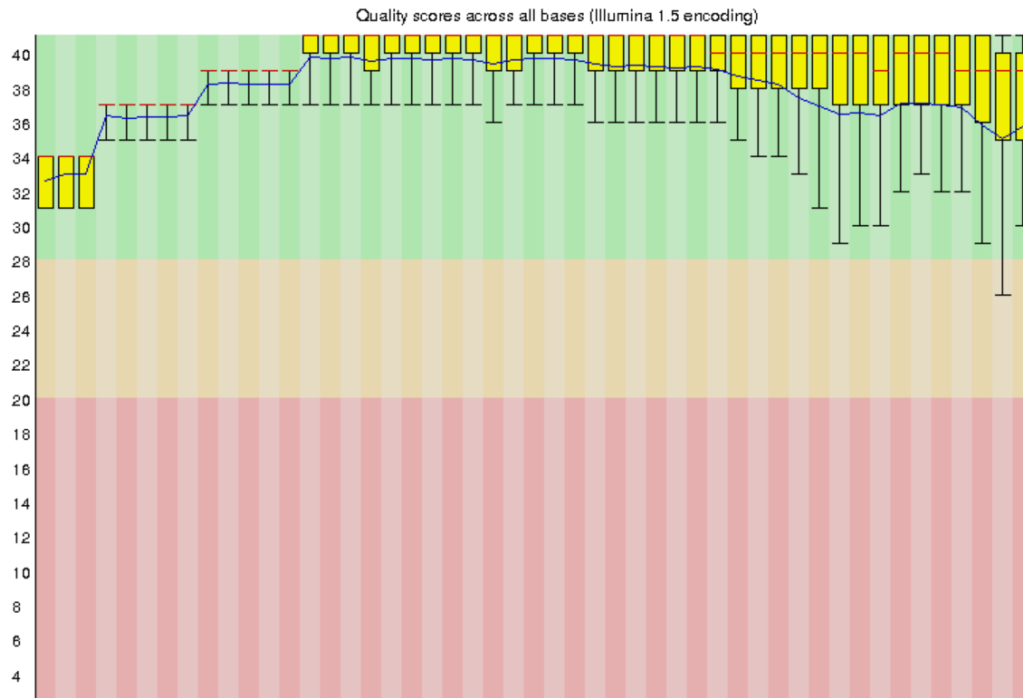


Part 2

- a. I ran FASTQC on the uploaded file and the boxplot of quality scores is shown below. Each quality score is in the green region, so I would say, in general, the quality of the data is good.



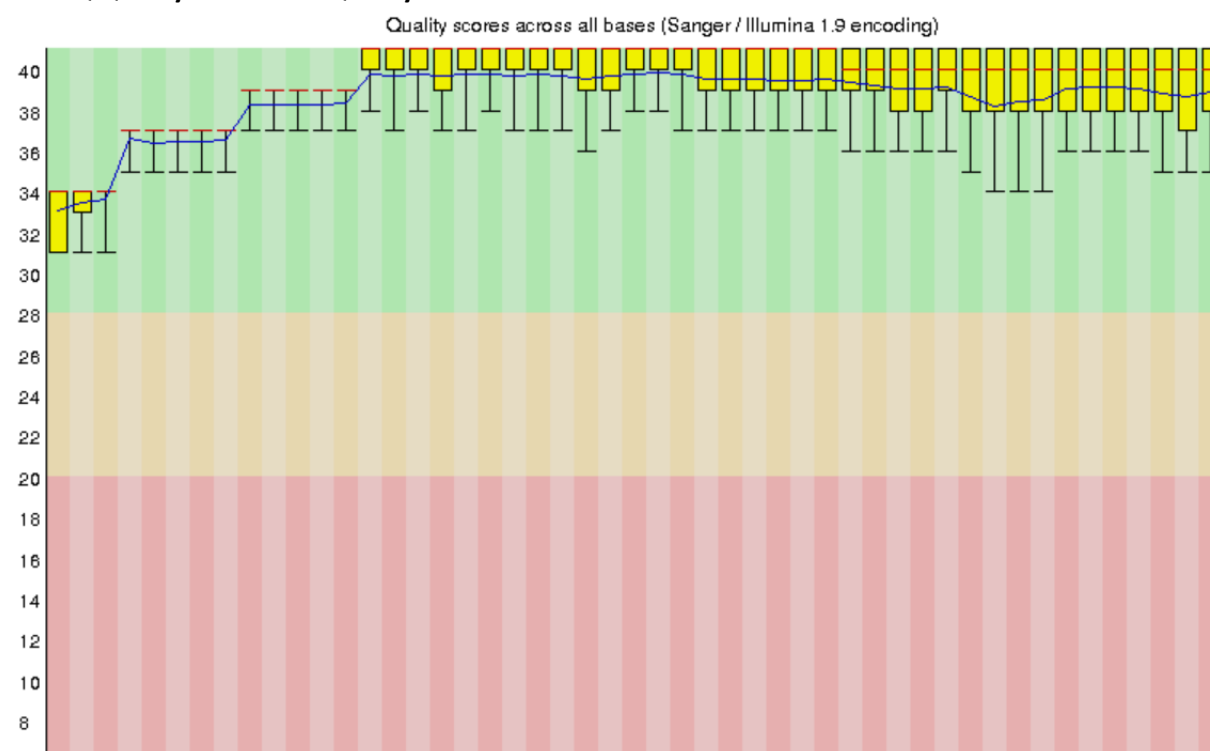
- b. This data uses the Illumina 1.5 phred encoding scheme, which is a previous version of Illumina. The read lengths are 49 and there are 20,000 reads in the file.
- c. Running FASTQ Groomer to convert the phred quality scores and then rerunning FASTQC shows that the phred encoding scheme was converted to Sanger / Illumina 1.9. The boxplot look almost exactly the same as the one in Part 2a.

