

Bioinformatics programming assignment

Write a NGS read simulator that randomly picks out reads from a genome and outputs them as a fastq file (use dummy quality values). (Use read length of 50 bp and generate 100,000 reads from the human genome)

Add a uniform error rate of 0.01 (1% of the time a base is randomly replaced with another base) to the fastq file.

Align the resulting fastq file with bwa and find the error rate (read aligned to a part of the genome other than where it originated from).

Validation of Introduced error rate from 'stats.csv' file generated:

Total number of reads : 100000bp

Read length : 50bp

Total errors introduced : 49954

$$\text{Sequencing Error Rate} = \frac{\text{Total errors introduced}}{\text{Total no. of reads} \times \text{Read length}}$$

$$\text{Sequencing Error Rate} = \frac{49954}{100000 \times 50} \times 100 = 0.99908$$

Which is approximately 1%

Error Rate due to mapping of reads on genome positions other than where it originated from:

Results from 'calc_error.py' script:

Reads mapped position other than its origin : 13653

Reads mapped to genome : 99192

Reads unmapped : 808

Results from 'samtools idxstats' command:

Reads mapped to genome : 99192

Reads unmapped : 808

$$\text{Error Rate} = \frac{\text{Reads Mapped to position other than its origin}}{\text{Total no. of mapped reads}} \times 100$$

$$\text{Error Rate} = \frac{13653}{99192} \times 100 = 13.76$$

Which is approximately 14%

Note : all analysis files are given in 'Analysis Files' folder.