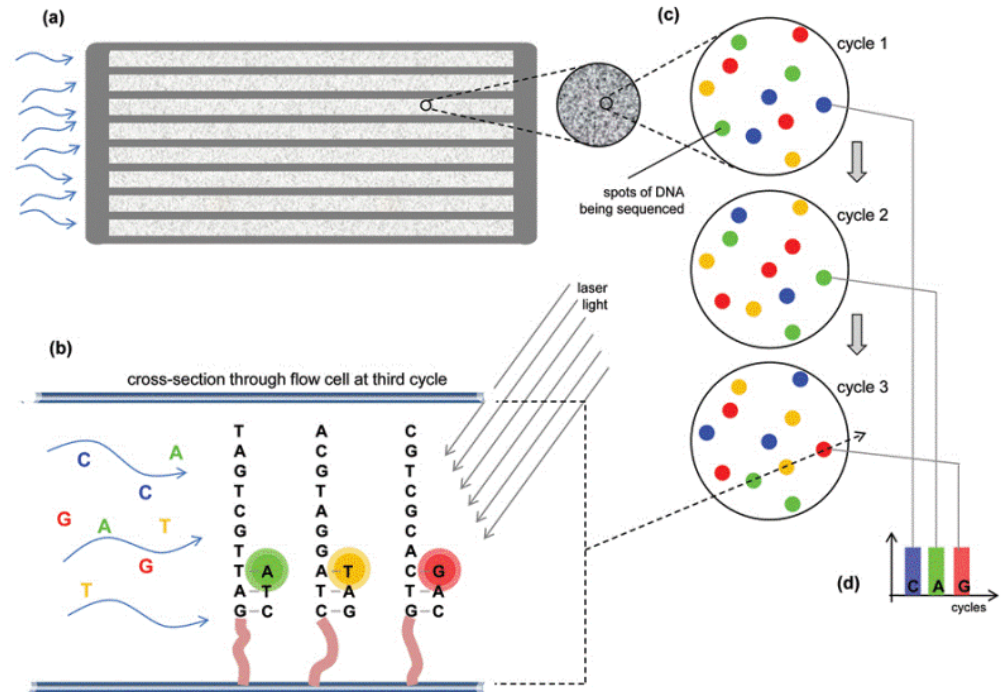
NGS Quality Control

Sandra Smit

RNAseq course, June 5-7

NBIC/EPS/WUR

NGS data



DATA: sequences (reads) and quality values

FASTQ

Many formats:

- FASTQ
- FASTA & QUAL
- SCARF
- SFF

Sequence ID

Sequence

```
@HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1  
TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTTNNNNNNNNNTAGTTTCTTGAGATTTGTTGGGGGAGACATTTTTGTGATTGCCTTGAT  
+HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1  
efcffffcfeeffcfffffdff`feed]`_Ba^__[YBBBBBBBBBRT\ ] [ ] dddd`ddd^dddadd^BBBBBBBBBBBBBBBBBBBBBBBB
```

Quality values

FASTQ quality values



- S - Sanger Phred+33, raw reads typically (0, 40)
- X - Solexa Solexa+64, raw reads typically (-5, 40)
- I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)
- J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)
with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
(Note: See discussion above).
- L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)

Quality Control (QC) of NGS data

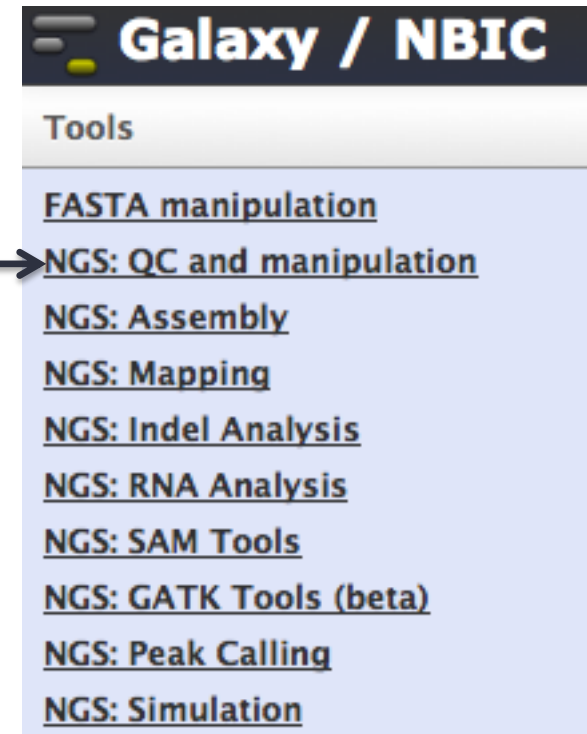
- Quality assessment of data
- Improvement of quality (trimming and filtering)
 - selection: remove problematic reads
 - trimming: discard regions of low fidelity
 - correction: replace improbable basecalls
- Quality assessment again, see improvement?

