

1 A mutation rate model at the basepair resolution identifies the mutagenic effect of 2 Polymerase III transcription

3
4
5 Vladimir Seplyarskiy^{1,2,*}, Daniel J. Lee^{1,2,*}, Evan M. Koch^{1,2,*}, Joshua S. Lichtman³, Harding H. Luan³, Shamil R.
6 Sunyaev^{1,2}

7
8 ¹Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

9 ²Brigham and Women's Hospital, Division of Genetics, Harvard Medical School, Boston, MA, USA

10 ³NGM Biopharmaceuticals, South San Francisco, CA, USA

11 *Contributed equally

12
13 ***De novo* mutations occur with substantially different rates depending on genomic location, sequence context
14 and DNA strand¹⁻⁴. The success of many human genetics techniques, especially when applied to large
15 population sequencing datasets with numerous recurrent mutations⁵⁻⁷, depends strongly on assumptions
16 about the local mutation rate. Such techniques include estimation of selection intensity⁸, inference of
17 demographic history⁹, and mapping of rare disease genes¹⁰. Here, we present Roulette, a genome-wide
18 mutation rate model at the basepair resolution that incorporates known determinants of local mutation rate
19 (<http://genetics.bwh.harvard.edu/downloads/Vova/Roulette/>). Roulette is shown to be more accurate than
20 existing models^{1,6}. Roulette has sufficient resolution at high mutation rate sites to model allele frequencies
21 under recurrent mutation. We use Roulette to refine estimates of population growth within Europe by
22 incorporating the full range of human mutation rates. The analysis of significant deviations from the model
23 predictions revealed a 10-fold increase in mutation rate in nearly all genes transcribed by Polymerase III,
24 suggesting a new mutagenic mechanism. We also detected an elevated mutation rate within transcription
25 factor binding sites restricted to sites actively utilized in testis and residing in promoters.**

26 The human single nucleotide mutation rate varies along the genome at different scales^{4,11,12}. Some of this
27 variation is explained by the combination of mutation type and immediately adjacent nucleotides,
28 conceptualized as the mutation spectra^{6,13}. The CpG di-nucleotide context induces by far the largest spectrum
29 effect because of the strongly mutagenic effect of methylation at cytosines followed by guanine¹⁴. Previous
30 studies demonstrated that the extended sequence context, well beyond the two adjacent bases, exerts an
31 additional effect on mutation rates^{1,15-17}. Mutation spectra also vary along the genome, indicating that rate
32 differences are not fully explained by the surrounding DNA sequence^{4,18}. Some of this variability tracks DNA
33 properties like replication timing and gene expression⁴. Other effects, such as spikes of multinucleotide
34 mutations in oocytes, lack obvious epigenetic correlates^{4,19,20}. In addition to regional variation, the rates of
35 many mutation types depend on the DNA strand^{4,21-23}. Transcription alters the mutation spectra between
36 transcribed and non-transcribed strands, while replication leads to differences between leading and lagging
37 strands.

38 We developed “Roulette” a mutation rate model that incorporates these factors and more (see Methods). Each
39 nucleotide has three potential mutations, and we hereafter refer to each of these potential mutations as a
40 site. The extended sequence context is included by estimating the effect of the 6 upstream and 6 downstream
41 nucleotides adjacent to each site (Figure 1a). Due to sparsity, it is impossible to accurately estimate the effect
42 of each unique 12-nucleotide context. To account for this, we estimated the effect of the central pentamer
43 (two nucleotides on either side) separately from the individual effects of the 8 more distant nucleotides, which
44 are included as covariates (Figure 1a,b). For epigenomic features, Roulette incorporates methylation level (for
45 both CpG transitions and CpG transversions), transcription direction, gene expression level in testis (for sites
46 within gene bodies), and quantitative estimates of replication direction (Figure 1a,c). The incorporation of
47 transcription and replication directions makes the model strand-dependent with unequal rates for mutations of
48 the same type on the two DNA strands. To our knowledge, strand-dependency has not been incorporated into

49 existing context-dependent and regional models^{1,6,24} of germline mutation rate, but was incorporated in the
50 context of cancer mutagenesis²⁵. The Roulette model accounts for local mutation rate variation by including
51 the observed mutability of each tri-nucleotide context in 50KB windows (Figure 1d). Known epigenetic factors
52 contribute to the regional variation of mutation rate at this scale, but the SNV density is a more direct proxy to
53 the local mutation rates than noisy epigenetic tracks. This approach also has a benefit of accounting for the
54 regional variation unexplained by existing epigenetic features^{17,18} (Supplementary Figure 1-3). Some DNA
55 repair pathways act differently in intergenic regions, gene bodies and promoters^{26,27}. We fit separate statistical
56 models for each genomic compartment and each pentamer, thereby allowing the effects of covariates to vary
57 independently among compartments and pentamers (64118 models total, see Methods). SNV probabilities
58 were modeled using logistic regression. We fit models with pairwise interactions (115 parameters) and without
59 pairwise interactions (25 parameters) between all covariates and selected the best performing model using
60 cross-validation based on a 50/50 test/train split (see Methods). Two simpler models were also analyzed to
61 prevent overfitting in pentamer-compartment pairs with too few mutations. Finally, we grouped the predicted
62 rates into 100 bins because discrete mutation rate classes facilitate many applications such as analyses of allele
63 frequency distributions.

64 It is impossible to fit parameter-rich mutation models to currently available *de novo* mutation datasets because
65 of data sparsity. To train Roulette, we collected all non-coding SNVs with frequency below 0.001 from gnomAD
66 v3 whole genomes⁶ (524M rare SNVs total). The derived allele of these rare SNVs correspond to mutational
67 events. We assume that rare alleles are always derived (as opposed to ancestral); simulations suggest that this
68 is violated at most for one in 33,000 SNVs (see Methods). The distribution of very rare non-coding SNVs along
69 the genome is primarily driven by mutation rate differences, with the effects of biased gene conversion, direct
70 and background selection being negligible⁴.

71
72 Due to the sample size of contemporary human sequencing data, many rare SNVs represent recurrent
73 mutations that have occurred multiple times in the genealogical history of the sequenced cohort. Because
74 Roulette only fits the density of monomorphic sites, we transformed SNV probabilities to the mutation rate
75 scale by assuming the probability a site remains monomorphic is given by the zero class of the Poisson
76 distribution for the expected number of variants per site. The expected number of variants is proportional to
77 mutation rate and the overall coalescent depth. We assume that the coalescent depth is approximately
78 constant for a very large sample from a growing population (see Methods).

