

1 **SuPreMo: a computational tool for streamlining *in silico* perturbation using sequence-**
2 **based predictive models**

3

4 Ketrin Gjoni^{1,2} and Katherine S. Pollard^{1,2,3}

5

6 ¹Gladstone Institute of Data Science and Biotechnology, San Francisco, CA 94158, USA

7 ²Department of Epidemiology & Biostatistics, University of California, San Francisco, CA 94158,

8 USA

9 ³Chan Zuckerberg Biohub, San Francisco, CA 94158, USA

10 **Summary**

11 Computationally editing genome sequences is a common bioinformatics task, but
12 current approaches have limitations, such as incompatibility with structural variants,
13 challenges in identifying responsible sequence perturbations, and the need for vcf file
14 inputs and phased data. To address these bottlenecks, we present Sequence Mutator
15 for Predictive Models (SuPreMo), a scalable and comprehensive tool for performing *in*
16 *silico* mutagenesis. We then demonstrate how pairs of reference and perturbed
17 sequences can be used with machine learning models to prioritize pathogenic variants
18 or discover new functional sequences.

19

20 **Availability and Implementation**

21 SuPreMo was written in Python, and can be run using only one line of code to generate
22 both sequences and 3D genome disruption scores. The codebase, instructions for
23 installation and use, and tutorials are on the Github page:

24 <https://github.com/ketringjoni/SuPreMo/tree/main>.

25

26 **Contact**

27 katherine.pollard@gladstone.ucsf.edu

28

29 **Supplementary information**

30 Supplementary data are available at Bioinformatics online.

31 **1 Introduction**

32 Many machine learning (ML) models have been developed that predict cellular profiles
33 from input DNA sequences (Supp Table 1). These sequence-to-profile models can
34 predict biological features—including gene expression (Enformer (Avsec *et al.* 2021a),
35 ExPecto (Zhou *et al.* 2018), Xpresso (Agarwal and Shendure 2020)), genome folding
36 (Akita (Fudenberg, Kelley and Pollard 2020), C.origami (Tan *et al.* 2023), DeepC
37 (Schwessinger *et al.* 2020), ORCA (Zhou 2022)), chromatin accessibility (Basenji
38 (Kelley 2020), Basset (Kelley, Snoek and Rinn 2016)), and epigenetic marks
39 (DeepFIGV (Hoffman *et al.* 2019), HyenaDNA (Nguyen *et al.* 2023), Sei (Chen *et al.*
40 2022a))—with incredible accuracy. These approaches are becoming increasingly
41 popular for exploring biological questions at lower cost and higher throughput than
42 experimental methods allow, and to address questions that are not possible to test
43 experimentally. One exciting potential is to use sequence-to-profile models in tandem
44 with *in silico* mutagenesis (ISM), in order to investigate how genomic alterations alter
45 cellular profiles. This strategy generates testable, causal hypotheses about genotype-
46 phenotype relationships (Chen *et al.* 2022b). ISM has been applied to the genomes of
47 modern humans, archaic hominins (McArthur, Rinker and Capra 2021), and other
48 species (Keough *et al.* 2023) to prioritize putative pathogenic variants for experimental
49 studies (Benegas, Batra and Song 2023), decode the grammar of noncoding DNA
50 sequences (Deng *et al.* 2023), discover new sequence motifs (Avsec *et al.* 2021b),
51 design tissue-specific enhancers (Gosai *et al.* 2023), and uncover novel roles of
52 sequence elements (Gunsalus, Keiser and Pollard 2023).