What could be wrong?

- Base calling errors
- Uncalled bases
- GC bias
- Homopolymers
- Low quality bases (3' end)
- Clonal duplicates
- Contamination (pathogens)
- Nuclear vs. organellar reads
- Sequencing artifacts (adaptors, vectors, clones, chimeric reads)

QC Tools

- FastQC
 - <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- PINSEQ
 - <http://prinseq.sourceforge.net/>
- FASTX Toolkit
 - http://hannonlab.cshl.edu/fastx_toolkit/
- NGS QC Toolkit
 - <http://www.nipgr.res.in/ngsqctoolkit.html>
- NGS QC in Galaxy

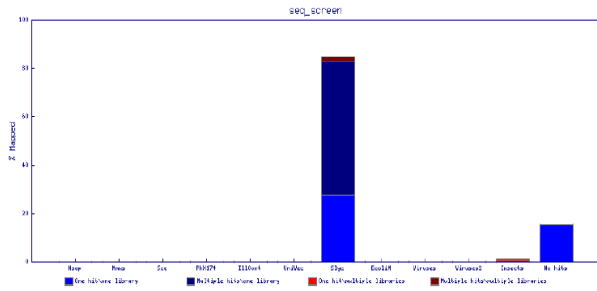


QC and manipulation in Galaxy

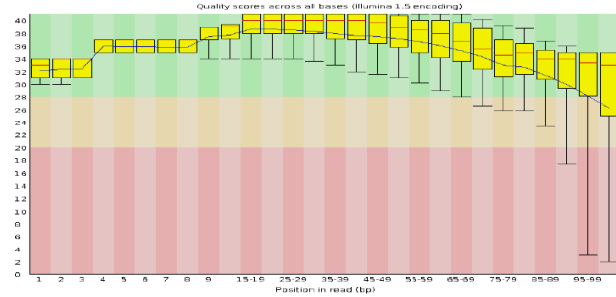
Manipulation of FASTQ data with Galaxy
Blankenberg et al.
Bioinformatics 2010