79
80 After estimation and rescaling, we found that Roulette captures expected genomic mutation rate variation
81 when applied to synonymous sites not used in the model training. For instance, nearly two-fold rate
82 differences between the transcribed and non-transcribed strands are predicted accurately (Figure 1c). Despite
83 not using replication timing, histone modifications, or recombination rate²⁸ as covariates, the direct inclusion of
84 the regional variation is able to capture associations between mutability and these epigenetic factors
85 (Supplementary Figure 2). The importance of this regional correction is illustrated by DNA segments that are
86 hypermutable in oocytes (sometimes called regions of maternal mutagenesis)^{19,20,29,30}. Maternal mutagenesis is
87 responsible for a localized increase in C>G mutations on the left arm of the chromosome 8 (Figure 1d,
88 Supplementary Figure 3).

89
90 As a second point of validation, we tested whether Roulette estimates resolve the old riddle of “cryptic
91 variation.” Early comparative genomics literature³¹⁻³³ observed that the frequency of triallelic SNVs is higher
92 than expected based on the probability of pairs of biallelic SNVs assuming independence and a three
93 nucleotide mutational model. Roulette accurately predicts the probability of triallelic SNVs (Supplementary
94 Figure 4), suggesting that previous observations of “cryptic variation” reflected residual mutation rate variance
95 in earlier models associated with extended nucleotide context and local genomic factors.

96
97 We compared Roulette with two existing mutation rate models to further validate its performance.

98 Karczewski et al. (2020)⁶ which used trinucleotide context and methylation levels to estimate rates for
99 gnomAD v2, and Carlson et al. (2018)¹ which used heptamer context along with several epigenetic features
100 including methylation levels for the BRIDGES study (see Supplementary Table 1 for a more detailed description
101 of model differences). We re-fit the model from Karczewski et al. (2020)⁶ on gnomAD v3, and we used the
102 publicly available estimates from Carlson et al. (2018). We hereafter refer to these models as gnomAD and
103 Carlson.

104
105 While previous studies evaluated goodness of fit of mutation rate models¹, none to our knowledge have
106 attempted to estimate the remaining residual variance. We used two novel site-by-site metrics to analyze each
107 model's ability to predict the rate and location of observed SNVs from separate datasets. The first metric is an
108 adjusted version of Nagelkerke's pseudo-R² for logistic models³⁴ that measures the residual variance between
109 the observed and expected likelihood given the inherently stochastic nature of mutational processes. Our
110 Pseudo-R² assumes that there is no variance among sites with the same predicted mutation rate, so that errors
111 result solely from misclassification among mutation rate bins. We therefore developed a second per-site metric
112 that estimates this additional variance within bins using observations of multiple mutations occurring at the
113 same site. We compare the rate of *de novo* mutations at sites where an SNV was observed to the *de novo* rate
114 at sites without SNVs. If the mutation rates corresponding to each bin are estimated without error, the *de novo*
115 mutation rate in both groups should be equal. This SNV-conditional method uses the difference in *de novo*
116 rates, depending on whether an SNV is observed or not, to estimate the within-bin variance. Both methods
117 necessarily require assumptions about the true distribution of mutation rates. For pseudo-R², we assume that
118 the full distribution is well-captured by the model even if per-site estimates are subject to error, and for the
119 SNV-conditional method, we assume that the distribution of true mutation rates within each bin is log-normal.

120
121 We first compared the models using synonymous variants from the gnomAD v2 whole exome dataset (~125K
122 individuals and ~1.9M synonymous SNVs). Since Roulette was trained on non-coding variants only, synonymous
123 variants are an independent dataset. Roulette predicts the rate of synonymous SNVs with higher accuracy than
124 the Carlson and gnomAD models, reaching a pseudo-R² of 0.86 compared to 0.81 and 0.78 respectively (Figure
125 2a). Next, we estimated pseudo-R² in a UK Biobank whole genome sequencing dataset (200K individuals) for
126 both synonymous (0.88, 0.83, 0.80) and non-coding sites (0.99, 0.94, 0.83) (Figure 2a, Supplementary Figure 5).
127 Performances increased for non-coding likely because these were used in model training. Roulette's Pseudo-R²
128 was also the highest for *de novo* synonymous mutations compiled from three independent trio-sequencing
129 studies (41,816 trios and 2,759 *de novo* synonymous mutations; Pseudo-R²: 0.93, 0.87, 0.85)^{29,35,36}.
130 We assessed Roulette's performance relative to the other mutational models using bootstrap samples of
131 synonymous sites (see Methods) and showed that Roulette provides similar improvements across all validation
132 sets ($p < 0.001$; Figure 2b).

133
134 As expected in the presence of residual mutation rate variation, sites that harbor SNVs in gnomAD had an
135 excess of *de novo* mutations even when predicted mutation rates were within the same bin. The mean excess
136 was 34% within Roulette bins, 47% within Carlson bins, and 94% within gnomAD bins (Supplementary Figure 6).
137 These result in estimated residual variances of 19%, 25%, and 51% for the Roulette, Carlson, and gnomAD
138 models (Figure 2c). While overall residual variances are larger for the SNV-conditional method, Roulette still
139 explains around 5% more of the variance in human mutation rates than the Carlson model.

140
141 Many population genetics applications rely on aggregated mutation rate estimates by gene or within a genomic
142 window. We evaluated the relevance of Roulette for these applications by aggregating synonymous sites by
143 gene for gnomAD v2 and predicting the number of SNVs. Aggregate estimates generated using Roulette are
144 more accurate than those for gnomAD or Carlson (Figure 2d). There are 1758 genes with a Z-score greater than
145 2 or less than -2 for Roulette rates, substantially fewer outlier genes than 2468 for Carlson or 2295 for the
146 gnomAD model. An area of population genetics inference with important applications in human disease
147 genetics is the estimation of selective constraints for protein truncating variants (PTVs). All methods to infer

148 strong selection rely on estimates of local mutation rate. We recomputed estimates of two measures of strong
149 heterozygous selection, s_{het} and LOEUF^{6,8}, using Roulette mutation rates. The new s_{het} estimates (available at
150 <http://genetics.bwh.harvard.edu/genescorers/selection.html>) showed a slight but statistically significant
151 improvement ($p < 0.001$) in detection of autosomal dominant disease genes annotated in DDG2P
152 (Supplementary Table 2), while updated LOEUF estimates showed no significant change.

153
154 We next investigated the utility of precise mutation rate estimates for the inference of demographic history
155 (specifically, historical changes in effective population size) from the site frequency spectrum (SFS, the number
156 of observed alleles at each frequency)^{9,37}. Most studies rely on the “infinite number of sites” model which
157 assumes that single mutation events contribute to each segregating site³⁸. Under this model and neutrality, the
158 relative distribution of allele frequencies only depends on genealogies and is independent of mutation rate,
159 while the overall level of variation is linearly dependent on mutation rate. The presence of recent recurrent
160 mutations breaks the key assumption of the infinite sites model and induces a dependency between the shape
161 of site frequency spectrum and mutation rate^{5,7,39} (Figure 3a, Supplementary Figure 7). Using a set of SFS
162 curves at different mutation rates can increase power and reduce biases due to recurrent mutations.

