

1

2

3 **simuG: a general-purpose genome simulator**

4

5

6

7

8 Jia-Xing Yue<sup>1\*</sup> and Gianni Liti<sup>1\*</sup>

9

10 <sup>1</sup> Université Côte d'Azur, CNRS, INSERM, IRCAN, Nice, France.

11

12 corresponding author: [yuejiaxing@gmail.com](mailto:yuejiaxing@gmail.com) and [gianni.liti@unice.fr](mailto:gianni.liti@unice.fr)

13

14

15

16

17

## 1 **Abstract**

2

### 3 **Summary:**

4 Simulated genomes with pre-defined and random genomic variants can be very useful for  
5 benchmarking genomic and bioinformatics analyses. Here we introduce simuG, a  
6 lightweight tool for simulating the full-spectrum of genomic variants (SNPs, INDELs, CNVs,  
7 inversions and translocations) for any organisms (including human). The simplicity and  
8 versatility of simuG makes it a unique general purpose genome simulator for a wide-range  
9 of simulation-based applications.

10

11 **Availability and implementation:** Code in Perl along with user manual and testing data is  
12 available at <https://github.com/yjx1217/simuG>. This software is free for use under the MIT  
13 license.

14

## 15 **1 Introduction**

16 Along with the rapid progress of genome sequencing technologies, many bioinformatics  
17 tools have been developed for characterizing genomic variants based on genome  
18 sequencing data. While there is an increasing availability of experimentally validated gold-  
19 standard genome sequencing data set from real biological samples, *in silico* simulation  
20 remains a powerful approach for gauging and comparing the performance of bioinformatics  
21 tools. Correspondingly, many read simulators have been developed for different sequencing  
22 technologies, such as ART (Huang *et al.*, 2012) for Illumina and 454, SimLoRD (Stöcker *et al.*,  
23 2016) for PacBio, and DeepSimulator (Li *et al.*, 2018) for Oxford Nanopore. However, when

1 it comes to tools for simulating genome sequences with embedded variants, the choices  
2 appear much more limited. The current available tools are either too simple or too  
3 specialized. For example, SInC (Pattnaik *et al.*, 2014) can introduce random single nucleotide  
4 polymorphisms (SNPs), Insertion/Deletions (INDELs), and copy number variants (CNVs) into  
5 a user-provided reference genome but lacks the ability to simulate known variants, which is  
6 actually highly relevant in some simulation applications. Simulome (Price *et al.*, 2017) is  
7 another random variant simulator that provides finer control options, but it is designed for  
8 prokaryote genomes only. More sophisticated tools exist, such as VarSim (Mu *et al.*, 2015)  
9 and Xome-Blender (Semeraro *et al.*, 2018), but these tools are mostly tailored for human  
10 cancer genome simulation and often require additional third-party databases. Therefore,  
11 we feel there is need for a genome simulator that strikes a balance between simplicity and  
12 versatility. With this in mind, we developed a general-purpose genome simulator simuG,  
13 which is versatile enough to simulate both small (i.e. SNPs and INDELs) and large (i.e. CNVs,  
14 inversions, and translocations) genomic variants while staying lightweight with no extra  
15 dependency and minimal input requirements. In addition, simuG provides a rich array of  
16 fine-grained controls, such as simulating SNPs in different coding partitions (e.g. coding sites,  
17 noncoding sites, 4-fold degenerate sites, or 2-fold degenerate sites); simulating CNVs with  
18 different formation mechanisms (e.g. segmental deletions, dispersed duplications, and  
19 tandem duplications); and simulating inversions and translocations with specific types of  
20 breakpoints. These features together make simuG highly amenable to a wide range of  
21 application scenarios.