FastQC Quality Checking Tool

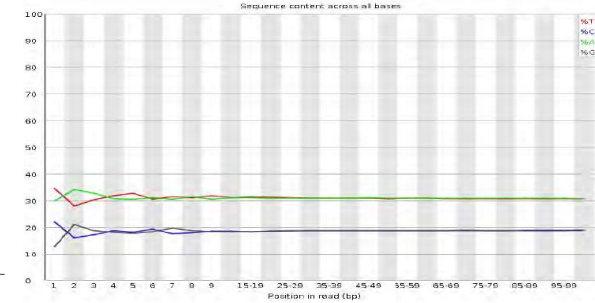
Contamination screen
fastq screen



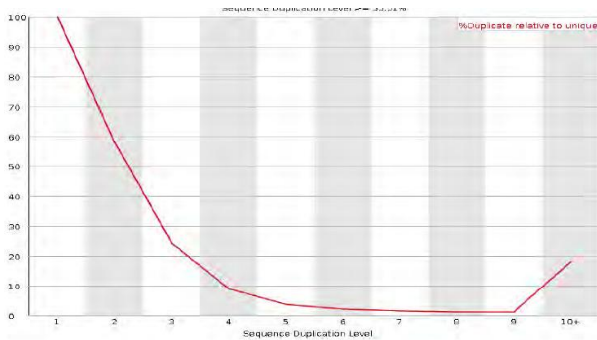
Per base sequence quality



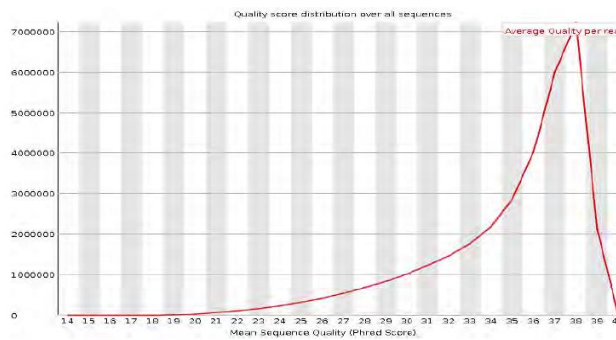
Per base sequence content



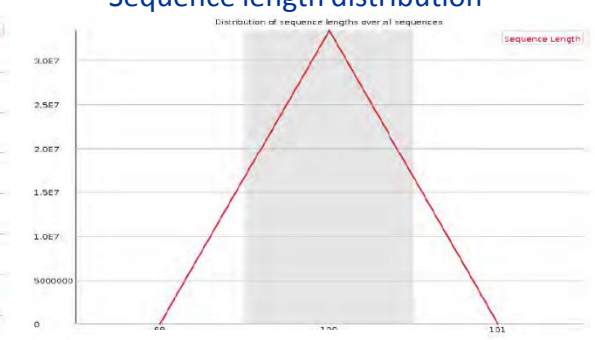
Sequence duplication



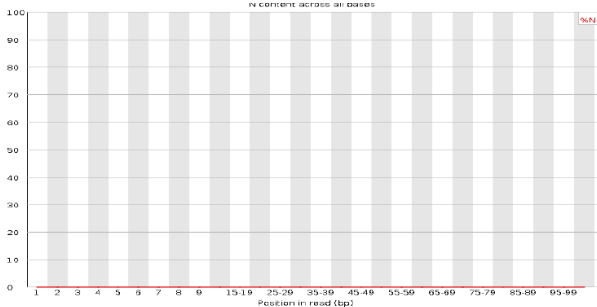
Per sequence quality



