

# Genomorientierte Bioinformatik

## Report

### Read Simulator

Malte Weyrich

NOVEMBER 2024

---

In bioinformatics, **Sequencing** is a term used to describe the process of gathering genomic data by reading the nucleotides of a DNA molecule. This is done by a *sequencer*. There are different types of sequencing techniques (e.g. *Illumina* (next generation sequencing), *Oxford Nanopore*, *Pacbio* (third generation sequencing), etc.) and different variants of *sequencing* (*ATAC-seq*, *scRNA-seq*, *ChIP-seq*, ...). In this report, we will focus on a *Read Simulator* which is a tool used to simulate the results of a sequencing experiment (in this case paired end sequencing using *Illumina*). The program was executed using "*Homo\_sapiens.GRCh37.75.dna.toplevel.fa*" as *Reference Genome*, its corresponding *fasta index* and its annotation in form of a *Gene Transfer Format-File* (GTF). The *Read Simulator* was written in *Java* and will be benchmarked and analyzed by its complexity and correctness. The results of simulator itself will also be discussed in this report.

---

## 1 — Introduction

The simplified process of *Illumina* sequencing is as follows: Several DNA target sequences get treated with ultra sound, in order to break it down into smaller fragments of a certain length with a certain margin of error (e.g. 200 bp +/-  $x$  bp). These fragments are labeled and then placed onto a flow cell, where they are amplified and sequenced. This generates an abundance of short reads, which often overlap with each other and contain mutations. The reads are then aligned to a *Reference Genome* via a *Mapper* like *STAR*, in order to determine the original sequence. The *Read Simulator* is a tool that simulates this process of generating fragments and reads for given transcript sequences. In our case we simulate a paired end sequencing experiment, where we generate two reads for each fragment, one for each end of the fragment. We can test the performance of the *Mapper* and other tools that are used in the analysis of sequencing data, since we know where our reads originated from. A downside of the *Read Simulator* is that it assumes a normal distribution of fragment lengths and their starting position inside the transcript and a constant mutation rate across the entire read sequences, which in reality is not the case. But it is still a useful tool for testing the performance of other tools.

## 2 — Java Implementation

### 2.1. Logic

The logic of the *Read Simulator* is split up into four main steps:

#### (A) Read Gene & Transcript IDs to Simulate:

The user specifies a *tsv* file containing the gene and transcript IDs of the sequences and the amount of reads to simulate for each transcript. The file is passed to the *ReadSimulator* via the `-readcounts` argument. The entries are stored inside a nested

*HashMap<String, HashMap<String,Integer>> readCounts* object where the first key corresponds to the gene ID and maps to a second map, which maps the transcript ID to the amount of reads to simulate. This way we avoid storing the same gene ID multiple times and can easily access the amount of reads to simulate for a given transcript.

#### (B) Initialize Genome using a GTF-File:

The *Genome* class is initialized by passing the path of the *Reference Genome* and the *Fasta Index* file to the constructor. The path of the *Reference Genome* is used to create a *RandomAccessFile* object, which is used to access the large fasta file in a more efficient way by utilizing the indices stored in the *Fasta Index*. This will later be used to extract the sequences of *Genes* containing the *Transcripts* we want to simulate reads for. The *GTF* file and *readCounts* object are then passed to the *Genome* object to initialize the gene and transcript coordinates. Each line of the *GTF-File* is filtered using *GenomeUtils.filterLine(line, readCounts)*, which basically checks if the *Gene ID* of the line is present in the *readCounts* object. The method works by counting the number of seen *<tabs>* in the current line and then extracting the *Gene ID*, which is located in between the 8th and 9th *<tab>* of the line. This way we don't call expensive *split()* operations on each line of the *GTF-File*. If a valid line was found, we only need to check if it is a *gene/transcript/exon* entry and either create a new *Gene/Transcript/Exon* object. A *Gene* can have several *Transcript*'s and a *Transcript* can have several *Exon*'s. Due to the filtering of *GenomeUtils.filterLine*, our *Genome* object will only contain *Genes* and *Transcripts* we want to simulate reads for.

