



Evidence for negative selection against somatic mutations induced in normal fibroblasts by *N*-ethyl-*N*-nitrosourea

Ronald Cutler, Johanna Heid, Shixiang Sun, et al.

Genome Res. published online May 15, 2026

Access the most recent version at doi:[10.1101/gr.281376.125](https://doi.org/10.1101/gr.281376.125)

P<P	Published online May 15, 2026 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Title**

2 **Evidence for negative selection against somatic mutations induced in**
3 **normal fibroblasts by *N*-ethyl-*N*-nitrosourea**

4
5 **Short Title**

6 Negative selection of mutations in normal cells

7
8 **Authors**

9 Ronald Cutler^{1,2,*,}, Johanna Heid^{1,=}, Shixiang Sun¹, Moonsook Lee¹, Alexander Y. Maslov^{1,3}, Lei
10 Zhang⁴, Simone Sidoli², Xiao Dong^{4, 5}, Jan Vijg^{1,*}

11
12 **Affiliations**

13 ¹ Department of Genetics, Albert Einstein College of Medicine, Bronx, NY 10461, USA

14 ² Department of Biochemistry, Albert Einstein College of Medicine, Bronx, NY 10461, USA

15 ³ Laboratory of Applied Genomic Technologies, Voronezh State University of Engineering
16 Technology, Voronezh 394000, Russia

17 ⁴ Institute on the Biology of Aging and Metabolism, University of Minnesota, Minneapolis, MN
18 55455, USA

19 ⁵ Department of Genetics, Cell Biology and Development, University of Minnesota, Minneapolis,
20 MN 55455, USA

21 ⁼ Authors contributed equally to this work

22 ^{*} Correspondence to J.V. (jan.vijg@einsteinmed.edu) & R.C. (ronald.cutler@einsteinmed.edu)

23

24

25

26

27 **Abstract**

28 Mutations accumulate with age in most human tissues. While some undergo clonal expansion and
29 contribute to disease, the mutational burden tolerated by a normal cell without functional decline
30 remains unknown. Here, we repeatedly treat proliferating human primary fibroblasts with the point
31 mutagen *N*-ethyl-*N*-nitrosourea, and analyze mutation burden by single-cell whole-genome
32 sequencing. Mutation burden increases linearly to ~56,000 single-nucleotide variants per cell, with
33 only a modest reduction in growth rate. We detect negative selection against potentially deleterious
34 coding and non-coding variants, including mutations affecting pathways important for cell growth
35 and maintenance. These findings suggest that selective depletion of harmful variants helps
36 proliferating cells maintain function despite an extreme mutation burden. Because most adult
37 tissues are largely non-dividing and cannot remove damaging mutations through a growth
38 disadvantage, somatic mutations that accumulate during aging may have pronounced functional
39 consequences *in vivo*.

40

41

42

43

44

45

46

47 **Introduction**

48 Aging is a complex process of continuous decline in cellular and tissue function that is
49 accompanied with increased disease risk (López-Otín et al., 2023). DNA damage is thought to play
50 a central role in this process and may be a universal cause of aging (Schumacher et al., 2021).
51 Among the major molecular consequences of DNA damage during aging are somatic DNA
52 mutations. Somatic mutations are the result of errors during repair or replication of a damaged
53 DNA template and can only be purged via the death of the cell or organism. The accumulation of
54 somatic mutations has since long been implicated as a major cause of aging (Failla, 1958; Szilard,
55 1959), which is strongly supported by the age-related exponential increase of cancer, a disease
56 known to be caused by mutations (Laconi et al., 2020). Recent single-cell and single-molecule
57 methods now allow accurate quantification of small somatic mutations in normal (non-cancerous)
58 tissues. Application of such methods has provided ample evidence that somatic mutation rate is
59 far higher than the germline mutation rate and that thousands of somatic mutations accumulate
60 with age in most, if not all, human tissues (Ren et al., 2022), including postmitotic tissues such as
61 the brain and heart (Choudhury et al., 2022; Lodato et al., 2018). In addition, evidence has emerged
62 that somatic mutations are a cause of a wide variety of human diseases other than cancer (Erickson,
63 2010; Li & Williams, 2013).

64 What thus far remains unknown is the possible functional impact, collectively, of thousands of
65 random mutations accumulating in a non-cancerous cell, for example, by adversely affecting gene-
66 coding or gene-regulatory sequences relevant for the given cell type (Vijg & Dong, 2020). If so,
67 one would expect an upper limit to the burden of random mutations a cell can tolerate or selection
68 against mutations in crucial functional regions. *N*-ethyl-*N*-nitrosourea (ENU) is well suited to test
69 this hypothesis because it is a highly efficient mutagen that induces mutations in a dose-dependent

70 manner while generating predominantly single-nucleotide variants (SNVs) rather than large-scale
71 alterations. Its mutational spectrum reflects the formation of ethyl adducts such as O6-ethylguanine
72 and O2/O4-ethylthymine and is typically enriched for C:G>T:A, T:A>C:G, and T:A>A:T
73 substitutions (Bronstein et al., 1992; Chen et al., 2000; Yang et al., 1994). Here, we treated actively
74 proliferating human primary fibroblasts with low, repeated doses of ENU to test if and how normal
75 cells can tolerate high levels of somatic mutations.

76

77 **Results**

78 **Repeated Mutagen Treatment of Normal Cells Only Slightly Affects Cell** 79 **Growth and Death**

80 To determine the impact of increased somatic mutation burden on primary human cells *in vitro*,
81 we treated low-passage fetal lung fibroblasts (IMR-90) multiple times with a sublethal dose of 50
82 $\mu\text{g/mL}$ of ENU. Each treatment was followed by a recovery period of 7 days, after which 1 million
83 cells (determined by cell counting) were replated and treated again. This process was repeated for
84 9 cycles (Figure 1A; Methods). We measured cell growth by calculating population doublings
85 (PDs), and found that repeated ENU treatment gradually reduced cell growth over the course of
86 the experiment, such that by cycle 9 control cells had undergone a total of 28.25 ± 1.87 (mean \pm
87 standard deviation) PDs compared to 24.1 ± 1.34 PDs for ENU-treated cells (Figure 1B). This
88 effect was more pronounced at the 3-day time point (Figure 1C) than at the 7-day time point
89 (Figure 1D), indicating that ENU causes an early slowing of proliferation that partially recovers
90 over time, consistent with a transient DNA damage response (Dogliotti et al., 1987). To assess
91 whether increased cell death also contributed to the overall reduction in growth, we measured

92 markers of early and late apoptosis at both time points. ENU-treated cells showed a slight but
93 consistent increase in both early and late apoptosis relative to controls at 3 and 7 days after
94 treatment (Figures 1E-F and 1G-H). Together, these results indicate that the modest decline in cell
95 growth after repeated ENU exposure reflects both an immediate slowing of proliferation and a
96 concurrent increase in apoptosis. To analyze possible effects of a DNA damage response induced
97 by ENU treatment in more detail, we performed bulk mRNA-sequencing analysis on samples
98 collected at cycles 1 and 9 for the control group and cycle 9 for the ENU-treated group (Figures
99 S1A-K; Table S1). As a positive control for the effect of cell passaging, we first compared the
100 control cycle 9 and the control cycle 1 groups, which resulted in 979 significantly differentially
101 expressed genes (DEGs) (Figures 2A; Table S2). We then performed targeted gene set enrichment
102 analysis (GSEA) which showed that cell cycle-related genes were significantly decreased, while
103 apoptosis-related genes were not significantly changed (Figures 2B-C). As cells are known to
104 become senescent during passaging, we also tested for an enrichment of senescence-related genes
105 using the SenMayo gene set (Saul et al., 2022), which were found to be significantly increased
106 (Figure 2D). Furthermore, we found the DEGs to be positively correlated with DEGs found in a
107 previous study on replicative senescence (Figure S1N) (Lackner et al., 2014).

108 To assess the effects of the repeated ENU treatment, we compared the ENU-treated cycle 9 and
109 the control cycle 9 groups. Unexpectedly, this resulted in only 2 significant DEGs (Figure 2E;
110 Table S2; Supplemental Note 1). However, after performing GSEA (which does not rely on *p*-
111 value thresholds), we found that repeated ENU treatment had a similar, albeit weaker, effect as
112 cell passaging – a significant decrease of cell cycle-related genes and a significant increase in
113 senescence-related genes (Figure 2F, 2H). In contrast to the effect of cell passaging, we found a
114 significant increase in the abundance of apoptosis-related genes in the ENU-treated cells relative

115 to the controls (Figure 2G), confirming the results of the direct apoptosis measurements (Figures
116 1G-H).

117 Taken together, these phenotypic and molecular data suggest that the observed slight decrease in
118 cell growth rate after repeated ENU treatment is, at least in part, due to DNA damage-induced
119 cell cycle inhibition and increased apoptosis. As the resulting effect sizes due to repeated ENU
120 treatment were relatively small, these results demonstrate the tolerance of normal primary cells
121 proliferating in culture to this mutagen, which is in agreement with previous literature (Barbaric
122 et al., 2007).

123

124 **Repeated Mutagen Treatment of Normal Cells Results in Extremely High** 125 **Mutation Burden**

126 Somatic DNA mutations in normal cells, which are largely unique to each cell, cannot be
127 distinguished from sequencing artifacts after bulk whole genome sequencing. Hence, we used our
128 established and highly accurate single-cell whole genome sequencing (SCWGS) assay, which has
129 been extensively validated by parallel analysis of kindred single-cell clones (Dong et al., 2017; L.
130 Zhang et al., 2024) (Supplemental Note 2), to analyze mutation burden in 21 single-cells from
131 control and ENU-treated groups collected in 3 independent batches at cycles 3, 6, and 9 from the
132 repeated ENU treatment experiment (Figure 1A, Table S3, Supplemental Methods). Cells for
133 mutation analysis were collected at the 7-day time point after each treatment cycle, as DNA
134 damage responses are abated by that time and all mutations fixed (Dogliotti et al., 1987). We
135 corrected for the less than 100% genome coverage and sensitivity, which provided estimated
136 mutation burdens higher than the observed burdens (Figures S2A-G; Supplemental Methods). We

137 also corrected for technical artifacts in the control group for cycle 1, which reduced mutation
138 burden somewhat but had no effect at all on the ENU-induced mutation burden (Figures S3A-D;
139 Supplemental Methods). We found that the number of estimated SNVs per cell in the control group
140 was initially $1,996 \pm 332$ at cycle 1, which then increased to $19,064 \pm 823$ following 3 ENU
141 treatment cycles, (~ 7.8 -fold increase compared to cycle 1, p -value = 1.54×10^{-6}), and then $33,098$
142 $\pm 1,785$ after 6 ENU treatment cycles (~ 1.7 -fold increase compared to cycle 3, p -value = $1.90 \times$
143 10^{-5}), and finally reached $55,954 \pm 4,066$ after the 9th ENU treatment cycle (~ 1.7 -fold increase
144 compared to the cycle 6, p -value = 3.33×10^{-9}) (Figures 3A and S2H). The most frequent types of
145 mutations were thymine-to-adenine (T>A) transversions, followed by thymine-to-cytosine (T>C)
146 transitions, both known to be induced by ENU (Barbaric et al., 2007) (Figures S2K-L).
147 Importantly, there was no significant difference between the mutation burden of control cells at
148 cycle 1 and at cycle 9, ruling out that the increase in SNVs in ENU-treated cells was due to
149 passaging (Supplemental Note 3). We also checked for the possibility of clonal expansion by
150 looking at shared mutations between cells within each group, but only detected 8 SNVs that were
151 shared between 2 cells in the ENU-treated group at cycle 6 (Figure S2F), indicating a shared
152 common ancestor, but not clonal expansion. This was expected as extensive clonal expansion
153 cannot occur over the relatively short number of cell divisions. We conclude that our treatment
154 paradigm increased mutational load to extremely high levels, with the increase in SNVs occurring
155 linearly with the treatment cycles (Figure 3B).

156 While ENU is known to mainly induce base substitution mutations through alkylation damage, a
157 slight increase in small insertions and deletions (INDELs) was also observed after ENU treatment,
158 starting from 123 ± 49 in the control group at cycle 1, to 246 ± 49 in the ENU-treated group at
159 cycle 3 (~ 1.8 -fold increase, p -value = 4.10×10^{-3}); however, there was no further increase at cycles

160 6 and 9 (Figures S2I-J). INDELs have been previously observed to be induced by ENU, albeit at
161 low frequencies (Watson et al., 1998), and could be caused by DNA polymerase slippage during
162 DNA replication at sites of alkylation damage. As with SNVs, there was no significant difference
163 in the amount of estimated INDELs between control cells from cycle 1 and cycle 9. While the lack
164 of an increase in INDELs after cycle 3 in the ENU-treated group may be suggestive of negative
165 selection acting to prevent a further increase, as was seen in a recent study (Zhang et al., 2025),
166 this finding remains inconclusive due to the low amount of observed INDELs per cell (~20
167 INDELs per cell).