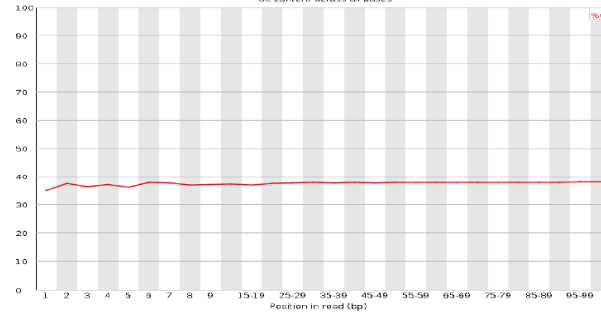
Sequence length distribution



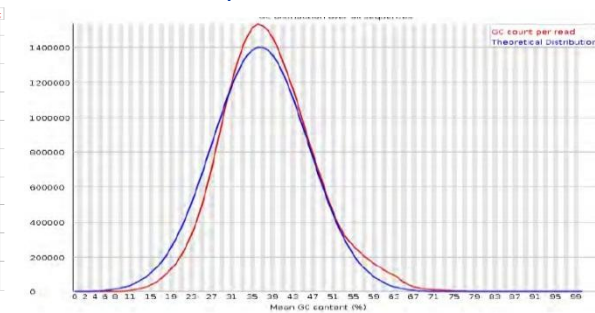
Per base N-content



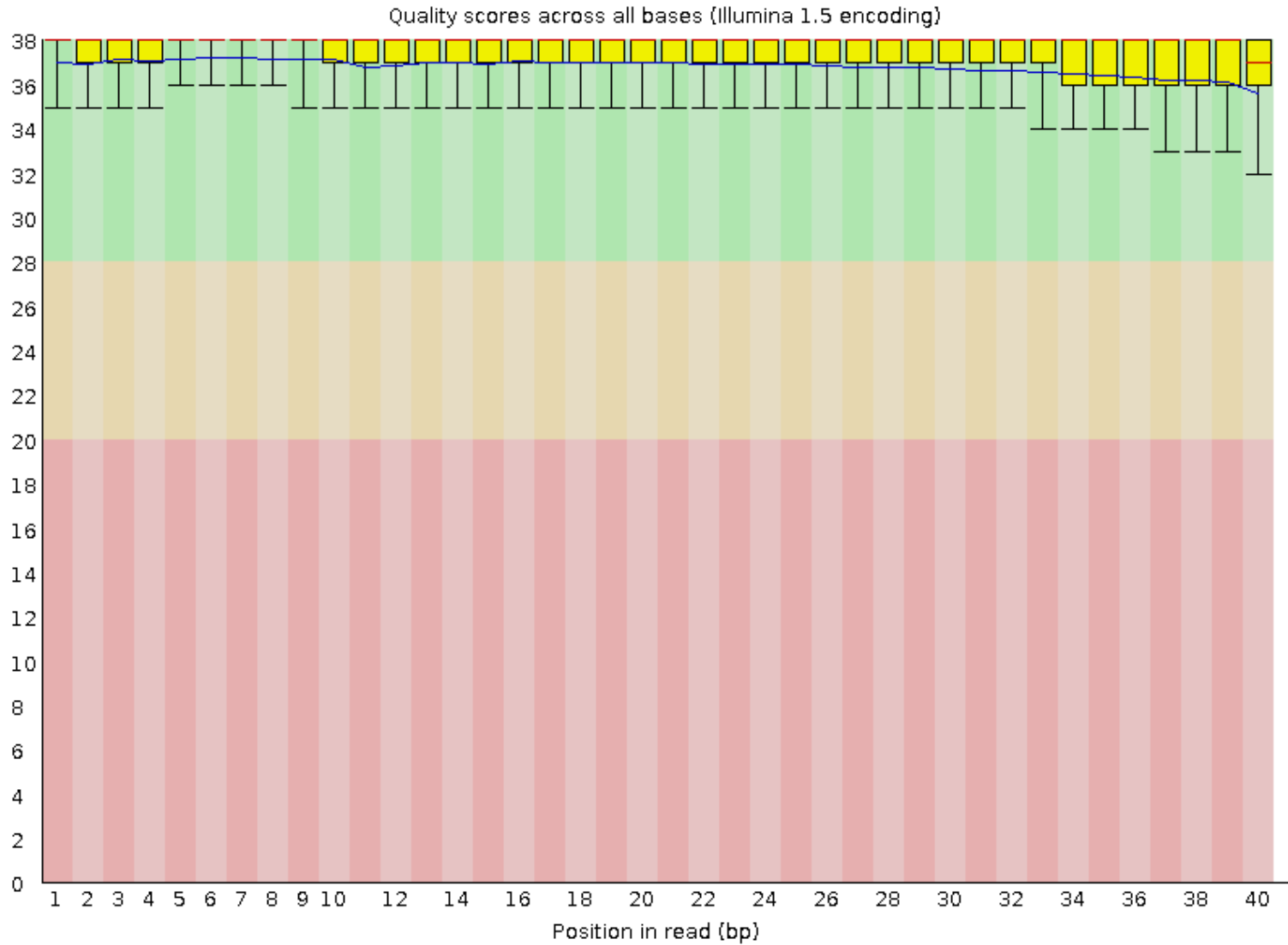
Per base GC content



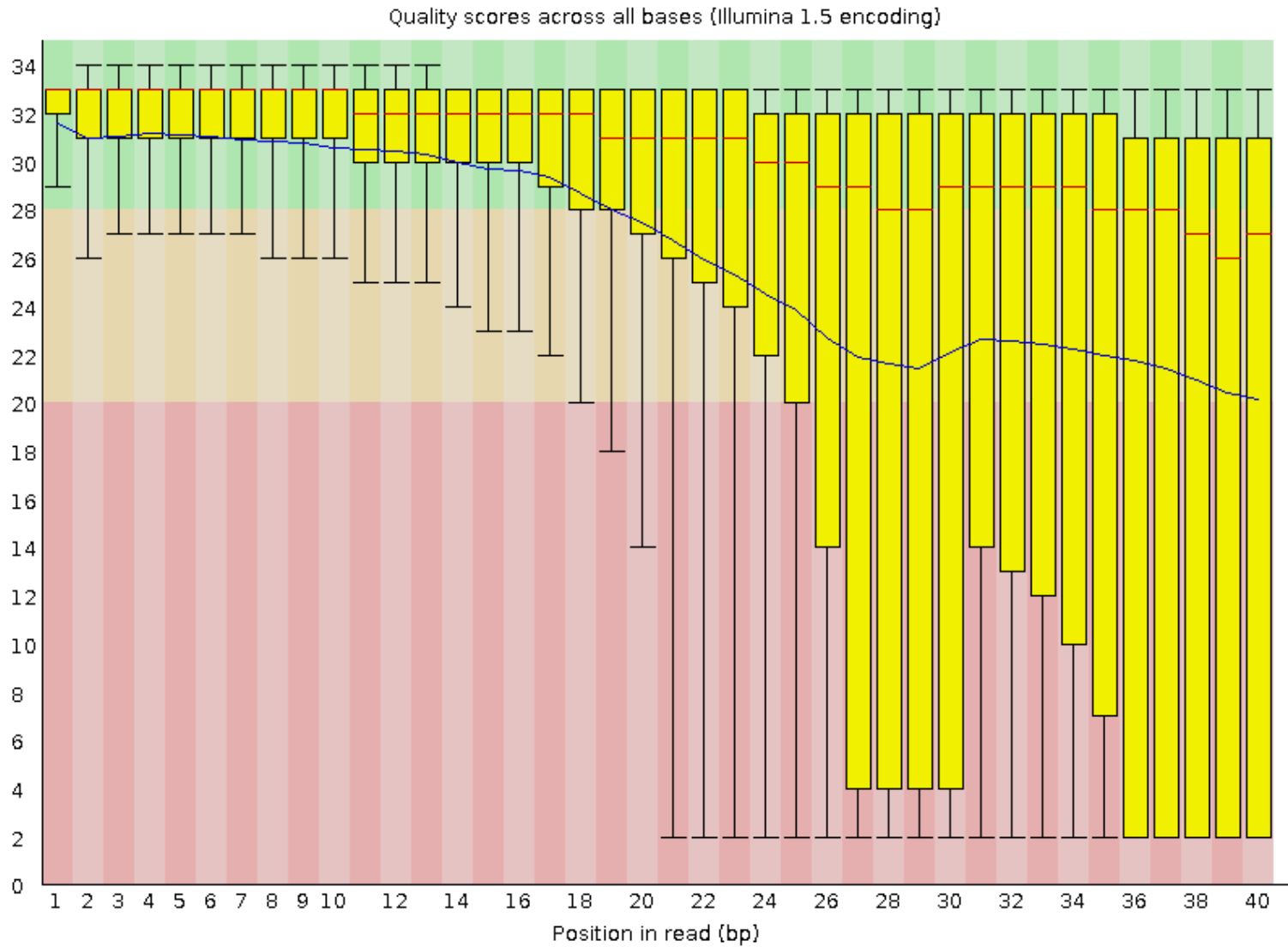
Per sequence GC content



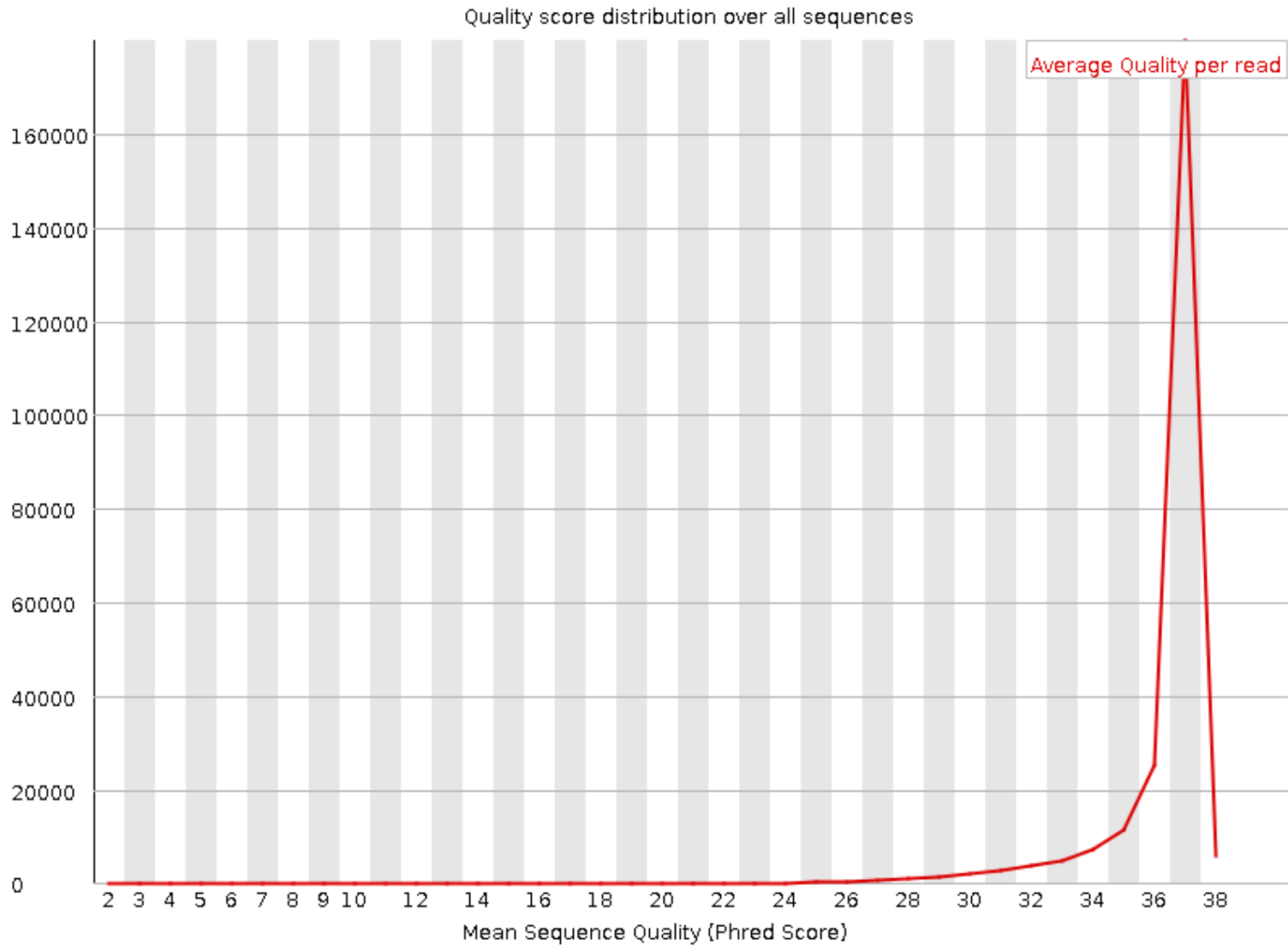
Good



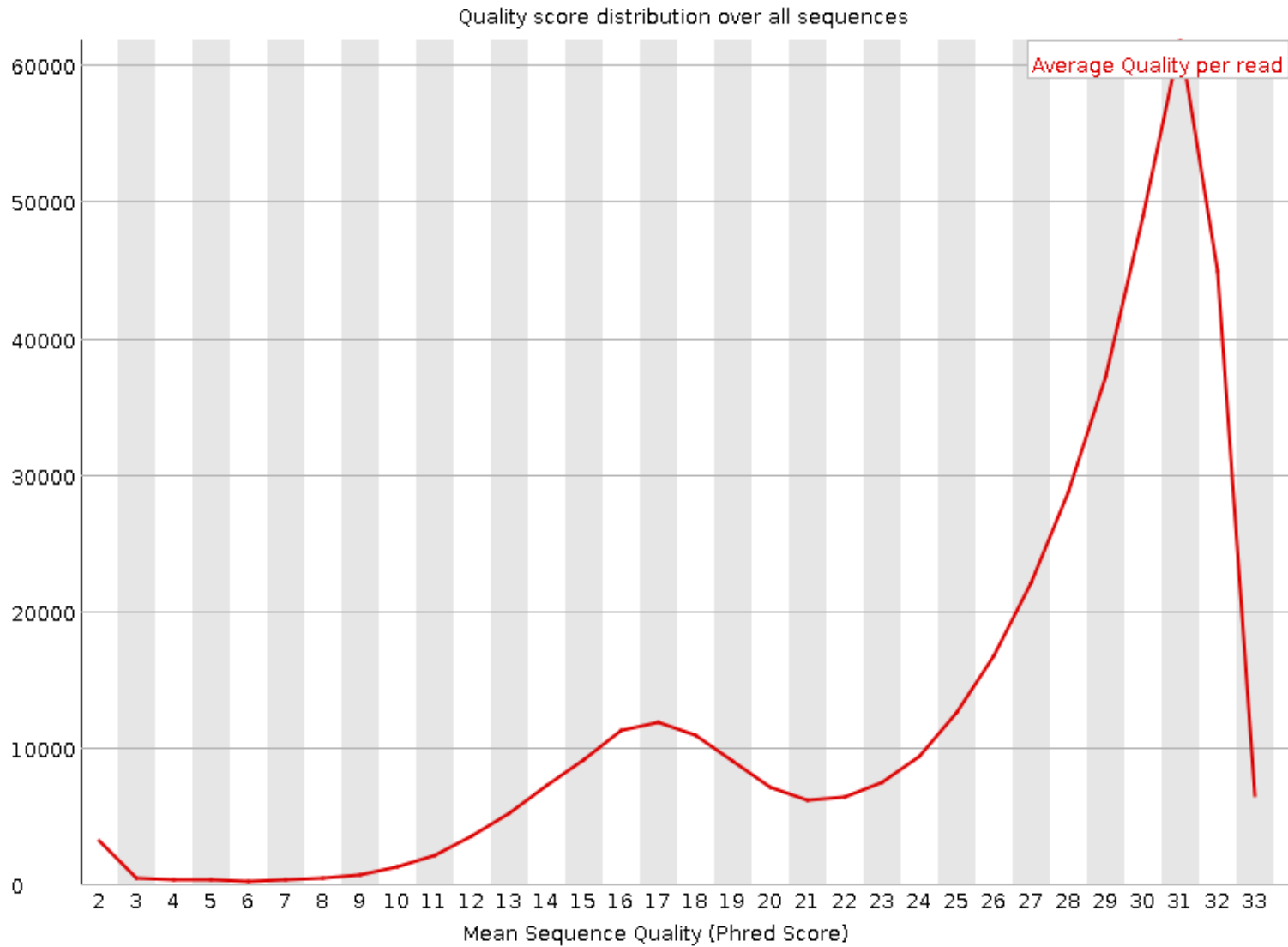
Bad



Good



Bad