163
164 We re-fit a model of European demographic history⁹ to evaluate the ability of Roulette to model the shape of
165 the SFS across the range of mutation rates. We used simulations that allowed for recurrent mutations⁴⁰ to fit to
166 the whole range of mutation rates. The inclusion of high mutation rate sites is meaningful because these are
167 more informative per-variant about population growth than low-rate sites (Figure 3c). The demographic model
168 allows for faster-than-exponential growth in the recent past, and we updated the acceleration parameter from
169 1.120 to 1.122 and the initial growth rate from 0.0050 to 0.0057 with a final population size estimated at 8.1
170 million compared to 2.5 million. This model fits the shape of the SFS well even as the mutation rate becomes
171 large enough that recurrent mutation substantially skews to shape towards less rare variants^{5,7,39} (Figure 3a).
172 The fine-scale mutation rate bins defined by Roulette provide a much better fit to the SFS shape than can be
173 achieved by dividing mutations into only two bins, one for low rates and one for high (Figure 3b). This is due to
174 sufficient recurrent mutation within the low-rate bin and sufficient rate variation within the high-rate bin to
175 make single-rate summaries inadequate to capture the shape of the SFS. While one solution is to filter high
176 mutation rates sites so that the infinite sites assumption remains reasonable, this removes sites that are more
177 informative on a per-variant level (Figure 3c). This utility extends to selection inference where it is possible to
178 identify individual strongly constrained sites when mutation rates are in the neighborhood of $1e-07$ per
179 generation⁴¹.

180
181 While much of mutation rate variation is adequately captured by Roulette (Figure 1, 2) including various
182 epigenetically active sites like enhancers and promoters (Supplementary Figure 8), strong local deviations can
183 be used to identify new mutagenic mechanisms in humans. Regional variation in mutation rates and spectra
184 have previously been characterized and biologically interpreted at scales exceeding 10kb.^{4,28} However, many
185 mutagenic mechanisms arise due to epigenetic factors acting at much shorter scales. Data sparsity prevents the
186 application of unsupervised statistical techniques to characterize variation at short scales⁴. We analyzed
187 extreme deviations from Roulette predictions at the 100bp scale genome-wide (Figure 4a). The choice of the
188 scale was determined by the need to balance resolution and statistical power.

189 While many likely represent unfiltered sequencing and mapping artifacts, the most striking observation is that
190 25.6% of 100bp genomic windows with extremely high SNV counts unexplained by the Roulette features lie
191 within RNA genes transcribed by polymerase III (Pol III). These outlier windows contain multiallelic variants and
192 overall harbor over 70 SNVs per 100bp, while some windows have more than 100 SNVs. The two most
193 prominent gene classes transcribed by Pol III are tRNA and small nuclear RNA genes (RNU) (Figure 4a, b,
194 Supplementary Figure 9,10). Analyses of allelic imbalance and other sequence quality metrics suggested that
195 these are true SNVs rather than sequencing artifacts (Supplementary Figure 11). We also found that the
196 number of *de novo* mutations increases with paternal age, as expected for real germline mutations

197 (Supplementary Figure 12). Elevated mutation rates in tRNA genes were recently noted by a comparative
198 genomics study⁴², although the magnitude of the effect was likely underestimated by not accounting for
199 recurrent mutations. Similarly, while we observe a 7-fold increase in SNV rate in RNU genes (Figure 4a,b), we
200 expect that recurrent mutation means this is an underestimation of the true hypermutability. Indeed, we
201 observed that *de novo* mutations in parent-child trio sequencing studies were detected at a 32-fold (19-50,
202 95% Poisson CI) higher rate in RNU genes.

203 To validate the link between Pol III transcription and elevated mutation rate, we compared mutation rates
204 between active RNU genes and pseudogenes. The increased mutation rate is almost exclusively limited to
205 active genes, suggesting that active transcription rather than genomic location or sequence context is
206 responsible (Supplementary Figure 13a-c). The few exceptions are pseudogenes that show H3K27ac chromatin
207 marks associated with active transcription (Supplementary Figure 13a,b), suggesting that apparent
208 hypermutable RNU pseudogenes are misannotated active genes. The association between Pol III transcription
209 and high SNV density extends to all other classes of non-coding RNAs (Supplementary Figure 13d) but not to
210 SINE repeats (Supplementary Figure 13f), which may also be transcribed by Pol III⁴³.

211 We next sought to further characterize the elevated mutation rates in tRNA and RNU genes. To this end, we
212 developed a statistical model to estimate the distribution of mutation rates among observed SNVs
213 (Supplementary Note 1) by taking advantage of the fact that recurrent mutations induce a shape-dependency
214 of the SFS on mutation rate^{5,7,39}. We observed that SFS in tRNAs and RNUs have the expected shift away from
215 very rare variants, though milder than observed for variants at the top range of Roulette estimates
216 (Supplementary Figure 14). By modeling SFS for a mixture of variants with different mutation rates, we
217 estimated that mutation rate within Pol III transcripts is highly variable. Both RNU and tRNA genes have a large
218 fraction of highly mutable sites, with mutability greatly exceeding the Roulette predictions (Figure 4c,d,
219 Supplementary Figure 14). To validate the SFS-based predictions we calculated *de novo* mutation rates for RNU
220 and tRNA sites, conditioning on the presence/absence of SNVs as well as transition/transversion status. There
221 is a stark difference between estimated *de novo* mutation rates in polymorphic and monomorphic positions
222 (Figure 4e), consistent with a high heterogeneity of mutation rate within RNU and tRNA genes. SFS-based
223 analysis and the analysis of *de novo* mutations in polymorphic sites estimated that the rate of RNU transitions
224 in highly mutable sites is higher than for any Roulette bin (Figure 4c,e). RNU genes contain some of the most
225 mutable positions in the human genome. The high mutation rate in Pol III transcripts masks the effect of
226 purifying selection and leads to an unrealistic selection inference^{44,45}.