168 We then characterized the underlying SNV mutational processes in the aggregated mutations from
169 the cells at all cycles within the control and ENU-treated groups. The ENU-treated group showed
170 an expected increase in the relative contribution of T>A transversions, followed by T>C transitions
171 (Figures 4A, S2K, S3A). This is in keeping with previously published ENU-induced mutational
172 spectra by us and others (Barbaric et al., 2007; Maslov et al., 2022). We then extracted SNV
173 mutational signatures *de novo* using non-negative matrix factorization with a reference set of
174 mutations from previous SCWGS studies (Alexandrov et al., 2020; Dong et al., 2017; Huang et
175 al., 2022; Miller et al., 2022) (Figures S3E-F). Two signatures were identified, labeled as SBSA
176 and SBSB, which were dominant in the cells of the ENU-treated and control groups, respectively
177 (Figures 4B and S3G). Importantly, we found that SBSA resembled the SIGNAL human signature
178 associated with ENU-treatment, validating the source of the induced mutations (Figure 4C) (Kucab
179 et al., 2019). We then compared our *de novo* signatures with those in the COSMIC database and
180 found that SBSB was most similar to SBS5, a ‘clock-like’ mutational signature that has been
181 previously found to correlate with cellular replication during aging (Alexandrov et al., 2020;

182 Hwang et al., 2025) (Figure 4D). Though relatively few INDELs were observed (369 in total), we
183 also provide a characterization of their mutational spectra (Figure S3H; Supplemental Note 4).

184

185 **Negative Selection Against ENU-Induced Accumulation of Damaging Coding** 186 **and Non-Coding Variants**

187 Based on the results above, we hypothesized that the sustained cell survival and growth in the face
188 of extremely high mutation burden could be explained by negative selection against damaging
189 mutations. Our dataset presented a unique opportunity to test this hypothesis as there were a total
190 of 108,600 SNVs observed collectively in all cells. Note that this was smaller than the *estimated*
191 number after corrections for sensitivity and genome coverage as mentioned above, but still
192 expected to provide sufficient power for statistical testing of selection pressure (Figure S2H).

193 To obtain a random background distribution of expected mutations against which to compare the
194 observed mutations in testing for selection pressure, we used the SigProfilerSimulator tool
195 (Bergstrom et al., 2020). This method controls for variation in genomic coverage, mutational
196 signatures, and sequence context of the observed mutations. With this, we simulated 10,000
197 instances of random mutations for each treatment cycle in both the control and ENU-treated cells
198 (Methods). Following this, both observed and expected SNVs in coding regions were annotated as
199 synonymous (non-protein-altering) or nonsynonymous (protein-altering) variants, using the
200 Ensembl Variant Effect Predictor (VEP) program (Table S4; Supplemental Methods) (McLaren et
201 al., 2016). The simulations of expected mutations performed well in predicting the number of
202 synonymous variants ($R = 0.94$, $p\text{-value} = 3.6 \times 10^{-10}$; Figure S4A), which should be under minimal
203 selection. We then measured selection pressure within each group by calculating the

204 $\text{Log}_2(\text{observed/expected})$ (written as $\text{Log}_2(\text{O/E})$ hereafter) ratios for each variant type using the
205 frequency of observed and expected variants (Lek et al., 2016). The direction of selection pressure
206 (i.e. neutral, positive, or negative) was ascertained by comparing the resulting $\text{Log}_2(\text{O/E})$ ratio to
207 0, where no difference from 0 is what would be expected under neutral selection (Greenman et al.,
208 2006) (Methods). We then used bootstrapping to determine the threshold for the frequency of
209 observed variants needed to obtain sufficient power for statistically significant results (p -value <
210 0.05) at various $\text{Log}_2(\text{O/E})$ ratios (Figures S4B; Table S5; Methods). Finally, to account for the
211 known effect of DNA repair in depleting the frequency of mutations within transcribed regions
212 through transcription-coupled repair (TCR) and mismatch repair (MMR) (Frigola et al., 2017;
213 Heilbrun et al., 2021), we normalized the $\text{Log}_2(\text{O/E})$ ratios of variants within transcribed regions
214 by the $\text{Log}_2(\text{O/E})$ ratio of mutations within intronic regions (assumed to be equally transcribed,
215 but under neutral selection).

216 With this approach, we found that while there were significant cycle-dependent increases in both
217 synonymous and nonsynonymous variant frequencies in the ENU-treated group relative to the
218 controls (Figure 5A), the $\text{Log}_2(\text{O/E})$ ratios of nonsynonymous variants at cycles 6 and 9 in the
219 ENU-treated group indicated that the frequency of these damaging variants was significantly lower
220 than what would be expected, consistent with negative selection (Figures 5B and S4C-E, right
221 panels). The negative $\text{Log}_2(\text{O/E})$ ratios of nonsynonymous variants at cycles 6 and 9 corresponded
222 to a ~20% reduction in mutation burden, equivalent to a decrease of 14.5 ± 9.1 and 19.9 ± 12.5
223 variants per cell at cycles 6 and 9, respectively (Figure 5C). In regard to synonymous variants, the
224 $\text{Log}_2(\text{O/E})$ ratios of these variants trended towards 0 across all cycles in the ENU-treated group,
225 which would suggest neutral selection (as would be expected for these benign variants). Of note,
226 statistically meaningful results could only be obtained with a sufficiently high number of

227 mutations, which were only present at the two highest treatment cycles. Indeed, for synonymous
228 mutations, which are generally lower in number than the nonsynonymous events, the frequency
229 even in the highest treatment cycle was slightly below the threshold set by the previously described
230 power analysis, which prevented us from making any definite conclusions (Figures 5B and S4C-
231 E, left panels; Table S5). See Supplemental Note 5 for discussion on the relationship of negative
232 selection strength and mutation burden.

233 We then sought to verify our finding of negative selection against nonsynonymous variants by
234 using a method which calculates dN/dS ratios (i.e. the *dNdScv* program), an alternative measure of
235 selection commonly used in cancer or evolutionary studies (Martincorena et al., 2017). This metric
236 compares the observed number of nonsynonymous mutations (i.e. missense, nonsense, and splice-
237 site) to the number expected under neutral evolution, using synonymous mutations as a neutral
238 reference. As this is effectively the same as calculating the O/E of nonsynonymous mutations over
239 the O/E of synonymous mutations, we calculated the $\text{Log}_2(\text{O/E})$ nonsynonymous/synonymous
240 (N/S) to facilitate the comparison. Synonymous variants, assumed to be functionally neutral, occur
241 in the same sequence contexts and at the same depth as nonsynonymous variants, so they serve as
242 an internal baseline for the true mutation rate introduced by ENU. A deficit of nonsynonymous
243 variants relative to this neutral baseline is therefore attributable to negative selection rather than to
244 sequencing bias or differences in mutagenesis. $\text{Log}_2(\text{O/E})$ (N/S) ratios were significantly below 0
245 at cycles 6 and 9 in the ENU-treated group (Figures 5D, S4F-G), indicating negative selection and
246 corroborating our previous result (Figure 5B). By contrast, using the *dNdScv* method, no evidence
247 for negative selection could be found (Figure S5A). The discrepancy is most likely due to
248 important differences between this method and our O/E method. Unlike our O/E method, which
249 empirically estimates expected mutations via genome-wide simulations accounting for coverage

250 and mutational signatures (Lek et al., 2016), the *dNdScv* method relies on probabilistic modeling
251 to estimate the expected mutations, assumes uniform genomic coverage, and expects abundant
252 mutations in coding regions. Therefore, the assumptions of this method make it less suitable for
253 our SCWGS data type as it is limited by variable genome coverage and a limited number of coding
254 mutations (Figures S2A and 5A). See Supplemental Note 6 for a detailed comparison between the
255 methods.

256 To identify the contribution of the various types of nonsynonymous variants to the negative
257 selection signal observed in the ENU-treated group, we stratified the nonsynonymous variants into
258 stop retained, missense, start lost, stop gained, and stop lost variants and calculated their respective
259 $\text{Log}_2(\text{O/E})$ ratios. This revealed that missense variants made up the majority of nonsynonymous
260 variants in the ENU-treated group (Figure 5E) and also that the negative selection signal was
261 driven almost exclusively by these variants (Figures 5F and S4H-J). See Supplemental Note 7 and
262 Table S6 for details on stop gained variants.

263 Next, we tested if variants in non-coding regions were under selection pressure. For this, we again
264 used the VEP program to annotate SNVs in non-coding regions, i.e., introns, intergenic, 5'
265 untranslated regions (UTR), 3' UTR, and splice sites, and then quantified their frequencies (Table
266 S4, Supplemental Methods) (McLaren et al., 2016). As with the coding variants, non-coding
267 variants showed cycle-dependent increases in the ENU-treated group relative to the controls, with
268 variants in intronic and intergenic regions accounting for 49.9% and 35.1% of all non-coding
269 variants, respectively (Figure 5G). For variants in intronic regions in the ENU-treated groups at all
270 cycles, $\text{Log}_2(\text{O/E})$ ratios were approximately -0.08, significantly less than what would be expected
271 by chance (Figure 5H, Table S5). The vast majority of introns (with the exception of splice sites
272 and other non-coding regulatory elements) do not contain functional sequences and are evolving

273 at neutrality in the germline. However, they are subject to TCR and MMR (Frigola et al., 2017;
274 Heilbrun et al., 2021). Hence, we interpret the negative $\text{Log}_2(\text{O/E})$ ratios in these regions as a
275 signal of TCR rather than negative selection. In intergenic regions for ENU-treated groups at all
276 cycles, we found that the $\text{Log}_2(\text{O/E})$ was approximately 0.13, significantly greater than what would
277 be expected by chance and indicative of mutation enrichment (Figure 5H, Table S5). Similar to
278 introns, the vast majority of intergenic regions lack functional sequences and are evolving at
279 neutrality in the germline. However, DNA repair activity in these regions is known to be
280 substantially reduced, which likely explains the positive $\text{Log}_2(\text{O/E})$ ratios (Dukler et al., 2022;
281 Yurchenko et al., 2023) (Supplemental Note 8).

282 Of all the non-coding variants, only $\text{Log}_2(\text{O/E})$ ratios of variants occurring in 3' UTR and splice
283 regions in the ENU-treated group at cycle 9 were significantly less than 0 (Figures 5H, S4K-M,
284 Table S5), which is consistent with negative selection and supported by previous studies on
285 germline selection pressures (Denisov et al., 2014; Dukler et al., 2022; Findlay et al., 2024). The
286 $\text{Log}_2(\text{O/E})$ ratios of 3' UTRs and splice sites in the ENU-treated group at cycle 9 corresponded to
287 ~20-25% reduction in mutation burden, a decrease of 19.3 ± 12.5 and 6.2 ± 2.7 variants per cell,
288 respectively (Figure 5I).

289 Taken together, these results show that potentially damaging variants in both coding and a subset
290 of non-coding regions did not accumulate to levels that would be expected, indicating negative
291 selection. This explains, in part, how a population of cells exposed to repeated mutagen treatments
292 was able to prevent a massive decline in growth rate, even as overall SNV burden continued to
293 increase linearly throughout the experiment. Of note, the selection pressure analysis was limited
294 to mutations that were *observed* (i.e. actually physically detected), which are ~20% of the ~56,000
295 estimated mutations at cycle 9 in the ENU-treated group (Figures 3A and S2H). Hence, these

296 results have relevance for the numbers of SNVs that are estimated to accumulate during normal
297 aging, e.g., from ~3,000 in human B lymphocytes to ~ 9,000 in human hepatocytes (Ren et al.,
298 2022).

299

300 **Pathways Supporting Fibroblast Growth and Survival are Protected from** 301 **Mutation Accumulation**

302 To gain insight into selective pressures at broader levels associated with cellular function, we
303 investigated if gene sets active in proliferating human fibroblasts were under negative selection.
304 To do this, we analyzed the SNVs across annotated functional genomic regions relevant to human
305 lung fibroblasts (i.e. exons, upstream promoters, and enhancers) and stratified these into various
306 gene sets (Boix, 2023; Boix et al., 2021; Carlson M, 2015). We then calculated the frequency of
307 SNVs within these regions as well as the $\text{Log}_2(\text{O/E})$ ratios in a similar manner to the previous
308 analyses.

