

1 **Nucleosome positioning stability is a significant modulator of germline**
2 **mutation rate variation across the human genome**

3

4 Cai Li^{1,#} and Nicholas M. Luscombe^{1,2,3}

5

6 ¹ The Francis Crick Institute, 1 Midland Road, London, NW1 1AT, UK

7 ² Okinawa Institute of Science & Technology Graduate University, Okinawa,
8 904-0495, Japan

9 ³ UCL Genetics Institute, University College London, Gower Street, London,
10 WC1E 6BT, UK

11 # Correspondence: cai.li@crick.ac.uk

12 **Summary**

13 Nucleosome organization is suggested to affect local mutation rates in a genome.
14 However, the lack of *de novo* mutation and high-resolution nucleosome data have
15 limited investigation. Further, analyses using indirect mutation rate measurements
16 have yielded contradictory and potentially confounded results. Combining >300,000
17 human *de novo* mutations with high-resolution nucleosome maps, we reveal
18 substantially elevated mutation rates around translationally stable ('strong')
19 nucleosomes. Translational stability is an under-appreciated nucleosomal property,
20 with greater impact than better-known factors like occupancy and histone
21 modifications. We show that the mutational mechanisms affected by strong
22 nucleosomes are low-fidelity replication, insufficient mismatch repair and increased
23 double-strand breaks. Strong nucleosomes preferentially locate within young
24 SINE/LINE transposons; subject to increased mutation rates, transposons are then
25 more rapidly inactivated. Depletion of strong nucleosomes in older transposons
26 suggests frequent re-positioning during evolution, thus resolving a debate about
27 selective pressure on nucleosome-positioning. The findings have important
28 implications for human genetics and genome evolution.

29

30 **1 Introduction**

31 Germline *de novo* mutations, which can be passed to offspring, are the primary
32 source of genetic variation in multicellular organisms, contributing substantially to
33 biological diversity and evolution. *De novo* mutations are also thought to play
34 significant roles in early-onset genetic disorders such as intellectual disability, autism
35 spectrum disorder, and developmental diseases (Veltman and Brunner 2012; Acuna-
36 Hidalgo et al. 2016). Thus, investigating the patterns and genesis of *de novo*
37 mutations in the germline is important for understanding genome evolution and
38 human diseases.

39 Germline and somatic mutation rates vary across the human genome at diverse
40 scales ranging from nucleotide to chromosomal resolution (Hodgkinson and Eyre-
41 Walker 2011; Segurel et al. 2014). Studies revealed factors linked to local mutation
42 rate variation, including sequence context (Michaelson et al. 2012), replication timing
43 (Stamatoyannopoulos et al. 2009), recombination rate (Francioli et al. 2015), DNA
44 accessibility (Sabarinathan et al. 2016) and histone modifications (Michaelson et al.
45 2012; Schuster-Bockler and Lehner 2012). However, genomic features identified so
46 far explain less than 40% of the observed germline mutation rate variation (at 100Kb
47 to 1Mb resolution) (Terekhanova et al. 2017; Smith et al. 2018). Therefore, important
48 factors remain to be found. Moreover, due to the limited availability of *de novo*
49 mutation datasets, studies focused on coarse-grained mutation rate variation
50 (typically $\geq 1\text{kb}$ windows for germline data), or used within-species polymorphisms
51 and inter-species divergence whose observations are potentially confounded by
52 natural selection and other evolutionary processes.

53 Moreover, the underlying mutational processes causing the observed mutation rate
54 variation are poorly understood, though recent studies have highlighted the
55 contributions of error-prone replicative processes (Harris and Nielsen 2014; Lujan et
56 al. 2014; Reijns et al. 2015; Seplyarskiy et al. 2017; Seplyarskiy et al. 2018) and
57 differential DNA repair efficiencies (Supek and Lehner 2015; Perera et al. 2016;
58 Sabarinathan et al. 2016; Frigola et al. 2017). Despite these advances, it remains a
59 challenge to understand the molecular mechanisms associated with mutation rate
60 variation, particularly in the germline.

61 Here, we focus on the role of nucleosomes in modulating germline mutation rates.
62 Chromatin is considered important because structural constraints could affect the
63 mutability of genomic sequences (Makova and Hardison 2015). Nucleosome
64 organization (including positioning and occupancy) has been reported as a significant

65 factor in humans and other eukaryotes (Sasaki et al. 2009; Tolstorukov et al. 2011;
66 Chen et al. 2012; Michaelson et al. 2012; Lujan et al. 2014; Pich et al. 2018). Studies
67 in different lineages (Sasaki et al. 2009; Tolstorukov et al. 2011; Lujan et al. 2014)
68 reported increased substitution rates around the centers of nucleosomal sequences
69 and increased insertion/deletion rates in linker DNA. However, there are also
70 disagreements between published studies. For example, Michaelson et al. (2012)
71 suggested that high nucleosome occupancy tends to suppress *de novo* mutations,
72 but Smith et al. (2018) found that a comparative analysis using datasets from
73 different studies resulted in opposing conclusions. Due to few available *de novo*
74 mutations for humans, analysis of many studies was based on variant data from
75 within-species polymorphisms or inter-species divergence, which can be affected by
76 natural selection and non-adaptive processes such as GC-biased gene conversion.
77 Furthermore, because of the limitation of available nucleosome maps, some previous
78 studies treated all annotated nucleosomes equally, ignoring the diverse contexts in
79 which they form. Therefore, combined with the scarcity of *de novo* mutation datasets,
80 the effects of nucleosome organization on germline mutation rate variation,
81 particularly at high resolution remain to be elucidated. Here we take advantage of the
82 rapid increase in the number of *de novo* mutation datasets and better understanding
83 of nucleosome organization in the human genome to perform a systematic analysis
84 of this topic.

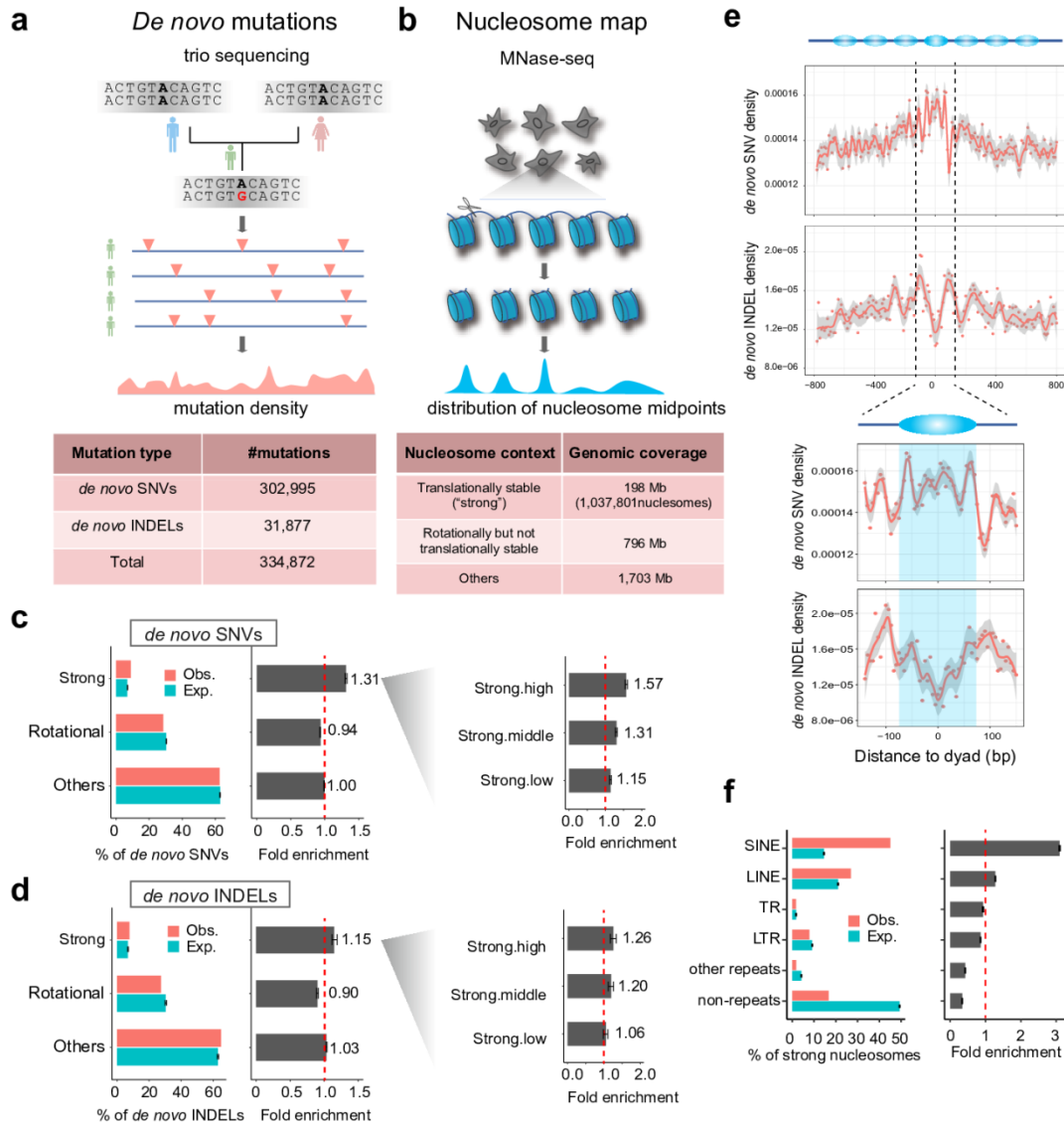
85

86 **2 Results**

87 **2.1 Datasets used for analysis**

88 We used >300,000 human *de novo* single-nucleotide variants (SNVs) and >30,000
89 short insertions/deletions (INDELs), having removed genomic regions that could
90 confound downstream analysis (**Fig. 1a, Supplementary Fig. 1a**; see Methods).
91 Most data come from three large-scale trio sequencing projects which contribute
92 about 100,000 mutations each (Jonsson et al. 2017; Turner et al. 2017a; Yuen et al.
93 2017). We also examined extremely rare variants (allele frequency ≤ 0.0001) from
94 the gnomAD database (Lek et al. 2016) which are approximated to *de novo*
95 mutations because they are thought to undergo limited selection and non-adaptive
96 evolutionary processes (Carlson et al. 2018).

97 Nucleosome positioning on the genome is described by the translational setting,
98 which defines the location of the nucleosomal midpoint (also called 'dyad') and the
99 rotational setting, which defines the orientation of the DNA helix on the histone
100 surface (Gaffney et al. 2012). Using MNase-seq measurements, Gaffney et al. (2012)
101 identified ~1 million 'strong' nucleosomes that adopt highly stable translational
102 positioning across seven lymphoblastoid cell lines. Rotationally stable nucleosomes
103 were previously identified from DNase-seq measurements across 43 cell types
104 (Winter et al. 2013), covering 892Mb of the genome. There is a ~50Mb overlap
105 between regions bound by strong nucleosomes and rotationally stable nucleosomes.
106 Using these data, we classified the genome into three groups of regions (**Fig. 1b**; sex
107 chromosomes excluded): i) those containing translationally stable, 'strong',
108 nucleosomes (198Mb); ii) those with rotationally but not translationally stable
109 nucleosomes (796Mb); and iii) all other non-N base genomic regions (1,703Mb).
110 West et al. (2014) reported that with the exception of a few specific loci such as
111 transcription start sites, overall nucleosome positioning varies little between cell types.
112 None of the nucleosomal datasets were produced using germ cells, therefore as a
113 precaution we excluded nucleosomes that differ in positioning between cell types
114 (~23Mb; see Methods).



115

116 **Fig. 1 *De novo* mutations are enriched in strong nucleosomes.** (a) Summary of
 117 germline *de novo* mutation data included in study. (b) Summary of nucleosome
 118 positioning data analysed in study. (c, d) Observed versus expected occurrence and fold
 119 enrichments of *de novo* (c) SNVs and (d) INDELS in the three different nucleosome
 120 contexts. Right-hand panel subdivides strong nucleosomes according to high, medium
 121 and low translational stabilities. Error bars depict 95% confidence intervals. (e) Top
 122 panels, meta-profiles of *de novo* SNV and INDEL densities relative to position of strong
 123 nucleosome dyads. Bottom panel, same meta-profiles zoomed into the middle
 124 nucleosome. (f) Fold enrichment of strong nucleosomes in different repeat elements:
 125 SINE (Short Interspersed Nuclear Element), LINE (Long Interspersed Nuclear Element),
 126 TR (Tandem Repeat) and LTR (Long Terminal Repeat).

127

2.2 *De novo* SNVs and INDELS are enriched in strong nucleosomes

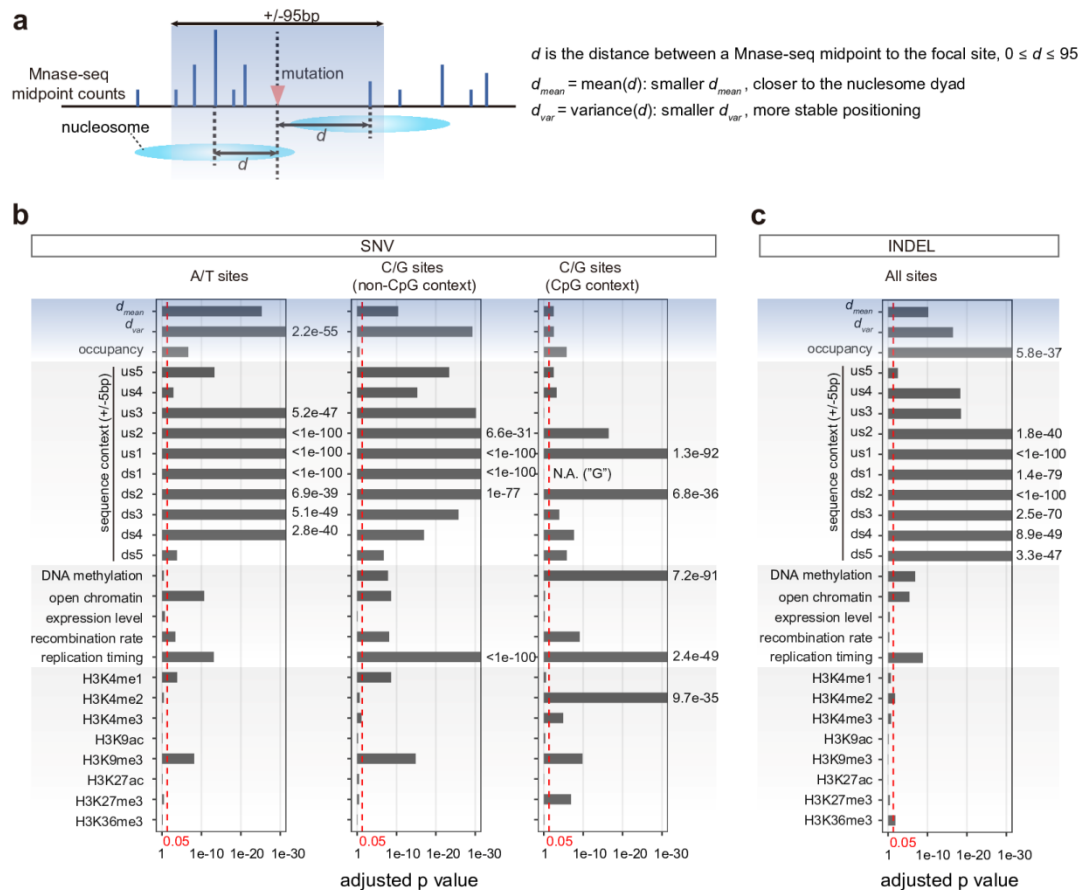
128 Genomic regions containing strong nucleosomes have ~30% more *de novo* SNVs
129 (**Fig. 1c**) and ~15% more *de novo* INDELS (**Fig. 1d**) than expected. Similar increases
130 are also apparent for extremely rare variants (**Supplementary Fig. 1b,c**), though
131 effect sizes are smaller than for *de novo* mutations, probably due to the fact that
132 highly mutable sites are under-represented among extremely rare variants (Harpak et
133 al. 2016). Restricting the analysis to strong nucleosomes, we found that those with
134 higher translational stability scores also exhibit higher mutation rates (**Fig. 1c,d**;
135 scores from Gaffney et al., 2012). These results suggest that translational stability is
136 associated with local variation in mutation rates across the genome, a previously
137 unappreciated aspect. Regions containing rotationally stable nucleosomes, in
138 contrast, are slightly depleted of both mutation types; we didn't perform further
139 analysis on this, as effect of rotational positioning has been comprehensively
140 discussed by Pich et al. (2018). A more detailed view with meta-profiles clearly
141 depicts increased SNV and reduced INDEL densities around dyad regions of strong
142 nucleosomes compared with flanking linker regions (**Fig. 1e**), in line with
143 observations made using polymorphism data (Tolstorukov et al. 2011).