53

54 In theory, ISM is very high-throughput, making it feasible to quantify the effects of a
55 large set of sequence perturbations, such as all variants in an individual's genome or a
56 cohort of patients. However, the application of ISM at scale is currently limited by the
57 process of generating sequences with and without perturbations. Existing tools
58 incorporate variants into a reference genome in ways that are not compatible with ISM
59 (bcftools consensus (Li 2011), GATK FastaAlternateReferenceMaker (Van der Auwera
60 and O'Connor 2020), perEditor (Rivas-Astroza *et al.* 2011), etc.). One of the biggest
61 limitations is that they incorporate all variants from an input variant call format (vcf)
62 (Danecek *et al.* 2011) file into a single output fasta file, making it very difficult to isolate
63 the effects of individual variants. Workarounds, such as generating an independent vcf
64 file for each variant (or variant combination) and looping over these or post-processing
65 the output fasta file to include one variant per locus, are extremely inefficient. Secondly,
66 existing tools are made largely for single nucleotide polymorphisms (SNPs) or small
67 insertions or deletions (indels), and cannot accommodate symbolic alleles—annotations
68 in vcf files for structural variants (SVs). A possible workaround is to convert symbolic
69 alleles into sequences by extracting them from a reference genome, but this becomes
70 infeasible with large structural variants due to limitations with both variant complexity
71 and memory allocation. One tool (perEditor(Rivas-Astroza *et al.* 2011)) is compatible
72 with some complex variants but is not comprehensive and has stringent requirements.
73 Finally, existing tools require the perturbations to be in a vcf format, which means that
74 pseudo input files must be generated if one wishes to apply ISM to custom or simulated
75 sequences (e.g., deleting all motifs for a given transcription factor or creating synthetic
76 enhancers).

77

78 Due to these limitations, it is common practice for ISM practitioners to write their own
79 code to generate input sequences for ISM studies. Indeed, the codebases for several
80 ML models include code examples or frameworks for performing ISM (Enformer, Sei
81 (Chen *et al.* 2022a), Basset), but these are restricted to simple variants (SNPs and
82 indels) and do not generate sequence files for input into other models. SVs make good
83 candidates for ISM since they span larger regions and are more likely to be damaging to
84 the genes, regulatory regions or active sites they overlap or neighbor. For example,
85 noncoding SVs have been shown to lead to cancer and developmental disorders by
86 disrupting genomic contacts of key genes (Paik, Maule and Gallo 2021). SVs also alter
87 more base pairs of the genome than any other type of genetic variation (1000 Genomes
88 Project Consortium *et al.* 2015). One major challenge with SVs is that, to adhere to the
89 fixed length input requirements of most ML models, input sequences must be padded,
90 and consequently, model outputs require un-padding and masking. Another
91 consideration is that—due to both biological effects and model artifacts related to
92 making predictions for fixed width genomic windows—models can be highly sensitive to
93 small changes in the input, such as masking, padding, and variant position in the
94 window. Therefore, it is important to make predictions for augmented input sequences
95 (shifted and/or reverse complement sequences) and evaluate them consistently across
96 perturbations. Thus, incorporating perturbations into a reference genome becomes
97 increasingly complicated and error-prone as variants get larger and more complex.

98

99 **2 Tool description**

100 To address these challenges, we developed SuPreMo, a framework for generating
101 perturbed sequences for input into predictive models that is scalable, flexible, and
102 comprehensive (Fig 1A). SuPreMo, which incorporates variants into the human
103 reference genome one at a time and generates model-ready sequences (Supp Fig 1A),
104 was extended to SuPreMo-Akita, which inputs those sequences into Akita (Fudenberg,
105 Kelley and Pollard 2020), an ML model that predicts chromatin contact maps, and
106 generates scores that measure variants' disruption to those maps (Supp Fig 1C).
107 Both tools accept a variety of variant files—vcf (version 4.1 and 4.2), txt, bed-like, and
108 tsv (generated from AnnotSV (Geoffroy *et al.* 2018))—making them flexible for use with
109 real or synthetic perturbations. The following variant types (marked by their Manta
110 (Chen *et al.* 2016) abbreviations) are supported: SNPs, indels, deletions (DEL),
111 duplications (DUP), inversions (INV), and complex rearrangements with breakends
112 (BNDs). Across a variety of datasets, including the 1K Genome Project, SuPreMo
113 makes it possible to analyze over 50% of SVs that would not be accessible with existing
114 tools (Fig 1B). In particular, symbolic alleles are now easily and uniformly processed
115 (Fig 1B, orange). On the other hand, insertions, which make up <20% of SVs, remain
116 inaccessible for sequence-based models since the precise inserted sequence is not
117 provided by SV calling methods (Fig 1B, gray).