309 First, we analyzed nonexpressed and expressed genes as determined by our bulk mRNA-
310 sequencing data (Figures S1A-B; Table S7). In keeping with the previous analyses of other
311 genomic regions, we found cycle-dependent increases of SNV frequency in nonexpressed and
312 expressed regions for the ENU-treated group relative to the controls (Figure 6A). $\text{Log}_2(\text{O/E})$ ratios
313 of SNVs in expressed genes were significantly below 0 at all cycles in the ENU-treated group even
314 after normalization with intronic regions, consistent with negative selection (Figures 6B and S6A-
315 C, right panels). In contrast, the $\text{Log}_2(\text{O/E})$ ratios of SNVs in nonexpressed genes did not
316 significantly deviate from the expectation in the ENU-treated group at cycle 9, consistent with
317 neutral selection (Figures 6B and S6A-C, left panels). Next, we subdivided the expressed genes

318 into nonessential and essential gene sets based on previous *in vitro* screens and repeated the same
319 analysis as above (Table S7) (Hart et al., 2015; Wang et al., 2015). This showed increased SNV
320 frequency in the ENU-treated group relative to the controls (Figure 6C), along with a $\text{Log}_2(\text{O/E})$
321 ratio significantly below 0 at cycle 9 in the ENU-treated group, consistent with negative selection
322 (Figures 6D and S6D-F, right panels). Of note, this is in agreement with global dN/dS ratios on
323 truncating mutations across 17 genes essential from a recent large-scale study sampling the normal
324 oral epithelium (Lawson et al., 2025).

325 As negative selection is expected to be generally weak due to the diploid nature of the human
326 genome (Martincorena et al., 2017; I. E. Vorontsov et al., 2016), but stronger in haploinsufficient
327 regions (i.e. regions where both alleles are required for function), we subdivided expressed genes
328 into haplosufficient ($\text{pLI} < 0.1$) and haploinsufficient ($\text{pLI} > 0.9$) gene sets based on a large scale
329 study of protein-coding genetic variation (Table S7) (Lek et al., 2016). As expected, analysis of
330 these gene sets showed increased SNV frequency in the ENU-treated group relative to the controls
331 (Figure 6E). $\text{Log}_2(\text{O/E})$ ratios were significantly below 0 in both haplosufficient and
332 haploinsufficient gene sets for the ENU-treated cells at cycle 9. However, only the ENU-treated
333 group at cycle 6 showed $\text{Log}_2(\text{O/E})$ ratios significantly below 0 for the haploinsufficient gene set
334 and not the haplosufficient gene set, suggesting overall stronger negative selection of
335 haploinsufficient genes, as would be expected (Figures 6F and S6G-I). In further support of this,
336 when the $\text{Log}_2(\text{O/E})$ was compared between haplosufficient and haploinsufficient gene sets for
337 ENU-treated cells at cycle 9, we found that the $\text{Log}_2(\text{O/E})$ ratio was considerably lower for the
338 haploinsufficient genes (-0.37 ± 0.16) as compared to the haplosufficient genes (-0.28 ± 0.04). In
339 summary, these results suggest that damaging SNVs are selectively purged from functional

340 genomic regions based on expression status, essentiality, and dosage sensitivity, extending our
341 previous results on negative selection of damaging SNVs in coding and non-coding regions.

342 Moving to the pathway level, we identified 399 Gene Ontology (GO) biological processes as
343 significantly enriched in the set of expressed genes from our bulk RNA-seq data (Figure S6J, Table
344 S8, Methods). As expected, the mean frequency of SNVs within each pathway showed significant
345 cycle-dependent increases of SNV frequency in the ENU-treated group relative to the controls
346 (Figure 6G). At the global level, $\text{Log}_2(\text{O/E})$ ratios of the identified pathways in the ENU-treated
347 group showed a cycle-dependent decrease, suggesting widespread negative selection at the
348 pathway level (Figures 6H, S6K-L, Table S9). At cycle 9 in the ENU-treated group, 6 pathways
349 showed $\text{Log}_2(\text{O/E})$ ratios significantly below 0 after multiple hypothesis correction, consistent
350 with negative selection. The pathways under negative selection were characterized as those related
351 to basic cellular processes supporting lung fibroblast cell growth *in vitro* and survival, i.e.,
352 ‘response to oxidative stress’ (GO:0006979), ‘respiratory tube development’ (GO:0030323),
353 ‘respiratory system development’ (GO:0006979), ‘regulation of response to DNA damage
354 stimulus’ (GO:2001020), ‘proteasome-mediated ubiquitin-dependent protein catabolic process’
355 (GO:0043161), and ‘cellular response to chemical stress’ (GO:0062197) (Figures 6I-K, Table S9).
356 The ‘respiratory system development’ GO term included the genes *FGF10*, *FOXF1*, and *COL3A1*,
357 which are known to be crucial for lung fibroblast identity and function (Bellusci et al., 1997;
358 Melboucy-Belkhir et al., 2014; Zoppi et al., 2004). To ensure that the pathways we had identified
359 to be under negative selection were not false positives, we repeated the pathway analysis but
360 employed bootstrapping to randomly shuffle genes amongst pathways before calculating the
361 $\text{Log}_2(\text{O/E})$ ratios (Figures S7A-C, Methods). These results had no significant overlap with the
362 pathways we identified to be under significant negative selection pressure (Figures S7D-E).

363 In summary, these results suggest that pathways supporting cell growth and survival in
364 proliferating cells *in vitro*, are protected from the accumulation of SNVs by negative selection.
365 Importantly, these pathways are biologically significant, as they appear to be relevant to the
366 conditions of *in vitro* cell culture and to the unique characteristics of lung fibroblasts,
367 highlighting the functional importance of preventing somatic mutation accumulation in
368 functional genomic regions. These analyses show how negative selection allows a population of
369 dividing cells to survive extremely high mutation burden and underscore the impact of random
370 somatic SNVs on essential cellular functions.

371

372 **Discussion**

373 Although somatic mutations accumulate with age in normal tissues, their collective functional
374 impact outside clonal expansion remains unclear (Franco et al., 2022; Ren et al., 2022). Here, we
375 tested that question directly, showing that repeated low-dose ENU treatment increased somatic
376 mutation burden in cultured fibroblasts from ~2,000 to ~56,000 SNVs per cell, more than 30-fold
377 above controls, with only modest effects on growth and death rates. One possible interpretation is
378 that fibroblasts tolerate large numbers of point mutations, which would be supported by results
379 from a recent study in which as many as ~30,000 SNVs were observed in intestinal crypts of
380 patients with genetic defects in DNA mismatch repair (Robinson et al., 2021).

381 An alternative explanation for the high tolerance of actively proliferating cells to mutation
382 accumulation is negative selection against mutations in genomic regions important for cell
383 survival. Evidence for negative selection against somatic mutations has indeed been found in
384 tumors (Bányai et al., 2021; Ilya E. Vorontsov et al., 2016; Zapata et al., 2018), although positive

385 selection of driver mutations is the dominant force in cancer (Martincorena et al., 2017). Evidence
386 for negative selection of somatic mutations in normal cells has so far been lacking, largely because
387 the smaller numbers of mutations restrict statistical power (Ren et al., 2022). Our untreated control
388 cells were similarly underpowered and did not show clear evidence of negative selection (Figure
389 S4D, Table S5), whereas ENU-treatment increased power sufficiently to gain meaningful results
390 (Figure S2H). Indeed, based on the high numbers of SNVs found after repeated treatment with
391 ENU we found evidence for negative selection acting on potentially deleterious coding and non-
392 coding variants. Although both synonymous and nonsynonymous variants increased linearly with
393 ENU treatment (Figure 5A), O/E analyses showed significant depletion of high-impact
394 nonsynonymous variants, especially missense variants, relative to the null model (Figures 5B-C).
395 This suggests that proliferating fibroblasts can selectively purge the most harmful coding
396 mutations, consistent with the selective constraint seen in germline datasets (Gudkov et al., 2024;
397 Karczewski et al., 2020; Lek et al., 2016; X. Zhang et al., 2024). Our findings are in agreement
398 with those from an *Msh2*^{-/-} mouse model, which showed negative selection against somatic stop
399 lost variants (Zhang et al., 2025). Additionally, our data showed negative selection acting on
400 variants in non-coding regions, such as 3' UTR and splice sites (Figures 5G-I), which is in
401 agreement with recent literature that has also analyzed selective pressures in these regions (Findlay
402 et al., 2024).

403 Our gene-set and pathway analyses further indicated negative selection at a broader functional
404 level, i.e., in expressed, essential, and haploinsufficient genes, which is consistent with a recent
405 population-scale study of somatic mutations in oral epithelium (Figures 6A-F) (Lawson et al.,
406 2025). We also identified signals of negative selection in actively expressed pathways, including
407 oxidative stress, respiratory system development, and response to DNA damage (Figures 6I-K),

408 all of which are relevant to our experimental conditions and cell type. These findings are consistent
409 with the idea that random somatic mutations affecting pathways important for cellular survival are
410 detrimental for lung fibroblasts proliferating *in vitro*.

411 Our work provides strong evidence that high somatic mutation burdens can affect normal cells
412 *independent* of clonally amplified mutations. Although negative selection was detected only at
413 extremely high mutation burdens, it should be noted that the number of SNVs directly observed
414 was no more than ~12,000, well below the ~56,000 estimate based on genome coverage and
415 sensitivity (Figures 3A and S2H). This is only moderately above the ~2,000-5,000 SNVs per cell
416 estimated at old age, depending on cell type (Ren et al., 2022). We therefore interpret the extreme
417 mutation burden in our system as providing the statistical power needed to detect negative
418 selection, rather than as evidence that this process operates only at extremely high mutation loads.
419 Consistent with this, we observed no association between mutation burden and the strength of
420 negative selection (Figure 5B), suggesting that negative selection may represent a general
421 mechanism for eliminating somatic mutations that impair cell function even at physiological
422 mutation burdens.

423 Another important consideration is how negative selection in our *in vitro* model compares with the
424 *in vivo* context. *In vivo*, additional mechanisms such as immune surveillance may contribute to the
425 elimination of damaged fibroblasts. More broadly, cell division rates are low in most adult tissues,
426 with notable exceptions including the lymphoid and intestinal systems (Sender & Milo, 2021). In
427 proliferating cells such as human B lymphocytes, we previously found that SNVs accumulate more
428 slowly with age in functional genomic regions than genome-wide, consistent with negative
429 selection during B-cell aging (Zhang et al., 2019). By contrast, somatic mutations in non-dividing
430 cells such as neurons and cardiomyocytes cannot be removed through a growth disadvantage.

431 Nevertheless, these cells show robust accumulation of damaging mutations (Choudhury et al.,
432 2022; Ganz et al., 2024). This suggests that age-related somatic mutation accumulation may have
433 its greatest impact in non-proliferative cells, as proposed by the disposable soma theory
434 (Kirkwood, 1977).

435 A key limitation of our study is that we mainly analyzed SNVs and the relatively small number of
436 INDELs. These variant classes are generally less likely than structural variants (SVs) (e.g. large
437 deletions, translocations, and copy number variants) to have major functional effects. Consistent
438 with this, INDELs did not accumulate linearly during ENU treatment (Figures S2I-J), and in a
439 recent study of *Msh2*^{-/-} fibroblasts we observed an upper limit of ~16,000 INDELs per cell (Zhang
440 et al., 2025). By contrast, somatic SVs remain poorly characterized because they are difficult to
441 detect accurately with current single-cell and single-molecule methods (Hård et al., 2023).
442 Although less frequent than SNVs, SVs can have major phenotypic consequences, as exemplified
443 by germline SVs in human disease (Yang, 2020). SVs may therefore have a stronger impact on
444 cell function than SNVs, but how they are selected against in non-dividing cells during aging
445 remains unclear.

446 Future efforts should extend the framework presented in this study to the *in vivo* context, especially
447 post-mitotic populations, and incorporate methods that capture SVs. Integrating mutational
448 landscapes with functional readouts will be crucial for determining how the balance between
449 information loss and selection shapes organismal aging and disease.

450

451

452

453 **Methods**

454 Additional methods can be found in the Supplemental Methods section in the Supplementary
455 Materials.

456

457 **Cell culture and treatment**

458 At the beginning of each cycle, 1×10^6 human primary lung fibroblasts (IMR-90) cells were
459 counted using a Countess 3 (Fisher Scientific), plated, and allowed to expand for 24 hours. Cells
460 were then treated with 50 $\mu\text{g}/\text{mL}$ *N*-ethyl-*N*-nitrosourea (ENU) and allowed to recover for 3 days.
461 Cells were then collected and assayed for cell number and viability for the 3-day time point. Of
462 the collected cells, 1×10^6 cells were re-plated and then allowed another 4 days to recover and
463 expand until they were collected to assess cell number and viability for the 7-day time point. This
464 procedure was repeated for 9 consecutive cycles over the course of ~10 weeks. This experiment
465 was repeated 3 separate times to create 3 independent biological replicates. Cells were isolated at
466 the 7-day time point at cycles 1, 3, 6, and 9 for somatic mutation analysis and at cycles 1 and 9 for
467 transcriptome analysis.