Part 3

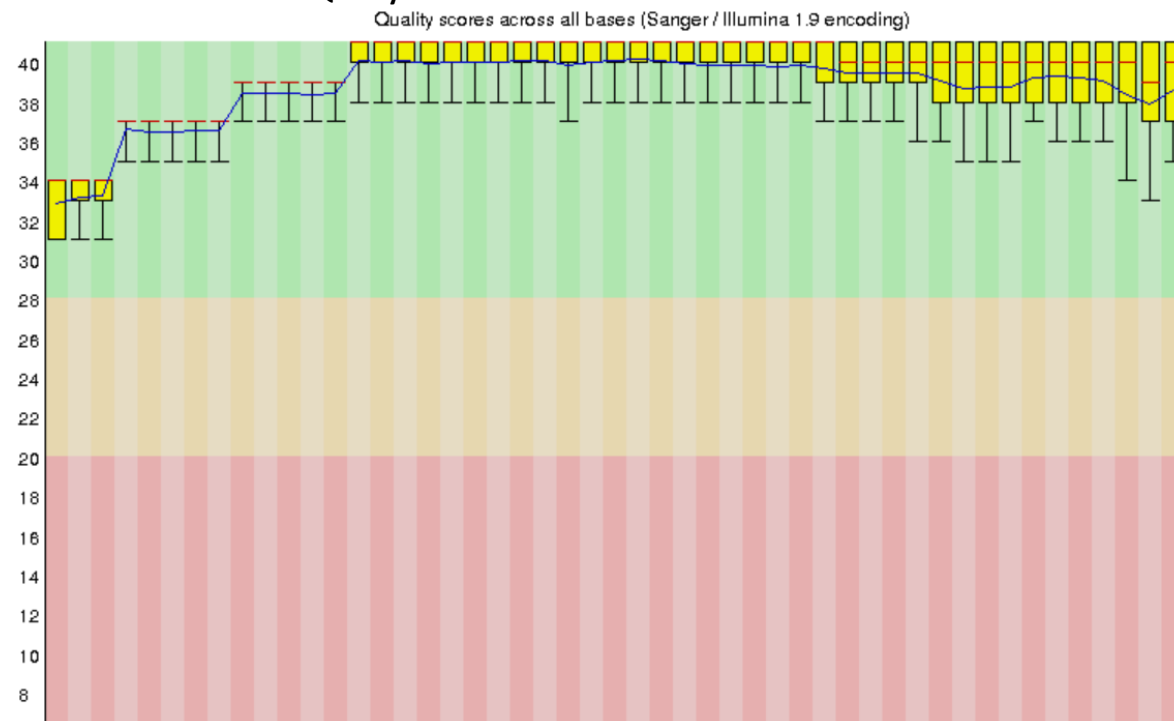
In order to compare the two trimmer tools in Galaxy, I ran FASTQ Quality Trimmer to trim the data with a sliding window of 4 bases and trimmed the reads until the average quality score of the window is greater than 30. I then ran the Trimmomatic tool on the groomed data and used the same parameters.

