

SomatoSim User Manual

Marwan A. Hawari, Celine S. Hong and Leslie G. Biesecker

Contact: marwan.hawari@nih.gov or celine.hong@nih.gov

October 22, 2020

Contents

1	INTRODUCTION	2
2	USAGE.....	3
3	INPUT	6
4	OUTPUT	7
5	TEST DATA	8
6	REFERENCES.....	10

1 Introduction

1.1 Software dependencies

1. Python version 3.6.8 (<https://www.python.org/downloads/>)
2. Numpy version 1.16.2 (<https://numpy.org/>)
3. Pandas version 0.25.1 (<https://pandas.pydata.org/>)
4. Matplotlib version 3.1.1 (<https://matplotlib.org/>)
5. Pysam version 0.15.0 (<https://pysam.readthedocs.io/en/latest/>)
6. SAMtools version 1.9 (<http://www.htslib.org/>)

1.2 Download

Visit <https://github.com/BieseckerLab/SomatoSim> for scripts and test data.

1.3 Installation

Clone the github repository and execute the following commands:

```
$ git clone https://github.com/BieseckerLab/SomatoSim.git
$ cd SomatoSim
$ python3 -m pip install .
```

1.4 Attribution

This software was developed by Marwan A. Hawari, Celine S. Hong, and Leslie G. Biesecker at the National Human Genome Research Institute (NHGRI), National Institutes of Health (NIH). Please include proper attribution of the NHGRI as the developer of this program and include a link to the following [<https://github.com/BieseckerLab/SomatoSim>] in all publications and other public disclosures that reference the program and/or include data or research results that were generated utilizing the program.

1.5 Public Domain Notice

This software is a United States Government Work. Anyone may use the software on a worldwide and royalty-free basis for any purpose and anyone may reproduce and prepare derivative works without restriction. Although all reasonable efforts have been taken to ensure the accuracy and reliability of the software, the National Human Genome Research Institute (NHGRI), National Institutes of Health (NIH) and the U.S. Government do not and cannot warrant the performance or any results that may be obtained by using this software. NHGRI, NIH and the U.S. Government disclaim all warranties as to performance, merchantability or fitness for any particular purpose. No indemnification is intended or provided by the US government.

2 Usage

2.1 Command line execution

```
$ somatosim \
-i <input bam file path> \
-b <input bed file path> \
-o <output directory path> \
--vaf-low <number> \
--vaf-high <number> \
--output-prefix <output files prefix> \
--number-snv <number> \
--random-seed <number> \
--down-sample <number> \
--target-coverage <number> \
--coverage-tolerance <number> \
--coverage-MQ <number> \
--coverage-BQ <number> \
--minimum-separation <number> \
<--verbose> \
--read-min-MQ <number> \
--position-min-BQ <number>
```

2.2 Command line options descriptions

-h, --help

Show the SomatoSim help message and exit

--version

Show the SomatoSim version number and exit

Required arguments:

-i, --input-bam-file

The full path of the input BAM file. See Section 3.1 for details.

-b, --input-bed-file

The full path of the input BED file. See Section 3.2 for details.

-o, --output-directory

The full path where the SomatoSim output files will be deposited. If the directory does not already exist, it will automatically be created.

Optional arguments:

--vaf-low

The lower bound of the desired VAF range. Input this value as a decimal, not as a percentage. This is ignored if the BED file input contains pre-defined VAF values.
[Default: 0.01]

--vaf-high

The upper bound of the desired VAF range. Input this value as a decimal, not as a percentage. This is ignored if the BED file input contains pre-defined VAF values. [Default: 0.20]

--output-prefix

The desired prefix for the output files generated by SomatoSim. [Default: input BAM file prefix]

--number-snv

The desired number of SNVs to simulate. [Default: Number of lines in the input BED file]

--random-seed

The desired random seed. Choosing a seed will help with reproducibility. [Default: random number between 0 and 100]

--down-sample

The desired average coverage to which the input BAM file will be down-sampled. [Default: No down-sampling]

--target-coverage

The desired depth of coverage to target when selecting for genomic positions. [Default: None]

--coverage-tolerance

The tolerance allowed when targeting a specific depth of coverage. Input this value as a decimal, not as a percentage. This option is ignored if the --target-coverage value is not used. [Default: 0.10]