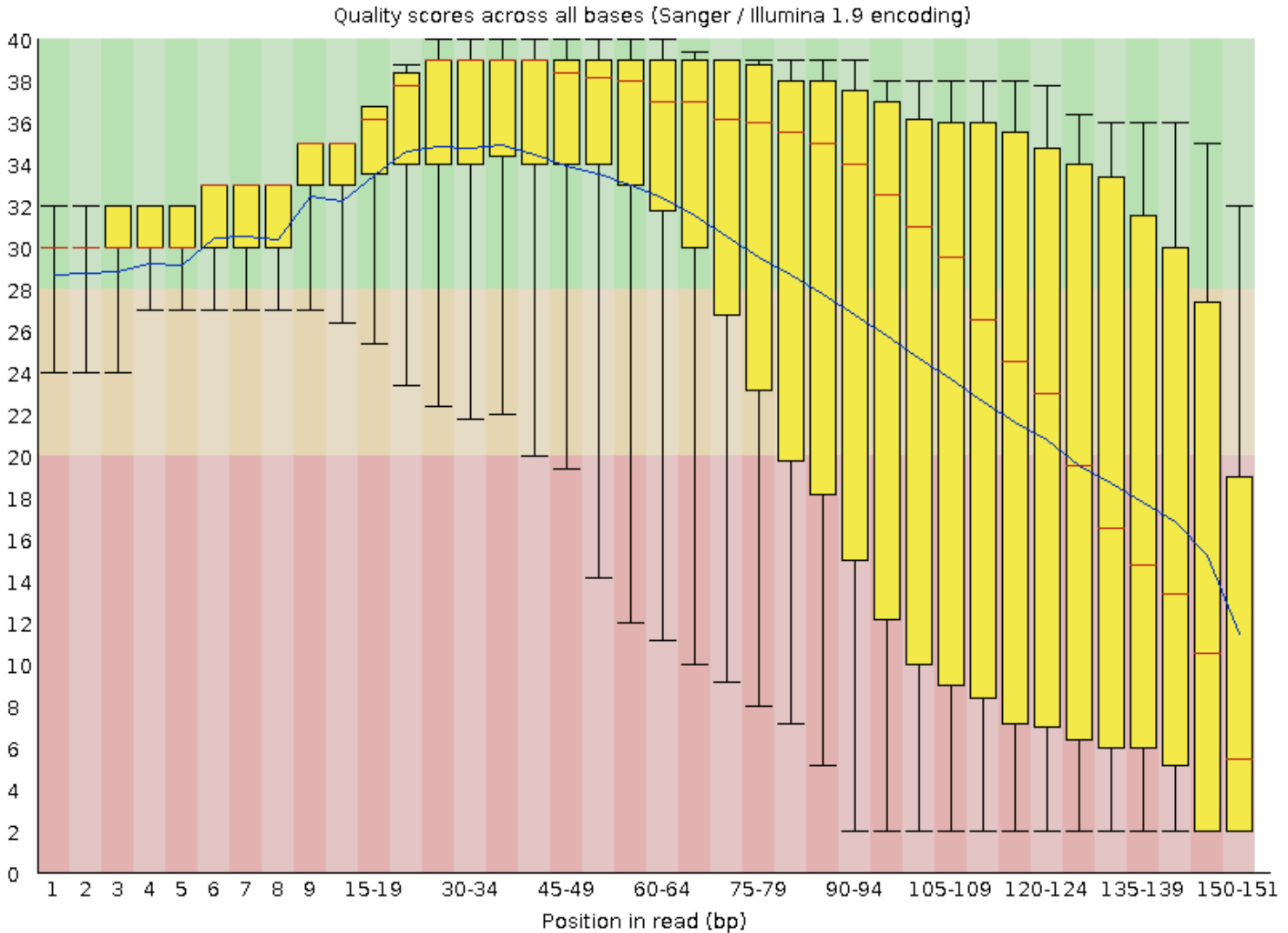
Tools in galaxy

- FASTQ groomer
 - Verify and convert between the known FASTQ variants
- Quality statistics
- Read Trimmer
- Quality filter
- Dealing with paired-end data

Tutorial

- Assess the quality in a good and bad data set
- Clean up the bad data set
- Examine the improvements

bad_MiSeq_r10000.fq (100000 reads)



Strategy