- a. To access the quality of each dataset, I ran FASTQC on both of the created datasets above. The box plots of quality scores are shown here:

FASTQ Quality Trimmer – Quality Score Box Plot



Trimmomatic Box Plot – Quality Score Box Plot



- b. The data quality scores after trimmer are better than the ones that we saw in the boxplot from Part 2. For example, as the reads got longer in the untrimmed data, you can see the quality of the scores went down. The trimmed data has a more consistent quality score across the top of the boxplot. This does make sense to me, because we are trimming the low quality portion of the reads and leaving only high quality data.
- c. The tool trimming tools seemed to give almost the same result, however quality scores seemed to be slightly better in the mid-length reads from Trimmomatic, while the longer length reads seemed to have better quality of FASTQ Quality Trimmer. Based on these boxplots alone, I would use either tool, but if I had to choose one it would be Trimmomatic because the the quality scores seemed to be slightly more consistent across all read lengths.

Part 4

To identify SNPs in two FASTQ files with paired-end sequencing data from the 1000 Genomes Project (reference genome hg19), I read the following data into Galaxy in the 'fastq' format with genome build hg19.

Forward

reads: ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG00117/sequence_read/SRR044234_1.filt.fastq.gz

Reverse

reads: ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG00117/sequence_read/SRR044234_2.filt.fastq.gz

To determine quality encoding I ran FASTQC on each file. The encoding said Sanger/Illumina 1.9, but when I tried running Trimmomatic, it could not find any fastsanger datasets, so I ran FASTA Groomer on both files and then ran FASTQC on both groomed files.

Next, I trimmed the data using Trimmomatic to remove the low-quality bases from each file. I used a window size of 4 bases and a required average quality score of > 20. Then I reran FASTQC on the trimmed data to ensure that low quality bases were removed.

On the trimmed data, I ran HISAT2 to align the files to the hg19 reference genome, using paired-end reads and outputting in sorted BAM format. I attached the BAM file to my submission. I also wanted to note that I tried Bowtie2, and it took a lot longer, as suspected, since HISAT2 is advertised as using a fast algorithm.

I then ran the FreeBayes variant calling program, limiting the output to chr22:0-51304566. This algorithm showed 329 variants in the resulting VCF file, which I attached to my submission. I filtered the variants down to those that show heterozygosity with estimated allele frequency of 0.5 and read depth > 10 with VCFfilter. This narrowed the results down to 16. The filtered VCF file is also attached.

I downloaded RefSeq known genes from the UCSC Table Browser in the region of chr22:0-51304566. The BED file is attached to my submission. I uploaded the BED file to Galaxy and ran VCFannotate to intersect this BED file with the filtered VCF file. The goal was to locate variants within genes on chromosome 22. The final filtered, annotated VCF file was is attached to this submission.

-
- chr22 (q13.1) 22p13 22p12 22p11.2 22q11.21 22q12.1 22q12.2 22q12.3 22q13.1 22q13.2 22q13.31
- Scale chr22: 39,040,541 39,040,545 39,040,549 39,040,551
- CTCF
- Overlap summary of Ensembl ChIPSeq binding peaks across available datasets
- Chromosome Bands Localized by FISH Mapping Clones
- UCSC Genes (RefSeq, GenBank, CCDS, Rfam, tRNAs & Comparative Genomics)
- Simple Nucleotide Polymorphisms (dbSNP 150) Found in >= 1% of Samples
- K562 CTCF ChIA-PET Interactions Rep 1 from ENCODE/GIS-Ruan
- K562 CTCF ChIA-PET Signal Rep 1 from ENCODE/GIS-Ruan
- K562 Pol2 ChIA-PET Interactions Rep 1 from ENCODE/GIS-Ruan
- Chromatin Interaction Analysis Paired-End Tags (ChIA-PET) from ENCODE/GIS-Ruan
- HeLa-S3 Pol2 ChIA-PET Interactions Rep 1 from ENCODE/GIS-Ruan
- Chromatin Interaction Analysis Paired-End Tags (ChIA-PET) from ENCODE/GIS-Ruan
- Chromatin Interactions by SC from ENCODE/Dekker Univ. Mass.
- FAM227A
- Common SNPs(150)
- 13909..39514502,2
- 52673..39053304,2
- 30 -
- K562 CTCF Sig 1
- 0 -
- 78218..39081150,2
- 80649..39102943,2
- 92906..39095105,2
- 95908..39095535,2
- 77676..39079385,3
- 74132..39076609,2
- 76471..39078147,2
- 96553..39099995,3
- K562 Pol2 Sig 1
- 89477..46418124,2
- HeLaS3 Pol2 Sig 1
- GM12878 FK
- H1-hESC FK