468

469 **Single-cell isolation, whole-genome amplification, library preparation, and** 470 **sequencing**

471 Single-cells were isolated using the CellRaft AIR system (Cell Microsystems). The CellRaft array
472 was prepared following the manufacturer's instructions and washed three times for 3 minutes with
473 warm PBS before addition of cell suspension. 1,000 cells were seeded in 1 mL for one CellRaft

474 array. Cells were allowed to settle for 3-4 hours at 37°C before isolation. Then, the raft was
475 carefully washed with warm medium to remove potential debris and dead cells. Individual rafts
476 containing one cell were deposited in PCR tubes with 2.5 µL of PBS and stored at -80°C until
477 further processing. Single-cell whole-genome amplification and library preparation for WGS was
478 performed as previously described (L. Zhang et al., 2024). In brief, single-cell whole genomes
479 were amplified using SCMDA and subject to quality control using a locus dropout test (Gundry et
480 al., 2012). Libraries were prepared using a NEBNext Ultra II FS kit (NEB). Library quantity and
481 size were assessed with a Qubit kit and TapeStation, respectively. Libraries were sequenced using
482 the 150 paired-end mode on an Illumina NovoSeq platform by Novogene Corp Inc., CA.

483

484 **Selection analysis of annotated variants**

485 Selection analysis was performed by first simulating 10,000 instances of each single-cell in our
486 dataset to be used as the null expectation (i.e. neutral selection) where mutations are randomly re-
487 distributed throughout the genome. This was done using the SigProfilerSimulator program
488 (Bergstrom et al., 2020), where for each single-cell, observed mutations were randomly re-
489 distributed across the genome while accounting for cell-specific sequencing coverage, mutational
490 signatures, and nucleotide contexts.

491 $\text{Log}_2(\text{observed/expected})$ ratios (written hereafter as $\text{Log}_2(\text{O/E})$) were calculated for each single-
492 cell or by group (mutations from all cells within a group are pooled), by first summing the
493 frequency of annotated consequences for observed and expected mutations, respectively. The
494 $\text{Log}_2(\text{O/E})$ ratio was then calculated with respect to each simulation instance (resulting in 10,000
495 $\text{Log}_2(\text{O/E})$ ratios per cell or group) using the following equation:

$$\text{Log}_2(O/E) = \text{Log}_2 \frac{(\text{observed consequence freq.} + 1)}{(\text{simulated consequence freq.} + 1)}$$

A permutation test was then used to determine the p -value of the observed consequence frequency where the simulated consequence frequency served as the null distribution and is as follows:

$$p = \frac{1}{n} \sum_{i=1}^n I(\text{simulated consequence freq.}_i \geq \text{observed consequence freq.})$$

Where p is the permutation test p -value, n is the number of permutations (i.e. 10,000), i is the i -th permutation, and I is the indicator function, which equals 1 if the condition inside is true, and 0 otherwise.

We performed a power analysis to determine how many observed variants are required to obtain reliable statistical results at varying $\text{Log}_2(O/E)$ ratios. To do this, we used the mutation data from ENU-treated cells to obtain the relationship between the standard deviation in the number of expected variants given the number of observed variants for common types of variants that were identified. We fit a gamma GLM with a log link to model standard deviation as a function of log-transformed observed variant count. The model showed excellent fit (residual deviance = 3.19 on 144 df vs. null deviance = 147.37 on 145 df; dispersion = 0.020), with log-transformed observed variant count as a highly significant predictor irrespective of the variant type (p -value < 2×10^{-16}). With this, we then calculated $\text{Log}_2(O/E)$ ratios by simulating 10,000 instances of expected variant counts while varying the observed variant counts from 1-1,000. We then varied the $\text{Log}_2(O/E)$ ratios from 0-1 by multiplying the observed variant count by a scaling factor before calculating the $\text{Log}_2(O/E)$ ratio. Finally, we calculated p -values for the $\text{Log}_2(O/E)$ ratios using a permutation test, as described above.

516 The power analysis demonstrated how varying the observed variant count affects the p -value of
517 the $\text{Log}_2(\text{O/E})$ ratio at different levels. Importantly, this analysis provided us with estimates of how
518 many observed variants were required to achieve sufficient power (p -value <0.05) (i.e. to obtain
519 robust statistical results) given a specific $\text{Log}_2(\text{O/E})$ (Figure S4B, solid line). This threshold was
520 then applied to the $\text{Log}_2(\text{O/E})$ results, at both the levels of variant types and pathways, to determine
521 which of these were amenable to further statistical testing where the results of this can be seen in
522 the figure legends (e.g. Figure 5B; Table S5). For the $\text{Log}_2(\text{O/E})$ nonsynonymous/synonymous
523 (N/S) calculation, we determined if these were amenable to further statistical testing if either the
524 nonsynonymous or synonymous $\text{Log}_2(\text{O/E})$ results in each group reached the established observed
525 variant threshold (Table S5).

526 When assessing for selection, it is also necessary to account for the cell-intrinsic bias in mutation
527 frequency due to DNA repair mechanisms (i.e. transcription coupled repair), which is known to
528 reduce mutation frequency in transcribed regions and is agnostic to damage that would result in
529 synonymous or nonsynonymous variants (Selby et al., 2023). This bias in mutation frequency was
530 indeed observed in our data and consistent with a depletion of mutations due to DNA repair (Figure
531 5H). Thus, we normalized the $\text{Log}_2(\text{O/E})$ ratio of coding variants (e.g. synonymous, missense, etc.)
532 and non-coding variants (5'/3' UTR and splice region) by the $\text{Log}_2(\text{O/E})$ ratio of mutations which
533 occurred within intronic regions and not in splice regions. We chose to use this normalization
534 because the occurrence of both exonic and intronic mutations is dependent on transcription coupled
535 repair (Heilbrun et al., 2021) however, since intronic mutations do not alter the translated protein
536 (with the exception of splice site variants), they are not subject to selective mechanisms. In
537 addition, the majority of mutations in our data occur within introns (~46.7%, Figure 5G), allowing
538 us to make more confident $\text{Log}_2(\text{O/E})$ ratio for these mutations as compared to using the $\text{Log}_2(\text{O/E})$

539 ratio of synonymous mutations which only made up only ~0.23% of all mutations. A comparison
540 of unnormalized and normalized $\text{Log}_2(\text{O/E})$ ratio for nonsynonymous and synonymous variants is
541 shown for example by comparing Figures 5B to S4C, where the effect of normalization can be
542 seen to slightly increase the $\text{Log}_2(\text{O/E})$ ratio.

543 To estimate the number of variants that are protected against resulting from negative selection, we
544 subtracted the average number of variants from simulated samples from the number of variants in
545 the respective observed sample.

546

547 **Selection analysis of gene sets and pathways**

548 Nonexpressed and expressed genes were determined using bulk mRNA transcriptome data from
549 cycle 9 control and ENU-treated groups using the zPKM program (Figure S2A; Table S7) (Hart
550 *et al.*, 2013). Nonessential and essential control gene sets were obtained from the DepMap database
551 (depmap.org) (Broad, 2024), where the nonessential gene set was the negative controls from Hart
552 *et al.* (2015) and the essential gene set was the intersection of the essential gene sets resulting from
553 the Hart *et al.* (2015) and Blomen *et al.* (2014) studies (Blomen *et al.*, 2015; Hart *et al.*, 2015). These
554 gene sets were additionally filtered for expressed genes by overlapping them with the list of
555 expressed genes we determined previously (Table S7). Haplosufficient and haploinsufficient gene
556 sets were obtained from the ExAC database where haplosufficient genes were defined as $\text{pLI} < 0.1$
557 and haploinsufficient genes were defined as $\text{pLI} > 0.9$ (Lek *et al.*, 2016). These gene sets were also
558 filtered for expressed genes in the same manner as essential genes (Table S7). Pathway level gene
559 sets were obtained by performing over representation analysis for Gene Ontology biological
560 processes using the clusterProfiler program (Wu *et al.*, 2021) on the list of expressed genes that

561 we determined previously (Table S8). For all gene sets, genomic coordinates of exons, introns,
562 promoters (5KB upstream of transcription start site), and annotated enhancers (Boix et al., 2021)
563 were retrieved using the GenomicRanges software (version 1.60.0) (Lawrence et al., 2013).

564 For each gene set, two $\text{Log}_2(\text{O/E})$ ratios were calculated. One for the version of the gene set where
565 the intronic regions were included, and one for the version of the gene set intronic regions were
566 excluded. To account for DNA repair (as described above for $\text{Log}_2(\text{O/E})$ calculation of coding and
567 non-coding variants), the $\text{Log}_2(\text{O/E})$ ratios of each gene set were normalized by subtracting
568 $\text{Log}_2(\text{O/E})$ ratio of the version that included introns from the version that excluded introns. *P*-
569 values were estimated, as described above, and were adjusted for multiple hypothesis testing using
570 the Benjamini-Hochberg method.

571 To serve as a negative control for the pathway analysis, for each enriched biological process gene
572 set, we randomly shuffled the genes between all other identified biological processes gene sets and
573 created 10,000 randomly permuted gene sets. Average $\text{Log}_2(\text{O/E})$ ratios were then calculated for
574 the 10,000 permutations of each gene set followed by *p*-value estimation, as described above.

575

576 **Data Access**

577 All raw and processed sequencing data generated in this study have been submitted to the NCBI
578 Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession number
579 GSE271867. The WGS and RNA-seq data generated in this study have been submitted to the
580 NCBI BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number

581 PRJNA1129131. Scripts used for processing and analyzing data can be found in the Supplemental
582 Material and have been deposited at:

583 [figshare.com/articles/software/Negative_selection_allows_human_primary_fibroblasts_to_tolera](https://figshare.com/articles/software/Negative_selection_allows_human_primary_fibroblasts_to_tolerate_high_somatic_mutation_loads_induced_by_N-ethyl-N-nitrosourea/26110453)
584 [te_high_somatic_mutation_loads_induced_by_N-ethyl-N-nitrosourea/26110453](https://figshare.com/articles/software/Negative_selection_allows_human_primary_fibroblasts_to_tolerate_high_somatic_mutation_loads_induced_by_N-ethyl-N-nitrosourea/26110453)

585

586 **Competing Interest Statement**

587 M.L., A.Y.M., X.D., and J.V. are co-founders and shareholders of SingulOmics Corp. A.Y.M.,
588 and J.V. are co-founders of MutagenTech Inc. All other authors declare no conflict of interest.

589

590 **Acknowledgements**

591 We would like to acknowledge the members of the Vijg and Sidoli labs for helpful discussions
592 during the preparation of the manuscript.

593

594 **Funding**

595 This study was supported by National Institutes of Health grants: T32GM007491 (R.C.)
596 T32AG023475 (R.C.), U19AG056278 (J.V.), U01HL145560 (J.V.), U01ES029519 (J.V.),
597 P01AG017242 (J.V.), P01AG047200 (J.V.), RF1AG068908 (J.V.), and P30AG038072 (J.V.).
598 Other funding sources: US Department of Defense grant BC180689P1 (J.V.), Glenn Foundation
599 for Medical Research Postdoctoral Fellowships in Aging Research (J.H.), The Michael J. Fox

600 Foundation (J.V.), and The Paul F. Glenn Center for the Biology of Human Aging at the Albert
 601 Einstein College of Medicine (J.V.).

602 **Author Contributions**

603 A.Y.M., and J.V. conceived this study and designed the experiments. J.H., M.S., and A.Y.M
 604 performed the experiments. R.C., J.H., and S.S. analyzed the data. X.D. and L.Z. provided
 605 guidance on the data analysis. R.C., J.H., S.S. and J.V. wrote the manuscript.

