

23 **ABSTRACT**

24 Transitions of cytosine to thymine in CpG dinucleotides are the most frequent type of mutations
25 observed in cancer. This increased mutability is commonly explained by the presence of 5-
26 methylcytosine (5mC) and its spontaneous hydrolytic deamination into thymine. Here, we describe
27 observations that question whether spontaneous deamination alone causes the elevated
28 mutagenicity of 5mC. Tumours with somatic mutations in DNA mismatch-repair genes or in the
29 proofreading domain of DNA polymerase ϵ (Pol ϵ) exhibit more 5mC to T transitions than would be
30 expected, given the kinetics of hydrolytic deamination. This enrichment is asymmetrical around
31 replication origins with a preference for the leading strand template, in particular in methylated
32 cytosines flanked by guanines (GCG). Notably, GCG to GTG mutations also exhibit strand asymmetry
33 in mismatch-repair and Pol ϵ wild-type tumours. Together, these findings suggest that mis-
34 incorporation of A opposite 5mC during replication of the leading strand might be a contributing
35 factor in the mutagenesis of methylated cytosine.

36 **KEYWORDS**

37 Mutagenesis; DNA methylation; DNA Replication; Cancer genomics

38 **1. INTRODUCTION**

39 C to T transitions in a CpG context (CpG>TpG) are the most frequently observed mutations in cancer
40 and genetic disorders [1,2]. Two independent observations link these mutations to 5-methylcytosine
41 (5mC), an epigenetic modification of cytosine. First, most cytosines in CpG dinucleotides are
42 methylated in humans [3]. Moreover, the increased C>T mutagenicity in CpG dinucleotides is
43 present specifically in cytosines that are methylated, compared to unmodified or hydroxymethylated
44 cytosines [4]. Second, it was shown *in vitro* that methylated cytosine (5mC) has a four-fold higher
45 rate of spontaneous deamination than unmodified cytosine [5]. The products of deamination can be
46 repaired by base excision repair (BER). DNA glycosylases involved in BER of T:G mismatches (*TDG*
47 and *MBD4*) excise T from the mismatch, leading to the restoration of C:G [6,7]. Notably, while the
48 deamination of 5mC produces thymine, leading to a T:G mismatch, C and 5-hydroxymethylcytosine

49 (5hmC) deaminate into uracil and 5-hydroxymethyluracil, respectively. Since these bases do not
50 normally occur in DNA, they are potentially more efficiently recognised and replaced by BER [8].
51 Moreover, deamination of 5hmC does not contribute to the steady-state levels of 5hmU in mouse
52 embryonic stem cells, suggesting either infrequent deamination or very fast repair [9]. Failure to
53 correct the T:G mismatch before replication results in a mutation in one daughter cell, due to the
54 semiconservative nature of DNA replication. Thus replication of a T:G mismatch leads to a C:G>T:A
55 mutation.

56 Mutations can also arise through mis-incorporation of bases during cell division. The fidelity of DNA
57 replication relies on proofreading by the major replicative polymerases Pol ϵ and Pol δ , and on post-
58 replicative DNA mismatch-repair (MMR) which removes errors from the newly synthesised DNA
59 strand [10]. Deficiency in any of these protective mechanisms leads to an increase in the number of
60 mutations. In particular, defects in MMR genes lead to “hypermutable” (10^4 - 10^5 mutations per
61 Gbp), and mutations in the proofreading domain of Pol ϵ lead to “ultra-hypermutable”, often
62 exceeding 10^5 mutations per Gbp [11,12]. Moreover, defects in Pol ϵ and Pol δ proofreading cause
63 tumours in mice [13] and germline mutations in *POLE* and *POLD1* (encoding the catalytic subunits of
64 Pol ϵ and δ , respectively) and genes of the MMR pathway predispose to cancer in humans [10].

65 DNA polymerase proofreading and post-replicative MMR (in their canonical, replication-linked
66 functions) are highly unlikely to play a role in the repair of 5mC deamination induced mutations, as
67 they operate *after* parental strands have been separated during replication, at which point a 5mC to
68 T deamination event is indistinguishable from other thymines. Therefore, although the total
69 frequency of mutations due to unrepaired errors introduced during replication increases drastically
70 in polymerase proofreading/MMR deficient samples, it would be expected that the frequency of
71 CpG>TpG mutations should only increase by a small amount.