#### (C) Initialize Gene Sequences of Interest:

In order to simulate reads, we first need to extract the exonic sequences of the *Transcripts*.

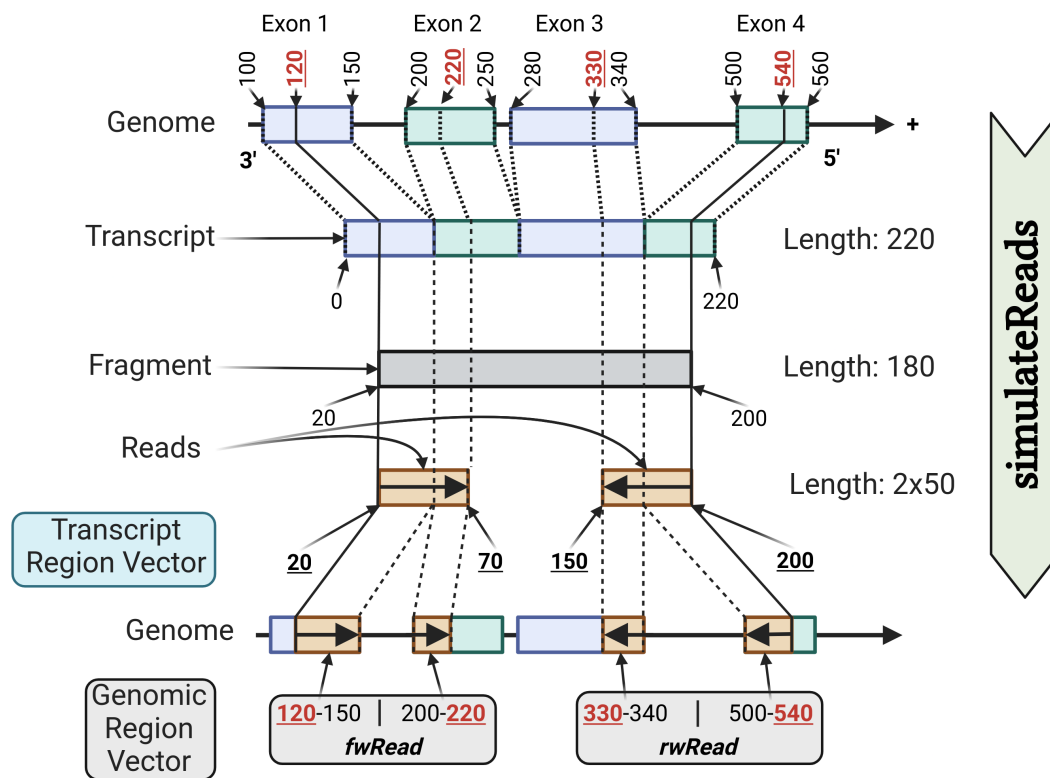
This is done by calling *Genome.initTargetSequences(readCounts)*, which iterates over all *Genes* and *Transcripts* in the *readCounts* object. For each *Gene* we extract its sequence using the *RandomAccessFile* object and store it in a *String seq* object. This is done by utilizing the start/end coordinates of the *Gene* together with the *Fasta Index* to calculate the byte offset of the *Gene* in the *Reference Genome* file and read in the sequence. The *seq* object can now be used to concatenate the *Transcript* sequences by cutting out the exonic sequences of the corresponding *Transcript* based on their start and end coordinates. This way, we avoid repeatedly accessing the *Reference Genome* file which is time inefficient due to the large size of the file and the vast amount of *Exons* we need to extract. If a *Gene* is located on the reverse strand, we go through the *Exons* in reverse order.