606

607 **References**

- 608
 609 Alexandrov, L. B., Kim, J., Haradhvala, N. J., Huang, M. N., Tian Ng, A. W., Wu, Y., Boot, A.,
 610 Covington, K. R., Gordenin, D. A., Bergstrom, E. N., Islam, S. M. A., Lopez-Bigas, N.,
 611 Klimczak, L. J., McPherson, J. R., Morganello, S., Sabarinathan, R., Wheeler, D. A.,
 612 Mustonen, V., Getz, G.,...Stratton, M. R. (2020). The repertoire of mutational signatures
 613 in human cancer. *Nature*, 578(7793), 94-101. <https://doi.org/10.1038/s41586-020-1943-3>
 614 Bányai, L., Trexler, M., Kerekes, K., Csuka, O., & Patthy, L. (2021). Use of signals of positive
 615 and negative selection to distinguish cancer genes and passenger genes. *Elife*, 10.
 616 <https://doi.org/10.7554/eLife.59629>
 617 Barbaric, I., Wells, S., Russ, A., & Dear, T. N. (2007). Spectrum of ENU-induced mutations in
 618 phenotype-driven and gene-driven screens in the mouse. *Environ Mol Mutagen*, 48(2),
 619 124-142. <https://doi.org/10.1002/em.20286>
 620 Bellusci, S., Grindley, J., Emoto, H., Itoh, N., & Hogan, B. L. (1997). Fibroblast growth factor 10
 621 (FGF10) and branching morphogenesis in the embryonic mouse lung. *Development*,
 622 124(23), 4867-4878. <https://doi.org/10.1242/dev.124.23.4867>
 623 Bergstrom, E. N., Barnes, M., Martincorena, I., & Alexandrov, L. B. (2020). Generating realistic
 624 null hypothesis of cancer mutational landscapes using SigProfilerSimulator. *BMC*
 625 *Bioinformatics*, 21(1), 438. <https://doi.org/10.1186/s12859-020-03772-3>
 626 Blomen, V. A., Májek, P., Jae, L. T., Bigenzahn, J. W., Nieuwenhuis, J., Staring, J., Sacco, R., van
 627 Diemen, F. R., Olk, N., Stukalov, A., Marceau, C., Janssen, H., Carette, J. E., Bennett, K.
 628 L., Colinge, J., Superti-Furga, G., & Brummelkamp, T. R. (2015). Gene essentiality and
 629 synthetic lethality in haploid human cells. *Science*, 350(6264), 1092-1096.
 630 <https://doi.org/10.1126/science.aac7557>
 631 Boix, C. A. (2023). *EpiMap Repository*. <https://compbio.mit.edu/epimap/>
 632 Boix, C. A., James, B. T., Park, Y. P., Meuleman, W., & Kellis, M. (2021). Regulatory genomic
 633 circuitry of human disease loci by integrative epigenomics. *Nature*, 590(7845), 300-307.
 634 <https://doi.org/10.1038/s41586-020-03145-z>

- 635 Broad, D. (2024). *DepMap 24Q4 Public* (Figshare+).
 636 <https://doi.org/10.25452/figshare.plus.27993248.v1>
- 637 Bronstein, S. M., Skopek, T. R., & Swenberg, J. A. (1992). Efficient Repair of O6-Ethylguanine,
 638 but not O4-Ethylthymine or O2-Ethylthymine, Is Dependent upon O6-Alkylguanine-DNA
 639 Alkyltransferase and Nucleotide Excision Repair Activities in Human Cells. *Cancer*
 640 *Research*, 52(7), 2008-2011.
- 641 Carlson M, M. B. (2015). *TxDb.Hsapiens.UCSC.hg19.knownGene: Annotation package for TxDb*
 642 *object(s)*.
 643 <https://bioconductor.org/packages/release/data/annotation/html/TxDb.Hsapiens.UCSC.hg>
 644 [19.knownGene.html](https://doi.org/10.1038/s43587-022-00261-5)
- 645 Chen, Y., Yee, D., Dains, K., Chatterjee, A., Cavalcoli, J., Schneider, E., Om, J., Woychik, R. P.,
 646 & Magnuson, T. (2000). Genotype-based screen for ENU-induced mutations in mouse
 647 embryonic stem cells. *Nature Genetics*, 24(3), 314-317. <https://doi.org/10.1038/73557>
- 648 Choudhury, S., Huang, A. Y., Kim, J., Zhou, Z., Morillo, K., Maury, E. A., Tsai, J. W., Miller, M.
 649 B., Lodato, M. A., Araten, S., Hilal, N., Lee, E. A., Chen, M. H., & Walsh, C. A. (2022).
 650 Somatic mutations in single human cardiomyocytes reveal age-associated DNA damage
 651 and widespread oxidative genotoxicity. *Nature Aging*, 2(8), 714-725.
 652 <https://doi.org/10.1038/s43587-022-00261-5>
- 653 Denisov, S. V., Bazykin, G. A., Sutormin, R., Favorov, A. V., Mironov, A. A., Gelfand, M. S., &
 654 Kondrashov, A. S. (2014). Weak negative and positive selection and the drift load at splice
 655 sites. *Genome Biol Evol*, 6(6), 1437-1447. <https://doi.org/10.1093/gbe/evu100>
- 656 Dogliotti, E., Vitelli, A., Terlizese, M., Muccio, A. D., Calcagnile, A., Saffiotti, U., & Bignami,
 657 M. (1987). Induction kinetics of mutations at two genetic loci, DNA damage and repair in
 658 CHO cells after different exposure times to N -ethyl- N -nitrosourea. *Carcinogenesis*, 8(1),
 659 25-31. <https://doi.org/10.1093/carcin/8.1.25>
- 660 Dong, X., Zhang, L., Milholland, B., Lee, M., Maslov, A. Y., Wang, T., & Vijg, J. (2017). Accurate
 661 identification of single-nucleotide variants in whole-genome-amplified single cells. *Nat*
 662 *Methods*, 14(5), 491-493. <https://doi.org/10.1038/nmeth.4227>
- 663 Dukler, N., Mughal, M. R., Ramani, R., Huang, Y. F., & Siepel, A. (2022). Extreme purifying
 664 selection against point mutations in the human genome. *Nat Commun*, 13(1), 4312.
 665 <https://doi.org/10.1038/s41467-022-31872-6>
- 666 Erickson, R. P. (2010). Somatic gene mutation and human disease other than cancer: an update
 667 [Research Support, Non-U.S. Gov't
 668 Review]. *Mutat Res*, 705(2), 96-106. <https://doi.org/10.1016/j.mrrev.2010.04.002>
- 669 Failla, G. (1958). The aging process and cancerogenesis. *Ann N Y Acad Sci*, 71(6), 1124-1140.
 670 <https://doi.org/10.1111/j.1749-6632.1958.tb46828.x>
- 671 Findlay, S. D., Romo, L., & Burge, C. B. (2024). Quantifying negative selection in human 3' UTRs
 672 uncovers constrained targets of RNA-binding proteins. *Nat Commun*, 15(1), 85.
 673 <https://doi.org/10.1038/s41467-023-44456-9>
- 674 Franco, I., & Eriksson, M. (2022). Reverting to old theories of ageing with new evidence for the
 675 role of somatic mutations. *Nat Rev Genet*, 23(11), 645-646.
 676 <https://doi.org/10.1038/s41576-022-00513-5>
- 677 Franco, I., Helgadottir, H. T., Moggio, A., Larsson, M., Vrtačnik, P., Johansson, A., Norgren, N.,
 678 Lundin, P., Mas-Ponte, D., Nordström, J., Lundgren, T., Stenvinkel, P., Wennberg, L.,
 679 Supek, F., & Eriksson, M. (2019). Whole genome DNA sequencing provides an atlas of

- 680 somatic mutagenesis in healthy human cells and identifies a tumor-prone cell type. *Genome*
681 *Biol*, 20(1), 285. <https://doi.org/10.1186/s13059-019-1892-z>
- 682 Franco, I., Johansson, A., Olsson, K., Vrtačnik, P., Lundin, P., Helgadottir, H. T., Larsson, M.,
683 Revêchon, G., Bosia, C., Pagnani, A., Provero, P., Gustafsson, T., Fischer, H., & Eriksson,
684 M. (2018). Somatic mutagenesis in satellite cells associates with human skeletal muscle
685 aging. *Nat Commun*, 9(1), 800. <https://doi.org/10.1038/s41467-018-03244-6>
- 686 Franco, I., Revêchon, G., & Eriksson, M. (2022). Challenges of proving a causal role of somatic
687 mutations in the aging process. *Aging Cell*, 21(5), e13613.
688 <https://doi.org/10.1111/ace1.13613>
- 689 Frigola, J., Sabarinathan, R., Mularoni, L., Muiños, F., Gonzalez-Perez, A., & López-Bigas, N.
690 (2017). Reduced mutation rate in exons due to differential mismatch repair. *Nat Genet*,
691 49(12), 1684-1692. <https://doi.org/10.1038/ng.3991>
- 692 Ganz, J., Luquette, L. J., Bizzotto, S., Miller, M. B., Zhou, Z., Bohrson, C. L., Jin, H., Tran, A. V.,
693 Viswanadham, V. V., McDonough, G., Brown, K., Chahine, Y., Chhouk, B., Galor, A.,
694 Park, P. J., & Walsh, C. A. (2024). Contrasting somatic mutation patterns in aging human
695 neurons and oligodendrocytes. *Cell*, 187(8), 1955-1970.e1923.
696 <https://doi.org/10.1016/j.cell.2024.02.025>
- 697 Greenman, C., Wooster, R., Futreal, P. A., Stratton, M. R., & Easton, D. F. (2006). Statistical
698 analysis of pathogenicity of somatic mutations in cancer. *Genetics*, 173(4), 2187-2198.
699 <https://doi.org/10.1534/genetics.105.044677>
- 700 Gudkov, M., Thibaut, L., & Giannoulatou, E. (2024). Context-adjusted proportion of singletons
701 (CAPS): a novel metric for assessing negative selection in the human genome. *NAR Genom*
702 *Bioinform*, 6(3), lqae111. <https://doi.org/10.1093/nargab/lqae111>
- 703 Gundry, M., Li, W., Maqbool, S. B., & Vijg, J. (2012). Direct, genome-wide assessment of DNA
704 mutations in single cells. *Nucleic Acids Res*, 40(5), 2032-2040.
705 <https://doi.org/10.1093/nar/gkr949>
- 706 Hård, J., Mold, J. E., Eisfeldt, J., Tellgren-Roth, C., Häggqvist, S., Bunikis, I., Contreras-Lopez,
707 O., Chin, C.-S., Nordlund, J., Rubin, C.-J., Feuk, L., Michaëlsson, J., & Ameer, A. (2023).
708 Long-read whole-genome analysis of human single cells. *Nature Communications*, 14(1),
709 5164. <https://doi.org/10.1038/s41467-023-40898-3>
- 710 Hart, T., Chandrashekhar, M., Aregger, M., Steinhart, Z., Brown, K. R., MacLeod, G., Mis, M.,
711 Zimmermann, M., Fradet-Turcotte, A., Sun, S., Mero, P., Dirks, P., Sidhu, S., Roth, F. P.,
712 Rissland, O. S., Durocher, D., Angers, S., & Moffat, J. (2015). High-Resolution CRISPR
713 Screens Reveal Fitness Genes and Genotype-Specific Cancer Liabilities. *Cell*, 163(6),
714 1515-1526. <https://doi.org/10.1016/j.cell.2015.11.015>
- 715 Hart, T., Komori, H. K., LaMere, S., Podshivalova, K., & Salomon, D. R. (2013). Finding the
716 active genes in deep RNA-seq gene expression studies. *BMC Genomics*, 14, 778.
717 <https://doi.org/10.1186/1471-2164-14-778>
- 718 Heilbrun, E. E., Merav, M., & Adar, S. (2021). Exons and introns exhibit transcriptional strand
719 asymmetry of dinucleotide distribution, damage formation and DNA repair. *NAR Genom*
720 *Bioinform*, 3(1), lqab020. <https://doi.org/10.1093/nargab/lqab020>
- 721 Huang, Z., Sun, S., Lee, M., Maslov, A. Y., Shi, M., Waldman, S., Marsh, A., Siddiqui, T., Dong,
722 X., Peter, Y., Sadoughi, A., Shah, C., Ye, K., Spivack, S. D., & Vijg, J. (2022). Single-cell
723 analysis of somatic mutations in human bronchial epithelial cells in relation to aging and
724 smoking. *Nat Genet*, 54(4), 492-498. <https://doi.org/10.1038/s41588-022-01035-w>