72 Contrary to this expectation, we observe that the frequency of CpG>TpG mutations in tumours with
73 defective Pol ϵ or MMR is approximately six-fold higher than for other types of mutations. We show
74 that the increased CpG>TpG mutation rate in Pol ϵ or MMR mutant cancers is linked to DNA

75 methylation, has a clear replication strand asymmetry, being enriched on the leading strand, with a
76 preference for a GCG sequence context. We also detect weaker but consistent replication strand
77 asymmetry of GCG>GTG mutations in Pol ϵ and MMR proficient samples. Together, our results
78 suggest that a substantial fraction of C>T mutations at methylated cytosines is independent of
79 spontaneous deamination, instead arising during DNA replication.

80 **2. MATERIALS AND METHODS**

81 **2.1. Somatic mutations**

82 Cancer somatic mutations in 3442 whole-genome sequencing samples (Supplementary Table 1)
83 were obtained from the data portal of The Cancer Genome Atlas (TCGA), the data portal of the
84 International Cancer Genome Consortium (ICGC), and previously published data in peer-review
85 journals [1,12,14–16]. MSI and *POLE*-MUT samples were combined from previous studies [11,12,17].
86 For the TCGA samples, aligned reads of paired tumour and normal samples were downloaded from
87 the UCSC CGHub website under TCGA access request #10140 and somatic variants were called using
88 Strelka (version 1.0.14) [18] with default parameters. Somatic mutations in autosomes only were
89 taken into account.

90 **2.2. DNA modification maps**

91 Maps of cytosine modifications (Supplementary Table 2) were obtained from BS-Seq data sets from
92 the data portals of The Cancer Genome Atlas (TCGA), Roadmap Epigenome, Blueprint, and from
93 previously published data in peer-review journals [19–22] and where needed converted to hg19
94 using liftover tool. For brain, kidney, and prostate maps, raw reads were processed with Trim galore,
95 Bismark[23] and Mark duplicates from Picard tools; and only sites covered with at least 5 reads were
96 taken into account.

97 **2.3. Mutation frequency with respect to modification levels**

98 All cytosines in the CpG context were divided into 10 right-open intervals according to their
99 modification levels (the number of unconverted reads divided by the number of all reads in BS-Seq):

100 [0-0.1), [0.1-0.2), ..., [0.9-1]. In each bin, the frequency of mutations was computed and plotted for
101 each sample. A linear regression model was fitted to the data (function `fitlm` in MatLab) and the
102 offset, slope, and last value, and fold-change from first to last value were measured. When
103 comparing CpG sites with low vs. intermediate vs. high modification levels, the thresholds (0.8 and
104 0.95) were chosen such that the three groups have approximately similar numbers of CpG sites in
105 most tissues.

106 **2.4. Direction of replication**

107 Left- and right-replicating domains were taken from [17]. Each domain (called territory in the
108 original source code and data) is 20kbp wide and annotated with the direction of replication and
109 with replication timing.

110 **2.5. Mutation frequency with respect to the direction of replication**

111 First, transitions between left- and right-replicated domains were computed as in [17]. These
112 transitions represent regions rich for replication origins. We computed the CpG>TpG mutation
113 frequency in the 20kbp domains distant 0 to 1Mbp from the closest left-/right- transition, with
114 respect to the strand (plus=Watson vs. minus=Crick) of the cytosine of the CpG. Template for the
115 leading strand then corresponds to the plus strand in the left direction and minus strand in the right
116 direction and *vice versa* for the lagging strand template. Finally, we annotated all cytosines in a CpG
117 context whether they are on the leading or lagging strand, and computed CpG>TpG mutation
118 frequency for the leading and lagging strand separately. `signtest` was used for evaluating
119 significance of CpG>TpG mutation frequency difference between the two strands.

120 **2.6. Spontaneous deamination estimates**

121 The number of years needed to reach the observed number of C>T mutations in methylated CpGs
122 observed in *POLE*-MUT and MSI samples was based on the spontaneous deamination rate of 5mC in
123 double-stranded DNA ($5.8 \cdot 10^{-13} \text{ s}^{-1}$) reported by Shen et al. [24], the number of seconds in a year
124 (31556736), the observed frequency of GCG>GTG mutations (*i.e.*, GmCG>T/GmCG; for mC with a

125 modification level of at least 0.9) in MSI ($5.133 \cdot 10^{-4}$) and *POLE*-MUT ($1.785 \cdot 10^{-3}$) samples, and
126 computed as:

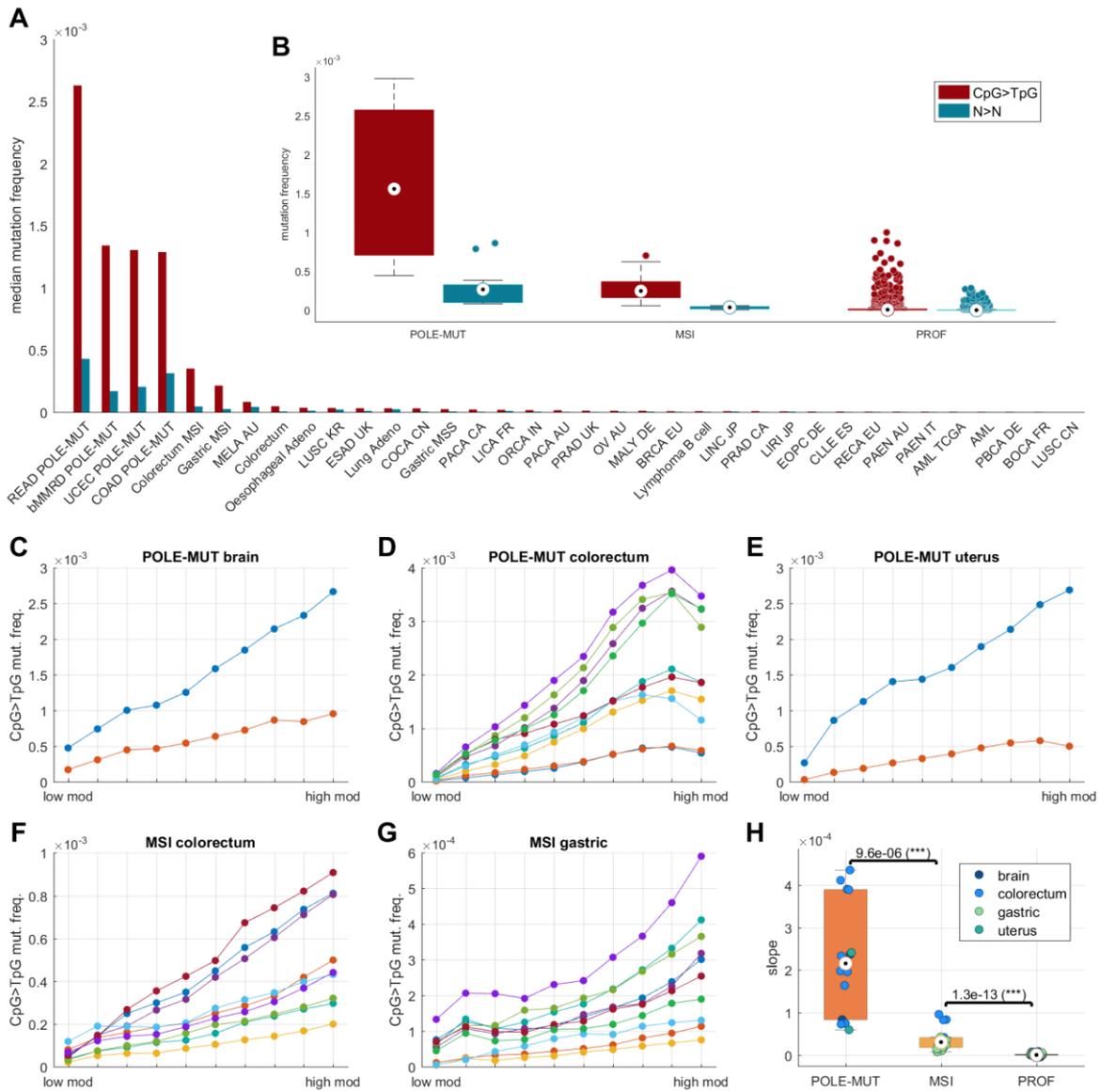
127
$$\text{MSI: } \frac{5.133 \cdot 10^{-4}}{5.8 \cdot 10^{-13} \cdot 31556736} = 28.05 \text{ years}$$

128
$$\text{POLE-MUT: } \frac{1.785 \cdot 10^{-3}}{5.8 \cdot 10^{-13} \cdot 31556736} = 97.53 \text{ years}$$