227 There are multiple, not necessarily exclusive explanations as to why Pol III transcription is strongly mutagenic.
228 First, unlike RNA polymerase II (Pol II), Pol III does not have the ability to recruit transcription coupled repair
229 (TCR). However, TCR only removes mutations on one of the two strands and thus cannot reduce the mutation
230 rate by more than a half. TCR alone is therefore insufficient to explain the observed 32-fold effect. Second,
231 transcription associated mutagenesis (TAM), a well-described phenomenon in yeasts⁴⁶, is attributed primarily
232 to ribonucleotide incorporation into DNA during transcription. A third possibility could involve an as-of-yet
233 uncharacterized transcription-associated mechanism specific to Pol III, because the biological machinery
234 transcribing Pol III-dependent genes differs substantially from the machinery for Pol II-dependent genes⁴⁷.
235 Interestingly, it was recently shown that damage-induced mutations can accumulate on the non-transcribed
236 strand outside of replication^{4,48}. However, this mechanism creates a very strong mutational asymmetry that is
237 absent for Pol III transcripts. Finally, transcription initiation by the transcription factor (TF) IIIB triggers
238 restructuring of the DNA-bound Pol III. This restructuring can be mutagenic by itself and create mutational
239 hotspots upstream of RNU genes.

240
241 Immunoglobulin kappa genes also exhibit long stretches of extreme hypermutability. In contrast to Pol III
242 transcripts, however, sequencing quality metrics raise concerns about the reliability of SNVs in these genes
243 (Supplementary Figure 11).

244

245 Transcription factor binding occurs at short scales and has been shown to be highly mutagenic in yeast and
246 human cancers either because of blocked resection of ribonucleotide primers introduced by polymerase alpha,
247 interference with the access of nucleotide excision repair, or altered DNA conformation^{27,49–51}. First, we
248 attributed TFBS activity to specific tissues by overlapping ChIP-seq signals with regions of open chromatin
249 measured by DNase I hypersensitivity. In the majority of TFBS, Roulette predicts mutation rates accurately,
250 confirming that the observed mutation rate elevation within TFBS is due to sequence context and regional
251 features⁵² included in the Roulette model (Figure 5a). TFBS active in testis are a notable exception
252 characterized by increases in the germline mutation rate over the background for most mutation types
253 (Supplementary Figure 15a, b), with the strongest effect for T>G mutations (median increase across TFs is 1.59-
254 fold, Figure 5a). This observation suggests a direct mutagenic effect of transcription factor binding.
255 Interestingly, binding of SNPC4, the factor responsible for RNU transcription, has the strongest (6-fold) impact
256 on mutation rate.

257 Furthermore, we found that the higher mutation rates are almost exclusively restricted to testis-active TFBS in
258 promoters (Figure 5b), and that TFBS overlapping multiple promoters have higher mutation rates than TFBS
259 overlapping a single promoter (Supplementary Figure 16). To allow the application of Roulette to these
260 genomic regions, we also provide mutation rates corrected for this TFBS effect (see methods, Supplementary
261 Figure 17). Interestingly, similarly to germline mutations, UV-induced mutations at TFBS in melanoma also have
262 different rates in and outside of promoters (Supplementary Figure 18)⁵¹.

263 As shown above, Roulette offers a significantly more accurate human mutational model and has demonstrated
264 utility across different biological fields. Mutation rate estimates from the three analyzed models are made
265 available here: <http://genetics.bwh.harvard.edu/downloads/Vova/Roulette/>. Future work may explain the
266 sources of the demonstrated residual mutation rate variation, some of which may derive from evolving rates
267 through time and variability between populations^{53–55}.

268 **Conflict of interest**

269 Joshua S. Lichtman and Harding H. Luan are employed by NGM Biopharmaceuticals. Vladimir Seplyarskiy,
270 Evan M. Koch and Shamil R. Sunyaev are partially funded by NGM Biopharmaceuticals.

271 **Code availability**

272 Code used to perform the analysis is available at <https://github.com/vseplyarskiy/Roulette>.

273 **Data availability**

274 Polymorphism data used in the study is freely available at <https://gnomad.broadinstitute.org/>.

275 De novo mutations have been aggregated from supplementary materials to the refs 18 and 30.

276 Mutation rate estimates for autosomes <http://genetics.bwh.harvard.edu/downloads/Vova/Roulette/>.

277 Shet values re-calculated with the help of Roulette could be found here
278 <http://genetics.bwh.harvard.edu/genescorers/selection.html>.

279 **Acknowledgements**

280 We thank J. Wakeley and L. Fan for helpful suggestions on population genetics theory. We thank D.J. Balick for
281 providing a forward Wright-Fisher simulator. This research was supported by National Institutes of Health
282 grants R35-GM127131, R01-MH101244, U01-HG012009, R01-HG010372, and R01-HG010372 along with
283 funding from NGM Biopharmaceuticals. D.J.L. was supported by NLM T15LM007092.

284

285

286 **References**

- 287 1. Carlson, J. *et al.* Extremely rare variants reveal patterns of germline mutation rate heterogeneity
288 in humans. *Nature Communications* **9**, 3753 (2018).
- 289 2. Carlson, J., DeWitt, W. S. & Harris, K. Inferring evolutionary dynamics of mutation rates through
290 the lens of mutation spectrum variation. *Current Opinion in Genetics & Development* **62**, 50–57
291 (2020).
- 292 3. Seplyarskiy, V. B. & Sunyaev, S. The origin of human mutation in light of genomic data. *Nat Rev*
293 *Genet* (2021) doi:10.1038/s41576-021-00376-2.
- 294 4. Seplyarskiy, V. B. *et al.* Population sequencing data reveal a compendium of mutational processes
295 in the human germ line. *Science* **373**, 1030–1035 (2021).
- 296 5. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291
297 (2016).
- 298 6. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456
299 humans. *Nature* **581**, 434–443 (2020).
- 300 7. Harpak, A., Bhaskar, A. & Pritchard, J. K. Mutation Rate Variation is a Primary Determinant of the
301 Distribution of Allele Frequencies in Humans. *PLOS Genetics* **12**, e1006489 (2016).
- 302 8. Cassa, C. A. *et al.* Estimating the selective effects of heterozygous protein-truncating variants from
303 human exome data. *Nat Genet* **49**, 806–810 (2017).
- 304 9. Gao, F. & Keinan, A. Explosive genetic evidence for explosive human population growth. *Curr Opin*
305 *Genet Dev* **41**, 130–139 (2016).
- 306 10. Samocha, K. E. *et al.* A framework for the interpretation of de novo mutation in human disease.
307 *Nat Genet* **46**, 944–950 (2014).