- 725 Hwang, T., Sitko, L. K., Khoirunnisa, R., Navarro-Aguad, F., Samuel, D. M., Park, H., Cheon, B.,
 726 Mutsnaini, L., Lee, J., Otlu, B., Takeda, S., Lee, S., Ivanov, D., & Gartner, A. (2025).
 727 Comprehensive whole-genome sequencing reveals origins of mutational signatures
 728 associated with aging, mismatch repair deficiency and temozolomide chemotherapy.
 729 *Nucleic Acids Res*, 53(1). <https://doi.org/10.1093/nar/gkae1122>
- 730 Karczewski, K. J., Francioli, L. C., Tiao, G., Cummings, B. B., Alföldi, J., Wang, Q., Collins, R.
 731 L., Laricchia, K. M., Ganna, A., Birnbaum, D. P., Gauthier, L. D., Brand, H., Solomonson,
 732 M., Watts, N. A., Rhodes, D., Singer-Berk, M., England, E. M., Seaby, E. G., Kosmicki, J.
 733 A.,...Genome Aggregation Database, C. (2020). The mutational constraint spectrum
 734 quantified from variation in 141,456 humans. *Nature*, 581(7809), 434-443.
 735 <https://doi.org/10.1038/s41586-020-2308-7>
- 736 Kirkwood, T. B. L. (1977). Evolution of ageing. *Nature*, 270(5635), 301-304.
 737 <https://doi.org/10.1038/270301a0>
- 738 Kucab, J. E., Zou, X., Morganella, S., Joel, M., Nanda, A. S., Nagy, E., Gomez, C., Degasperi, A.,
 739 Harris, R., Jackson, S. P., Arlt, V. M., Phillips, D. H., & Nik-Zainal, S. (2019). A
 740 Compendium of Mutational Signatures of Environmental Agents. *Cell*, 177(4), 821-836
 741 e816. <https://doi.org/10.1016/j.cell.2019.03.001>
- 742 Lackner, D. H., Hayashi, M. T., Cesare, A. J., & Karlseder, J. (2014). A genomics approach
 743 identifies senescence-specific gene expression regulation. *Aging Cell*, 13(5), 946-950.
 744 <https://doi.org/10.1111/accel.12234>
- 745 Laconi, E., Marongiu, F., & DeGregori, J. (2020). Cancer as a disease of old age: changing
 746 mutational and microenvironmental landscapes. *Br J Cancer*, 122(7), 943-952.
 747 <https://doi.org/10.1038/s41416-019-0721-1>
- 748 Lawrence, M., Huber, W., Pagès, H., Aboyoun, P., Carlson, M., Gentleman, R., Morgan, M. T., &
 749 Carey, V. J. (2013). Software for computing and annotating genomic ranges. *PLoS Comput*
 750 *Biol*, 9(8), e1003118. <https://doi.org/10.1371/journal.pcbi.1003118>
- 751 Lawson, A. R. J., Abascal, F., Nicola, P. A., Lensing, S. V., Roberts, A. L., Kalantzis, G., Baez-
 752 Ortega, A., Brzozowska, N., El-Sayed Moustafa, J. S., Vaitkute, D., Jakupovic, B., Nessa,
 753 A., Wadge, S., Österdahl, M. F., Paterson, A. L., Rassl, D. M., Alcantara, R. E., O'Neill,
 754 L., Widaa, S.,...Martincorena, I. (2025). Somatic mutation and selection at population
 755 scale. *Nature*, 647(8089), 411-420. <https://doi.org/10.1038/s41586-025-09584-w>
- 756 Lek, M., Karczewski, K. J., Minikel, E. V., Samocha, K. E., Banks, E., Fennell, T., O'Donnell-
 757 Luria, A. H., Ware, J. S., Hill, A. J., Cummings, B. B., Tukiainen, T., Birnbaum, D. P.,
 758 Kosmicki, J. A., Duncan, L. E., Estrada, K., Zhao, F., Zou, J., Pierce-Hoffman, E.,
 759 Berghout, J.,...MacArthur, D. G. (2016). Analysis of protein-coding genetic variation in
 760 60,706 humans. *Nature*, 536(7616), 285-291. <https://doi.org/10.1038/nature19057>
- 761 Li, C., & Williams, S. M. (2013). Human Somatic Variation: It's Not Just for Cancer Anymore.
 762 *Current Genetic Medicine Reports*, 1(4), 212-218. <https://doi.org/10.1007/s40142-013-0029-z>
- 763
 764 Lodato, M. A., Rodin, R. E., Bohrsen, C. L., Coulter, M. E., Barton, A. R., Kwon, M., Sherman,
 765 M. A., Vitzthum, C. M., Luquette, L. J., Yandava, C. N., Yang, P., Chittenden, T. W.,
 766 Hatem, N. E., Ryu, S. C., Woodworth, M. B., Park, P. J., & Walsh, C. A. (2018). Aging
 767 and neurodegeneration are associated with increased mutations in single human neurons.
 768 *Science*, 359(6375), 555-559. <https://doi.org/10.1126/science.aao4426>

- 769 López-Otín, C., Blasco, M. A., Partridge, L., Serrano, M., & Kroemer, G. (2023). Hallmarks of
 770 aging: An expanding universe. *Cell*, *186*(2), 243-278.
 771 <https://doi.org/10.1016/j.cell.2022.11.001>
- 772 Martincorena, I., Raine, K. M., Gerstung, M., Dawson, K. J., Haase, K., Van Loo, P., Davies, H.,
 773 Stratton, M. R., & Campbell, P. J. (2017). Universal Patterns of Selection in Cancer and
 774 Somatic Tissues. *Cell*, *171*(5), 1029-1041.e1021.
 775 <https://doi.org/10.1016/j.cell.2017.09.042>
- 776 Maslov, A. Y., Makhortov, S., Sun, S., Heid, J., Dong, X., Lee, M., & Vijg, J. (2022). Single-
 777 molecule, quantitative detection of low-abundance somatic mutations by high-throughput
 778 sequencing. *Sci Adv*, *8*(14), eabm3259. <https://doi.org/10.1126/sciadv.abm3259>
- 779 McLaren, W., Gil, L., Hunt, S. E., Riat, H. S., Ritchie, G. R., Thormann, A., Flicek, P., &
 780 Cunningham, F. (2016). The Ensembl Variant Effect Predictor. *Genome Biol*, *17*(1), 122.
 781 <https://doi.org/10.1186/s13059-016-0974-4>
- 782 Melboucy-Belkhir, S., Pradère, P., Tadbiri, S., Habib, S., Bacrot, A., Brayer, S., Mari, B., Besnard,
 783 V., Mailleux, A., Guenther, A., Castier, Y., Mal, H., Crestani, B., & Plantier, L. (2014).
 784 Forkhead Box F1 represses cell growth and inhibits COL1 and ARPC2 expression in lung
 785 fibroblasts in vitro. *Am J Physiol Lung Cell Mol Physiol*, *307*(11), L838-847.
 786 <https://doi.org/10.1152/ajplung.00012.2014>
- 787 Miller, M. B., Huang, A. Y., Kim, J., Zhou, Z., Kirkham, S. L., Maury, E. A., Ziegenfuss, J. S.,
 788 Reed, H. C., Neil, J. E., Rento, L., Ryu, S. C., Ma, C. C., Luquette, L. J., Ames, H. M.,
 789 Oakley, D. H., Frosch, M. P., Hyman, B. T., Lodato, M. A., Lee, E. A., & Walsh, C. A.
 790 (2022). Somatic genomic changes in single Alzheimer's disease neurons. *Nature*,
 791 *604*(7907), 714-722. <https://doi.org/10.1038/s41586-022-04640-1>
- 792 Ren, P., Dong, X., & Vijg, J. (2022). Age-related somatic mutation burden in human tissues. *Front*
 793 *Aging*, *3*, 1018119. <https://doi.org/10.3389/fragi.2022.1018119>
- 794 Robinson, P. S., Coorens, T. H. H., Palles, C., Mitchell, E., Abascal, F., Olafsson, S., Lee, B. C.
 795 H., Lawson, A. R. J., Lee-Six, H., Moore, L., Sanders, M. A., Hewinson, J., Martin, L.,
 796 Pinna, C. M. A., Galavotti, S., Rahbari, R., Campbell, P. J., Martincorena, I., Tomlinson,
 797 I., & Stratton, M. R. (2021). Increased somatic mutation burdens in normal human cells
 798 due to defective DNA polymerases. *Nat Genet*, *53*(10), 1434-1442.
 799 <https://doi.org/10.1038/s41588-021-00930-y>
- 800 Saul, D., Kosinsky, R. L., Atkinson, E. J., Doolittle, M. L., Zhang, X., LeBrasseur, N. K., Pignolo,
 801 R. J., Robbins, P. D., Niedernhofer, L. J., Ikeno, Y., Jurk, D., Passos, J. F., Hickson, L. J.,
 802 Xue, A., Monroe, D. G., Tchkonja, T., Kirkland, J. L., Farr, J. N., & Khosla, S. (2022). A
 803 new gene set identifies senescent cells and predicts senescence-associated pathways across
 804 tissues. *Nature Communications*, *13*(1), 4827. [https://doi.org/10.1038/s41467-022-32552-](https://doi.org/10.1038/s41467-022-32552-1)
 805 [1](https://doi.org/10.1038/s41467-022-32552-1)
- 806 Schumacher, B., Pothof, J., Vijg, J., & Hoeijmakers, J. H. J. (2021). The central role of DNA
 807 damage in the ageing process. *Nature*, *592*(7856), 695-703.
 808 <https://doi.org/10.1038/s41586-021-03307-7>
- 809 Selby, C. P., Lindsey-Boltz, L. A., Li, W., & Sancar, A. (2023). Molecular Mechanisms of
 810 Transcription-Coupled Repair. *Annu Rev Biochem*, *92*, 115-144.
 811 <https://doi.org/10.1146/annurev-biochem-041522-034232>
- 812 Sender, R., & Milo, R. (2021). The distribution of cellular turnover in the human body. *Nat Med*,
 813 *27*(1), 45-48. <https://doi.org/10.1038/s41591-020-01182-9>

- 814 Szilard, L. (1959). ON THE NATURE OF THE AGING PROCESS. *Proc Natl Acad Sci U S A*,
815 45(1), 30-45. <https://doi.org/10.1073/pnas.45.1.30>
- 816 Vijg, J., & Dong, X. (2020). Pathogenic Mechanisms of Somatic Mutation and Genome Mosaicism
817 in Aging. *Cell*, 182(1), 12-23. <https://doi.org/10.1016/j.cell.2020.06.024>
- 818 Vorontsov, I. E., Khimulya, G., Lukianova, E. N., Nikolaeva, D. D., Eliseeva, I. A., Kulakovskiy,
819 I. V., & Makeev, V. J. (2016). Negative selection maintains transcription factor binding
820 motifs in human cancer. *BMC Genomics*, 17 Suppl 2(Suppl 2), 395.
821 <https://doi.org/10.1186/s12864-016-2728-9>
- 822 Vorontsov, I. E., Khimulya, G., Lukianova, E. N., Nikolaeva, D. D., Eliseeva, I. A., Kulakovskiy,
823 I. V., & Makeev, V. J. (2016). Negative selection maintains transcription factor binding
824 motifs in human cancer. *BMC Genomics*, 17(2), 395. <https://doi.org/10.1186/s12864-016-2728-9>
- 825
- 826 Wang, T., Birsoy, K., Hughes, N. W., Krupczak, K. M., Post, Y., Wei, J. J., Lander, E. S., &
827 Sabatini, D. M. (2015). Identification and characterization of essential genes in the human
828 genome. *Science*, 350(6264), 1096-1101. <https://doi.org/10.1126/science.aac7041>
- 829 Watson, D. E., Cunningham, M. L., & Tindall, K. R. (1998). Spontaneous and ENU-induced
830 mutation spectra at the cII locus in Big Blue Rat2 embryonic fibroblasts. *Mutagenesis*,
831 13(5), 487-497. <https://doi.org/10.1093/mutage/13.5.487>
- 832 Wu, T., Hu, E., Xu, S., Chen, M., Guo, P., Dai, Z., Feng, T., Zhou, L., Tang, W., Zhan, L., Fu, X.,
833 Liu, S., Bo, X., & Yu, G. (2021). clusterProfiler 4.0: A universal enrichment tool for
834 interpreting omics data. *Innovation (Camb)*, 2(3), 100141.
835 <https://doi.org/10.1016/j.xinn.2021.100141>
- 836 Yang, J.-L., Lee, P.-C., Lin, S.-R., & Lin, J.-G. (1994). Comparison of mutation spectra induced
837 by N-ethyl-N-nitrosourea in the hprt gene of Mer⁺ and Mer⁻ diploid human fibroblasts.
838 *Carcinogenesis*, 15(5), 939-945. <https://doi.org/10.1093/carcin/15.5.939>
- 839 Yang, L. (2020). A Practical Guide for Structural Variation Detection in the Human Genome. *Curr*
840 *Protoc Hum Genet*, 107(1), e103. <https://doi.org/10.1002/cphg.103>
- 841 Yurchenko, A. A., Rajabi, F., Braz-Petta, T., Fassih, H., Lehmann, A., Nishigori, C., Wang, J.,
842 Padioleau, I., Gunbin, K., Panunzi, L., Morice-Picard, F., Laplante, P., Robert, C.,
843 Kannouche, P. L., Menck, C. F. M., Sarasin, A., & Nikolaev, S. I. (2023). Genomic
844 mutation landscape of skin cancers from DNA repair-deficient xeroderma pigmentosum
845 patients. *Nat Commun*, 14(1), 2561. <https://doi.org/10.1038/s41467-023-38311-0>
- 846 Zapata, L., Pich, O., Serrano, L., Kondrashov, F. A., Ossowski, S., & Schaefer, M. H. (2018).
847 Negative selection in tumor genome evolution acts on essential cellular functions and the
848 immunopeptidome. *Genome Biol*, 19(1), 67. <https://doi.org/10.1186/s13059-018-1434-0>
- 849 Zhang, L., Dong, X., Lee, M., Maslov, A. Y., Wang, T., & Vijg, J. (2019). Single-cell whole-
850 genome sequencing reveals the functional landscape of somatic mutations in B
851 lymphocytes across the human lifespan. *Proc Natl Acad Sci U S A*, 116(18), 9014-9019.
852 <https://doi.org/10.1073/pnas.1902510116>
- 853 Zhang, L., Lee, M., Hao, X., Ma, X., Xia, C., Zhao, Y., Ehlert, J., Chi, Z., Jin, B., Cutler, R.,
854 Maslov, A. Y., Barabási, A. L., Hoeijmakers, J. H. J., Edelman, W., Vijg, J., & Dong, X.
855 (2025). Divergent accumulation patterns of SNVs and INDELS reveal negative selection
856 in noncancerous cells. *Innovation (Camb)*, 6(10), 101008.
857 <https://doi.org/10.1016/j.xinn.2025.101008>