129

130 3. RESULTS

131 We explored the mutation spectra of 14 tumour samples with a mutation in Pol ϵ (*POLE*-MUT
132 samples), 19 samples with microsatellite-instability (MSI) deficient in MMR, and 3409 other cancer
133 samples (proficient; PROF). The median overall mutation frequency per base was $1.5 \cdot 10^{-6}$ (IQR
134 $0.6 \cdot 10^{-6}$ – $3.5 \cdot 10^{-6}$) in PROF samples, $36.9 \cdot 10^{-6}$ (IQR $18.0 \cdot 10^{-6}$ – $47.4 \cdot 10^{-6}$) in MSI samples, and $267.4 \cdot 10^{-6}$
135 (IQR $99.9 \cdot 10^{-6}$ – $300.5 \cdot 10^{-6}$) in *POLE*-MUT samples (N>N in Fig. 1A–B). In PROF samples, the median
136 CpG>TpG mutation frequency (i.e., the number of CpG>TpG mutations relative to the number of
137 CpGs in the genome) was $7.4 \cdot 10^{-6}$ (IQR $3.7 \cdot 10^{-6}$ – $16.8 \cdot 10^{-6}$), approximately 5 fold higher than the
138 overall mutation frequency (i.e., the number of all mutations relative to the number of all positions
139 in the genome). Notably, the CpG>TpG mutation frequency also increased in MSI and *POLE*-MUT
140 samples, compared to the overall mutation frequency (MSI: median $247.7 \cdot 10^{-6}$ per CpG, IQR
141 $162.7 \cdot 10^{-6}$ – $367.3 \cdot 10^{-6}$; *POLE*-MUT: median $1559.8 \cdot 10^{-6}$ per CpG, IQR $707.9 \cdot 10^{-6}$ – $2574.2 \cdot 10^{-6}$)
142 (CpG>TpG in Fig. 1A–B, Fig. 1-supplement 1). This observation is surprising, since neither MMR nor
143 proofreading during DNA replication by Pol ϵ are thought to be essential for effective repair of
144 deamination induced T:G mismatches [8].



145

146 **Fig. 1: Frequency of C to T mutations in a CpG context is unexpectedly high in *POLE*-MUT and MSI samples and correlates**
 147 **with DNA modification levels. A:** Median CpG>TpG and N>N (overall) mutation frequency in each cancer type separately.
 148 **B:** Distribution of CpG>TpG and N>N mutation frequency in *POLE*-MUT, MSI, and PROF (other) samples. The white circle
 149 with the black dot inside denotes the median. **C-G:** Fraction of mutated CpG sites as a function of modification levels. The
 150 x-axis represents CpG sites grouped into 10 bins by their modification levels (0-0.1, ..., 0.9-1.0). The y-axis represents C>T
 151 mutation frequency in each bin. Individual samples are plotted in different colours. **H:** Distribution of the slope of the linear
 152 relationship between DNA modification levels and CpG>TpG mutation frequency in four tissues (brain, colorectum, gastric,
 153 and uterus). The Wilcoxon ranksum test was used to evaluate differences between the groups (*POLE*-MUT, MSI, and PROF)
 154 of samples. See the distribution of offsets in Fig. 1-supplement 2.

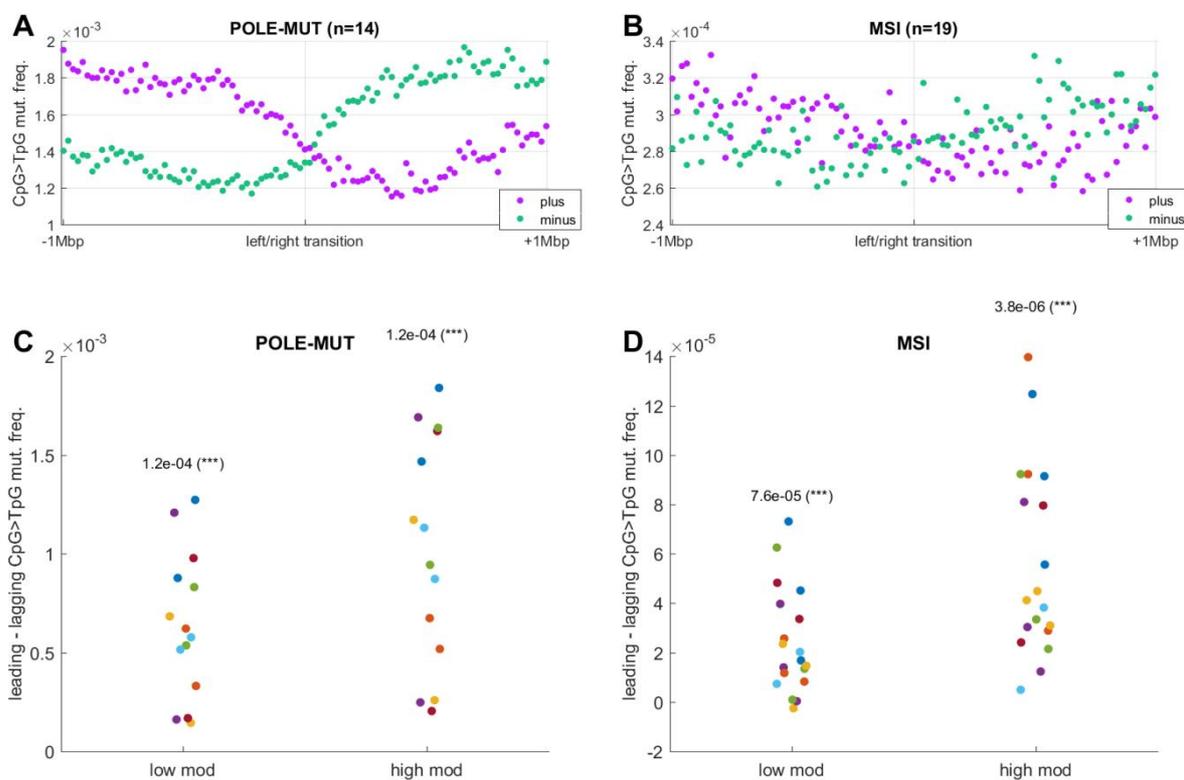
155 We next used bisulfite-sequencing (BS-seq) derived DNA modification maps from normal tissue of
 156 the same organ as each cancer sample to explore whether DNA modifications play a role in the
 157 occurrence of CpG>TpG mutations in *POLE*-MUT and MSI samples. These maps represent levels of