144 Interestingly, ~80% of strong nucleosomes overlap with repeats (**Fig. 1f**,
145 **Supplementary Fig. 1d**), especially SINE/Alu (~44%) and LINE/L1 elements (~26%).
146 Genetic variations in repeats are traditionally hard to detect because of poor
147 mappability and so analyses have tended to be cautious in calling variants, resulting
148 in many false negatives (though, few false positives; Lee and Schatz (2012)).
149 Therefore, the above observations probably underestimate the true enrichment of *de*
150 *novo* mutations in strong nucleosomes. We subdivided strong nucleosomes into
151 three groups: i) Alu-associated, ii) L1-associated and iii) others. Alu-associated
152 nucleosomes display increased SNV rates around the dyads, as seen in the
153 metaprofiles for all strong nucleosomes (**Supplementary Fig. 1e**), whereas non-Alu
154 nucleosomes show increased SNV rates ~60bp away from the dyads, close to the
155 nucleosome edges. Such differences may be due to the different local sequence
156 composition (discussed in next section). In contrast, the patterns of INDEL densities
157 are relatively similar among different groups (**Supplementary Fig. 1e**).

158 **2.3 Controlling for potential confounding factors**

159 Many factors are associated with mutation rate variation. One of the most important
160 is local sequence context - for example, CpG sites are known to be highly mutable
161 and CpG density profiles correlate well with mutation rate profiles in strong
162 nucleosomes (**Supplementary Fig. 1e**). Functional factors like DNA methylation,

163 histone modification, chromatin accessibility, replication timing and recombination
 164 rate are also relevant. Therefore, to systematically assess the contribution of
 165 nucleosomes to mutation rate variation, we used a logistic regression framework to
 166 control for potential confounding factors (**Fig. 2**).



167

168 **Fig. 2 Controlling for potential confounding factors in evaluating contribution of**
 169 **nucleosome organization to mutation rate variation.** (a) Schematic diagram
 170 describing two nucleosome positioning-related variables (d_{mean} and d_{var}) relative to a
 171 given genomic position. Lower d_{var} corresponds to higher translational stability. (b, c)
 172 Independent statistical significance of potential contributing factors to mutation rate
 173 variation, having controlled for other factors; (b) for SNVs and (c) INDELS. Tests for
 174 SNVs were performed separately at A/T and C/G sites (non-CpG and CpG contexts
 175 respectively). Vertical red lines indicate the threshold for statistical significance (0.05).
 176 'us', upstream; 'ds', downstream.

177 We defined three variables to quantify nucleosomal properties relative to a specific
 178 nucleotide position in the genome. Two relate to translational positioning: d_{mean} , the
 179 mean distance between the focal position and the midpoints of mapped MNase-seq
 180 fragments (maximum distance of 95 bp) and d_{var} , the variance of these distances (**Fig.**
 181 **2a**). A smaller d_{mean} means that a nucleotide position is closer to nucleosome dyads

182 and a smaller d_{var} indicates that the nucleosomes around it are more translationally
183 stable. As the relationship between between d_{mean} and SNV rates is non-linear, we
184 defined d_{mean} a categorical variable binned into five intervals (Methods; **Fig. 1e**,
185 **Supplementary Fig. 1e**). The third variable is nucleosome occupancy calculated as
186 a normalised per-base MNase-seq fragment coverage (see Methods). Other factors
187 considered are local nucleotide sequences (± 5 bp of the focal site) and functional
188 genomic measurements in human germ cells or other cell types if no available germ-
189 cell data (see Methods). d_{var} has a relatively weak but statistically significant
190 correlation with many of these factors, suggesting non-independence
191 (**Supplementary Fig. 2**).

192 To assess the contribution of each factor to local mutation rates, we compared a full
193 logistic regression model encompassing all variables against reduced models
194 missing individual variables; the reported p values indicate how significant a factor is
195 associated with mutation rate variation, having controlled for other factors (**Fig. 2b,c**;
196 Methods). For SNVs, we tested A/T (comprising A>C, A>G and A>T mutations), CpG
197 and non-CpG C/G sites separately (both C>A, C>G and C>T; **Fig. 2b**), whereas they
198 were pooled for INDELS.

199 Our statistical framework recapitulates reported observations (**Fig. 2b,c**,
200 **Supplementary Fig. 3**). In agreement with previous studies (Carlson et al. 2018),
201 local sequence context is the biggest contributor to local mutation rate variation (**Fig.**
202 **2b,c**), with effect sizes generally declining with increasing distance from the surveyed
203 site. DNA methylation and H3K9me3 are two common epigenetic marks associated
204 with mutation rate variation in general (Schuster-Bockler and Lehner 2012), whereas
205 H3K4me1, H3K4me2, H3K4me3 H3K27me3 and H3K36me3 are linked with specific
206 mutation types. Replication timing has highly statistically significant associations with
207 both SNVs and INDEL mutation types. Recombination rate and open chromatin
208 (measured by ATAC-seq) are also associated with many mutation types.
209 Transcription levels, however, lack any links with local mutation rates here.

210 Turning to nucleosomal properties, translational stability (d_{var}) is associated with
211 elevated mutation rates at A/T, non-CpG C/G and CpG sites, with the first two
212 showing the greatest effect sizes. INDELS also show similar effects, though the
213 higher p values compared with SNVs could partly be due to the smaller sample size.
214 Examining specific SNV mutation types, d_{var} is significantly associated with all A/T
215 and C/G mutations (**Supplementary Fig. 3**), except for CpG>TpG (adjusted p =
216 0.10).. The regression coefficients for d_{var} are always negative (i.e., nucleosome

217 variability is anti-correlated with mutation rate, see coefficients in **Supplementary**
218 **Table 1**), indicating that translational stability is positively associated with mutation
219 rates thus corroborating the patterns observed in **Fig. 1**. As expected from **Fig. 1**, the
220 mean distance to dyads, d_{mean} , also displays statistically significant associations with
221 mutations rates at A/T and C/G sites (**Fig. 2b,c**). Finally, nucleosome occupancy is
222 also statistically significant; in contrast to the positioning variables however, here the
223 effect is much larger for INDELS than SNVs (**Fig. 2b,c**; INDELS, adjusted $p = 5.8e-37$;
224 SNVs, adjusted $p = 0.21, 1.6e-6$ and $2.2e-7$). The regression coefficients of
225 occupancy are negative for SNVs at A/T sites, but positive for SNVs at CpG sites
226 (**Supplementary Table 1**), suggesting that occupancy can have opposing effects on
227 mutability depending on sequence context.

228 Nucleosome positioning stability is at least partly determined by the occupied DNA
229 sequence and thus its effects on mutation rates to some degree can be attributed to
230 the associated sequence (this also applies to other reported factors such as
231 replication timing). However, higher-order interactions among the long stretches of
232 nucleotides which guide nucleosome positioning are difficult to model properly.
233 Nonetheless, we achieved similar statistical significance for translational stability after
234 including non-additive two-way interaction effects for ± 5 nucleotides and the 7-mer
235 mutability estimates from Carlson et al. in regression models (Methods;
236 **Supplementary Fig. 4a,b**).

237 Since many strong nucleosomes are associated with repeat elements, we added
238 repeat status as a predictor in the regression models (Methods). We still achieved
239 strong statistical significance for translational stability after considering repeat status
240 (**Supplementary Fig. 4c**), suggesting that translational stability is independently
241 associated with mutation rate variation. We also tested repeat and non-repeat
242 regions separately, and in most tests (including those for non-repeat regions)
243 translational stability is a significant factor (**Supplementary Fig. 4d**).

244 Taken together, the logistic regression modeling analysis recapitulated known factors
245 and confirmed the independent contribution of nucleosome translational stability as a
246 new significant factor to local mutation rate variation.

247 **2.4 Mutational processes associated with elevated mutability around strong** 248 **nucleosomes**

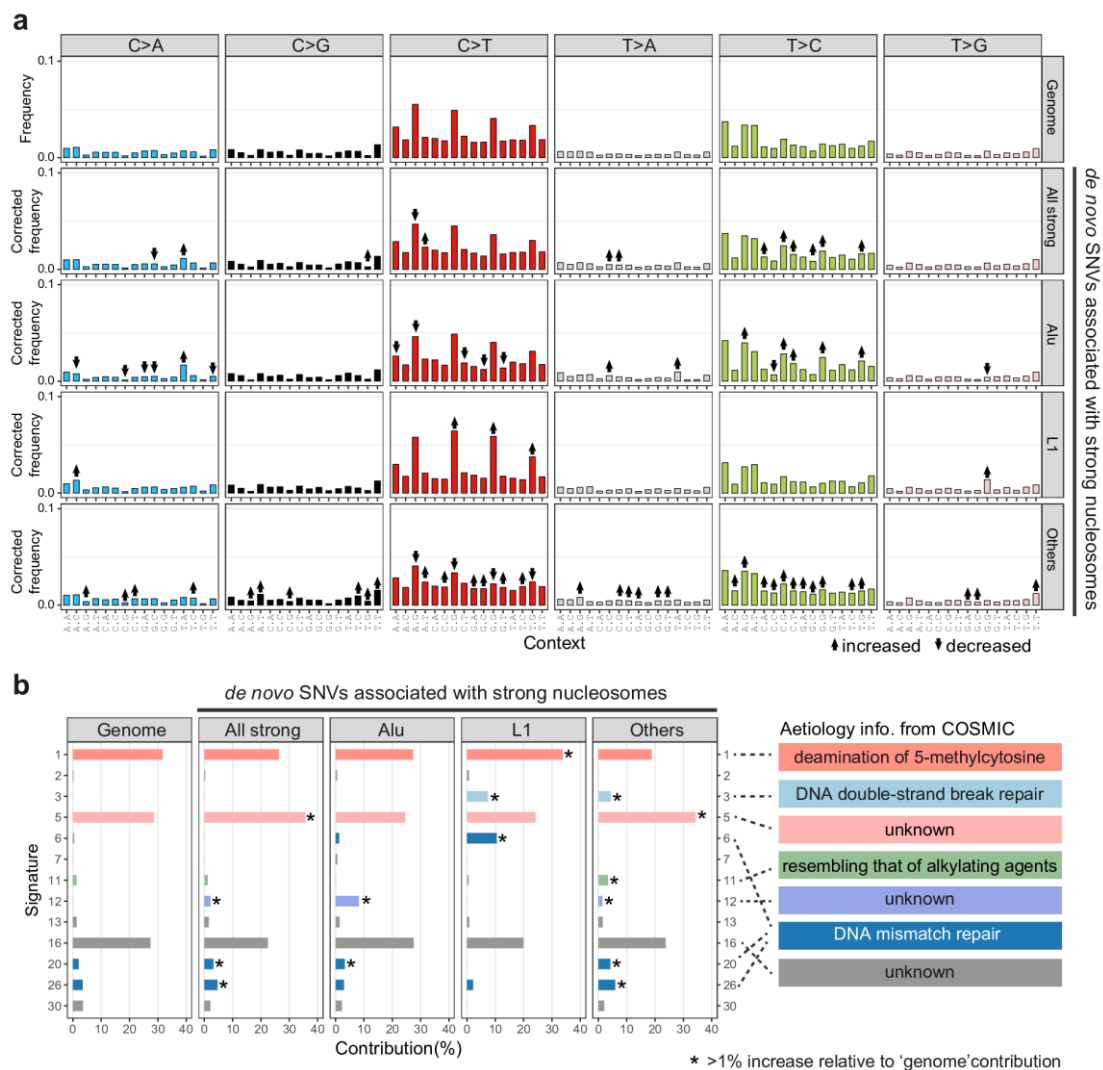
249 **2.4.1 Mutational signature analysis**

250 Having established an association between mutation rate and nucleosome
251 translational stability, we next sought to identify mutational mechanisms that might
252 explain it. As an initial screen, we compared the COSMIC mutational signatures for
253 *de novo* mutations within strong nucleosomes and those in genomic background.
254 Mutational signatures were originally developed to infer the mutational processes
255 underlying cancer progression by combining the relative frequencies of 96 possible
256 mutation types (six types of single nucleotide substitutions C>A, C>G, C>T, T>A,
257 T>C and T>G, each considered in the context of the bases immediately 5' and 3' to
258 each mutated base; Alexandrov et al. (2013)).

259 We first consider the relative frequencies of the 96 mutation types in the whole
260 genome and strong nucleosomes in different repeat contexts (**Fig. 3a**). The results
261 account for background differences in trinucleotide frequencies between these
262 regions (Methods). Several mutation types display distinct frequencies in strong
263 nucleosomes, suggesting differences in the underlying mutational processes. For
264 instance, 6 out of 16 T>C mutation types are more prevalent in strong nucleosomes
265 and different repeat-based subgroups display distinct C>T mutation frequencies. L1-
266 associated strong nucleosomes tend to show the most similar mutation frequencies
267 to genomic background, whereas the 'Others' group show the most changes,
268 perhaps reflecting the heterogeneity of constituent genomic regions.