- 858 Zhang, L., Lee, M., Maslov, A. Y., Montagna, C., Vijg, J., & Dong, X. (2024). Analyzing somatic
 859 mutations by single-cell whole-genome sequencing. *Nat Protoc*, *19*(2), 487-516.
 860 <https://doi.org/10.1038/s41596-023-00914-8>
- 861 Zhang, X., Theotokis, P. I., Li, N., Wright, C. F., Samocha, K. E., Whiffin, N., & Ware, J. S.
 862 (2024). Genetic constraint at single amino acid resolution in protein domains improves
 863 missense variant prioritisation and gene discovery. *Genome Med*, *16*(1), 88.
 864 <https://doi.org/10.1186/s13073-024-01358-9>
- 865 Zoppi, N., Gardella, R., De Paepe, A., Barlati, S., & Colombi, M. (2004). Human Fibroblasts with
 866 Mutations in COL5A1 and COL3A1 Genes Do Not Organize Collagens and Fibronectin
 867 in the Extracellular Matrix, Down-regulate $\alpha 2\beta 1$ Integrin, and Recruit $\alpha v\beta 3$ Instead of
 868 $\alpha 5\beta 1$ Integrin*. *Journal of Biological Chemistry*, *279*(18), 18157-18168.
 869 <https://doi.org/https://doi.org/10.1074/jbc.M312609200>

870

871 **Figure 1. Repeated Mutagen Treatment of Normal Cells Only Slightly Affects Cell Growth** 872 **and Death**

873 **(A)** Schematic depiction of the experimental design. Human fetal lung fibroblasts (IMR-90) were
 874 treated with a sublethal dose of 50 $\mu\text{g}/\text{mL}$ of *N*-ethyl-*N*-nitrosourea (ENU), which was followed
 875 by a recovery period of 7 days. At 3-day and 7-day time points, 1 million cells were passaged, and
 876 a sample was taken for further analysis. This process was repeated for 9 cycles (Methods). **(B)**
 877 Cumulative population doublings (PDs) of the control (CTRL) and ENU-treated (ENU) groups
 878 throughout the experiment. Statistics: $n=3$; Data represent mean \pm S.D.; P-value estimated using a
 879 two-way repeated measures ANOVA. **(C)** Cumulative PDs of the CTRL and ENU groups assessed
 880 at the 3-day time point. Statistics: Same as in (B). **(D)** Same as in (C), but at 7-day time point.
 881 Statistics: Same as in (B). **(E)** Percentage of cells detected in early apoptosis at the 3-day time
 882 point at each cycle (left) or with all cycles pooled (right). Statistics: $n=2$; Data represent mean \pm
 883 S.D.; P-value of data from each cycle was estimated using a two-way repeated measures ANOVA.
 884 P-value of data from all cycles pooled was estimated using a paired t-test; **: $p \leq 0.005$; ****:
 885 $p \leq 0.0001$. **(F)** Same as in (E), but for late apoptosis. Statistics: Same as in (E). **(G)** Same as in
 886 (E), but for 7-day time point. Statistics: Same as in (E). **(H)** Same as in F, but for 7-day time point.
 887 Statistics: Same as in (E).

888

889 **Figure 2. RNA-sequencing Reveals Increased Cell Cycle Inhibition and Apoptosis in ENU-** 890 **treated Cells**

891 **(A)** Differential gene expression (DGE) of control cycle 9 vs. control cycle 1 groups. Red points
 892 indicate Log_2 fold change > 2 and adjusted *p*-value < 0.05 . See Table S1 for sample details and
 893 Table S2 for DGE results. Statistics: $n=3$; DGE *p*-values estimated using a Wald test followed by
 894 Benjamini-Hochberg correction for multiple hypothesis testing. **(B-D)** Gene set enrichment
 895 analysis (GSEA) of DGE results from the control cycle 9 vs. control cycle 1 comparison for KEGG
 896 cell cycle (B), KEGG apoptosis (C), and SenMayo senescence (D) gene sets. NES: normalized
 897 enrichment score. Statistics: $n=3$; *p*-values estimated using a permutation test. **(E)** Same as in (A),
 898 but for the ENU-treated cycle 9 vs. control cycle 9 comparison. See Supplemental Note 1 for batch

899 correction note. Statistics: Same as in (A). **(F-H)** Same as in (B-D), but for the ENU-treated cycle
 900 9 vs. control cycle 9 comparison. Statistics: Same as in (B-D).

901

902 **Figure 3. Repeated Mutagen Treatment of Normal Cells Results in Extremely High Mutation**
 903 **Burden.**

904 **(A)** Estimated load of somatic single nucleotide variants (SNVs) per cell after correction for
 905 genome coverage, sensitivity, and technical artifacts (Methods). See Figure S2H for observed SNV
 906 load per cell and Figures S2K-L for ENU-specific mutation burden. See Table S3 for sample
 907 details and Table S10 for mutation catalogue. Statistics: n=3 cells for cycles 1, 3, and 6; n=6 cells
 908 for cycle 9; data represent mean \pm S.D. p-values of comparisons between cycles within the same
 909 condition were estimated with a two-sided linear model with estimated marginal means, Tukey's
 910 HSD. P-values of comparisons of control and ENU-treated groups were estimated using a two-
 911 way ANOVA and shown at very top of plot. p-value legend: ns: $p > 0.05$, ****: $p \leq 0.0001$. **(B)**
 912 Fit of linear model to estimated SNV load per cell shown in (A) vs. the treatment cycle cells were
 913 collected at. Note that cells from control cycle 1 are included and control cycle 9 is omitted. Shown
 914 at top is the fit of a linear equation. Statistics: data represent mean \pm S.D.; Pearson R correlation
 915 coefficient and p-value of the correlation was estimated using a t-test.

916

917 **Figure 4. Mutational Spectra and Signatures Resulting from Repeated Mutagen**
 918 **Treatment**

919 **(A)** Mutational spectra of the relative contribution of SNV types per cell grouped by condition.
 920 See Figure S3A for groups stratified by cycle. Note that cells from the control cycle 1 group were
 921 excluded from the mutational spectra and signatures analyses due to abnormally high C>T
 922 frequency suggestive of a sample preparation artifact (Figures S3A-D; Methods). Statistics: n=3
 923 cells for cycles 3 and 6; n=6 cells for cycle 9; Data represent mean \pm S.D. **(B)** Relative contribution
 924 of de novo mutational signatures SBSA and SBSB to individual cells (shown as rows) grouped by
 925 condition and cycle. See Figure S3F for determination of number of signatures to extract and
 926 Figure S3G for absolute contribution. **(C)** Comparison of de novo mutational signature SBSA to
 927 the ENU signature from the SIGNAL database. Statistics: cosine similarity 0.858, residual sum of
 928 squares= 7.83×10^{-3} . **(D)** Comparison of de novo mutational signature SBSB to most similar fitted
 929 signature SBS5 from the COSMIC database. Statistics: cosine similarity=0.955, residual sum of
 930 squares= 1.79×10^{-3} .

931

932 **Figure 5. Negative Selection Against Damaging Coding and Non-Coding Variants in ENU-**
 933 **treated Cells.**

934 **(A)** Observed frequency of synonymous and nonsynonymous variants per cell. See Table S3 for
 935 sample details and Table S5 for variant annotation. Statistics: n=3 cells for cycles 1, 3, and 6; n=6
 936 cells for cycle 9; data represent mean \pm S.D; P-values of comparisons of control and ENU-treated
 937 groups were estimated using a two-way ANOVA and shown at very top of plot; p-value legend:
 938 ***: $p \leq 0.001$. **(B)** $\text{Log}_2(\text{O/E})$ ratio of synonymous and nonsynonymous variants for ENU-
 939 treated cells. See figures S4C, S4D, and S4E for results of without normalization, results of control

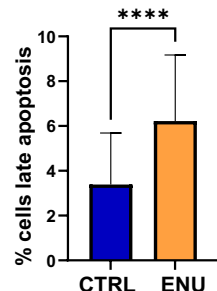
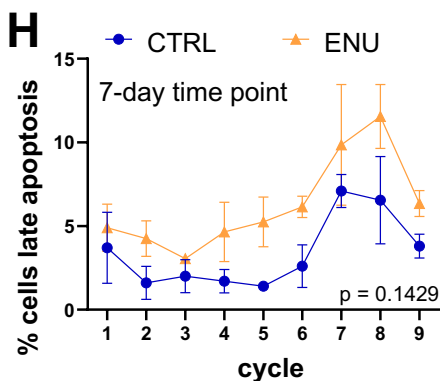
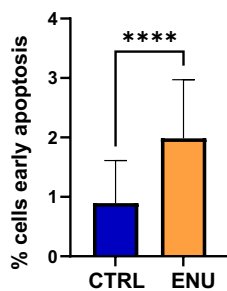
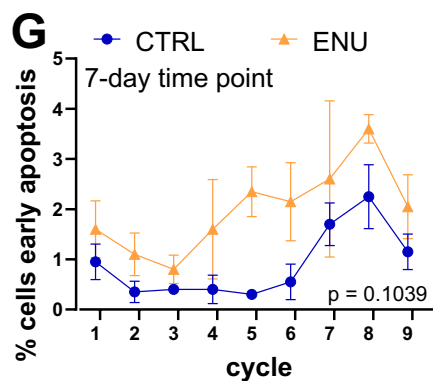
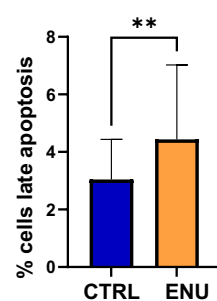
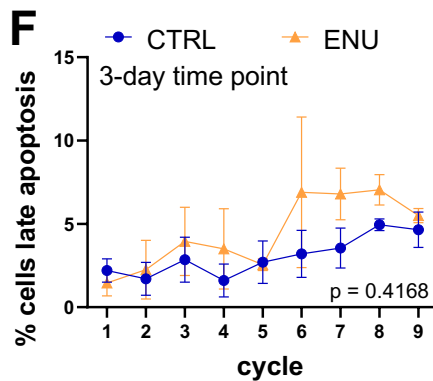
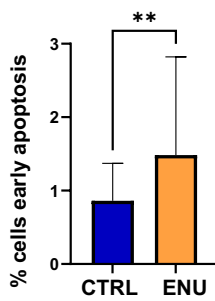
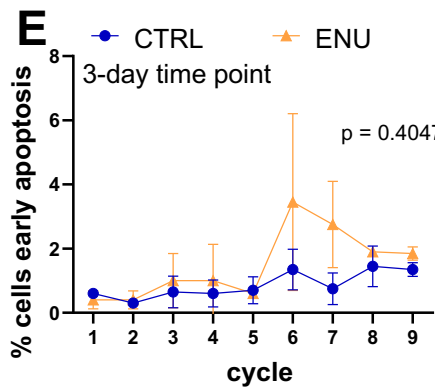
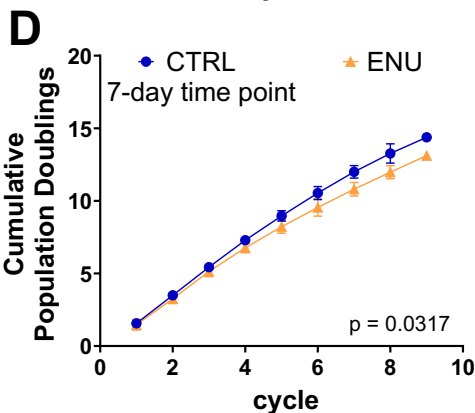
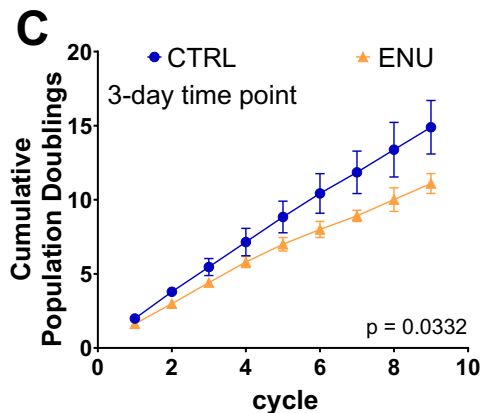
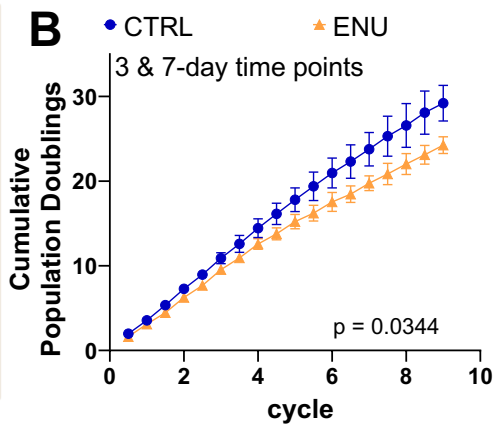
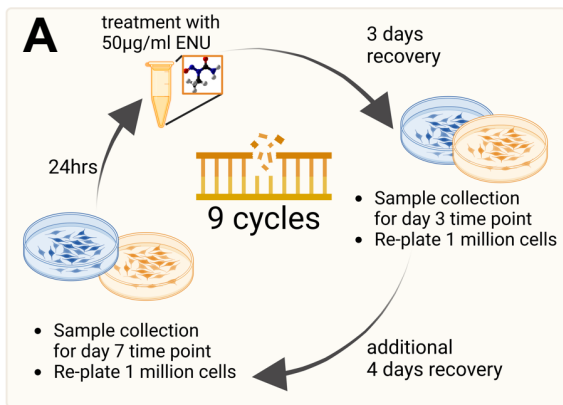
940 cells, and results per cell, respectively. See Table S5 for detailed results. Statistics: n=3 cells for
 941 cycles 3 and 6; n=6 cells for cycle 9; data represent mean \pm S.D; P-values to assess if the $\text{Log}_2(\text{O}/\text{E})$
 942 ratio significantly deviates from 0 were estimated using a permutation test (n=10,000) but were
 943 only estimated if the frequency of observed variants allowed for sufficiency statistical power
 944 (Figure S4B; Table S5); p-value legend: **: $p \leq 0.01$. **(C)** Estimated percent of nonsynonymous
 945 variants out of total coding variants that are reduced in the ENU-treated group due to negative
 946 selection. Statistics: n=3 cells for cycles 3 and 6; n=6 cells for cycle 9; data represent mean \pm S.D.
 947 **(D)** $\text{Log}_2\text{O}/\text{E}$ of the N/S ratio for ENU-treated cells. See figures S4I and S4J for results of control
 948 cells and per cell, respectively. Statistics: Same as in (B). p-value legend: *: $p \leq 0.05$; **: $p \leq$
 949 0.01 . **(E)** Observed frequency of various types of nonsynonymous variants per cell. Statistics:
 950 Same as in (A). p-value legend: ****: $p \leq 0.0001$. **(F)** $\text{Log}_2\text{O}/\text{E}$ of various types of
 951 nonsynonymous variants for ENU-treated cells. See figures S4F, S4G, and S4H for results of
 952 without normalization, results of control cells, and results per cell, respectively. Statistics: Same
 953 as in (B). p-value legend: *: $p \leq 0.05$; **: $p \leq 0.01$. **(G)** Observed frequency of various types of
 954 non-coding variants per cell. Statistics: Same as in (A). ****: $p \leq 0.0001$. **(H)** $\text{Log}_2\text{O}/\text{E}$ of various
 955 types of non-coding variants for ENU-treated cells. See figures S4K, S4L, and S4M for results of
 956 control cells, without normalization, and per cell, respectively. See Supplemental Note 5 for
 957 discussion on results of the intergenic region. Statistics: Same as in (B). ****: $p \leq 0.0001$; **: p
 958 ≤ 0.01 . **(I)** Estimated percent reduction of variants in 3' UTR and splice region variants that are
 959 protected against in the ENU-treated group due to negative selection. Statistics: Same as in (C).