158 both the more frequent 5mC as well as the less frequent 5hmC, since BS-seq alone cannot
159 distinguish between these two modifications. In all *POLE*-MUT and MSI samples, the CpG>TpG
160 mutation frequency was positively correlated with modification levels (Fig. 1C–G). Moreover, the
161 slope of the correlation was significantly higher in *POLE*-MUT than in MSI, and in MSI than in tissue-
162 matched PROF samples (Fig. 1H, 1-supplement 2). These results support the notion that the
163 mechanism responsible for the elevated mutation rate of CpGs in *POLE*-MUT and MSI samples is
164 linked to epigenetic DNA modifications.

165 It is unlikely that Pol ϵ or MMR, through their canonical, replication-linked activity, are used for the
166 repair of deamination-induced T:G mismatches that happened before replication. However, it is
167 possible that their non-canonical, replication unrelated, activity is involved in the repair of
168 deamination induced mismatches. Conversely, the CpG>TpG mutations could be replication related,
169 but independent of spontaneous deamination of 5mC. We therefore explored whether the CpG>TpG
170 mutagenicity in *POLE*-MUT and MSI samples shows any replication-linked characteristics, to
171 distinguish between the potential replication-unrelated repair of spontaneous deamination, and a –
172 yet undescribed – replication-related source of CpG>TpG mutations.

173 In eukaryotic cells, DNA replication is initiated around replication origins (ORI) from where it
174 proceeds in both directions, synthesizing the leading strand continuously and the lagging strand
175 discontinuously. As Pol ϵ is the main leading strand DNA polymerase [25,26], mutations in *POLE*-
176 MUT samples are distributed asymmetrically on the leading and lagging strands [11,17]. MSI samples
177 also display replication strand bias across several types of mutations [17], presumably because MMR
178 is involved in balancing the differences in fidelity of the leading and lagging polymerases [27]. In
179 order to determine whether CpG>TpG mutations in *POLE*-MUT and MSI samples happened during or
180 before replication, we computed the frequency of CpG>TpG mutations on the plus (Watson) and
181 minus (Crick) strand around transitions between left- and right-replicating regions, as defined in
182 [17]. The transitions correspond to regions enriched for replication origins.