- 308 11. Terekhanova, N. V., Seplyarskiy, V. B., Soldatov, R. A. & Bazykin, G. A. Evolution of Local Mutation
309 Rate and Its Determinants. *Mol. Biol. Evol.* **34**, 1100–1109 (2017).
- 310 12. Hodgkinson, A. & Eyre-Walker, A. Variation in the mutation rate across mammalian genomes. *Nat.*
311 *Rev. Genet.* **12**, 756–766 (2011).
- 312 13. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415–421
313 (2013).
- 314 14. Ehrlich, M., Norris, K. F., Wang, R. Y., Kuo, K. C. & Gehrke, C. W. DNA cytosine methylation and
315 heat-induced deamination. *Biosci. Rep.* **6**, 387–393 (1986).
- 316 15. Dietlein, F. *et al.* Identification of cancer driver genes based on nucleotide context. *Nature*
317 *Genetics* **52**, 208–218 (2020).
- 318 16. Rodriguez-Galindo, M., Casillas, S., Weghorn, D. & Barbadilla, A. Germline de novo mutation rates
319 on exons versus introns in humans. *Nat Commun* **11**, 3304 (2020).
- 320 17. Aggarwala, V. & Voight, B. F. An expanded sequence context model broadly explains variability in
321 polymorphism levels across the human genome. *Nat. Genet.* **48**, 349–355 (2016).
- 322 18. Vöhringer, H., Hoeck, A. V., Cuppen, E. & Gerstung, M. Learning mutational signatures and their
323 multidimensional genomic properties with TensorSignatures. *Nat Commun* **12**, 3628 (2021).
- 324 19. Goldmann, J. M. *et al.* Germline de novo mutation clusters arise during oocyte aging in genomic
325 regions with high double-strand-break incidence. *Nat. Genet.* **50**, 487–492 (2018).
- 326 20. Jónsson, H. *et al.* Parental influence on human germline *de novo* mutations in 1,548 trios from
327 Iceland. *Nature* **549**, 519–522 (2017).
- 328 21. Green, P., Ewing, B., Miller, W., Thomas, P. J. & Green, E. D. Transcription-associated mutational
329 asymmetry in mammalian evolution. *Nature Genetics* **33**, 514 (2003).

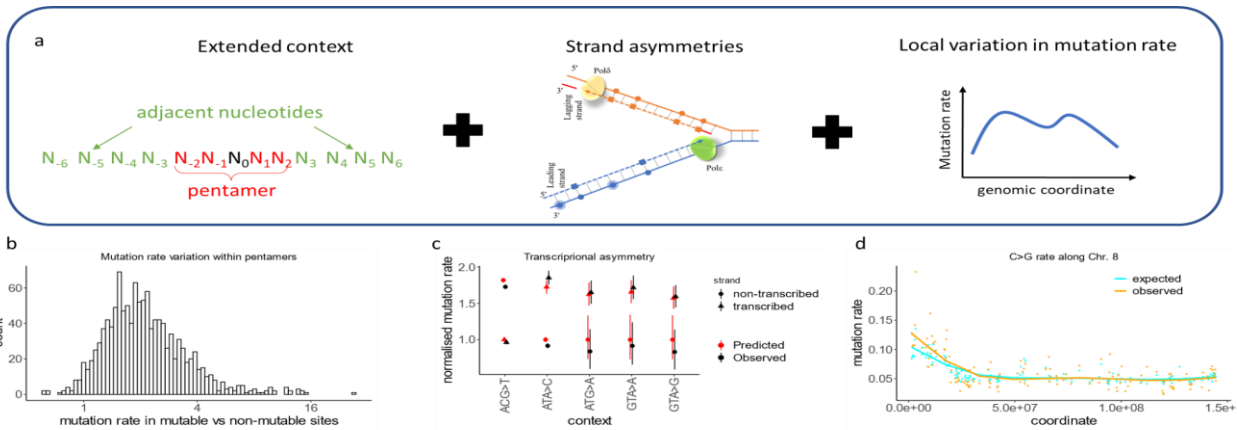
- 330 22. Chen, C.-L. *et al.* Replication-associated mutational asymmetry in the human genome. *Mol. Biol.*
331 *Evol.* **28**, 2327–2337 (2011).
- 332 23. Seplyarskiy, V. B. *et al.* Error-prone bypass of DNA lesions during lagging-strand replication is a
333 common source of germline and cancer mutations. *Nature Genetics* **51**, 36 (2019).
- 334 24. Bethune, J., Kleppe, A. & Besenbacher, S. A method to build extended sequence context models
335 of point mutations and indels. 2021.12.06.471476 Preprint at
336 <https://doi.org/10.1101/2021.12.06.471476> (2021).
- 337 25. Sherman, M. A. *et al.* Genome-wide mapping of somatic mutation rates uncovers drivers of
338 cancer. *Nat Biotechnol* **40**, 1634–1643 (2022).
- 339 26. Adar, S., Hu, J., Lieb, J. D. & Sancar, A. Genome-wide kinetics of DNA excision repair in relation to
340 chromatin state and mutagenesis. *Proc. Natl. Acad. Sci. U.S.A.* **113**, E2124–2133 (2016).
- 341 27. Perera, D. *et al.* Differential DNA repair underlies mutation hotspots at active promoters in cancer
342 genomes. *Nature* **532**, 259–263 (2016).
- 343 28. Agarwal, I. & Przeworski, M. Signatures of replication timing, recombination, and sex in the
344 spectrum of rare variants on the human X chromosome and autosomes. *PNAS* **116**, 17916–17924
345 (2019).
- 346 29. Halldorsson, B. V. *et al.* Characterizing mutagenic effects of recombination through a sequence-
347 level genetic map. *Science* **363**, eaau1043 (2019).
- 348 30. Wong, W. S. W. *et al.* New observations on maternal age effect on germline de novo mutations.
349 *Nat Commun* **7**, 10486 (2016).
- 350 31. Hodgkinson, A., Ladoukakis, E. & Eyre-Walker, A. Cryptic Variation in the Human Mutation Rate.
351 *PLOS Biology* **7**, e1000027 (2009).

- 352 32. Seplyarskiy, V. B., Kharchenko, P., Kondrashov, A. S. & Bazykin, G. A. Heterogeneity of the
353 transition/transversion ratio in *Drosophila* and Hominidae genomes. *Mol. Biol. Evol.* **29**, 1943–
354 1955 (2012).
- 355 33. Johnson, P. L. F. & Hellmann, I. Mutation Rate Distribution Inferred from Coincident SNPs and
356 Coincident Substitutions. *Genome Biol Evol* **3**, 842–850 (2011).
- 357 34. NAGELKERKE, N. J. D. A note on a general definition of the coefficient of determination.
358 *Biometrika* **78**, 691–692 (1991).
- 359 35. An, J.-Y. *et al.* Genome-wide de novo risk score implicates promoter variation in autism spectrum
360 disorder. *Science* **362**, eaat6576 (2018).
- 361 36. Satterstrom, F. K. *et al.* Large-Scale Exome Sequencing Study Implicates Both Developmental and
362 Functional Changes in the Neurobiology of Autism. *Cell* **180**, 568-584.e23 (2020).
- 363 37. Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H. & Bustamante, C. D. Inferring the Joint
364 Demographic History of Multiple Populations from Multidimensional SNP Frequency Data. *PLOS*
365 *Genetics* **5**, e1000695 (2009).
- 366 38. Crow, J. F. & Kimura, M. *An Introduction to Population Genetics Theory*. (The Blackburn Press,
367 2009).
- 368 39. Wakeley, J., Fan, W.-T. L., Koch, E. & Sunyaev, S. Recurrent mutation in the ancestry of a rare
369 variant. *Genetics* iyad049 (2023) doi:10.1093/genetics/iyad049.
- 370 40. Weghorn, D. *et al.* Applicability of the Mutation–Selection Balance Model to Population Genetics
371 of Heterozygous Protein-Truncating Variants in Humans. *Molecular Biology and Evolution* **36**,
372 1701–1710 (2019).
- 373 41. Agarwal, I. & Przeworski, M. Mutation saturation for fitness effects at human CpG sites. *Elife* **10**,
374 e71513 (2021).