269 Next, we applied the MutationalPatterns software (Blokzijl et al. 2018) to calculate the
270 contribution of COSMIC mutational signatures to different sets of *de novo* SNVs.
271 Three major signatures (Signatures 1, 5 and 16) are present in all tested groups
272 (contributing 87.7% for the whole-genome group, 77.0%~84.5% for strong-
273 nucleosome groups; **Fig. 3b**). Four signatures (Signatures 5, 12, 20 and 26) show
274 increased contribution (>1%) to the 'all strong-nucleosome' group relative to the
275 genomic background. The aetiologies of Signatures 5 (~7% increase in strong-
276 nucleosome regions) and 12 (2.2% increase) are currently unknown according to the
277 COSMIC website, but a recent study (Roy et al. 2018) suggested that Signature 5 is
278 likely associated with POL θ -mediated mutagenesis and double-strand break repair.
279 Signatures 20 (1.3% increase) and 26 (1.2% increase) are associated with DNA
280 mismatch repair. There are further differences in associated signatures among strong
281 nucleosome-associated SNVs in different repeat contexts ('Alu', 'L1' and 'Others'; **Fig.**
282 **3b**), such as signatures 1, 3, 5, 6, 11, 12, 20 and 26. Such differences between
283 different groups could be due to the heterogeneity of contributing mutational
284 processes and redundancy among some COSMIC signatures.

285 It is worth highlighting that COSMIC mutational signatures were designed for use
 286 with cancer genomes and so some germline mutational processes may not be well
 287 represented. Nevertheless, our analysis identified several candidate mutational
 288 processes associated with strong nucleosomes, such as the mutagenesis linked to
 289 DNA mismatch repair (Signatures 6, 20 and 26) and DNA double-strand repair
 290 (Signatures 3 and 5). Therefore, to gain deeper insights and to obtain independent
 291 evidence for these mutational processes, we examined multiple published genomic
 292 and functional genomic datasets below.



293

294 **Fig. 3 De novo SNVs in strong nucleosomes display distinct mutation type**
 295 **frequencies and COSMIC mutational signatures.** (a) Frequencies of 96 mutation
 296 types among *de novo* SNVs; 6 nucleotide substitutions in the context of the bases
 297 immediately 5' and 3' of the mutated site. SNVs are grouped into those overlapping
 298 strong nucleosomes and those elsewhere, and among the former into those overlapping
 299 with different classes of repeat elements. ↑ and ↓ indicate mutation types showing
 300 statistically significant differences relative to the genomic background SNV set (adjusted

301 $p < 0.05$, Fisher's exact test). **(b)** Percentage contribution of COSMIC mutational
302 signatures among different groups of SNVs; only signatures with non-zero values are
303 shown. * indicate mutational signatures displaying $>1\%$ increase relative to the genomic
304 background SNV set. Brief summaries of the aetiologies of affected signatures are
305 shown on the right (descriptions taken from the COSMIC website).

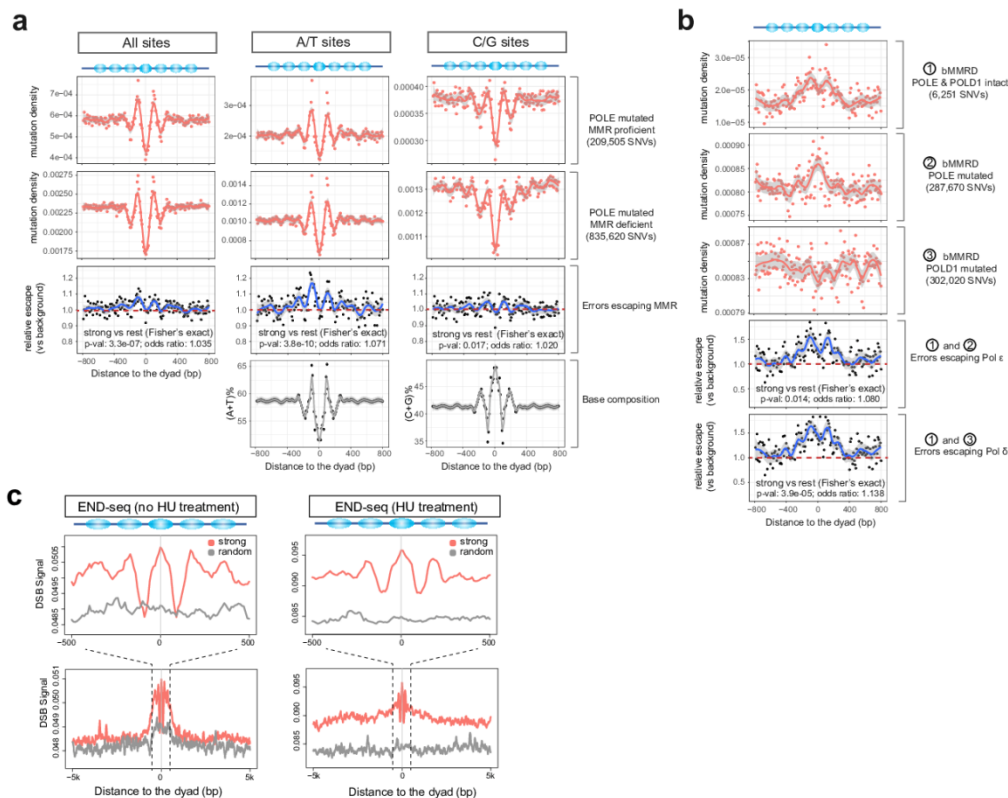
306 **2.4.2 Mismatch repair (Signatures 6, 20 and 26)**

307 DNA Mismatch repair (MMR) is a major pathway that is active during DNA replication:
308 it mainly repairs mismatches and short INDELS introduced by DNA synthesis that
309 have escaped polymerase proofreading. Mutations arising from inefficiencies in MMR
310 are represented by Signatures 6, 20 and 26, which show increased contribution to *de*
311 *novo* SNVs in the 'All strong nucleosomes' group (2% increase collectively) and three
312 repeat-based subgroups of mutations (1.6%, 6.7% and 4.3% increase for 'Alu', 'L1'
313 and 'Others', respectively).

314 We analyzed somatic mutations from two sets of ultra-hypermuted cancer
315 genomes (Campbell et al. 2017). The first comprised genomes with driver mutations
316 in the *POLE* gene encoding the catalytic subunit of DNA polymerase ϵ (Pol ϵ , the
317 major replicase for the leading strand) and in one or more of the core MMR genes
318 (*MLH1*, *MSH2*, *MSH6*, *PMS1* and *PMS2*). The second contained cancers with
319 mutated *POLE* but intact MMR. As it is even more challenging to detect somatic
320 mutations in tumor-derived data than re-sequencing of normal individuals, we
321 focused this analysis on strong nucleosomes found in high-mappability regions of the
322 genome (Methods).

323 We reasoned that differences in mutation distributions between the two sets of
324 genomes could be attributed to the MMR pathway. The overall mutation patterns are
325 similar in both cases, with much higher mutation rates at strong nucleosome
326 boundaries and adjacent linker DNA than the surrounding regions (**Fig. 4a**). This
327 implies that errors introduced during error-prone replication by a deficient Pol ϵ
328 escape repair by the MMR pathway when they coincide with strong nucleosomes.
329 Next, we calculated an 'MMR escape ratio' to quantify the relative amount of
330 replication errors that escapes MMR repair in the *POLE* only mutant cancers
331 compared with the *POLE* and MMR double mutants. Strong nucleosomal regions
332 (especially boundaries and adjacent linkers) display $\sim 10\%$ higher escape ratios than
333 the genome-wide background (**Fig. 4a**). Although A/T sites have higher escape ratios
334 than C/G sites around strong nucleosomes, both C/G and A/T sites exhibit similarly
335 elevated escape ratio profiles, suggesting independence of sequence context.

336 Moreover, the apparent ~200-bp periodicity in escape ratio and mutation density
 337 profiles are suggestive of associations with nucleosome positioning rather than
 338 sequence alone. Together, these observations strongly indicate a relationship
 339 between replication errors, MMR and strong nucleosomes in elevating mutation rates.



340

341 **Fig. 4 Mismatch repair (MMR), DNA polymerase fidelity and double strand breaks (DSB)**
 342 **explain increased mutation rates in strong nucleosomes.** (a) Mutation density profiles
 343 relative to strong nucleosome dyads in cancer genomes harboring driver mutations in the
 344 *POLE* and MMR pathway genes. Numbers of mutations used are indicated in the brackets.
 345 The MMR escape ratio compares the mutation densities in the MMR proficient and MMR
 346 deficient genomes. (b) Mutation density profiles relative to strong nucleosome dyads for
 347 bMMRD cancer genomes with different driver mutation statuses in the *POLE* and *POLD1*
 348 genes. The escape ratios compare the mutation densities for Pol ϵ -deficient and Pol δ -
 349 deficient cancers with the proficient ones. (c) END-seq signal indicating the density of DSBs
 350 relative to strong nucleosome dyads. HU, hydroxyurea. Fisher's exact test was used for
 351 testing the association of strong nucleosomal regions (dyad \pm 95bp) with differential
 352 MMR/polymerase performance.

353 2.4.3 DNA polymerase fidelity (Signatures 10 and possibly 12)

354 We also studied the effect of strong nucleosomes on replication fidelity by examining
 355 data from children with inherited biallelic mismatch repair deficiency (bMMRD);

356 (Shlien et al. 2015); these include ultra-hypermuted genomes arising from Pol ϵ
357 and polymerase δ defects (Pol δ , the major replicase for the lagging strand). We
358 estimated Pol δ and Pol ϵ escape ratios (escaping the proofreading correction of
359 polymerases) using the same reasoning as above (**Fig. 4b**). We found that strong
360 nucleosomes have higher escape ratios for both polymerases relative to the genomic
361 background (**Fig. 4b**), implying that they have lower replication fidelity in these
362 regions. The proofreading escape ratios for both polymerases are even higher than
363 that for MMR (**Fig. 4a,b**) and A/T sites display higher proofreading escape ratios than
364 C/G sites (**Supplementary Fig. 5a**). Again, the periodic pattern in the relative escape
365 profiles (**Fig. 4b, Supplementary Fig. 5a**) suggests that nucleosome positioning
366 contributes to the heterogeneity in replicase fidelity across the genome.

367 The aetiology of Signature 12 is currently unknown. Here, we found that it contributes
368 21.15%~21.99% to mutations in *POLD1*-mutant bMMRD genomes (inferred by
369 MutationalPatterns, **Supplementary Fig. 5b,c**), but much less for other bMMRD
370 samples (0~2.88% for *POLE*-mutant, and 3.32%~10.43% for *POLE/POLD1*-intact).
371 This suggests that Signature 12 is probably associated with Pol δ and that many *de*
372 *novo* mutations around strong nucleosomes arise from errors escaping Pol δ
373 proofreading. Surprisingly, Signature 10, known to be associated with Pol ϵ
374 deficiency, is absent from strong nucleosomal *de novo* SNVs (**Fig. 3b**). This
375 suggests that although both Pol ϵ and Pol δ have high proofreading escape ratios (i.e.
376 low fidelities) around strong nucleosomes (**Fig. 4b**), the majority of the replication
377 errors that are eventually converted to *de novo* mutations are derived from lagging
378 strand replicase Pol δ .

379 Reijns et al (2015) showed that in budding yeast, Okazaki junctions formed during
380 lagging strand replication tend to be near nucleosome dyads and display elevated
381 mutation rates (Reijns et al. 2015). We tested this by re-analyzing OK-seq data from
382 human lymphoblastoid cells (Petryk et al. 2016). Unlike yeast, Okazaki junctions in
383 humans are more frequently located in the linker regions (**Supplementary Fig. 6**)
384 rather than the dyads, suggesting that the mutagenic effects of Okazaki junctions are
385 different in the two organisms. This may partly be because yeast lacks the typical H1
386 histone found in human and other eukaryotes. However, the very short reads (single-
387 ended 50bp) of OK-seq data restricted our analysis to nucleosomes with high
388 mappability (~10% of strong nucleosomes), limiting the strength of the conclusions
389 here.

390 **2.4.4 Double-strand breaks (Signatures 3 and 5)**

391 Double-strand break (DSB) repair represented by Signatures 3 and 5 is another
392 potential mechanism involved in strong nucleosome-associated mutations (**Fig. 3b**).
393 Tubbs et al. (2018) studied the genome-wide distribution of DSBs using END-seq
394 and suggested that poly(dA:dT) tracts are recurrent sites of replication-associated
395 DSBs. Our analysis of this data revealed a higher frequency of DSBs around strong
396 nucleosomes compared with genomic background (**Fig. 4c**). The trend holds for
397 experiments with and without hydroxyurea treatment (HU, a replicative stress-
398 inducing agent), suggesting that strong nucleosomes are endogenous hotspots (i.e.
399 without HU treatment) of DSBs during replication. It is notable that young Alu and L1
400 elements harbor prominent poly(dA:dT) tracts, which are enriched at the boundary
401 and linker regions of strong nucleosomes (**Supplementary Fig. 7a**). The patterns of
402 high DSB frequency still hold true when looking at strong nucleosomes associated
403 with different repeats (**Supplementary Fig. 7b,c**). However, because the END-seq
404 data were sequenced with single-ended 75bp reads and majority of young Alu and
405 L1 elements cannot be assessed with such short reads, we could not pursue further
406 detailed analysis. Since DSB repair can be error-prone (Rodgers and McVey 2016),
407 even using high-fidelity homologous recombination, frequent DSB formation and
408 subsequent error-prone repair likely contribute to the elevated mutation rates around
409 strong nucleosomes.