Tool: Filter by quality

Name:	Filter by quality on data 5
Created:	Jun 04, 2013
Filesize:	19.2 MB
Dbkey:	?
Format:	fastqsanger
Galaxy Tool Version:	1.0.0
Tool Version:	
Tool Standard Output:	<u>stdout</u>
Tool Standard Error:	<u>stderr</u>
Tool Exit Code:	0
API ID:	5abf71dc00cff9f5

Input Parameter	Value	Note for rerun
Library to filter	5: FASTQ Groomer on data 2	
Quality cut-off value	20	
Percent of bases in sequence that must have quality equal to / higher than cut-off value	80	

Tool: FASTQ Quality Trimmer

Name:	FASTQ Quality Trimmer on data 7
Created:	Jun 04, 2013
Filesize:	18.9 MB
Dbkey:	?
Format:	fastqsanger
Galaxy Tool Version:	1.0.0
Tool Version:	
Tool Standard Output:	<u>stdout</u>
Tool Standard Error:	<u>stderr</u>
Tool Exit Code:	0
API ID:	af7fa13876122ee5

Input Parameter	Value	Note for rerun
FASTQ File	7: Filter by quality on data 5	
Keep reads with zero length	False	
Trim ends	3' only	
Window size	1	
Step Size	1	
Maximum number of bases to exclude from the window during aggregation	0	
Aggregate action for window	min score	
Trim until aggregate score is	>=	
Quality Score	20.0	

After cleaning (50989 reads)