118

119 SuPreMo provides flexibility through the following parameters: prediction window shift,
120 output sequence length, maximum variant length, and reverse complement. Variants

122 ends (Supp Fig 1B) or the shift parameter is used, which moves the prediction window
123 around the variant and therefore changes its position in the sequence (Supp Fig 2). The
124 length of the generated sequences can be customized to match the required input of the
125 model of interest (Supp Table 1). For SV inputs, the maximum SV length considered is
126 limited to two thirds of the input sequence, unless otherwise specified by the user.
127 Lastly, the reverse complement of the generated sequence can be outputted.
128 Generated sequences are accompanied by the relative position of the variant in each
129 sequence. Thus, SuPreMo is a flexible tool for performing ISM that can be applied
130 across sequence-based ML models.

131
132 SuPreMo-Akita generates an array of 3D genome disruption scores, predicted contact
133 frequency maps for reference (wild type) and alternate (perturbed) sequences, and
134 genomic tracks of disruption scores across the prediction window. Akita predicts contact
135 frequency maps for a ~1 megabase (Mb) input sequence at a ~2 kilobase (kb)
136 resolution. SuPreMo-Akita inputs variants as described above and optionally also takes
137 in already generated sequences. Since methods for scoring contact maps are biased
138 and sometimes only target certain features, we have made available 13 different
139 predefined metrics (Gunsalus *et al.* 2023) to use with this tool, with the defaults being
140 the most common measures: mean squared error (MSE) and Spearman's rank
141 correlation coefficient (referred to here as just correlation). To assess the robustness of
142 the generated disruption scores, the augmentation parameter optionally provides
143 averages of scores from standard sequences, sequences with -1 bp and +1 bp shifts,
144 and reverse complement sequences, or any other augmentations specified. Each

145 generated map will be accompanied by the start genomic position and the relative bin
146 that the variant lies in.

147

148 Lastly, we considered computational efficiency. To enable customization to different
149 hardware, the user can choose the number of rows to be processed at a time from the
150 input file and what outputs to request, keeping in mind storage and memory limitations.

151 We measured the run time, peak memory, and size of outputs on 3 GHz CPUs using a
152 set of 100-1000 SVs of different types from the reference cancer cell line in Figure 1B.

153 SuPreMo-Akita is fast and easily scaled up—with the augmentation parameter it takes
154 approximately 19 seconds and accumulates ~ 0.15 megabytes (MB) of memory per
155 variant (Supp Table 2).

156

157 We implemented SuPreMo using two models, although our framework is extendable to
158 any model utilizing genome sequences as input. First, we used SuPreMo with DeepSEA
159 (Zhou and Troyanskaya 2015) to rank a set of CTCF deletions based on their predicted
160 effect on epigenetic marks. Second, we used SuPreMo-Akita on cancer SVs (Supp Fig
161 3A). SVs were scored using MSE and correlation, and the top 3 scoring variants for
162 each SV type and scoring method were selected (Supp Fig 3B). We separately ranked
163 variants by their type because their 3D genome disruption scores vary, and by the
164 scoring method because each has unique biases. Using SuPreMo-Akita, contact
165 frequency maps and disruption tracks were generated for these selected SVs and the
166 most interesting variants, based on the structures they disrupt, were chosen (Supp Fig
167 3C). This method prioritized a deletion of an insulated site that is predicted to cause

168 increased contact frequency between neighboring regions (Supp Fig 3C, left panel).