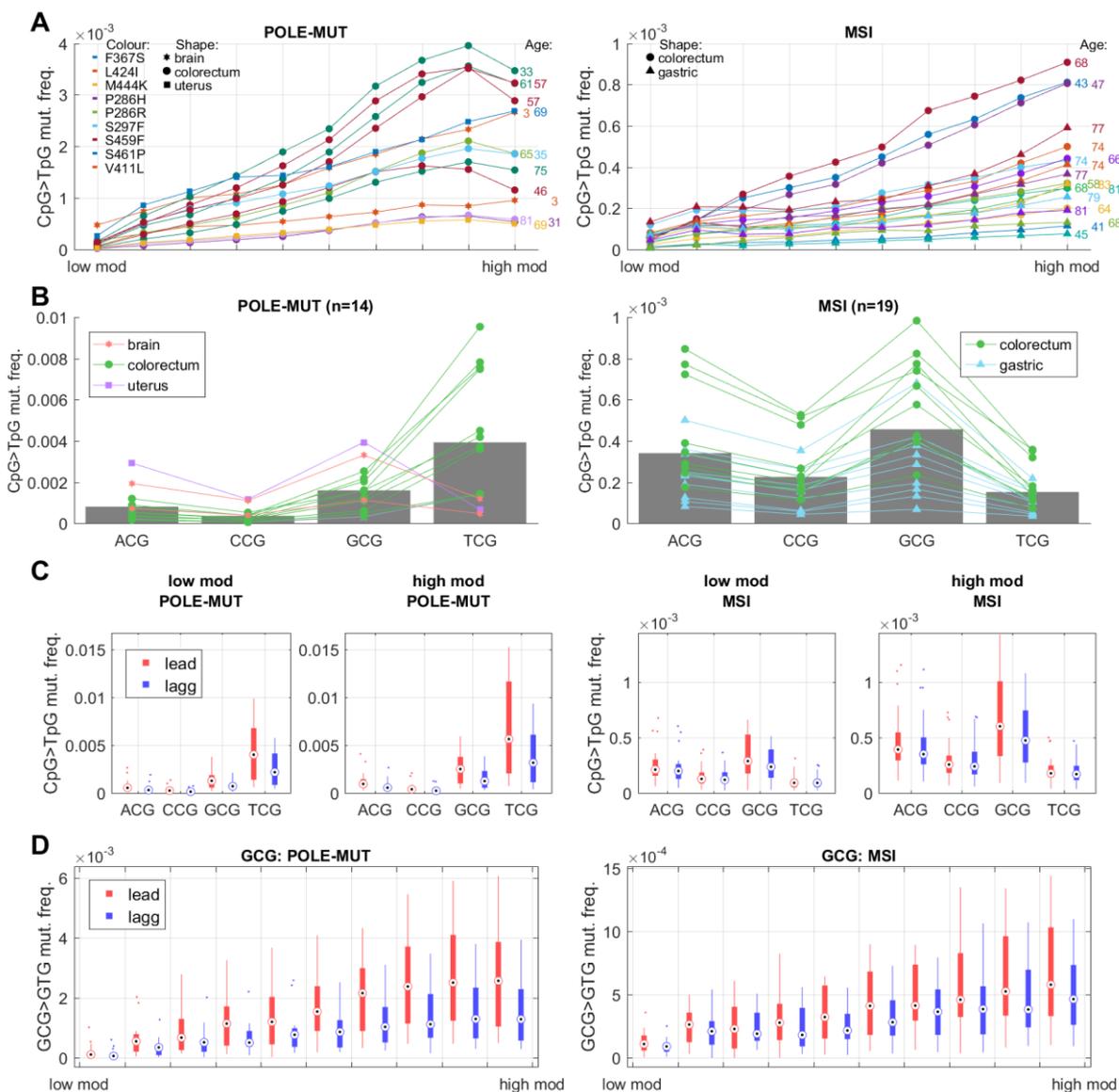
183 In the *POLE*-MUT and MSI samples, we observed a strong enrichment of CpG>TpG mutations on the
 184 leading strand template (plus strand in the left direction, minus strand in the right direction) (Fig. 2).
 185 Moreover, the strand asymmetry was at least as strong or stronger in highly modified CpGs (top
 186 tertile) than in lowly modified CpGs (bottom tertile) (Fig. 2C–D). This effect was furthermore
 187 observed across cancer types and across modification levels (Fig. 2 supplement 1). It thus appears
 188 that DNA repair deficient cells accumulate more CpG>TpG mutations in cytosines that were modified
 189 on the template for the leading strand, suggesting that they are related to replication.



190
 191 **Fig. 2: Frequency of C to T mutations in a CpG context in *POLE*-MUT and MSI samples is higher on the leading strand than**
 192 **on the lagging strand, especially in modified CpG sites. A-B:** Mean CpG>TpG mutation frequency on the plus (Watson) and
 193 minus (Crick) strand around transitions between left- and right-replicating regions. The transitions correspond to regions
 194 enriched for replication origins. The leading strand template corresponds to the plus strand in the left direction and the
 195 minus strand in the right direction, whereas the lagging strand template corresponds to the minus strand in the left
 196 direction and the plus strand in the right direction. **C-D:** Difference in the leading and lagging CpG>TpG mutation
 197 frequency in each sample (signtest was used for evaluating significance between leading and lagging strand).