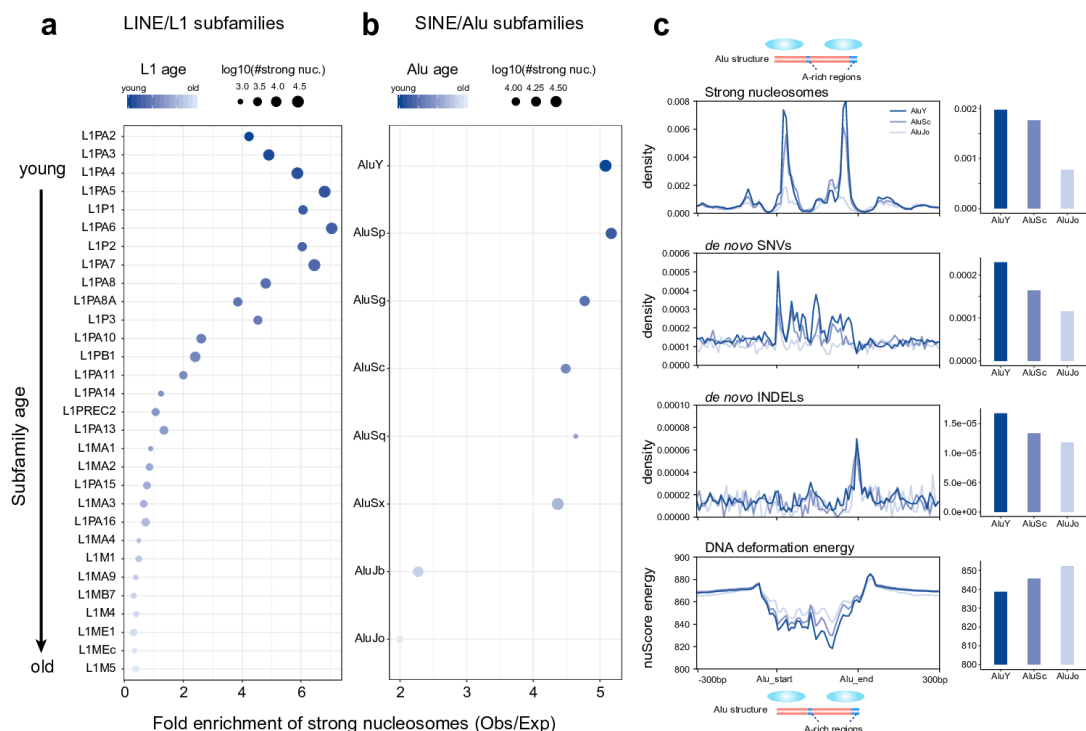
410 **2.5 Strong nucleosome positioning is mostly associated with young repeat** 411 **elements and undergoes frequent turnover**

412 Above, we highlighted that ~70% of strong nucleosomes are located in Alu and L1
413 retrotransposons (**Supplementary Fig. 1d**). Upon examination of the subfamilies
414 (**Fig. 5a,b**), we uncovered a strong enrichment for evolutionarily young L1s (e.g.
415 L1PA2 to L1PA11) and Alus (e.g. AluY to AluSx). Since younger repeats have poorer
416 mappability, these observations probably underestimate the true enrichment. This
417 may also explain why several of the youngest L1 subfamilies (L1PA2 to L1PA5) have
418 lower enrichments than the slightly older subfamilies (**Fig. 5a**).

419 The preference for nucleosomes to occupy specific sections of Alu elements is
420 supported by both *in vitro* and *in vivo* evidence (Englander et al. 1993; Englander and
421 Howard 1995; Salih et al. 2008; Tanaka et al. 2010). We recapitulated these
422 observations for strong nucleosomes using the Gaffney et al. MNase-seq data (**Fig.**
423 **5c**): there are two hotspots of strong nucleosomes in young Alus, which fade away in
424 older elements. We also observed that younger Alus exhibit elevated *de novo*
425 mutation rates compared with old ones (**Fig. 5c**), and the weaker translational

426 stability in older Alus is accompanied by reduced *de novo* mutation rates for both
 427 SNVs and INDELs (Fig. 5c). Thus, there is an intriguing interplay between Alus,
 428 strong nucleosomes and mutation rates.

429 The histone octamer is thought to preferentially bind DNA sequences presenting
 430 lower deformation energy costs (Tolstorukov et al. 2008). We estimated deformation
 431 energies using the nuScore software (Tolstorukov et al. 2008) based on the DNA
 432 sequence and nucleosome core particle structure and we found that Alus do indeed
 433 exhibit lower deformation energies than surrounding regions (Fig. 5c). Furthermore,
 434 the energies of Alu elements tend to increase with age, suggesting that the
 435 accumulated mutations in Alu sequences reduced their nucleosome-binding stability.
 436 This is also supported by comparing deformation energies of Alu consensus
 437 sequences (ancestral states) and those of current genomic sequences
 438 (Supplementary Fig. 8a). We further analyzed the 3' end sequences of L1 elements
 439 harboring strong nucleosomes and observed similar patterns (Supplementary Fig.
 440 8b,c).



441

442 **Fig. 5 Strong nucleosomes are frequently found inside evolutionarily young LINE**
 443 **and SINE elements.** (a) Fold enrichment of strong nucleosome occurrence in L1
 444 subfamilies. The top 30 abundant subfamilies are shown ordered by evolutionary age.
 445 Dot sizes depict the numbers of strong nucleosomes and color-scale indicates the
 446 subfamily age. (b) Same as (a) but for Alu elements. (c) Densities of strong nucleosome
 447 dyads, *de novo* SNVs and *de novo* INDELs along the Alu sequences and flanking

448 regions, grouped by Alu subfamilies of different ages. Bar plots show the average
449 densities for all Alus of different subfamilies on the right. The bottom panel shows the
450 average DNA deformation energies along Alu sequences estimated using nuScore.
451 Profiles were plotted using Alu elements ≥ 250 bp and all elements were scaled up to a
452 300bp region in the plots.

453 Studies have suggested that natural selection appears to preserve nucleosome
454 positioning during evolution (Prendergast and Semple 2011; Tolstorukov et al. 2011;
455 Drillon et al. 2016), but they had differing views about the effects of selection on the
456 underlying sequence. In contrast, Warnecke et al. (2013) suggested that the
457 observed sequence divergence patterns around nucleosomes can be explained by
458 frequent nucleosome re-positioning after mutation, rather than by natural selection.
459 Since these results were mainly based on human polymorphisms or inter-species
460 divergence, indirect mutation rate measurements were potentially confounded by
461 selection and non-adaptive processes. The use of *de novo* mutations helps resolve
462 this debate to some extent.

463 As we showed above, there is considerable *de novo* mutation rate variation around
464 strong nucleosomes (**Fig. 1e, Supplementary Fig. 1**), which cannot be ignored in
465 any selection analysis. Furthermore, strong nucleosomes are clearly preferentially
466 present in young SINE/LINE elements and the strength of translational stability
467 decays substantially over time (**Fig. 5**). These observations support the re-positioning
468 model over a long evolutionary scale. Since a large majority of strong nucleosomes
469 associated with SINE/LINE elements are expected to become non-strong ones in
470 future, selection for preserving positioning might not be as widespread as previously
471 suggested, though it may happen at some particular regions or within a short
472 evolutionary scale.

473 **3 Discussion**

474 Though the involvement of nucleosome organization in DNA damage/repair
475 processes was recognised nearly 30 years ago (Smerdon 1991), its genome-wide
476 effects on germline mutation rates (particularly in higher eukaryotes) have remained
477 poorly understood. Our analysis combining large-scale *de novo* mutation and
478 nucleosome datasets in human provides several important insights into this topic.

479 A major finding is that strong translational positioning of nucleosomes is associated
480 with elevated *de novo* mutation rates, which is also supported by observations using
481 extremely rare variants in polymorphism data. The ability to use *de novo* mutations

482 here allowed us to bypass confounding evolutionary factors such as selection, thus
483 allowing direct assessment of the impact on background mutation rates. Importantly,
484 our statistical tests controlling for nucleosome occupancy and other related factors
485 confirmed the significant contribution of translational stability to mutation rate
486 variation. Therefore, we have discovered a novel factor that significantly modulate
487 germline mutation rate variation.

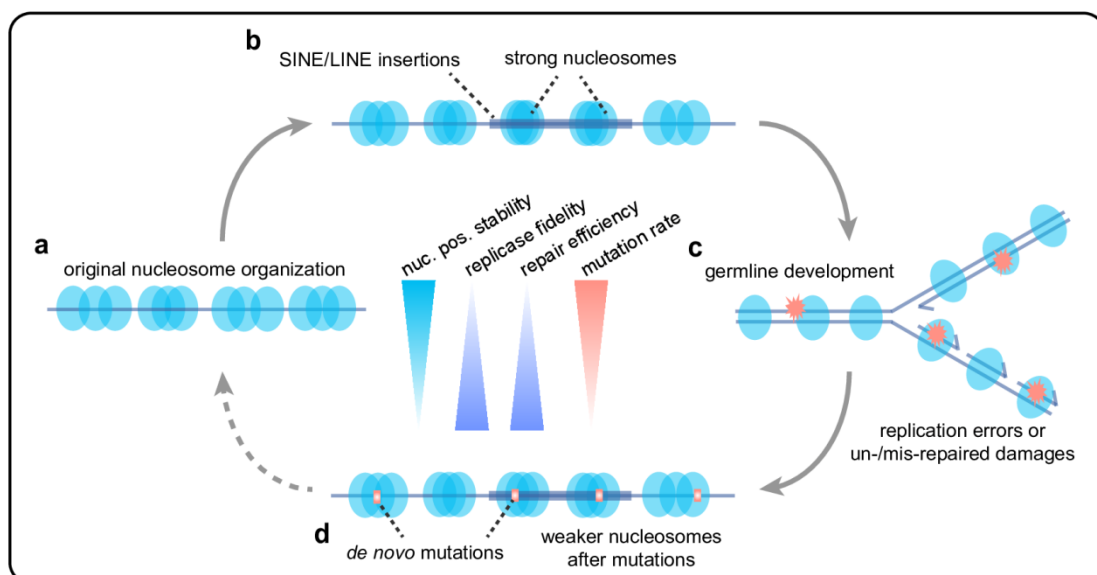
488 Investigating the underlying mutational processes responsible for this association
489 remains challenging. Nevertheless, we obtained several informative results regarding
490 potential mechanisms by leveraging published omics data related to DNA damage
491 and repair. In doing so, we revealed that MMR, replicase fidelity and DSB contribute
492 significantly to elevated mutation rates around strong nucleosomes. In particular,
493 multiple sets of ultra-hypermuted cancer data allowed us to quantify the
494 performance of MMR and replicases by calculating the repair escape ratios. The
495 results probably apply to germ cells because i) they agree nicely with the
496 observations from our mutational signature analysis with *de novo* mutations and ii)
497 recent studies suggested that replicative errors account for majority of mutations
498 arising in both somatic and germ cells (Tomasetti and Vogelstein 2015; Tomasetti et
499 al. 2017). The precise molecular interactions determining the relationships between
500 strong nucleosome positioning, replicase fidelity and DNA repair are still not clear.
501 However, based on the evidence from our analysis with the omics data and previous
502 studies (Li et al. 2009; Reijns et al. 2015; Tubbs et al. 2018), we speculate that
503 strong nucleosomes may act as particularly strong barriers which impair the
504 performance of the replication and repair machineries. There may be additional,
505 unexamined effects on DNA damage/repair processes related to germline
506 development, but many published genomic datasets about DNA damage/repair were
507 generated in non-germ cells and with very short sequencing reads (e.g. <100bp),
508 which hinder accurate analysis. Improved sequencing strategies such as long-read
509 sequencing and direct measurement in germ cells would benefit future related
510 studies.

511 Interestingly, we found that strong nucleosomes are preferentially located within
512 young LINE and SINE elements, two of the most common retrotransposons in the
513 human and other mammalian genomes. Owing to their potentially deleterious effects,
514 newly inserted retrotransposons are tightly repressed by multiple regulatory
515 mechanisms, such as DNA methylation and H3K9me3 (Slotkin and Martienssen
516 2007). Strong nucleosome positioning, which may mask access to the transcription
517 machinery, could be another layer of the repressive system. Furthermore, the

518 hypermutation in young SINEs/LINEs, partly contributed by associated strong
519 nucleosomes, could lead to the rapid reduction of retrotransposition capacity.
520 Therefore, the combination of strong nucleosome positioning and hypermutation in
521 SINEs/LINEs might have facilitated their expansion across the genome.

522 The decreasing numbers of strong nucleosomes in older LINE/SINE elements imply
523 widespread nucleosome re-positioning during evolution. Since nucleosome
524 positioning is strongly affected by the underlying DNA sequence, their re-positioning
525 probably arises from the accumulation of mutations. Our data largely disagree with
526 the previous hypothesis of widespread selection for maintaining nucleosome
527 positioning in the human genome (Prendergast and Semple 2011). Another reason
528 for favoring the re-positioning model is that most genomic regions do not employ
529 strong positioning, possibly due to its relatively high mutagenic potential.

530 Finally, we summarized our major findings in a proposed model in **Fig. 6**, which
531 demonstrates the relationship among nucleosome positioning, mutation rate variation,
532 retrotransposons and evolution. Given the importance of germline *de novo* mutations
533 in evolution and human diseases and the universal roles of nucleosomes in
534 eukaryotic genome organization and regulation, our work should have profound
535 implications in related research areas.



536

537 **Fig. 6 Proposed model of the interplay between nucleosome translational stability,**
538 **mutation rate and transposable elements.** (a) Most genomic regions are occupied by
539 nucleosomes lacking strong translational stability. (b) Strong nucleosomes are
540 preferentially associated with newly inserted SINE/LINE elements. (c) Strong
541 nucleosomal regions are subject to high mutation rates during germline development,

542 caused by mutational processes such as low replicase fidelity, inefficient MMR and DSB
543 repair. (d) Accumulation of mutations reduces translational stability of strong
544 nucleosomes and reduces transposition capacity of transposable elements.

545

546

547

548 **Methods**

549 **Mutation datasets**

550 *De novo* mutations identified in multiple large-scale trio sequencing project were
551 downloaded from denovo-db v1.6.1 (Turner et al. 2017b). Seven studies with >1000
552 *de novo* mutations (Genome of the Netherlands 2014; Turner et al. 2016; Yuen et al.
553 2016; Jonsson et al. 2017; Turner et al. 2017a; Yuen et al. 2017; Werling et al. 2018)
554 were considered in our analysis (**Supplementary Fig. 1a**). Extremely rare variants
555 (derived allele frequency ≤ 0.0001) were obtained from Genome Aggregation
556 Database (gnomAD, release 2.0.2) (Lek et al. 2016).

557 **Nucleosome datasets**

558 We used the 1,037,801 strong nucleosomes (i.e. translationally stable nucleosomes)
559 identified based on MNase-seq data of sequenced seven lymphoblastoid cell lines
560 from Gaffney et al. (Gaffney et al. 2012). The original hg18-based coordinates of
561 annotated nucleosomes were converted to hg19 using the 'liftOver' tool from UCSC
562 genome browser. The rotationally stable nucleosomes identified based on 49 DNase-
563 seq samples (43 distinct cell types) were from Winter et al. (Winter et al. 2013). We
564 classified the human genome into three groups based on the nucleosome contexts
565 (**Fig. 1b**): i) regions covered by translationally stable ('strong') nucleosomes; ii)
566 regions covered by rotationally but not stable translationally nucleosomes; and iii) the
567 remaining genomic regions. Chromosomes X and Y were excluded from analysis as
568 some other datasets used in our work lacked data for these chromosomes. As the
569 nucleosome maps we used were not derived from germ cells, for downstream
570 analysis we excluded the genomic regions in which nucleosome positioning were
571 found to differ between human embryonic stem cells and differentiated fibroblasts
572 (West et al. 2014). Based on the positioning stability scores defined in Gaffney et al.,
573 we divided the one million strong nucleosomes into three categories of equal sizes
574 with different levels of stability – 'high', 'middle' and 'low', which were used for
575 analysis shown in **Fig. 1** and **Supplementary Fig. 1**.