- 375 42. Thornlow, B. P. *et al.* Transfer RNA genes experience exceptionally elevated mutation rates. *Proc*
376 *Natl Acad Sci U S A* **115**, 8996–9001 (2018).
- 377 43. Zhang, X.-O., Gingeras, T. R. & Weng, Z. Genome-wide analysis of polymerase III–transcribed Alu
378 elements suggests cell-type–specific enhancer function. *Genome Res.* **29**, 1402–1414 (2019).
- 379 44. Dukler, N., Mughal, M. R., Ramani, R., Huang, Y.-F. & Siepel, A. Extreme purifying selection against
380 point mutations in the human genome. *Nat Commun* **13**, 4312 (2022).
- 381 45. Zhang, X., Fang, B. & Huang, Y.-F. Transcription factor binding sites are frequently under
382 accelerated evolution in primates. *Nat Commun* **14**, 783 (2023).
- 383 46. Jinks-Robertson, S. & Bhagwat, A. S. Transcription-associated mutagenesis. *Annu. Rev. Genet.* **48**,
384 341–359 (2014).
- 385 47. Abascal-Palacios, G., Ramsay, E. P., Beuron, F., Morris, E. & Vannini, A. Structural basis of RNA
386 polymerase III transcription initiation. *Nature* **553**, 301–306 (2018).
- 387 48. Anderson, C. J. *et al.* Strand-resolved mutagenicity of DNA damage and repair. 2022.06.10.495644
388 Preprint at <https://doi.org/10.1101/2022.06.10.495644> (2022).
- 389 49. Reijns, M. A. M. *et al.* Lagging strand replication shapes the mutational landscape of the genome.
390 *Nature* **518**, 502–506 (2015).
- 391 50. Sabarinathan, R., Mularoni, L., Deu-Pons, J., Gonzalez-Perez, A. & López-Bigas, N. Nucleotide
392 excision repair is impaired by binding of transcription factors to DNA. *Nature* **532**, 264–267
393 (2016).
- 394 51. Mao, P. *et al.* ETS transcription factors induce a unique UV damage signature that drives recurrent
395 mutagenesis in melanoma. *Nature Communications* **9**, 2626 (2018).
- 396 52. Vierstra, J. *et al.* Global reference mapping of human transcription factor footprints. *Nature* **583**,
397 729–736 (2020).

- 398 53. Harris, K. Evidence for recent, population-specific evolution of the human mutation rate. *Proc.*
399 *Natl. Acad. Sci. U.S.A.* **112**, 3439–3444 (2015).
- 400 54. Harris, K. & Pritchard, J. K. Rapid evolution of the human mutation spectrum. *eLife* **6**, e24284
401 (2017).
- 402 55. Narasimhan, V. M. *et al.* Estimating the human mutation rate from autozygous segments reveals
403 population differences in human mutational processes. *Nature Communications* **8**, 303 (2017).
- 404
- 405
- 406
- 407

408 **Figures**



409 **Figure 1. Roulette accounts for extended nucleotide context, strand asymmetries and local variation in**
 410 **mutation rate.**

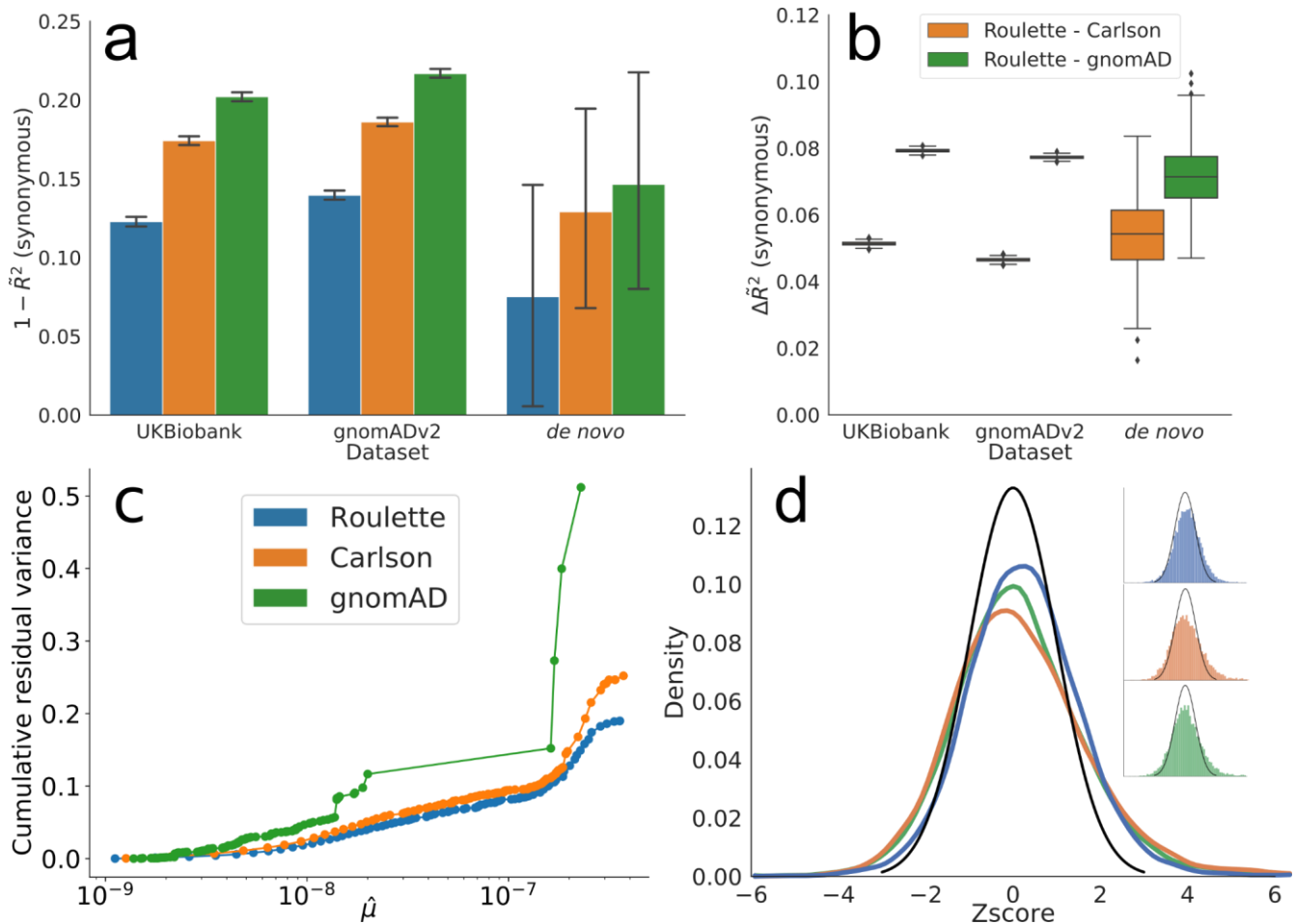
411 a) Roulette is implemented as logistic regression with pairwise interactions (see Methods). For each pentamer,
 412 we model the effect of eight surrounding nucleotides (left), strand specific information (middle), and context-
 413 specific variation along the genome (right). b) Ratio of observed *de novo* mutation rates between the Roulette
 414 predicted most and least mutable deciles for each pentamer shows large variation unexplained by the
 415 pentamer context alone. c) Effect of transcriptional asymmetry on the rate of rare synonymous SNVs in the
 416 genes with high expression in testis (top quartile). Mutation rate is relative to the least mutable strand. d) Spike
 417 of the density of rare synonymous SNVs on the left arm of chromosome 8. This region is known to be affected
 418 by increased maternal mutagenesis^{4,17,23,24}.