22

## 23 **2 Description and feature highlights**

1 simuG is a command-line tool written in Perl and supports all mainstream operating systems.  
2 It takes the user-supplied reference genome (in FASTA format) as the working template to  
3 introduce non-overlapping genomic variants of all major types (i.e. SNPs, INDELS, CNVs,  
4 inversions, and translocations). SNP and INDELS can be introduced simultaneously, whereas  
5 CNVs (implemented as segmental duplications and deletions), inversions, and translocations  
6 can be introduced with separated runs. For each variant type, simuG can simulate pre-  
7 defined or random variants depending on specified options. For pre-defined variants, a  
8 user-supplied VCF file that specifies all desired variants is needed, based on which simuG  
9 will operate on the input reference genome to introduce the corresponding variants. For  
10 random variants, simuG supports a wide-spectrum of fine control options, such as '-  
11 titv\_ratio' for specifying the transition/transversion ratio of SNPs, '-  
12 indel\_size\_powerlaw\_alpha' and '-indel\_size\_powerlaw\_constant' for specifying the size  
13 distribution of INDELS, '-cnv\_gain\_loss\_ratio' for specifying the ratio of segmental  
14 duplication versus segmental deletion, "-duplication\_tandem\_dispersed\_ratio" for  
15 specifying the ratio of tandem versus dispersed duplications, and '-centromere\_gff' for  
16 specifying the location of centromeres so that simulated random CNVs, inversions, and  
17 translocations will not disrupt the specified centromeres. An ancillary script vcf2model.pl is  
18 further provided to directly calculate the best parameter combinations for the random  
19 SNP/INDEL simulation based on real data. Moreover, given the strong association between  
20 gross chromosomal rearrangement breakpoints and repetitive sequences (e.g. transposable  
21 elements) observed in empirical studies (Zhang *et al.*, 2011; Yue *et al.*, 2017), simuG can  
22 restrict random inversions and translocations to only use user-defined breakpoints (by  
23 specifying the '-inversion\_breakpoint\_gff' or  
24 '-translocation\_breakpoint\_gff' option). The specific feature type and strand information of

1 these user-defined breakpoints will be considered during the breakpoint sampling. For  
2 example, the breakpoint pairs that can trigger inversion should belong to the same feature  
3 type but from opposite strands (e.g. inverted repeats). Also, when specified, centromeres  
4 will be given special consideration in random translocation simulation so that translocations  
5 leading to dicentric chromosomes will not be sampled. Finally, when needed, users can also  
6 define a list of chromosomes (e.g. mtDNA) to be excluded from variant introduction. Upon  
7 the completion of the simulation, three files will be produced: 1) a simulated genome  
8 bearing introduced variants in FASTA format, 2) a tabular file showing the genomic locations  
9 of all introduced variants relative to both the reference genome and the simulated genome,  
10 3) a VCF file showing the genomic locations of all introduced variants relative to the  
11 reference genome. Since simuG's major input/output formats (e.g. FASTA, VCF, and GFF3)  
12 are all widely used in the field, it should be fairly straightforward to connect simuG with  
13 other computational tools both upstream and downstream. Please note that when  
14 comparing the VCF outputs from simuG and other tools, all VCF files used for the such  
15 comparison should be normalized by tools like vt (Tan *et al.*, 2015) beforehand.

16

### 17 **3 Application demonstration**

18 To demonstrate the application of simuG in a real case scenario, we ran simuG with the  
19 budding yeast *Saccharomyces cerevisiae* (version R64-2-1) and human (version GRCh38)  
20 reference genomes to generate nine simulated genomes for each organism: A) with 10000  
21 SNPs, B) with 1000 random INDELS, C) with 10 random CNV due to segmental deletions, D)  
22 with 10 random CNV due to dispersed duplications, E) with 10 random CNV due to tandem  
23 duplications, F) with 5 random inversions, G) with 5 random inversions triggered by  
24 breakpoints sampled from pre-specified transposable elements (TEs), H) with 5 random