--coverage-MQ

The mapping quality used by SAMtools mpileup when calculating coverage. [Default: 20]

--coverage-BQ

The base quality used by SAMtools mpileup when calculating coverage. [Default: 20]

--minimum-separation

Set the minimum number of bases that must separate simulated variants. For example, if this value is 0, then variants can be simulated directly adjacent to each other. [Default: 1]

--verbose

Displays the specific positions not selected during variant selection stage. Typing --verbose will turn the option on, and not typing anything will leave it off. [Default: Off]

--read-min-MQ

The minimum mapping quality used when selecting reads to introduce simulated variants into. [Default: 20]

--position-min-BQ

The minimum base quality used when selected reads to introduce simulated variants into. [Default: 20]

2.3 Run time considerations

Table 1. Summary of how each input of SomatoSim can affect the run time

Argument	Change	Run time	Rationale
-b/--input-bam-file	Increase BAM file size	↑	More time necessary to compute average coverage and iterate through the reads
--vaf-low	Increase value	↑	This will result in more reads needing to be selected and mutated
--vaf-high	Increase value	↑	This will result in more reads needing to be selected and mutated
--number-snv	Increase value	↑	Selecting for more positions and reads means searching longer for those positions and reads
--random-seed	Change value	—	Changing random seed only changes random values
--down-sample	Selecting a value	↑	SomatoSim needs to down-sample the BAM file
--target-coverage	Selecting a value	↑	SomatoSim needs to iterate through more bases to select bases at the target coverage
--coverage-tolerance	Increase value	↓	Increasing the coverage tolerance allows more positions to be selected in each loop
--coverage-MQ	Change value	—	Changes the SAMtools parameter when the depth of coverage is calculated
--coverage-BQ	Change value	—	Changes the SAMtools parameter when the depth of coverage is calculated
--minimum-separation	Increase value	↑	Increasing the minimum number of bases in between spike-ins may make selecting suitable positions more difficult
--read-min-MQ	Increase value	↑	Increasing the minimum read MQ may make selecting suitable reads more difficult
--position-min-BQ	Increase value	↑	Increasing the minimum BQ for a position may make selecting suitable reads more difficult

3 Input

3.1 BAM file

The BAM file input for SomatoSim should be processed and analysis ready. This means that any pre-processing practices such as marking duplicates or base quality score recalibration should be completed prior to using SomatoSim for simulating SNVs. Additionally, the input BAM file should contain only paired end reads and be karyotypically sorted and indexed.

3.2 BED file

Random variant selection

With random variant selection, the input BED file needs only three fields: the chromosome, the genomic range start position, and the genomic range end position.

```
1 14042036 14042109
1 63048855 63048925
1 53796840 53796923
1 156127861 156127923
```

User-specified variant selection

With user-specified variant selection, the input BED file can contain specific single positions, VAF values, and variant alleles.

If the user wishes to simulate a variant at a very specific single genomic position, they can do so using the 0-based genomic coordinate format. In this format, the VAF and variant allele will still be randomly assigned to each position.

```
1 63048906 63048907
1 53796891 53796892
1 156127917 156127918
1 92729200 92729201
```

If you are specifying *only* specific VAF values, those values should be in the fourth field:

```
1 63048906 63048907 0.06
1 53796891 53796892 0.06
1 156127917 156127918 0.03
1 92729200 92729201 0.14
```

If specifying *only* specific variant alleles, those values should be in the fourth field:

```
1 63048906 63048907 G
1 53796891 53796892 T
1 156127917 156127918 G
1 92729200 92729201 G
```

If specifying *both* specific VAF values and variant alleles, the VAF values should be in the fourth field and the variant alleles should be in the fifth field, or vice versa.

```
1 63048906 63048907 0.06 G
1 53796891 53796892 0.06 T
1 156127917 156127918 0.03 G
1 92729200 92729201 0.14 G
```

4 Output

1. Simulated BAM file

A file with ".somasim.bam" suffix and the associated ".somasim.bam.bai" index file. The BAM file contains the mutated reads, while the original BAM file is unchanged.