#### (D) Generate Reads and Write to File:

For each *Transcript* we want to simulate reads for, we sample a random fragment length and starting position from a normal distribution. These two values and the specified *read length* are used to extract two substrings of the *Transcript* sequence:

```
do {
    fragmentLength = (int) Math.round(normalDist.sample());
} while (fragmentLength < length || fragmentLength > transcriptSeq.length());
int maxStartPos = transcriptSeq.length() - fragmentLength;
int randomStartPos = splittableRandom.nextInt(maxStartPos + 1);
String fwSeqRead = transcriptSeq.substring(
    randomStartPos,
    randomStartPos + length
);
String rwSeqRead = GenomeUtils.revComplement(
    transcriptSeq.substring(
        randomStartPos + fragmentLength - length,
        randomStartPos + fragmentLength
    ));
```

These two *Strings* resemble the unmutated forward (*fw*) and reverse (*rw*) sequence of the *Read*. The *start/end* positions inside the *Transcript*, the *Read* sequence, the *Read ID* and a *boolean* indicating if the *Read* is the forward or reverse *Read*, are then collected in a *Read* object (*fwRead* or *rwRead*). We iterate over each nucleotide of the *Read* and use the *mutationRate* to determine if we should mutate the nucleotide of the current position. Mutations to the original nucleotide are not considered, so we only mutate to one of the other three nucleotides. The next step is to derive the *Genomic Region Vector* of both *Reads*. In figure 1 we can see the mapping of the *Reads* to the *Reference Genome* using the prior knowledge of the *Exon* coordinates.



**Figure 1 – Mapping of Reads to the Reference Genome. This figure was created using BioRender.com 2024**

The method *getGenomicRegion* is used to annotate the generated *Reads* with the corresponding *Genomic Region Vector*. We start by extracting the list of *Exons* from the provided *Transcript* object and getting the *Read*'s *start* and *end* positions inside the *Transcript*. For each *Exon*, we calculate its genomic *start*, *end*, and *length*. By checking if the *Read* overlaps with the current *Exon* in *Transcript* coordinates, we identify the overlapping region. This region is then mapped

to genomic coordinates, accounting for the strand direction. Finally, the genomic regions are formatted and stored in a list, ensuring proper coordinate ordering for the reverse strand, before proceeding to the next *Exon*.

After both *Reads* have been generated and annotated with their *Genomic Region Vectors*, we write them to two separate *fastq* files ("*fw.fastq*" and "*rw.fastq*") in the following format:

```
@0  
CTAAAAGGCGCAAGAGAATGGATGATAGTAGTGTCCTCGAGGCCACACGGGTT...  
+0  
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII..
```

The *Read ID* is followed by the *Read* sequence and a + sign. The quality scores are represented by the *ASCII* characters *I* in this case. The quality scores are not considered in the *Read Simulator* and are set to the same value for all *Reads*. A summary of the generated *Reads* is also created and written to a *tsv* file ("*read.mappinginfo*") in the following format (table 1):

**Table 1—Format of the Summary File created by the Read Simulator. Each row corresponds to a Read Pair (*fwRead* and *rwRead*).**

[illegible]

The *Read* objects are discarded after writing them to the files, in order to save memory. Writing is done by three separate *BufferedWriter* objects, one for each file. The entries are constructed using three different *StringBuilder* objects, which are then written to the files and reset afterwards.

## 2.2. Complexity

The complexity of the *Read Simulator* can be broken down into the parts described above:

**(A) Read Gene & Transcript IDs to Simulate:**

Section (A) of the *ReadSimulator* scales linearly with the amount of lines ( $c$ ) in the *tsv* file:

$\mathcal{O}(c).$

**(B) Initialize Genome using a GTF-File:**

Section (B) also has a linear complexity with the amount of lines ( $l$ ) in the *GTF* file:

$$\mathcal{O}(l).$$

**(C) Initialize Gene Sequences of Interest:**

**(D) Generate Reads and Write to File:**

*2.3. Correctness*

*2.4. Benchmarking*

## **3 — Results**

## **A — Appendix Section**

hm

Text goes here