1 translocation, 1) with 5 random translocation triggered by breakpoints sampled from pre-  
 2 specified TEs. Based on each simulated genome, 50X 150-bp Illumina paired-end reads and  
 3 25X PacBio reads were simulated with ART (Huang *et al.*, 2012) and SimLoRd (Stöcker *et al.*,  
 4 2016) respectively and subsequently mapped to the yeast and human reference genomes.  
 5 The read mapping was performed by BWA (Li and Durbin, 2009) for Illumina reads and by  
 6 minimap2 (Li, 2018) for PacBio reads. With this setup, we evaluated the performance of  
 7 different variant callers for both small and large variants (Table 1 and Supplementary Note).  
 8 For small-variants (i.e. SNP and INDELS), we found freebayes (Garrison and Marth, 2012) and  
 9 the GATK4 HaplotypeCaller (Poplin *et al.*, 2018) both performed well, with the latter one  
 10 marginally won out in INDEL calling. For large structural variants like CNVs, inversions, and  
 11 translocations, we found both the short-read-based callers Delly (Rausch *et al.*, 2012) and  
 12 Manta (Chen *et al.*, 2016) and the long-read-based caller Sniffles (Sedlazeck *et al.*, 2018)  
 13 were able to identify most simulated events, especially when no TEs were associated with  
 14 the breakpoints. The long-read caller Sniffles showed superior accuracy in resolving the  
 15 exact breakpoints to the basepair resolution than short-read-based callers by taking  
 16 advantage of the long-reads, even with half of the sequencing coverage. Between the two  
 17 short-read-based callers, Manta outperformed Delly in terms of breakpoint accuracy at the  
 18 basepair level.  
 19

Variant type	Variant caller	Yeast			Human		
		Precision	Recall	F <sub>1</sub> score	Precision	Recall	F <sub>1</sub> score
SNP (n = 10000)	freebayes	1.000	0.971	0.985	0.999	0.981	0.990
	GATK4	1.000	0.970	0.985	1.000	0.977	0.988
INDEL	freebayes	0.954	0.931	0.942	0.939	0.930	0.935

(n = 1000)	GATK4	1.000	0.969	0.984	1.000	0.976	0.988
CNV:	Delly	1.000	1.000	1.000	1.000	1.000	1.000
segmental deletion	Manta	1.000	1.000	1.000	1.000	1.000	1.000
(n = 10)	Sniffles	1.000	1.000	1.000	1.000	1.000	1.000
CNV:	Delly	1.000	0.875	0.933	1.000	0.906	0.951
dispersed duplication	Manta	1.000	1.906	0.951	1.000	0.906	0.951
(n = 10)	Sniffles	1.000	0.875	0.933	1.000	0.906	0.951
CNV:	Delly	1.000	1.000	1.000	1.000	0.700	0.824
tandem duplication	Manta	1.000	1.000	1.000	1.000	0.700	0.824
(n = 10)	Sniffles	1.000	1.000	1.000	1.000	0.800	0.889
Inversion	Delly	1.000	1.000	1.000	1.000	1.000	1.000
(n = 5)	Manta	1.000	1.000	1.000	1.000	1.000	1.000
	Sniffles	1.000	1.000	1.000	1.000	1.000	1.000
Inversion with TE	Delly	1.000	0.200	0.333	1.000	1.000	1.000
breakpoints	Manta	1.000	0.200	0.333	1.000	1.000	1.000
(n = 5)	Sniffles	1.000	0.200	0.333	1.000	1.000	1.000
Translocation	Delly	1.000	1.000	1.000	0.800	0.800	0.800
(n = 5)	Manta	1.000	1.000	1.000	1.000	1.000	1.000
	Sniffles	1.000	1.000	1.000	1.000	1.000	1.000
Translocation with TE	Delly	NA	0.000	NA	1.000	1.000	1.000
breakpoints	Manta	NA	0.000	NA	1.000	1.000	1.000
(n = 5)	Sniffles	NA	0.000	NA	1.000	1.000	1.000

1

2 **Table 1. Benchmarking popular variant callers with the small and large genomic variants simulated by**  
3 **simuG.** For each variant type, number of introduced variants are shown in parentheses. TE: transposable  
4 elements (full-length Ty1 for *S. cerevisiae* and full-length intact L1 for human). Precision = true  
5 positive/(true positive + false positive). Recall = true positive/(true positive + false negative).  $F_1$  score = 2  
6 \* (recall \* precision)/(recall + precision). For a single CNV derived from dispersed duplication, there could  
7 be multiple duplicated copies inserted to different genomic locations, making it tricky to calculate

1 accuracy, precision, and  $F_1$  score by measuring the number of recovered CNV events. Therefore, we  
2 calculated these values based on the number of recovered breakpoints instead in this case.