169 Step-by-step instructions for both implementations are available on Github.

170

171 **3 Conclusion**

172 SuPreMo is a software tool that facilitates ISM with predictive models and extends this
173 principle with Akita to predict scores for 3D genome folding disruption. Potential use
174 cases include scoring all variants in an individual or cohort for disruption to genome
175 folding, generating predicted contact frequency maps to explore the effects of
176 noncoding variants on regulatory interactions, performing ISM to evaluate or discover
177 sequence motifs using Akita, and, more broadly, generating sequences for input into
178 predictive models of interest to evaluate variant effects. SuPreMo is scalable to a large
179 number of variants and only limited by the storage capacity the user has for the
180 expected outputs. Overall, SuPreMo allows for easy, fast and broadly applicable
181 analysis of simple variants, SVs, and chromosomal rearrangements in the context of
182 sequence-based predictive models.

183

184 **Acknowledgements**

185 We thank Shu Zhang for helpful discussion and feedback and for reviewing the manuscript. We
186 thank Maureen Pittman, Laura Gunsalus, and Shuzhen Kuang for helpful insight on using Akita.

187

188 **Funding**

189 This work was supported by the National Institutes of Health [grant number R03OD034499 and
190 grant number U01HL157989] and Gladstone Institutes.

- 191 1000 Genomes Project Consortium, Auton A, Brooks LD *et al.* A global reference for
192 human genetic variation. *Nature* 2015;**526**:68–74.
- 193 Agarwal V, Shendure J. Predicting mRNA Abundance Directly from Genomic Sequence
194 Using Deep Convolutional Neural Networks. *Cell Rep* 2020;**31**:107663.
- 195 Avsec Ž, Agarwal V, Visentin D *et al.* Effective gene expression prediction from
196 sequence by integrating long-range interactions. *Nat Methods* 2021a;**18**:1196–203.
- 197 Avsec Ž, Weilert M, Shrikumar A *et al.* Base-resolution models of transcription-factor
198 binding reveal soft motif syntax. *Nat Genet* 2021b;**53**:354–66.
- 199 Benegas G, Batra SS, Song YS. DNA language models are powerful predictors of
200 genome-wide variant effects. *bioRxiv* 2023:2022.08.22.504706.
- 201 Chen KM, Wong AK, Troyanskaya OG *et al.* A sequence-based global map of
202 regulatory activity for deciphering human genetics. *Nat Genet* 2022a;**54**:940–9.
- 203 Chen V, Yang M, Cui W *et al.* Best Practices for Interpretable Machine Learning in
204 Computational Biology. *bioRxiv* 2022b:2022.10.28.513978.
- 205 Chen X, Schulz-Trieglaff O, Shaw R *et al.* Manta: rapid detection of structural variants
206 and indels for germline and cancer sequencing applications. *Bioinformatics*
207 2016;**32**:1220–2.
- 208 Danecek P, Auton A, Abecasis G *et al.* The variant call format and VCFtools.
209 *Bioinformatics* 2011;**27**:2156–8.
- 210 Deng C, Whalen S, Steyert M *et al.* Massively parallel characterization of psychiatric

- 211 disorder-associated and cell-type-specific regulatory elements in the developing
212 human cortex. *bioRxiv* 2023, DOI: 10.1101/2023.02.15.528663.
- 213 Fudenberg G, Kelley DR, Pollard KS. Predicting 3D genome folding from DNA
214 sequence with Akita. *Nat Methods* 2020;**17**:1111–7.
- 215 Geoffroy V, Herenger Y, Kress A *et al.* AnnotSV: an integrated tool for structural
216 variations annotation. *Bioinformatics* 2018;**34**:3572–4.
- 217 Gosai SJ, Castro RI, Fuentes N *et al.* Machine-guided design of synthetic cell type-
218 specific cis -regulatory elements. *bioRxiv* 2023, DOI: 10.1101/2023.08.08.552077.
- 219 Gunsalus LM, Keiser MJ, Pollard KS. In silico discovery of repetitive elements as key
220 sequence determinants of 3D genome folding. *Cell Genomics* 2023;**3**:100410.
- 221 Gunsalus LM, McArthur E, Gjoni K *et al.* Comparing chromatin contact maps at scale:
222 methods and insights. *bioRxiv* 2023, DOI: 10.1101/2023.04.04.535480.
- 223 Hoffman GE, Bendl J, Girdhar K *et al.* Functional interpretation of genetic variants using
224 deep learning predicts impact on chromatin accessibility and histone modification.
225 *Nucleic Acids Res* 2019;**47**:10597–611.
- 226 Kaivola K, Chia R, Ding J *et al.* Genome-wide structural variant analysis identifies risk
227 loci for non-Alzheimer’s dementias. *Cell Genom* 2023;**3**:100316.
- 228 Kelley DR. Cross-species regulatory sequence activity prediction. *PLoS Comput Biol*
229 2020;**16**:e1008050.
- 230 Kelley DR, Snoek J, Rinn JL. Basset: learning the regulatory code of the accessible