576 **Accounting for mappability**

577 Sequencing read mappability can significantly affect variant calling results and other
578 aligned read-depth based measurements (e.g. nucleosome occupancy). The
579 sequencing reads for detecting *de novo* mutations used in our analysis were mainly
580 150bp paired-end reads, with fragment sizes ranging from 300-700bp
581 (**Supplementary Fig. 1**). We used the Genome Mappability Analyzer (GMA) (Lee

582 and Schatz 2012) to generate the mappability scores for simulated paired-end 150
583 reads with fragment sizes set to be 400bp. Only the regions with GMA mappability
584 scores of ≥ 90 (~2.59Gb) were considered for most analyses, unless specified
585 otherwise. We did not use the mappability tracks from ENCODE for the *de novo*
586 mutation data, because those tracks were only for single-ended reads. For some
587 analyses, additional filtering were applied if other associated datasets suffered from
588 more severe mappability issues. For measuring nucleosome occupancy, we used the
589 method described in the Gaffney et al. to simulate paired-end 25bp reads matching
590 the base compositions of MNase-seq data in the human genome, and then
591 calculated per-base coverage depth by the simulated fragments. The 10bp-bin ratios
592 between the MNase-seq read coverage and the simulated read coverage were used
593 for measuring the occupancy.

594 **Enrichment analysis for *de novo* mutations in different nucleosome contexts**

595 Genomic association tester (GAT) (Heger et al. 2013), a tool for computing the
596 significance of overlap between multiple sets of genomic intervals, was used to
597 estimate the expected numbers of mutations in different contexts (sampling ≥ 1000
598 times), which were then compared with the observed numbers. Low-mappability
599 regions were excluded from analysis. A similar analysis was also done for the
600 extremely rare variants of gnomAD. Analysis of meta-profiles along strong
601 nucleosomes was done using deepTools (Ramirez et al. 2014).

602 **Statistical modelling of the contribution of different factors to mutation rate** 603 **variation**

604 As described in the main text, for a given genomic position, we defined two variables
605 regarding the translational positioning of nearby nucleosomes (**Fig. 2a**):

$$d_{mean} = \frac{\sum_{i=1}^n d_i}{n}, 0 \leq d \leq 95,$$
$$d_{var} = \frac{\sum_{i=1}^n (d_i - d_{mean})^2}{n}$$

606 where d is the distance between a MNase-seq midpoint to the focal site. We
607 considered MNase-seq midpoints within ± 95 bp of the focal site, because genome-
608 wide nucleosome repeat length was estimated to be 191.4bp for the Gaffney et al.
609 data (Gaffney et al. 2012). Genomic sites without any MNase-seq midpoint within
610 ± 95 bp were excluded from analysis (123Mb out of 2.59Gb excluded). The
611 measurements for nucleosome occupancy were 10bp-bin ratios between the MNase-

612 seq read coverage and the simulated read coverage. We did not use the positioning
613 score $S(i)$ defined in Gaffney et al. to measure positioning stability in our modelling
614 analysis, because $S(i)$ was designed for identifying the stable dyads and so for non-
615 dyad positions it does not represent the positioning stability properly.

616 RNA expression, DNA methylation and chromatin accessibility (ATAC-seq) data from
617 human spermatogonial stem cells were from Guo et al. (Guo et al. 2017). For the
618 RNA-seq and ATAC-seq data from Guo et al., because the genome-wide read signal
619 tracks were not available, we downloaded, processed and mapped the raw reads to
620 generate the genome-wide tracks. Since suitable data for histone modifications in
621 human germ cells were not available, we used the ChIP-seq data of human
622 embryonic stem cells from ENCODE (ENCODE Consortium 2012). Replication timing
623 data (Repli-seq of GM12878) were also from ENCODE. The data of recombination
624 rates were from the HapMap project (International HapMap Consortium et al. 2007).

625 A binary logistic regression framework was used to assess the contribution of
626 different factors to mutation rate variation across the genome systematically. The
627 logistic regression model is described as below:

$$\begin{aligned}\mu = \Pr(y = 1) &= \frac{\exp(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)}{1 + \exp(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)} \\ &= \frac{\exp(\mathbf{X}\beta)}{1 + \exp(\mathbf{X}\beta)}\end{aligned}$$

$$\text{logit}(\mu) = \log\left(\frac{\mu}{1 - \mu}\right) = \mathbf{X}\beta$$

628 where $\mu = \Pr(y = 1)$ denotes the probability that a genomic position is mutated (for
629 testing individual SNV mutation types, e.g. A>T, μ is the probability that a site is
630 mutated to a specific nucleotide), \mathbf{X} represents the observations for the considered
631 variables (categorical or continuous, e.g. d_{mean} , d_{var} , adjacent nucleotides, etc.), and β
632 is the vector of parameters to be estimated.

633 We used the Bayesian logistic regression model implemented in the 'bayesglm'
634 (Gelman et al. 2008) of the R package 'arm', which was reported to perform well in
635 handling the complete separation issue in logistic regression models (Gelman et al.
636 2008). The complete separation issue is common when one class is rare relative to
637 the other and (or) there are many regressors in a model. As we had only ~300,000
638 *de novo* mutations, the probability for a given site to be mutated in our data is
639 ~1/10,000, which is a rare event.

640 Within the logistic regression framework, we compared the full model with all
641 considered variables to a reduced model without one specific variable by performing
642 likelihood-ratio tests in R ('anova' function) to evaluate the significance for each
643 variable. The resulting p values of a set of likelihood-ratio tests were adjusted for
644 multiple testing with Benjamini–Hochberg correction.

645 To perform the regression analysis, we generated the data of all variables for the *de*
646 *novo* mutation sites and subsampled a fraction of the non-mutated sites as the
647 control sites. We did not use all the non-mutated sites in the genome as it would lead
648 to a large imbalance in the sizes of two classes ('mutated' and 'non-mutated') and
649 much larger computational burden. For *de novo* SNVs, we randomly generated
650 2,561,953 non-mutated sites (about 1/1000 of the accessible genome, about 10
651 times as many as *de novo* SNVs) and 256,337 non-mutated sites (about 1/10,000 of
652 the accessible genome, about 10 times as many as *de novo* INDELs) for INDELs.
653 For *de novo* INDELs, we used the INDELs of ≤ 5 bp for regression analysis, because
654 long INDELs were rare and may have high false positive/negative rates. For RNA
655 expression, DNA methylation, chromatin accessibility, replication timing,
656 recombination rate and histone modifications data, we used the average value of the
657 ± 10 bp of a focal site for each specific feature based on the genome-wide signal
658 tracks. We also assessed different window sizes (± 5 bp and ± 20 bp), which led to
659 similar results.

660 For SNVs, we performed logistic regression tests for mutation types at A/T sites and
661 C/G sites separately and distinguished C/G sites in CpG and non-CpG contexts. We
662 also tested for nine individual SNV mutation types (three for A/T sites, three for C/G
663 sites at CpG contexts, and three for non-CpG contexts, **Supplementary Fig. 3**). The
664 regression coefficients for the full model of each test are given in **Supplementary**
665 **Table 1**.

666 Since the variable d_{mean} has a non-monotonic relationship with mutation rates, we
667 binned the values into five categories: [0,18], [19, 36], [37, 54], [55, 73] and [74, 95]
668 (first four bins implying nucleosome-bound regions, and the last bin implying close to
669 the linker).

670 In the regression models mentioned above, we did not consider the non-additive
671 effects of adjacent nucleotides (± 5 bp). When we tried adding non-additive effects for
672 ± 5 nucleotides (considering only two-way interactions; taking a much longer running
673 time), we got similar results regarding the association of translational stability (d_{var})

674 and mutation rates (**Supplementary Fig. 4**). We also tried using the 7-mer mutability
675 estimates from Carlson et al. (Carlson et al. 2018), which incorporated non-additive
676 effects among ± 3 nucleotides, as predictors in the regression models.

677 To evaluate how the sequence repeat status affects the effects of translational
678 stability on mutation rates, We added the repeat status ('Alu', 'L1', 'other repeat' or
679 'non-repeat') as a predictor in the regression models, and also ran the regression
680 tests for different repeat/non-repeat regions separately.

681 **Analysis of mutational processes**

682 COSMIC mutational signatures are based on frequencies of mutations in tri-
683 nucleotide contexts. Since the regions associated with strong nucleosomes have
684 different tri-nucleotide composition relative to genome background, we first
685 normalized the mutation type frequencies in regions associated with strong
686 nucleosomes as this: set $F_{i,strong}$ for the occurrence of a specific mutation type (e.g.
687 T[T>C]T), $N_{i,strong}$ for the occurrence of the considered tri-nucleotide context (e.g.
688 TTT) in strong-nucleosome regions and $N_{i,genome}$ for the occurrence of the
689 considered tri-nucleotide context in the whole-genome background, then the
690 corrected occurrence of a the mutation type for strong nucleosomes is $N'_{i,strong} =$
691 $F_{i,strong} \div N_{i,strong} \times N_{i,genome}$. Fisher's exact tests were performed to identify
692 mutation types that show significant increase or decrease in strong-nucleosome
693 regions relative to genome background. The contingency table used for running
694 'fisher.test' in R for a specific mutation type is
695 $matrix\left(c(F_{i,strong}, N_{i,strong} - F_{i,strong}, F_{i,genome} - F_{i,strong}, (N_{i,genome} - N_{i,strong}) -\right.$
696 $\left.(F_{i,genome} - F_{i,strong})\right), ncol = 2)$, where $F_{i,strong}$. And $F_{i,genome}$ are the
697 occurrences of the considered mutation type and $N_{i,strong}$ and $N_{i,genome}$ for the
698 occurrences of the considered tri-nucleotide context. Benjamini-Hochberg method
699 was used for multiple testing correction.

700 The contribution of COSMIC mutational signatures (Alexandrov et al. 2013) to
701 different sets of mutations (*de novo* SNVs and somatic mutations from bMMRD
702 samples) was predicted using the 'fit_to_signatures' function in the R package
703 'MutationalPatterns' (Blokzijl et al. 2018). For the sets of *de novo* SNVs associated
704 with strong nucleosomes, the corrected frequencies described above were used for
705 running 'fit_to_signatures'.

706 Mutations in *POLE* in cancers can lead to reduced base selectivity and/or deficient
707 proofreading during replication, producing unusually large numbers of mutations (so
708 called ‘ultra-hypermutation’) which facilitated our analysis. *POLE* mutated genomes
709 from PCAWG project (Campbell et al. 2017) were used to evaluate the differential
710 MMR efficiency between strong and non-strong nucleosome regions. We compared
711 the mutation densities in cancer genomes with *POLE* mutated and a deficient MMR
712 (4 individual samples) to those with *POLE* mutated and a proficient MMR (6 samples).
713 The MMR pathway was considered deficient if a driver mutation (annotated by the
714 PCAWG consortium) was found in one of five MMR core genes - *MLH1*, *MSH2*,
715 *MSH6*, *PMS1* and *PMS2*.

716 For a given bin (10bp-size) in the meta-profile, we calculated the relative MMR
717 escape ratio relative to genomic background around strong nucleosomes as
718 described in the following formula,

$$R_i^{escape} = \frac{\frac{m_i^{POLE^*, MMR^{WT}}}{m_i^{POLE^*, MMR^*}}}{\frac{\bar{m}^{POLE^*, MMR^{WT}}}{\bar{m}^{POLE^*, MMR^*}}}$$

719 where m_i is the mutation density for the i th bin (observed number of mutations in the
720 i th bin divided by the bin size), and \bar{m} is the genome-wide average mutation density
721 of a specific sample group (observed number of mutations in the simulated windows
722 divided by the total window size), estimated by simulating random windows in the
723 genome. A similar logic was used when evaluating relative proofreading escape
724 ratios of Pol ϵ (mutated *POLE*) and Pol δ (mutated *POLD1*) using the somatic
725 mutation data from the bMMRD project (Shlien et al. 2015).

726 When analyzing PCAWG and bMMRD data, to account for potential mappability
727 issues, we focused on the highly mappable regions based on the CrgMapability
728 scores from ENCODE. We used CrgMapability scores here, which are more stringent
729 than GMA ones, because detecting somatic mutations in tumors is more difficult than
730 for ordinary individual re-sequencing data. We considered the strong nucleosomes
731 which have a 100mer CrgMapability score of 1 (meaning any 100-bp read from these
732 regions can be mapped uniquely in the genome) within ± 800 bp of the dyads. We
733 then simulated a same number of 1600bp-sized regions from the genome that satisfy
734 the mappability requirement to calculate the background mutation density. Note that
735 in theory the mappability issue in the relative escape ratios should be very small

736 because the two sets of samples have the same mappability for a given bin and the
737 ratio calculation normalizes the effects of different mappability among regions.

738 The raw reads of OK-seq data (Petryk et al. 2016) were downloaded from NCBI and
739 mapped to the human genome. We kept only the uniquely mapped reads for inferring
740 Okazaki junctions. The very 5' end sites of aligned reads (separating reads mapped
741 to Watson and Crick strands) were considered putative Okazaki junction signals.

742 To investigate DSBs around strong nucleosomes, we downloaded the genome-wide
743 tracks of human END-seq data (GSM3227951 and GSM3227952) (Tubbs et al.
744 2018). Because the reads of END-seq data were single-ended 75bp, we considered
745 the strong nucleosomes which have a 75mer CrgMapability score of 1 within ± 500 bp
746 of the strong nucleosome dyads for analysis.

747 **Enrichment analysis for strong nucleosomes in different repeat contexts**

748 GAT (Heger et al. 2013) was used to estimate the expected numbers of strong
749 nucleosomes in different contexts (sampling ≥ 1000 times), which were compared to
750 the observed numbers. The annotations of repeat elements (Feb 2009, Repeat
751 Library 20140131) were downloaded from RepeatMasker (Tempel 2012). We also
752 did GAT analysis for LINE-1(L1) and Alu subfamilies of different ages. The age
753 information of repeat families was from Giordano et al. (Giordano et al. 2007). For
754 generating the MNase-seq midpoints along the repeat consensus sequences, we
755 made use of the alignment information in the RepeatMasker result files
756 ('hg19.fa.align.gz') and mapped the hg19-based coordinates to the coordinates in the
757 consensus sequences. Strong nucleosomes appear to be under-detected in very
758 young L1 elements, which we think is due to difficulties in mapping short MNase-seq
759 reads (Alus are easier to map because they are much smaller).