3

#### 4 **4 Conclusions**

5 We developed simuG, a simple, flexible, and powerful tool to simulate genome sequences  
6 with both pre-defined and random genomic variants. Simple as it is, simuG is highly versatile  
7 to handle the full spectrum of genomic variants, which makes it very useful to serve the  
8 purpose of various simulation studies.

9

#### 10 **Funding**

11 This work was supported by Agence Nationale de la Recherche (ANR-16- CE12-0019 and  
12 ANR-15-IDEX-01). J.-X. Yue was supported by a postdoctoral fellowship from Fondation ARC  
13 pour la Recherche sur le Cancer (PDF20150602803). Part of computation involved in this  
14 work was performed via the Extreme Science and Engineering Discovery Environment  
15 (XSEDE) (TG-BIO170065).

16

17

18 Conflict of Interest: none declared.

19



## 1 **References**

- 2 Chen,X. *et al.* (2016) Manta: Rapid detection of structural variants and indels for germline  
3 and cancer sequencing applications. *Bioinformatics*, **32**, 1220–1222.
- 4 Garrison,E. and Marth,G. (2012) Haplotype-based variant detection from short-read  
5 sequencing. *arXiv Prepr. arXiv1207.3907*, 9.
- 6 Huang,W. *et al.* (2012) ART: A next-generation sequencing read simulator. *Bioinformatics*,  
7 **15**, 593–594.
- 8 Li,H. (2018) Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics*, **34**,  
9 3094–3100.
- 10 Li,H. and Durbin,R. (2009) Fast and accurate short read alignment with Burrows-Wheeler  
11 transform. *Bioinformatics*, **25**, 1754–1760.
- 12 Li,Y. *et al.* (2018) DeepSimulator: a deep simulator for Nanopore sequencing. *Bioinformatics*,  
13 **34**, 2899–2908.
- 14 Mu,J.C. *et al.* (2015) VarSim: a high-fidelity simulation and validation framework for high-  
15 throughput genome sequencing with cancer applications. *Bioinformatics*, **31**, 1469–  
16 1471.
- 17 Pattnaik,S. *et al.* (2014) SInC: an accurate and fast error-model based simulator for SNPs,  
18 Indels and CNVs coupled with a read generator for short-read sequence data. *BMC*  
19 *Bioinformatics*, **15**, 40.
- 20 Poplin,R. *et al.* (2018) Scaling accurate genetic variant discovery to tens of thousands of  
21 samples. *bioRxiv*, 201178.
- 22 Price,A. *et al.* (2017) Simulome: a genome sequence and variant simulator. *Bioinformatics*,  
23 **33**, 1876–1878.
- 24 Rausch,T. *et al.* (2012) DELLY: Structural variant discovery by integrated paired-end and

- 1 split-read analysis. *Bioinformatics*, **28**, i333–i339.
- 2 Sedlazeck,F.J. *et al.* (2018) Accurate detection of complex structural variations using single-  
3 molecule sequencing. *Nat. Methods*, **15**, 461–468.
- 4 Semeraro,R. *et al.* (2018) Xome-Blender: A novel cancer genome simulator. *PLoS One*, **13**,  
5 e0194472.
- 6 Stöcker,B.K. *et al.* (2016) SimLoRD: Simulation of Long Read Data. *Bioinformatics*, **32**, 2704–  
7 2706.
- 8 Tan,A. *et al.* (2015) Unified representation of genetic variants. *Bioinformatics*, **31**, 2202–  
9 2204.
- 10 Yue,J.-X. *et al.* (2017) Contrasting evolutionary genome dynamics between domesticated  
11 and wild yeasts. *Nat. Genet.*, **49**.
- 12 Zhang,J. *et al.* (2011) Transposable Elements as Catalysts for Chromosome Rearrangements.  
13 In, *Plant Chromosome Engineering*. Humana Press, Totowa, NJ, pp. 315–326.
- 14