198 The link between C>T mutagenicity in methylated CpG sites and replication could either be a unique
 199 feature of *POLE*-MUT and MSI samples, or it could be present in all samples, but normally be
 200 suppressed by a combination of Pol ϵ proofreading and MMR. To explore the first option, we tested

201 the observed *POLE* and MMR mutations for signs of a “gain of function” mutation. A range of 9
 202 different variants in the proofreading domain of *POLE* were present in the 14 *POLE*-MUT samples, all
 203 of them showing an increase of CpG>TpG mutations in modified cytosine (Fig. 3A). The positive
 204 correlation of CpG>TpG mutagenicity with methylation seems to be independent of the type of *POLE*
 205 mutation, cancer type or age at diagnosis, and is present in both *POLE*-MUT and MSI samples (Fig
 206 3A). A gain-of-function mutation therefore seems unlikely.

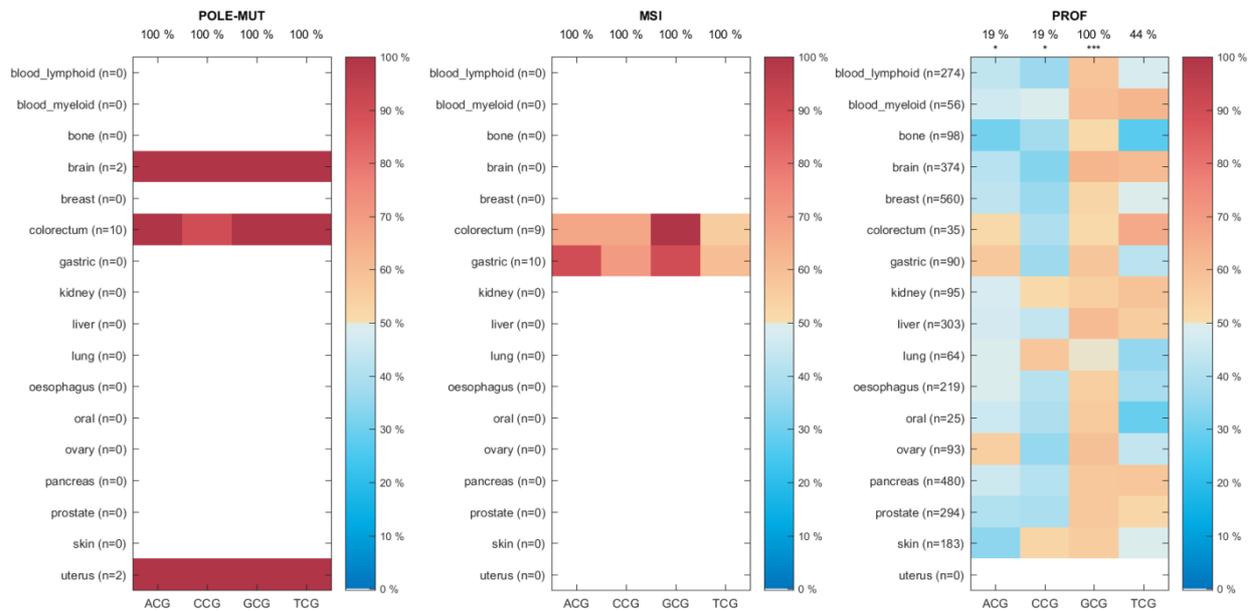


207
 208 **Fig. 3: Increase of C to T mutations in modified cytosine on the leading strand is most consistent in a GCG sequence**
 209 **context in *POLE*-MUT and MSI samples.** **A:** C>T mutation frequency in CpG context binned by the tissue-matched
 210 modification levels (0-0.1, ..., 0.9-1.0). Individual samples are plotted as separate traces. In *POLE*-MUT samples, the colour
 211 represents different variants of the *POLE* mutation. In both *POLE*-MUT and MSI samples, the shape of the marker
 212 represents different tissues. The age at diagnosis is shown next to the last value of the sample. **B:** CpG>TpG mutation

213 frequency stratified by the 5' flanking sequence context. The bars denote mean over samples and individual samples are
214 shown as markers with shape and colour distinguishing the tissue type. C: C>T mutation frequency in CpG sites in the
215 leading and lagging strands, in low mod (≤ 0.8) vs high mod (> 0.95), and stratified by the 5' sequence context: ACG, CCG,
216 GCG, and TCG. D: C>T mutation frequency in GCG context in leading and lagging strand binned by the tissue-matched
217 modification levels (0-0.1, ..., 0.9-1.0).