760 Nucleosome deformation energies of all sites in the human genome were estimated
761 using nuScore (Tolstorukov et al. 2008). We also used nuScore to estimate the
762 deformation energies of Alu/L1 subfamily consensus sequences. For the L1 analysis
763 shown in **Supplementary Fig. 8**, we only considered the 3' end regions of L1
764 subfamilies, because 5' end regions of L1 elements are usually truncated in the
765 genome and their subfamily identities are difficult to be determined.

766

767

768 References

- 769 Acuna-Hidalgo R, Veltman JA, Hoischen A. 2016. New insights into the generation
770 and role of de novo mutations in health and disease. *Genome biology* **17**(1):
771 241.
- 772 Alexandrov LB, Nik-Zainal S, Wedge DC, Aparicio SA, Behjati S, Biankin AV, Bignell
773 GR, Bolli N, Borg A, Borresen-Dale AL et al. 2013. Signatures of mutational
774 processes in human cancer. *Nature* **500**(7463): 415-421.
- 775 Blokzijl F, Janssen R, van Boxtel R, Cuppen E. 2018. MutationalPatterns:
776 comprehensive genome-wide analysis of mutational processes. *Genome*
777 *medicine* **10**(1): 33.
- 778 Campbell PJ, Getz G, Stuart JM, Korbelt JO, Stein LD. 2017. Pan-cancer analysis of
779 whole genomes. *BioRxiv*: 162784.
- 780 Carlson J, Locke AE, Flickinger M, Zawistowski M, Levy S, Myers RM, Boehnke M,
781 Kang HM, Scott LJ, Li JZ et al. 2018. Extremely rare variants reveal patterns
782 of germline mutation rate heterogeneity in humans. *Nature communications*
783 **9**(1): 3753.
- 784 Chen X, Chen Z, Chen H, Su Z, Yang J, Lin F, Shi S, He X. 2012. Nucleosomes
785 suppress spontaneous mutations base-specifically in eukaryotes. *Science*
786 **335**(6073): 1235-1238.
- 787 Drillon G, Audit B, Argoul F, Arneodo A. 2016. Evidence of selection for an accessible
788 nucleosomal array in human. *BMC genomics* **17**: 526.
- 789 ENCODE Consortium. 2012. An integrated encyclopedia of DNA elements in the
790 human genome. *Nature* **489**(7414): 57-74.
- 791 Englander EW, Howard BH. 1995. Nucleosome positioning by human Alu elements in
792 chromatin. *The Journal of biological chemistry* **270**(17): 10091-10096.
- 793 Englander EW, Wolffe AP, Howard BH. 1993. Nucleosome interactions with a human
794 Alu element. Transcriptional repression and effects of template methylation.
795 *The Journal of biological chemistry* **268**(26): 19565-19573.
- 796 Francioli LC, Polak PP, Koren A, Menelaou A, Chun S, Renkens I, Genome of the
797 Netherlands C, van Duijn CM, Swertz M, Wijmenga C et al. 2015. Genome-
798 wide patterns and properties of de novo mutations in humans. *Nature*
799 *genetics* **47**(7): 822-826.
- 800 Frigola J, Sabarinathan R, Mularoni L, Muinos F, Gonzalez-Perez A, Lopez-Bigas N.
801 2017. Reduced mutation rate in exons due to differential mismatch repair.
802 *Nature genetics* **49**(12): 1684-1692.
- 803 Gaffney DJ, McVicker G, Pai AA, Fondufe-Mittendorf YN, Lewellen N, Michelini K,
804 Widom J, Gilad Y, Pritchard JK. 2012. Controls of nucleosome positioning in
805 the human genome. *PLoS genetics* **8**(11): e1003036.
- 806 Gelman A, Jakulin A, Pittau MG, Su Y-S. 2008. A weakly informative default prior
807 distribution for logistic and other regression models. *The Annals of Applied*
808 *Statistics* **2**(4): 1360-1383.
- 809 Genome of the Netherlands C. 2014. Whole-genome sequence variation, population
810 structure and demographic history of the Dutch population. *Nature genetics*
811 **46**(8): 818-825.
- 812 Giordano J, Ge Y, Gelfand Y, Abrusan G, Benson G, Warburton PE. 2007.
813 Evolutionary history of mammalian transposons determined by genome-wide
814 defragmentation. *PLoS computational biology* **3**(7): e137.
- 815 Guo J, Grow EJ, Yi C, Mlcochova H, Maher GJ, Lindskog C, Murphy PJ, Wike CL,
816 Carrell DT, Goriely A et al. 2017. Chromatin and Single-Cell RNA-Seq Profiling
817 Reveal Dynamic Signaling and Metabolic Transitions during Human
818 Spermatogonial Stem Cell Development. *Cell stem cell* **21**(4): 533-546 e536.

- 819 Harpak A, Bhaskar A, Pritchard JK. 2016. Mutation Rate Variation is a Primary
820 Determinant of the Distribution of Allele Frequencies in Humans. *PLoS*
821 *genetics* **12**(12): e1006489.
- 822 Harris K, Nielsen R. 2014. Error-prone polymerase activity causes multinucleotide
823 mutations in humans. *Genome research* **24**(9): 1445-1454.
- 824 Heger A, Webber C, Goodson M, Ponting CP, Lunter G. 2013. GAT: a simulation
825 framework for testing the association of genomic intervals. *Bioinformatics*
826 **29**(16): 2046-2048.
- 827 Hodgkinson A, Eyre-Walker A. 2011. Variation in the mutation rate across
828 mammalian genomes. *Nature reviews Genetics* **12**(11): 756-766.
- 829 International HapMap Consortium Frazer KA Ballinger DG Cox DR Hinds DA Stuve LL
830 Gibbs RA Belmont JW Boudreau A Hardenbol P et al. 2007. A second
831 generation human haplotype map of over 3.1 million SNPs. *Nature*
832 **449**(7164): 851-861.
- 833 Jonsson H, Sulem P, Kehr B, Kristmundsdottir S, Zink F, Hjartarson E, Hardarson MT,
834 Hjorleifsson KE, Eggertsson HP, Gudjonsson SA et al. 2017. Parental influence
835 on human germline de novo mutations in 1,548 trios from Iceland. *Nature*
836 **549**(7673): 519-522.
- 837 Lee H, Schatz MC. 2012. Genomic dark matter: the reliability of short read mapping
838 illustrated by the genome mappability score. *Bioinformatics* **28**(16): 2097-
839 2105.
- 840 Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, O'Donnell-Luria
841 AH, Ware JS, Hill AJ, Cummings BB et al. 2016. Analysis of protein-coding
842 genetic variation in 60,706 humans. *Nature* **536**(7616): 285-291.
- 843 Li F, Tian L, Gu L, Li GM. 2009. Evidence that nucleosomes inhibit mismatch repair in
844 eukaryotic cells. *The Journal of biological chemistry* **284**(48): 33056-33061.
- 845 Lujan SA, Clausen AR, Clark AB, MacAlpine HK, MacAlpine DM, Malc EP, Mieczkowski
846 PA, Burkholder AB, Fargo DC, Gordenin DA et al. 2014. Heterogeneous
847 polymerase fidelity and mismatch repair bias genome variation and
848 composition. *Genome research* **24**(11): 1751-1764.
- 849 Makova KD, Hardison RC. 2015. The effects of chromatin organization on variation in
850 mutation rates in the genome. *Nature reviews Genetics* **16**(4): 213-223.
- 851 Michaelson JJ, Shi Y, Gujral M, Zheng H, Malhotra D, Jin X, Jian M, Liu G, Greer D,
852 Bhandari A et al. 2012. Whole-genome sequencing in autism identifies hot
853 spots for de novo germline mutation. *Cell* **151**(7): 1431-1442.
- 854 Perera D, Poulos RC, Shah A, Beck D, Pimanda JE, Wong JW. 2016. Differential DNA
855 repair underlies mutation hotspots at active promoters in cancer genomes.
856 *Nature* **532**(7598): 259-263.
- 857 Petryk N, Kahli M, d'Aubenton-Carafa Y, Jaszczyszyn Y, Shen Y, Silvain M, Thermes C,
858 Chen CL, Hyrien O. 2016. Replication landscape of the human genome.
859 *Nature communications* **7**: 10208.
- 860 Pich O, Muinos F, Sabarinathan R, Reyes-Salazar I, Gonzalez-Perez A, Lopez-Bigas N.
861 2018. Somatic and Germline Mutation Periodicity Follow the Orientation of the
862 DNA Minor Groove around Nucleosomes. *Cell* **175**(4): 1074-1087 e1018.
- 863 Prendergast JG, Semple CA. 2011. Widespread signatures of recent selection linked
864 to nucleosome positioning in the human lineage. *Genome research* **21**(11):
865 1777-1787.
- 866 Ramirez F, Dundar F, Diehl S, Gruning BA, Manke T. 2014. deepTools: a flexible
867 platform for exploring deep-sequencing data. *Nucleic acids research* **42**(Web
868 Server issue): W187-191.

- 869 Reijns MAM, Kemp H, Ding J, de Proce SM, Jackson AP, Taylor MS. 2015. Lagging-
870 strand replication shapes the mutational landscape of the genome. *Nature*
871 **518**(7540): 502-506.
- 872 Rodgers K, McVey M. 2016. Error-Prone Repair of DNA Double-Strand Breaks.
873 *Journal of cellular physiology* **231**(1): 15-24.
- 874 Roy S, Tomaszowski KH, Luzwick JW, Park S, Li J, Murphy M, Schlacher K. 2018. p53
875 orchestrates DNA replication restart homeostasis by suppressing mutagenic
876 RAD52 and POLtheta pathways. *eLife* **7**.
- 877 Sabarinathan R, Mularoni L, Deu-Pons J, Gonzalez-Perez A, Lopez-Bigas N. 2016.
878 Nucleotide excision repair is impaired by binding of transcription factors to
879 DNA. *Nature* **532**(7598): 264-267.
- 880 Salih F, Salih B, Kogan S, Trifonov EN. 2008. Epigenetic nucleosomes: Alu sequences
881 and CG as nucleosome positioning element. *Journal of biomolecular structure*
882 *& dynamics* **26**(1): 9-16.
- 883 Sasaki S, Mello CC, Shimada A, Nakatani Y, Hashimoto S, Ogawa M, Matsushima K,
884 Gu SG, Kasahara M, Ahsan B et al. 2009. Chromatin-associated periodicity in
885 genetic variation downstream of transcriptional start sites. *Science*
886 **323**(5912): 401-404.
- 887 Schuster-Bockler B, Lehner B. 2012. Chromatin organization is a major influence on
888 regional mutation rates in human cancer cells. *Nature* **488**(7412): 504-507.
- 889 Segurel L, Wyman MJ, Przeworski M. 2014. Determinants of mutation rate variation
890 in the human germline. *Annual review of genomics and human genetics* **15**:
891 47-70.
- 892 Seplyarskiy VB, Akkuratov EE, Akkuratova N, Andrianova MA, Nikolaev SI, Bazykin
893 GA, Adameyko I, Sunyaev SR. 2018. Error-prone bypass of DNA lesions
894 during lagging-strand replication is a common source of germline and cancer
895 mutations. *Nature genetics*.
- 896 Seplyarskiy VB, Andrianova MA, Bazykin GA. 2017. APOBEC3A/B-induced
897 mutagenesis is responsible for 20% of heritable mutations in the TpCpW
898 context. *Genome research* **27**(2): 175-184.
- 899 Shlien A, Campbell BB, de Borja R, Alexandrov LB, Merico D, Wedge D, Van Loo P,
900 Tarpey PS, Coupland P, Behjati S et al. 2015. Combined hereditary and
901 somatic mutations of replication error repair genes result in rapid onset of
902 ultra-hypermuted cancers. *Nature genetics* **47**(3): 257-262.
- 903 Slotkin RK, Martienssen R. 2007. Transposable elements and the epigenetic
904 regulation of the genome. *Nature reviews Genetics* **8**(4): 272-285.
- 905 Smerdon MJ. 1991. DNA repair and the role of chromatin structure. *Current opinion*
906 *in cell biology* **3**(3): 422-428.
- 907 Smith TCA, Arndt PF, Eyre-Walker A. 2018. Large scale variation in the rate of germ-
908 line de novo mutation, base composition, divergence and diversity in humans.
909 *PLoS genetics* **14**(3): e1007254.
- 910 Stamatoyannopoulos JA, Adzhubei I, Thurman RE, Kryukov GV, Mirkin SM, Sunyaev
911 SR. 2009. Human mutation rate associated with DNA replication timing.
912 *Nature genetics* **41**(4): 393-395.
- 913 Supek F, Lehner B. 2015. Differential DNA mismatch repair underlies mutation rate
914 variation across the human genome. *Nature* **521**(7550): 81-84.
- 915 Tanaka Y, Yamashita R, Suzuki Y, Nakai K. 2010. Effects of Alu elements on global
916 nucleosome positioning in the human genome. *BMC genomics* **11**: 309.
- 917 Tempel S. 2012. Using and understanding RepeatMasker. *Methods in molecular*
918 *biology* **859**: 29-51.