- 231 genome with deep convolutional neural networks. *Genome Res* 2016;**26**:990–9.
- 232 Keough KC, Whalen S, Inoue F *et al.* Three-dimensional genome rewiring in loci with
233 human accelerated regions. *Science* 2023;**380**:eabm1696.
- 234 Li H. A statistical framework for SNP calling, mutation discovery, association mapping
235 and population genetical parameter estimation from sequencing data.
236 *Bioinformatics* 2011;**27**:2987–93.
- 237 Mahmoud M, Gobet N, Cruz-Dávalos DI *et al.* Structural variant calling: the long and the
238 short of it. *Genome Biol* 2019;**20**:246.
- 239 McArthur E, Rinker DC, Capra JA. Quantifying the contribution of Neanderthal
240 introgression to the heritability of complex traits. *Nat Commun* 2021;**12**:4481.
- 241 Nguyen E, Poli M, Faizi M *et al.* HyenaDNA: Long-Range Genomic Sequence Modeling
242 at Single Nucleotide Resolution. *ArXiv* 2023.
- 243 Paik S, Maule F, Gallo M. Dysregulation of chromatin organization in pediatric and adult
244 brain tumors: oncoepigenomic contributions to tumorigenesis and cancer stem cell
245 properties. *Genome* 2021;**64**:326–36.
- 246 Rivas-Astroza M, Xie D, Cao X *et al.* Mapping personal functional data to personal
247 genomes. *Bioinformatics* 2011;**27**:3427–9.
- 248 Schwessinger R, Gosden M, Downes D *et al.* DeepC: predicting 3D genome folding
249 using megabase-scale transfer learning. *Nat Methods* 2020;**17**:1118–24.
- 250 Talsania K, Shen T-W, Chen X *et al.* Structural variant analysis of a cancer reference

- 251 cell line sample using multiple sequencing technologies. *Genome Biol* 2022;**23**:255.
- 252 Tan J, Shenker-Tauris N, Rodriguez-Hernaez J *et al*. Cell-type-specific prediction of 3D
253 chromatin organization enables high-throughput in silico genetic screening. *Nat*
254 *Biotechnol* 2023;**41**:1140–50.
- 255 Van der Auwera GA, O'Connor BD. *Genomics in the Cloud: Using Docker, GATK, and*
256 *WDL in Terra*. "O'Reilly Media, Inc.," 2020.
- 257 Zhou J. Sequence-based modeling of three-dimensional genome architecture from
258 kilobase to chromosome scale. *Nat Genet* 2022;**54**:725–34.
- 259 Zhou J, Theesfeld CL, Yao K *et al*. Deep learning sequence-based ab initio prediction of
260 variant effects on expression and disease risk. *Nat Genet* 2018;**50**:1171–9.
- 261 Zhou J, Troyanskaya OG. Predicting effects of noncoding variants with deep learning–
262 based sequence model. *Nat Methods* 2015;**12**:931–4.
- 263 Zook JM, Hansen NF, Olson ND *et al*. A robust benchmark for detection of germline
264 large deletions and insertions. *Nat Biotechnol* 2020;**38**:1347–55.