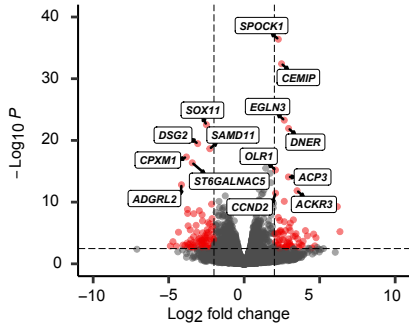
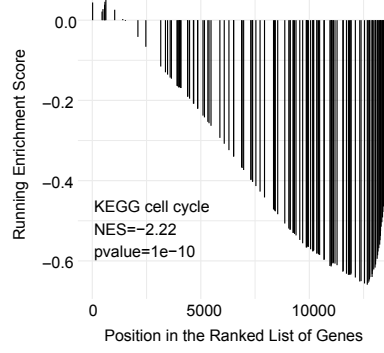
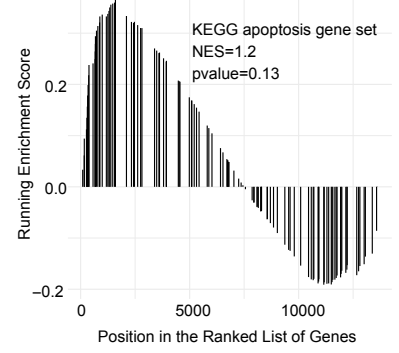
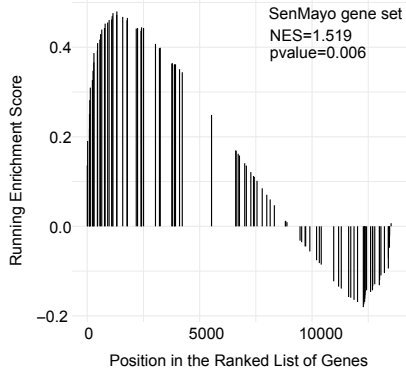
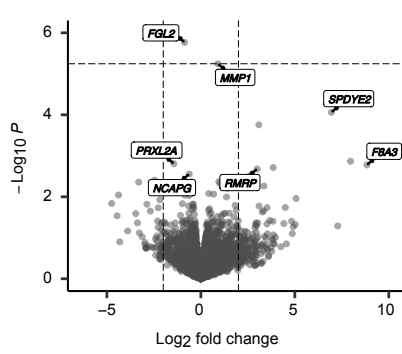
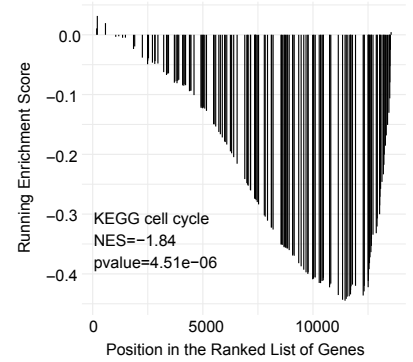
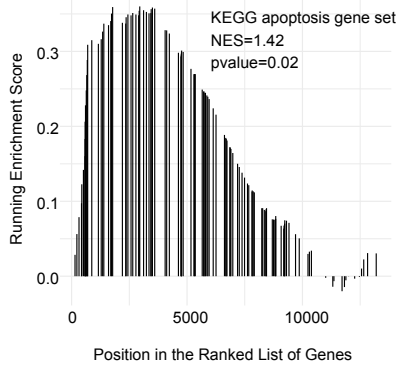
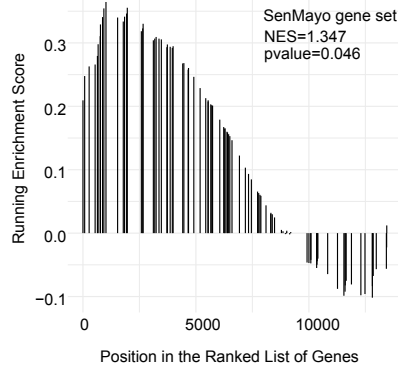
960

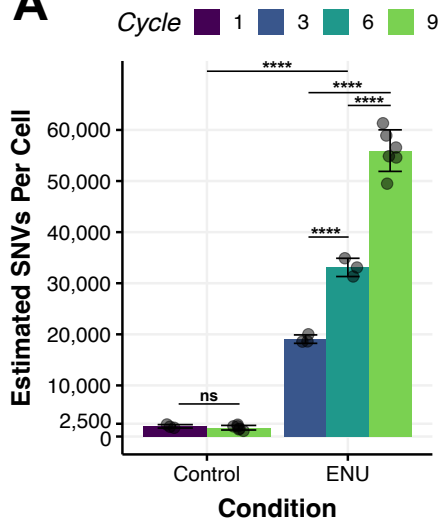
961 **Figure 6. Pathways Supporting Fibroblast Growth and Cell Identity are Protected from** 962 **Mutation Accumulation**

963 **(A)** Observed frequency of SNVs per cell in functional regions associated with genes that were
 964 determined to be non-expressed or expressed from bulk mRNA transcriptome data (Table S7).
 965 Statistics: n=3 cells for cycles 1, 3, and 6; n=6 cells for cycle 9; data represent mean \pm S.D; P-
 966 values of comparisons of control and ENU-treated groups were estimated using a two-way
 967 ANOVA; p-value legend: ****: $p \leq 0.0001$. **(B)** $\text{Log}_2(\text{O}/\text{E})$ ratio of SNVs in functional regions
 968 associated with genes determined to be non-expressed or expressed in bulk mRNA transcriptome
 969 data (Table S7). See figures S6A, S6B, and S6C for results without normalization, results of
 970 control cells, and results per cell, respectively. Statistics: n=3 cells for cycles 3 and 6; n=6 cells
 971 for cycle 9; data represent mean \pm S.D; P-values to assess if the $\text{Log}_2(\text{O}/\text{E})$ ratio significantly
 972 deviates from 0 were estimated using a permutation test (n=10,000) but were only estimated if the
 973 frequency of observed variants allowed for sufficiency statistical power (Figure S4B; Table S5);
 974 p-value legend: ****: $p \leq 0.0001$. **(C)** Observed frequency of SNVs per cell in functional regions
 975 associated with nonessential or essential gene sets expressed in bulk mRNA transcriptome data
 976 (Table S7). Statistics: Same in in (A); p-value legend: ***: $p \leq 0.001$. **(D)** $\text{Log}_2(\text{O}/\text{E})$ ratio of
 977 SNVs in functional regions associated with nonessential or essential gene sets expressed in bulk
 978 mRNA transcriptome data (Table S7). See figures S6D, S6E, and S6F for results without
 979 normalization, results of control cells, and results per cell, respectively. Statistics: Same as in (B);
 980 p-value legend: ****: $p \leq 0.0001$. **(E)** Observed frequency of SNVs per cell in functional regions
 981 associated with haplosufficient or haploinsufficient gene sets expressed in bulk mRNA
 982 transcriptome data (Table S7). Statistics: Same in in (A); p-value legend: ***: $p \leq 0.001$. **(F)**
 983 $\text{Log}_2(\text{O}/\text{E})$ ratio of SNVs in functional regions associated with haplosufficient or haploinsufficient
 984 gene sets expressed in bulk mRNA transcriptome data (Table S7). See figures S6G, S6H, and S6I

985 results without normalization, results of control cells, and results per cell, respectively. Statistics:
986 Same as in (B); p-value legend: **: $p \leq 0.01$; ****: $p \leq 0.0001$. **(G)** Mean observed frequency
987 of SNVs within the functional regions of pathways expressed in bulk mRNA transcriptome data
988 (Table S8). See figure S6J for top 20 most enriched pathways. Statistics: $n=399$ pathways; P-values
989 of comparisons of control and ENU-treated groups were estimated using a two-way ANOVA; p-
990 value legend: ****: $p \leq 0.0001$. **(H)** Mean $\text{Log}_2(\text{O/E})$ ratio of SNVs within the functional regions
991 of pathways expressed in bulk mRNA transcriptome data (Table S8). See figures S6K and S6L for
992 results of control cell and results without normalization, respectively. See table S9 for detailed
993 results. Statistics: $n=399$ pathways; Otherwise, same as in (B), with the addition of Benjamini-
994 Hochberg correction for multiple hypothesis testing of pathways. **(I-K)** Mean $\text{Log}_2(\text{O/E})$ ratio (I),
995 estimated percent reduction in SNVs (J), and estimated reduction in SNV frequency (K) within
996 functional regions of pathways found to be under significant negative selection ($p \leq 0.05$) at cycle
997 9 in the ENU-treated group. Statistics: Same as in (B); p-value legend: *: $p \leq 0.05$; **: $p \leq 0.01$;
998 ***: $p \leq 0.001$.



A Control Cycle 9 vs Control Cycle 1**B** Control Cycle 9 vs Control Cycle 1**C** Control Cycle 9 vs Control Cycle 1**D** Control Cycle 9 vs Control Cycle 1**E** ENU Cycle 9 vs Control Cycle 9**F** ENU Cycle 9 vs Control Cycle 9**G** ENU Cycle 9 vs Control Cycle 9**H** ENU Cycle 9 vs Control Cycle 9

A**B**