- 919 Terekhanova NV, Seplyarskiy VB, Soldatov RA, Bazykin GA. 2017. Evolution of Local
920 Mutation Rate and Its Determinants. *Molecular biology and evolution* **34**(5):
921 1100-1109.
- 922 Tolstorukov MY, Choudhary V, Olson WK, Zhurkin VB, Park PJ. 2008. nuScore: a
923 web-interface for nucleosome positioning predictions. *Bioinformatics* **24**(12):
924 1456-1458.
- 925 Tolstorukov MY, Volfovsky N, Stephens RM, Park PJ. 2011. Impact of chromatin
926 structure on sequence variability in the human genome. *Nature structural &
927 molecular biology* **18**(4): 510-515.
- 928 Tomasetti C, Li L, Vogelstein B. 2017. Stem cell divisions, somatic mutations, cancer
929 etiology, and cancer prevention. *Science* **355**(6331): 1330-1334.
- 930 Tomasetti C, Vogelstein B. 2015. Variation in cancer risk among tissues can be
931 explained by the number of stem cell divisions. *Science* **347**(6217): 78-81.
- 932 Tubbs A, Sridharan S, van Wietmarschen N, Maman Y, Callen E, Stanlie A, Wu W,
933 Wu X, Day A, Wong N et al. 2018. Dual Roles of Poly(dA:dT) Tracts in
934 Replication Initiation and Fork Collapse. *Cell* **174**(5): 1127-1142 e1119.
- 935 Turner TN, Coe BP, Dickel DE, Hoekzema K, Nelson BJ, Zody MC, Kronenberg ZN,
936 Hormozdiari F, Raja A, Pennacchio LA et al. 2017a. Genomic Patterns of De
937 Novo Mutation in Simplex Autism. *Cell* **171**(3): 710-722 e712.
- 938 Turner TN, Hormozdiari F, Duyzend MH, McClymont SA, Hook PW, Iossifov I, Raja A,
939 Baker C, Hoekzema K, Stessman HA et al. 2016. Genome Sequencing of
940 Autism-Affected Families Reveals Disruption of Putative Noncoding Regulatory
941 DNA. *American journal of human genetics* **98**(1): 58-74.
- 942 Turner TN, Yi Q, Krumm N, Huddleston J, Hoekzema K, HA FS, Doebley AL, Bernier
943 RA, Nickerson DA, Eichler EE. 2017b. denovo-db: a compendium of human de
944 novo variants. *Nucleic acids research* **45**(D1): D804-D811.
- 945 Veltman JA, Brunner HG. 2012. De novo mutations in human genetic disease. *Nature
946 reviews Genetics* **13**(8): 565-575.
- 947 Warnecke T, Becker EA, Facciotti MT, Nislow C, Lehner B. 2013. Conserved
948 substitution patterns around nucleosome footprints in eukaryotes and
949 Archaea derive from frequent nucleosome repositioning through evolution.
950 *PLoS computational biology* **9**(11): e1003373.
- 951 Werling DM, Brand H, An JY, Stone MR, Zhu L, Glessner JT, Collins RL, Dong S, Layer
952 RM, Markenscoff-Papadimitriou E et al. 2018. An analytical framework for
953 whole-genome sequence association studies and its implications for autism
954 spectrum disorder. *Nature genetics* **50**(5): 727-736.
- 955 West JA, Cook A, Alver BH, Stadtfeld M, Deaton AM, Hochedlinger K, Park PJ,
956 Tolstorukov MY, Kingston RE. 2014. Nucleosomal occupancy changes locally
957 over key regulatory regions during cell differentiation and reprogramming.
958 *Nature communications* **5**: 4719.
- 959 Winter DR, Song L, Mukherjee S, Furey TS, Crawford GE. 2013. DNase-seq predicts
960 regions of rotational nucleosome stability across diverse human cell types.
961 *Genome research* **23**(7): 1118-1129.
- 962 Yuen R, Merico D, Bookman M, J LH, Thiruvahindrapuram B, Patel RV, Whitney J,
963 Deflaux N, Bingham J, Wang Z et al. 2017. Whole genome sequencing
964 resource identifies 18 new candidate genes for autism spectrum disorder.
965 *Nature neuroscience* **20**(4): 602-611.
- 966 Yuen RK, Merico D, Cao H, Pellecchia G, Alipanahi B, Thiruvahindrapuram B, Tong X,
967 Sun Y, Cao D, Zhang T et al. 2016. Genome-wide characteristics of de novo
968 mutations in autism. *NPJ genomic medicine* **1**: 160271-1602710.

969

970 **Data availability**

971 All the analyses in this study were based on published datasets.

972 **Acknowledgments**

973 We are grateful to Tobias Warnecke, John Diffley, Anob Chakrabarti and Sara
974 Rohban for insightful comments. We thank Peter Van Loo, Jonas Demeulemeester
975 and Maxime Tarabichi for assistance in accessing the PCAWG genomic data. We
976 appreciate obtaining access to the *de novo* mutation data on SFARI Base. This work
977 is supported by the Francis Crick Institute which receives its core funding from
978 Cancer Research UK (FC001110), the UK Medical Research Council (FC001110),
979 and the Wellcome Trust (FC001110) (N.M.L.). N.M.L. is also supported by a
980 Wellcome Trust Investigator Award and core funding from the Okinawa Institute of
981 Science & Technology. C.L. is funded by an EMBO long-term postdoctoral fellowship
982 (ALTF 1499-2016).

983 **Author contributions**

984 C.L. conceived the project, performed the analyses and drafted the manuscript;
985 N.M.L. supervised the project and co-wrote the manuscript.

986 **Competing financial interests**

987 The authors declare no competing financial interests.

988

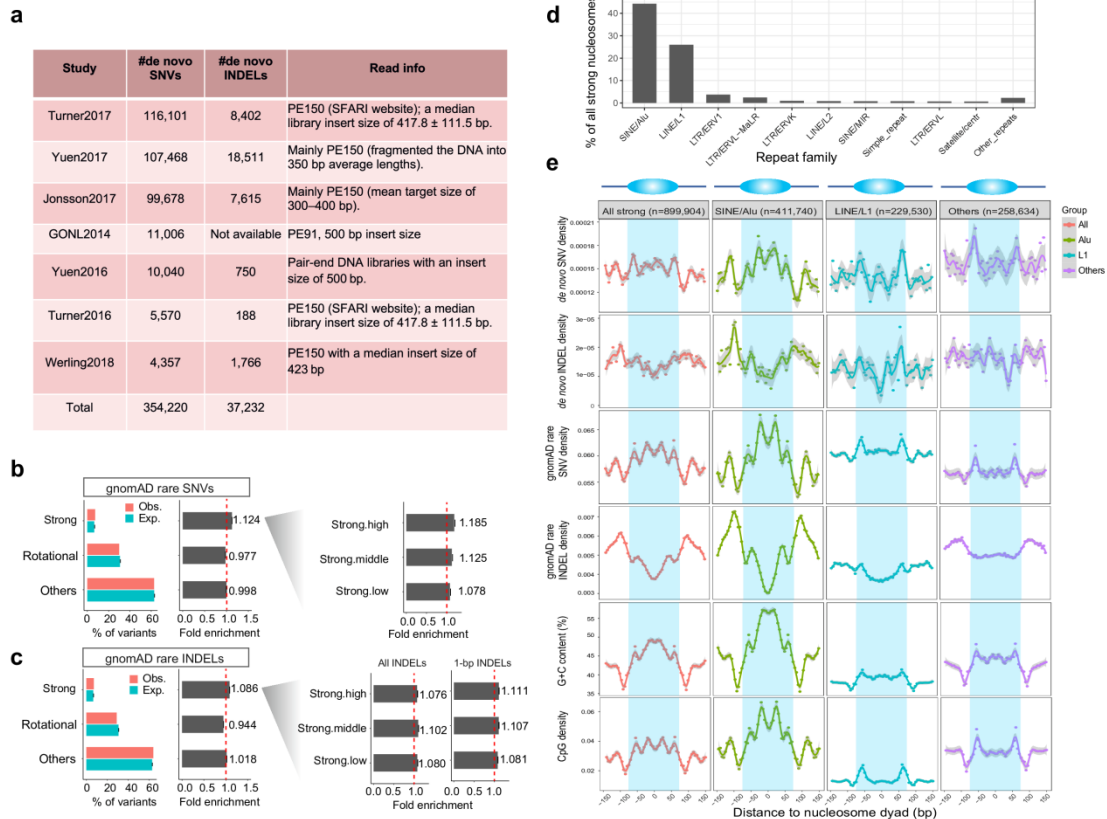
989 **Supplementary Tables and Figures**

990

991 **Supplementary Table 1 Coefficients of variables and other information from the**
992 **full regression models for different mutation types (in a separate Excel file).**

993 Note that for each of the categorical variables, the first category was used by the
994 regression model as reference category (other categories were compared with the
995 reference category) and thus there is no coefficient for that category. The statistics
996 (4th column) and p-values (5th column) in the table were from Wald tests defaultly
997 produced by 'bayesglm' (shown for reference), which are different from the likelihood
998 ratio test-based p-values and were not used in our discussion.

999

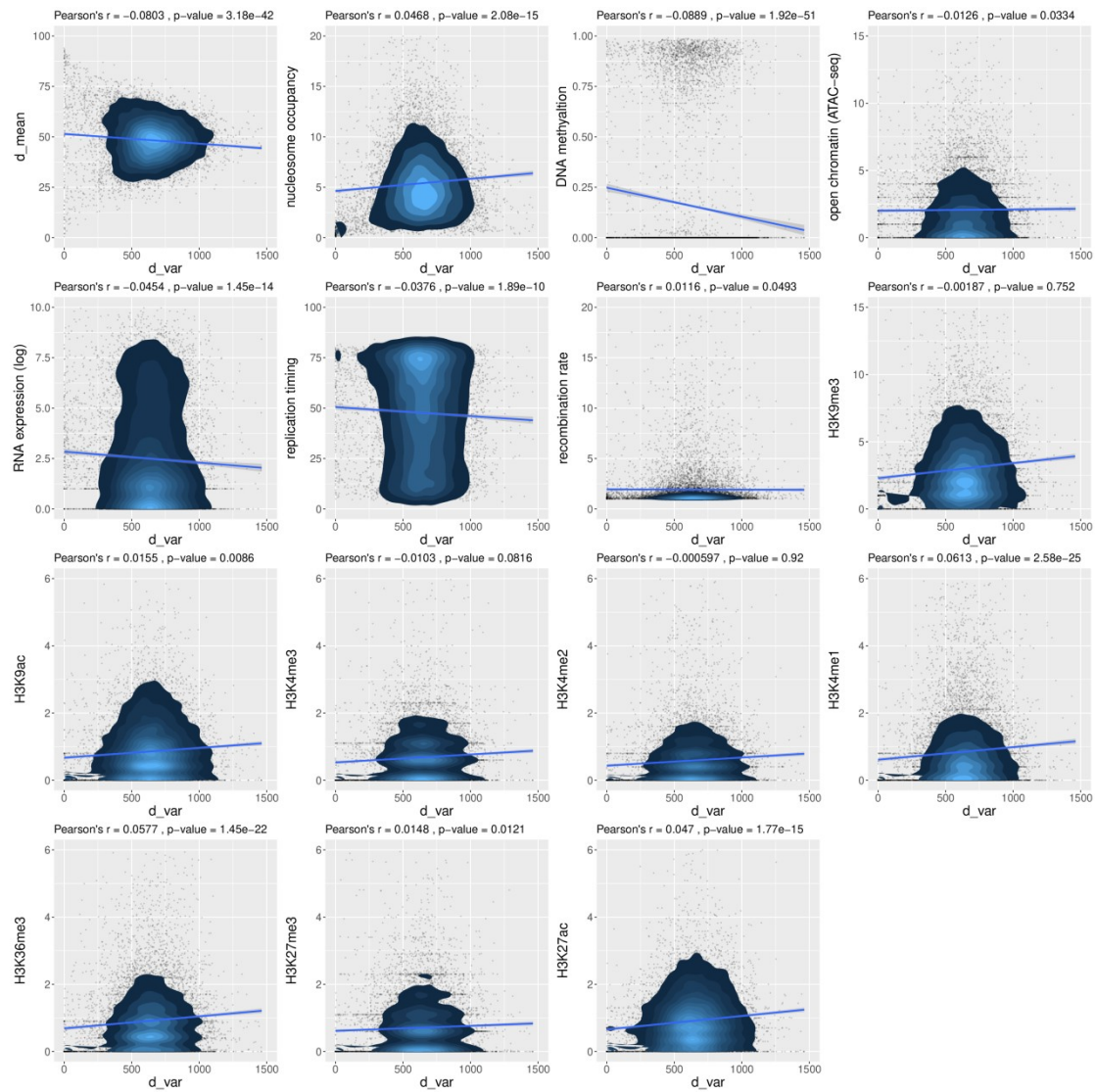


1000

1001 **Supplementary Figure 1 Mutations in different nucleosome contexts.** (a)
 1002 Information of the *de novo* mutation datasets from seven studies used in analysis. (b)
 1003 Fold enrichment/depletion of gnomAD extremely rare SNVs in different nucleosome
 1004 contexts. ‘Strong’, translationally stable positioning; ‘Rotational’, rotationally but not
 1005 translationally stable positioning; ‘Others’, the remaining genomic regions. On the left
 1006 is the fold enrichment for three subgroups of strong nucleosomes with different
 1007 stabilities. Error bars depict 95% confidence intervals. (c) Fold enrichment/depletion
 1008 of gnomAD INDELS in different nucleosome contexts. When using all INDELS the
 1009 ‘strong.high’ group does not have a higher mutation rate than other two groups, but if
 1010 using the 1-bp INDELS ‘strong.high’ does have the highest mutation rate among the
 1011 three groups. We speculated that there may be more false negatives of longer
 1012 INDELS in the ‘strong.high’ group. (d) Top 10 repeat families that are associated with
 1013 strong nucleosomes. (e) Meta-profiles of SNV/INDEL densities (*de novo* or extremely
 1014 rare variants) around all strong nucleosomes, or in different repeat-associated
 1015 subgroups. At the bottom are the G+C content and CpG content profiles.

1016

1017



1018

1019 **Supplementary Figure 2 Correlation analysis between nucleosome positioning**
1020 **stability (d_{var}) and other factors.** On the top of each panel are the Pearson's
1021 correlation coefficients and the corresponding p-values.

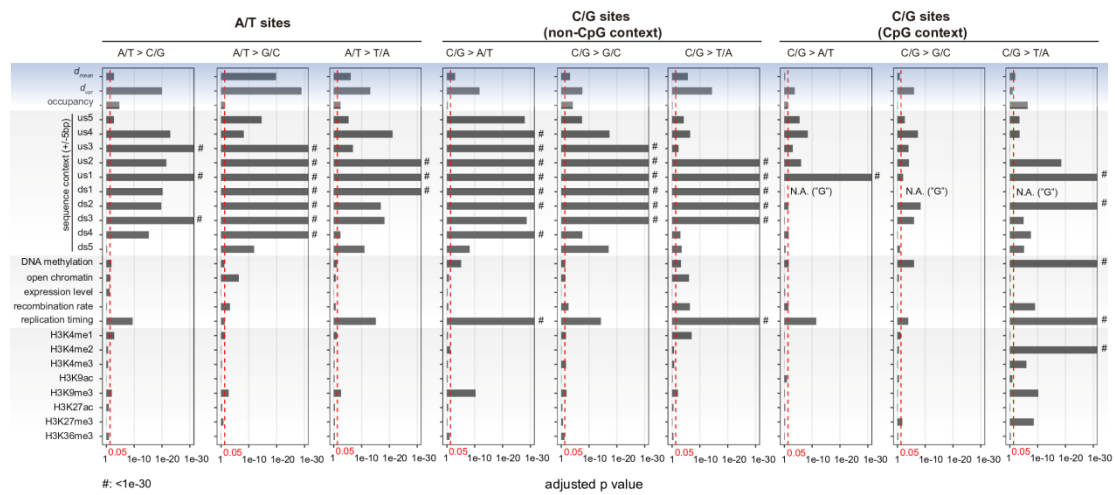
1022

1023

1024

1025

1026



1027

1028 **Supplementary Figure 3 Results of statistical tests for nine individual SNV**
 1029 **mutation types.** C/G sites in non-CpG contexts and C/G sites in CpG contexts were
 1030 tested separately. The red vertical lines represent the significance cut-off (0.05) for
 1031 the adjusted p values (Benjamini–Hochberg correction). ‘us’, upstream; ‘ds’,
 1032 downstream. ‘#’ means adjusted p < 1e-30.