2. Simulation output text file

A text file in a 0-based BED format containing information on the final output of the simulation. These are the final positions that were successfully mutated. The file has the following fields:

```
[chromosome/ position/ position/ input_VAF/ input_coverage/
output_VAF/ output_coverage/ ref_allele/ alt_allele/]
```

3. Simulation log file

This log file contains metrics from the variant selection stage, variant simulation stage, and the variant evaluation stage. Notably, the log file will report the user's input parameters, the version numbers of SomatoSim and its dependencies, the run time, the number of positions selected for variant simulation, the number of reads altered, the final number of positions where variants were simulated, the VAF distributions, and the variant allele distributions among other metrics.

4. Simulation failed mutation text file

A text file in a 0-based BED format containing positions with failed mutations. This file is not generated if there are no failed mutations. The file has the following fields:

```
[chromosome/ position/ position/ input_VAF/ coverage/
variant_coverage/ count_T_in/ count_G_in/ count_A_in/ count_C_in/
count_total_in/ monoallelic_in/ ref_allele/ alt_allele/ count_T_out/
count_G_out/ count_A_out/ count_C_out/ count_total_out/]
```

5 Test data

5.1 Data generation

The test data BAM file is derived from the National Institute of Standards and Technology (NIST) Genome in a Bottle (GIAB) NA12878 HiSeq 300X BAM file (Zook *et al.*, 2016).

<https://www.nist.gov/programs-projects/genome-bottle>

ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/NA12878/NIST_NA12878_HG001_HiSeq_300x/NHGRI_II_lumina300X_novoalign_bams/HG001.hs37d5.300x.bam

The exon regions used in the test BED file were derived from the GENCODE Release 27 (GRCh37) comprehensive gene annotation gff3 file (Frankish *et al.*, 2019).

https://www.gencodegenes.org/human/release_27lift37.html

ftp://ftp.ebi.ac.uk/pub/databases/gencode/Gencode_human/release_27/GRCh37_mapping/gencode.v27lift37.annotation.gff3.gz

To create the test BED file, 12 exonic regions for each chromosome (including both X and Y) were randomly selected from the GENCODE exon annotation. These selected exonic regions were intersected with the GIAB BAM file to create a sub-sampled BAM file known to contain the exonic regions in the test BED file. The advantage of this is that high depth of coverage in the original GIAB data was preserved but the BAM file size was reduced by several orders of magnitude. The BAM file was then re-aligned to the hs37d5 reference genome using SAMtools fastq and BWA-mem and had duplicate reads marked by Picard MarkDuplicates.

The test_BED_user.bed file was created by randomly selecting a single genomic position from each genomic range in the BED file and assigning it a VAF and variant allele.

5.2 Example simulation

Using the test data BAM and the BED file, the below command was run to simulate 1,000 SNVs with VAF range of 0.01 – 0.10 with a 100X target coverage and output BAM average coverage of 100X. During the read selection process, a total of 5,981 positions were checked as potential simulation positions before all 1,000 desired positions that passed the user input criteria were selected. The `--verbose` option details which positions were not selected and for what reason. For this example, 282 positions already contained an existing variant, 4,255 positions did not meet the position coverage criteria of 100X coverage with a 10% tolerance, and 441 positions were not sufficiently separated from previously selected positions (Table 2). Since SomatoSim checks if the calculated number of reads to mutate for a position is greater than zero, three positions that would have required zero reads to be mutated (after rounding) were excluded.

Users can use this report to determine the stringency of the input parameters. For this example simulation, the number of positions that do not pass the criteria can be greatly decreased by increasing the coverage tolerance. The coverage tolerance can be changed using *--coverage-tolerance* option.

```
$ somatosim \
-i test_BAM.bam \
-b test_BED.bed \
-o out_dir \
--vaf-low 0.01 \
--vaf-high 0.10 \
--number-snv 1000 \
--random-seed 38 \
--down-sample 100 \
--target-coverage 100 \
--verbose
```

Table 2. Number of positions that failed the criteria for position selection

Failed criterion	Number of positions
No existing variant	282
Within target coverage	4,255
Minimum SNV separation distance	441
Number of reads to mutate is greater than zero	3

6 References

- Frankish,A. *et al.* (2019). GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.*, **47**, D766-D773.
- Zook,J.M. *et al.* (2016). Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci Data*, **3**, 160025.