419

420

421

422



423

424

Figure 2 Roulette outperforms existing mutational models, under both per-gene and per-site metrics

425

426

427

428

429

430

431

432

433

434

435

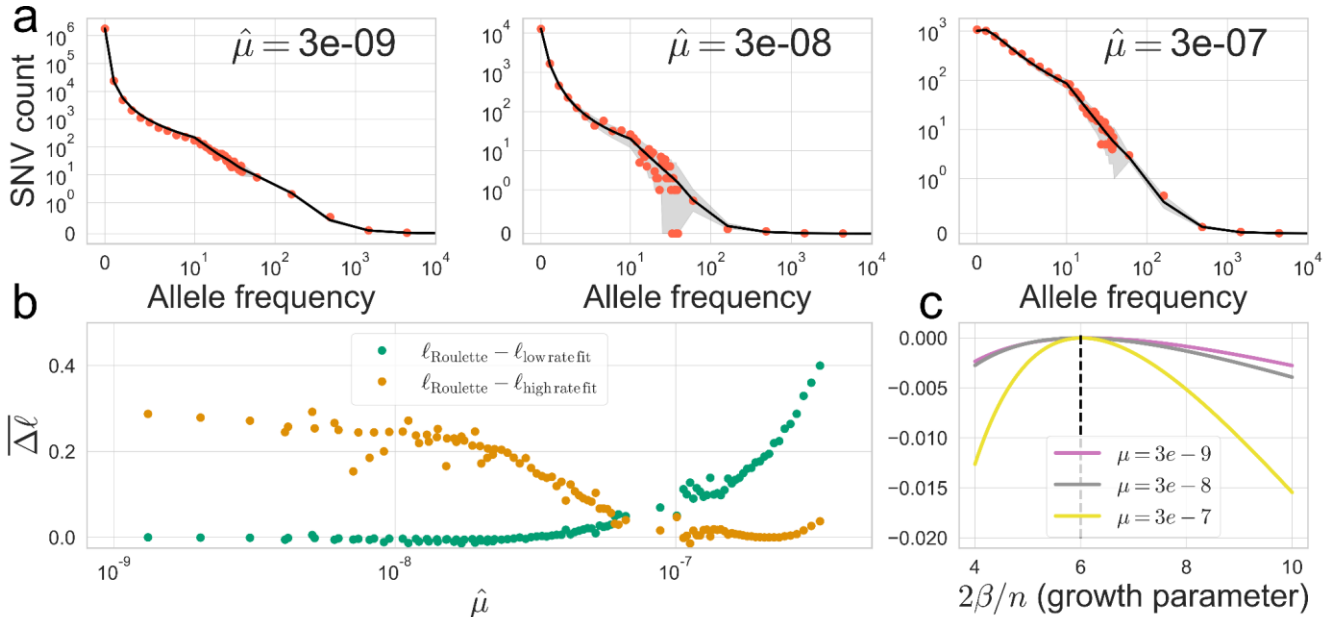
436

437

a) $1 - \text{pseudo-R}^2$ of the three mutational models on synonymous variants observed in population sequencing data (gnomAD v2.1.1 and UK Biobank) and *de novo* mutation datasets^{18,27,28}. A pseudo-R² of 0 is equivalent to using genome-wide mean mutation rate for every site. A pseudo-R² of 1 is the best per-site mutation rate estimate we can achieve, under the constraint that the mutation rates of synonymous sites follow the predicted genome-wide distribution. Error bars represent 95% confidence intervals estimated by bootstrap samples of synonymous sites. b) Difference in pseudo-R² between Roulette and the two other models. The difference was calculated over each bootstrapped sample and whiskers represent estimated 95% confidence intervals. c) The estimated cumulative residual variance for the Carlson, gnomAD and Roulette models after binning mutation rate estimates. Within-bin variance is scaled by the total variance estimated for Roulette. The x-axis gives the estimated mean in each mutation rate bin scaled to the observed per-generation *de novo* rate observed in trio data. d) Error distributions on the Z-scale for predicted counts of synonymous mutations within genes in gnomAD v2. The standard normal density is shown in black to provide a reference for the expected error distribution if mutation rates were known without error.