1033

1034

1035

1036

1037

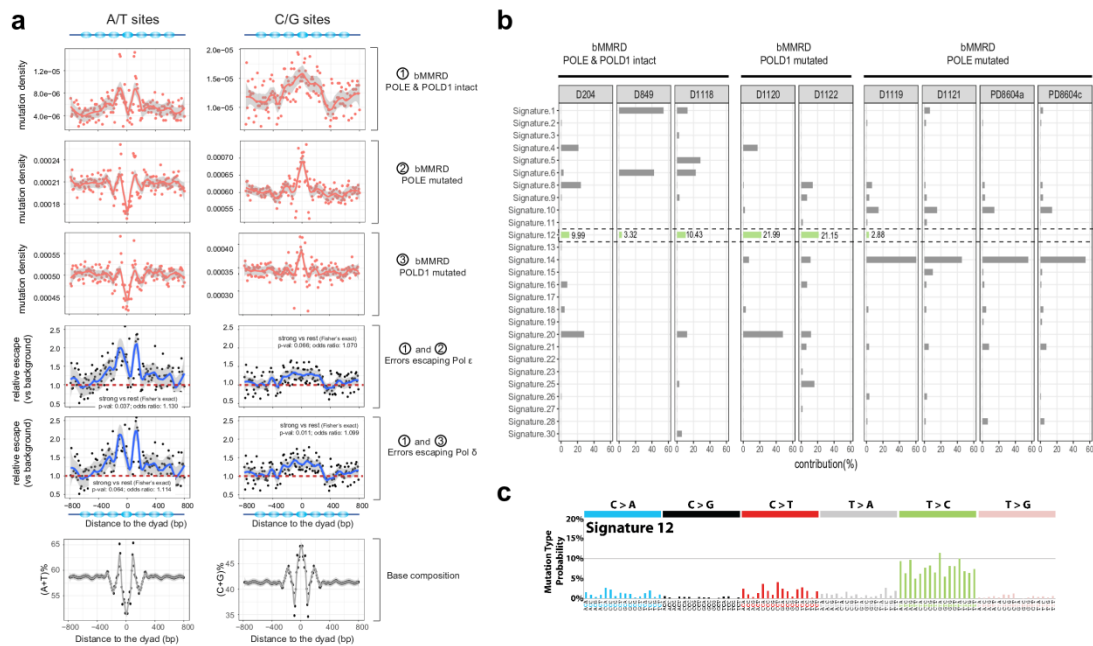


1038

1039 **Supplementary Figure 4 Results of statistical tests when considering two-way**
 1040 **interactions of adjacent nucleotides, 7-mer mutability estimates from Carlson**
 1041 **et al. and repeat status. (a) Adding the two-way interactions for ±5 nucleotides in**
 1042 **the regression models. (b) Adding the 7-mer mutability estimates from Carlson et al.**
 1043 **as predictors in the regression models. (c) Adding repeat status as a predictor in the**
 1044 **regression models. (d) Running regression models for regions associated with**
 1045 **different repeat contexts separately. We tested SNVs at A/T sites, C/G sites in non-**
 1046 **CpG context and C/G sites in CpG context separately. The red vertical lines**
 1047 **represent the significance cut-off (0.05) for the adjusted p values (Benjamini–**
 1048 **Hochberg correction). ‘us’, upstream; ‘ds’, downstream. ‘#’ means adjusted p < 1e-30.**

1049

1050



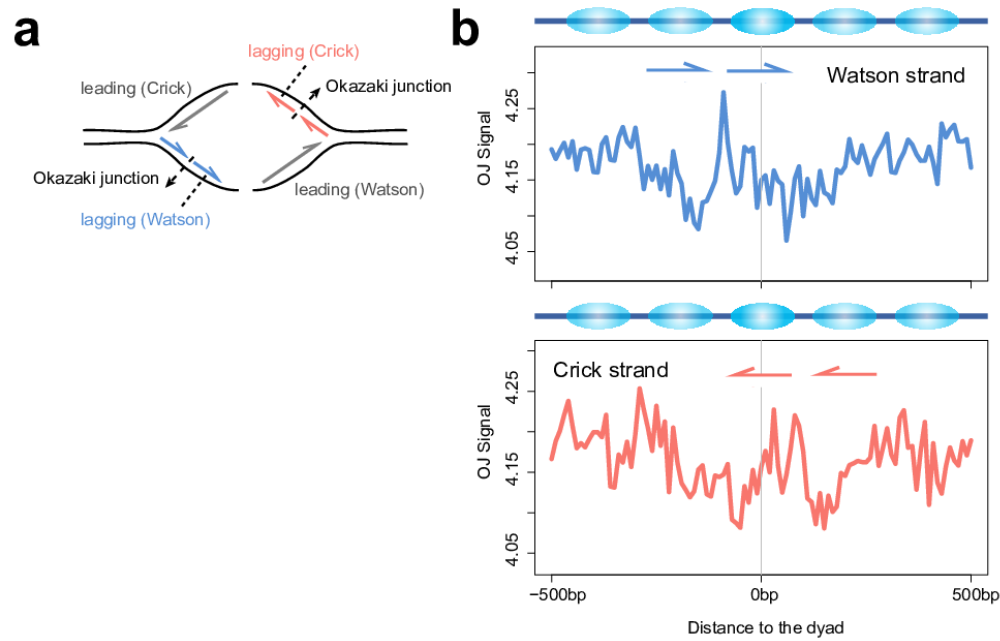
1051

1052 **Supplementary Figure 5 Analysis of related mutational processes using**
 1053 **bMMRD data. (a)** Mutation profiles around strong nucleosomes for bMMRD cancer
 1054 genomes and the estimated relative escape ratios of Pol ϵ or Pol δ , for mutations at
 1055 A/T sites and C/G sites respectively. Fisher's exact test was used for testing the
 1056 association of strong-nucleosome regions (dyad \pm 95bp) with differential polymerase
 1057 performance. (b) Comparison of the contribution of COSMIC mutational signatures
 1058 predicted by MutationalPatterns in different bMMRD genomes. Highlighted is
 1059 Signature 12, which shows a particularly high contribution in POLD1-mutated bMMRD
 1060 samples. (c) the tri-nucleotide mutational profile of Signature 12, obtained from
 1061 COSMIC website.

1062

1063

1064

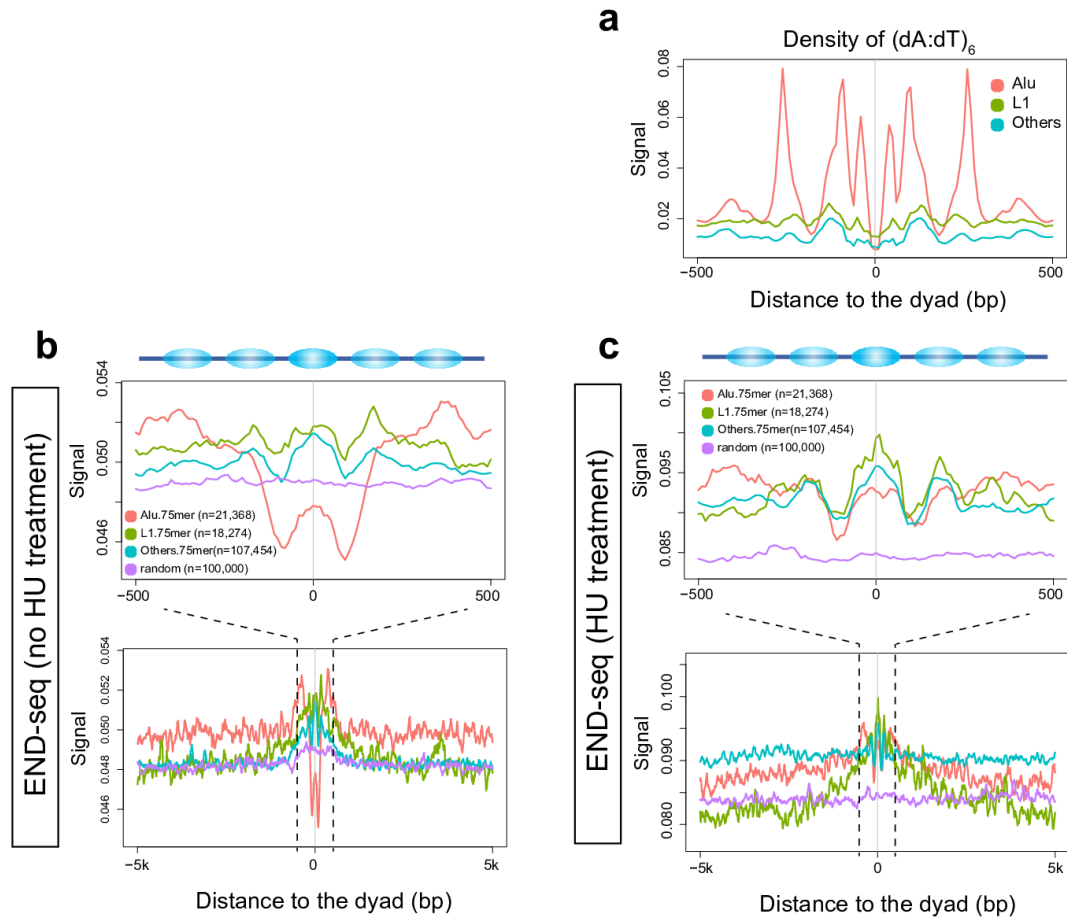


1065

1066 **Supplementary Figure 6 Analysis with OK-seq data.** (a) Schematic illustrating
1067 replication strands and Okazaki junctions (OJs). (b) Meta-profile of the density of
1068 Okazaki junctions inferred from alignments of OK-seq reads around strong
1069 nucleosomes (high-mappability). OJ signals for Watson strand and Crick strand
1070 were plotted separately. Replication directions of Okazaki fragments are shown by arrows.

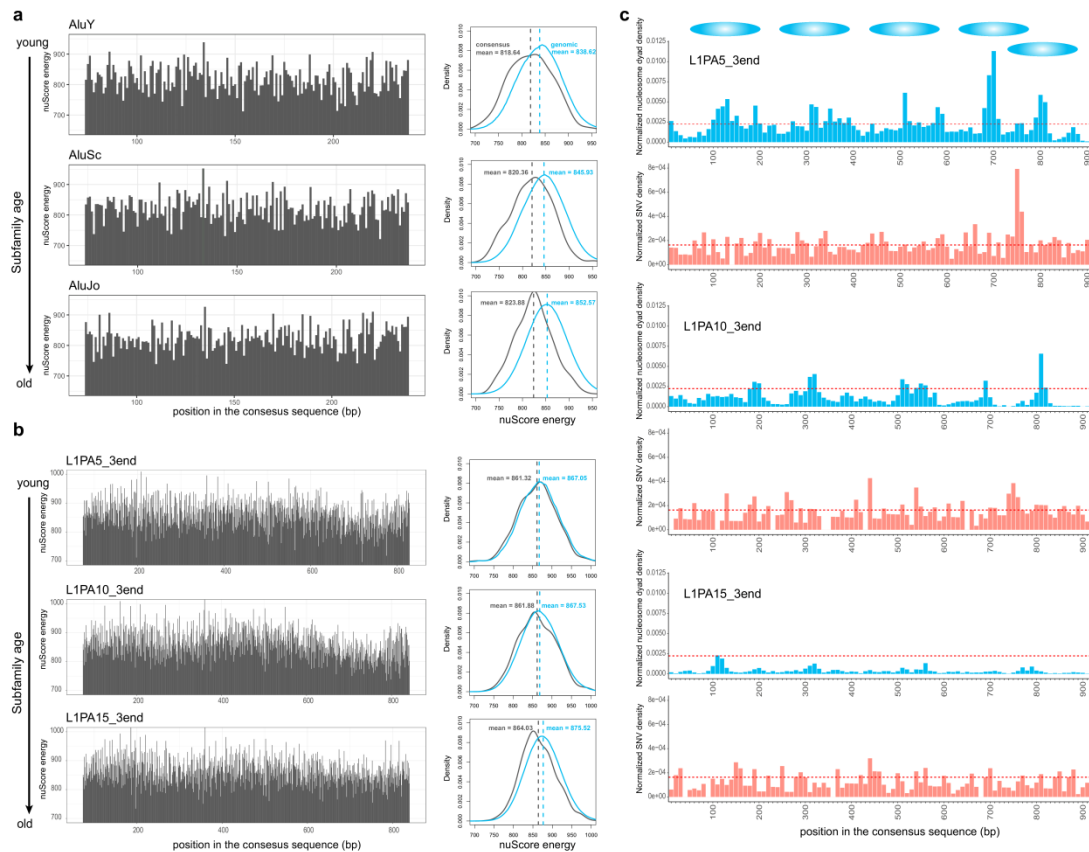
1071

1072



1073

1074 **Supplementary Figure 7 Analysis related to the DSBs around strong**
1075 **nucleosomes.** (a) Density of poly(dA:dT)₆ motifs) around strong nucleosomes. (b-c) Signal of DSBs based on the END-seq
1076 data around strong nucleosomes associated with different repeat elements. Only the strong nucleosomes of high 75-mer mappability within ±500bp were considered.
1077 Numbers of usable strong nucleosomes for each group are given in the brackets. HU
1078 (hydroxyurea) is a replicative stress-inducing agent.
1079
1080



1081

1082 **Supplementary Figure 8 Additional analysis about repeat subfamily ages and**
 1083 **strong nucleosomes. (a)** nuScore-estimated per-base nucleosome deformation
 1084 energies along three Alu subfamily consensus sequences. On the right are the
 1085 comparisons of deformation energy distributions of the consensus sequences
 1086 (ancestral states) and those of current genomic regions for the three subfamilies
 1087 respectively. The deformation energy profiles of the consensus sequences are
 1088 similar, but the average deformation energies increase over time, with older Alu
 1089 subfamilies displaying larger differences relative to the consensus. **(b)** Similar to (a),
 1090 but for three example L1 subfamilies. **(c)** Barplots for normalized densities of strong
 1091 nucleosome dyads and *de novo* SNVs along the consensus sequences of three L1
 1092 subfamilies, using 10-bp bins. Several loci that are enriched for dyads of strong
 1093 nucleosomes are shown on the top with ellipses. The red dash lines represent the
 1094 average densities for the L1PA5 subfamily. The densities of strong nucleosome
 1095 dyads and *de novo* SNVs appear to decrease over evolutionary time.

1096