438

439



440

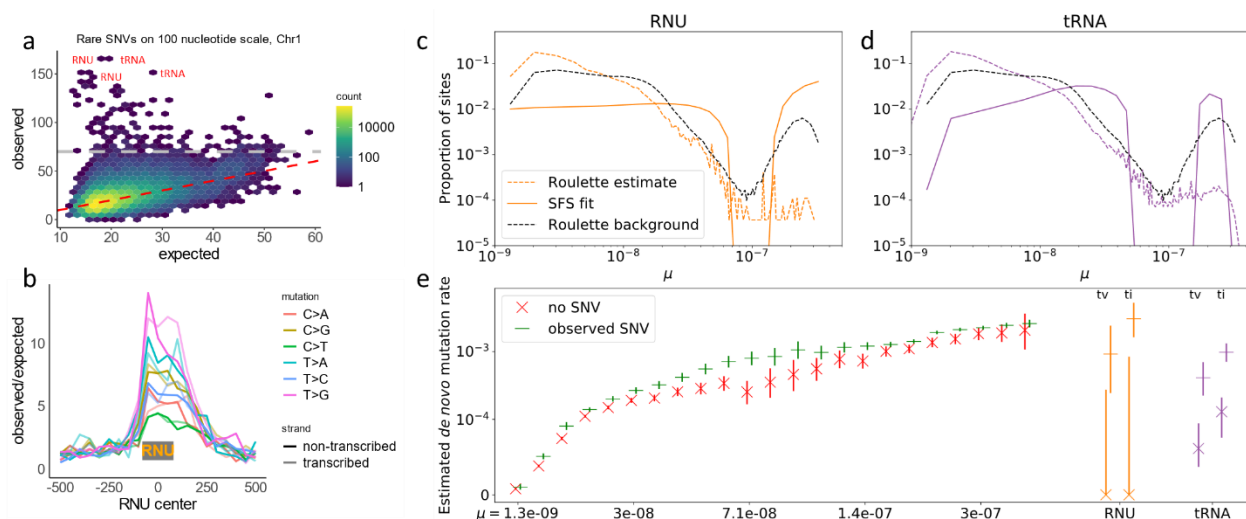
441

Figure 3 Accurate per-site mutation rate estimates improve population genetics inference

442

a) Estimated demographic history fits the SFS with mutation rate bins at different orders of magnitude. Red dots show the observed SFS at synonymous sites in gnomAD and black lines show the expected SFS under the inferred demographic model. Shaded areas correspond to 95% binomial confidence intervals. The observed SFS (red dots) shows the observed numbers of SNVs at allele counts 0-40. For more common alleles with counts above 40, red dots show numbers of SNVs for logarithmically (base 3) spaced bins. Allele counts are out of a total sample size of about 57K non-Finnish European individuals. b) Roulette bins improve fits to the shape of the SFS compared to demographic model predictions scaled to either low ($1e-09 - 3.3e-09$) or high rate ($1e-07 - 3e-07$) bins. Average log-likelihoods (per-SNV) are higher for Roulette after subtracting one to account for the additional parameter used to refit the mutation rate within each bin. Roulette improves over the model trained on sites with low mutation rate (mostly non-recurrent sites) because recurrent mutations change the shape of the SFS. It also improves over the high-rate model as one moves away from the mean mutation rate within the high-rate bin. c) High mutation rate SNVs are more informative about population growth parameters. The expected per-SNV log-likelihood relative to the maximum is shown using rare SNVs (1-40 allele counts). The compound population growth-rate / sample size parameter was chosen to approximate the observed synonymous SFS in gnomAD v2.

457



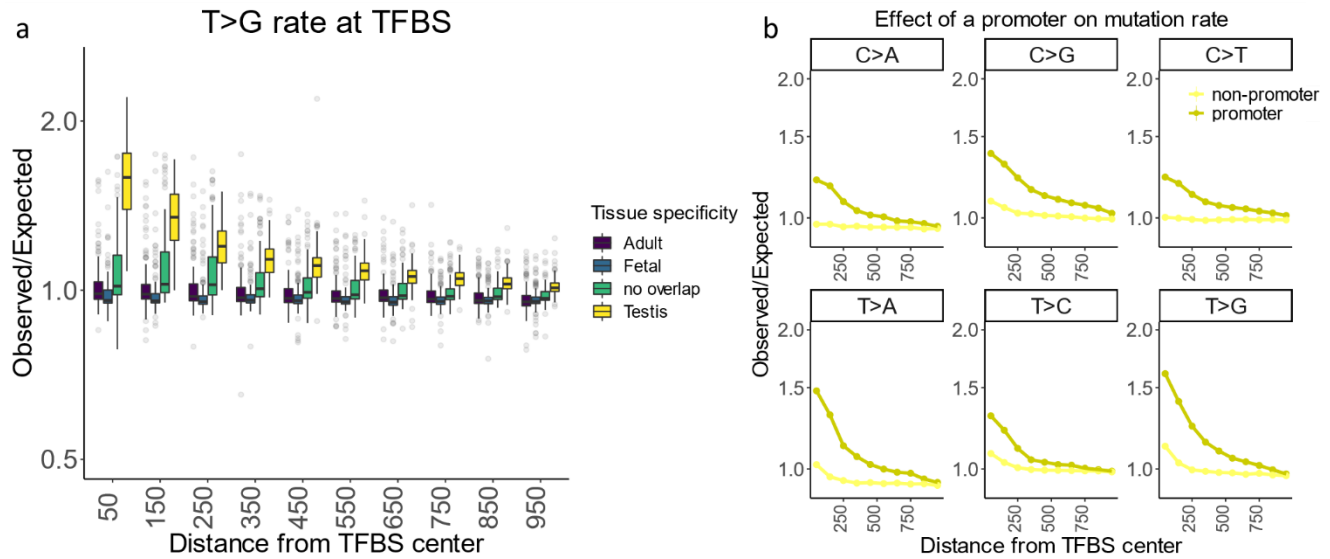
458

459 **Figure 4 Polymerase III transcripts and transcription binding sites are mutational hotspots**

460 a) Number of rare SNVs in 100 nucleotide non-overlapping windows. Expectation is calculated with Roulette.
461 While mutation counts in most regions show minor deviations from the prediction, a few loci have much higher
462 mutation rates (>70 SNVs, above the gray line). These loci are heavily enriched with Polymerase III transcripts.
463 b) Mutation rate at and around small nuclear RNAs (RNU); the median RNU size is depicted as a gray rectangle.
464 The mutation rate distributions for observed SNVs in c) RNU and d) tRNA genes was estimated by fitting the SFS
465 in these genes as a mixture of SFS shapes observed in Roulette bins. Fits are compared to the original Roulette
466 estimates and to the background rate distribution. e) SFS-based mutation rate predictions are validated by
467 estimating the *de novo* rate for mutations with and without observed SNVs in gnomAD v3. Mutations are
468 separated into transversions and transitions.

469

470



471

472

473 **Figure 5. TFBS are prone to high mutation rate.**

474 a) Box plot for the observed to expected rate of rare T>G mutations across different transcription factors.
475 Positions occupied with TF were annotated with chip-seq data. Tissues where TFBSs are active were
476 determined through overlap with tissue specific DHS peaks. b) mutagenic effect of TFBS active in testis
477 overlapping promoter (- 2 kb upstream of transcription start site, dark yellow) or not (light yellow). c) and d)
478 strand resolved observed to expected mutation rates at 100 nucleotide windows around TFBS centers in
479 promoters.

480