218 Interestingly, the frequency of C>T mutations was not only affected by the 3' sequence context, but
219 also the 5' base of cytosine. We noticed that, while C>T mutations in a TCG context (TCG>TTG)
220 dominate in colorectal *POLE*-MUT samples, both tissues with MSI samples and all tissues with *POLE*-
221 MUT samples exhibited high levels of C>T mutations in a GCG context (GCG>GTG) (Fig.3B, Fig3-
222 supplement 1). GCG>GTG mutations also showed particularly strong strand asymmetry and
223 correlation with modification levels in all MSI and *POLE*-MUT samples (Fig. 3C, D, 3-supplement 2).

224 Our observations could be explained by a model of CpG>TpG mutagenesis in which 5mC is
225 occasionally incorrectly paired with adenine by Pol ϵ during replication of the leading strand,
226 potentially due to the structural similarity of 5mC and thymine. If such mismatches were not
227 detected by the polymerase proofreading machinery or MMR, they would result in CpG>TpG
228 mutations most frequently where 5mC occurred in the leading strand template. Under this model of
229 decreased fidelity of wildtype Pol ϵ in replication of 5mC, we would expect that such errors could
230 sometimes escape the polymerase proofreading and MMR even in *POLE*-WT and MMR proficient
231 samples, resulting in a slight strand asymmetry of CpG>TpG mutations. To test this, we grouped
232 PROF samples by tissue, and in each tissue measured the percentage of samples with a higher
233 CpG>TpG mutation frequency on the leading than the lagging strand, while also distinguishing
234 between all four sequence contexts. The majority of samples exhibited leading strand bias for
235 GCG>GTG mutations in 13 out of 16 tissue types in lowly and intermediately modified CpGs (Fig. 4-
236 supplement 1). This effect was even stronger (16 out of 16 tissues) when restricting the analysis to
237 highly modified CpGs only (Fig. 4), supporting the hypothesis that CpG>TpG mutations can also be
238 caused by errors during the replication of methylated cytosine by Pol ϵ .



239
240

Fig. 4: GCG>GTG mutations are more frequent on the leading strand than on the lagging strand, even in Pol ϵ and MMR

241

proficient samples. The heatmap shows the percentage of samples with higher C>T mutation frequency on the leading

242

strand than on the lagging strand (only C>T mutations in highly modified (>0.95) CpG sites, using tissue-matched

243

modification maps): white colour denotes no data, blue colour denotes more frequent lagging bias, and red denotes more

244

frequent leading bias. The number above each column represents the percentage of cancer types with a leading strand

245

bias in a majority of samples. Asterisks represent significance of the bias in each column (signtest; ***P < 0.001; **P < 0.01;

246

*P < 0.05).

247

4. DISCUSSION

248

The increased rate of C>T mutations at CpG dinucleotides across tissue types has been thought to

249

primarily stem from spontaneous deamination of methylated cytosine. The fact that *POLE*-MUT and

250

MSI samples exhibit high CpG>TpG mutation frequency is therefore surprising, since neither MMR

251

nor proofreading by Pol ϵ are thought to be required for the repair of deamination damage. A similar

252

increase of CpG>TpG mutations in MSI and *POLE*-MUT colorectal cancer samples has also been

253

observed in another study that was published during the preparation of this manuscript [28], but the

254

correlation of these mutations with methylation levels was not explored in much detail.

255

Three theoretical models could explain this observation. In the first model, MMR and Pol ϵ —

256

through a non-canonical, replication-unrelated mechanism— are in fact essential for the repair of

257

T:G mismatches created by spontaneous deamination of 5mC. For MMR, this is the model proposed

258 in a recent study [28]. However, the observed number of CpG>TpG mutations in MSI and *POLE*-MUT
259 samples is difficult to reconcile with the known deamination kinetics of methylated cytosine in
260 double-stranded DNA, even under the unrealistic assumption that no repair mechanisms at all are
261 active in these samples. At 5.8×10^{-13} mutations per 5mC per second [24], it would take 28 years to
262 reach the observed C>T mutation frequency in modified GCG sites of MSI samples, and 98 years for
263 *POLE*-MUT samples (see Methods for calculations). These timescales are unlikely to represent the
264 real time between the acquisition of the MMR or Pol ϵ mutation and the collection of the sample.
265 Moreover, if spontaneous deamination was the source of CpG>TpG mutagenicity in MMR and Pol ϵ
266 deficient samples, one would not expect to see replication strand asymmetry. However, CpG>TpG
267 mutations are highly enriched on the leading strand in all these samples and therefore do not
268 support this first model.

269 The second possible explanation is that the Pol ϵ and MMR mutations are gain of function
270 mutations, causing a mutator phenotype that actively increases CpG>TpG mutagenicity during
271 replication. This mechanism has been suggested by Poulos et al. [28] for the *POLE*-MUT samples and
272 by Kane et al. [29] in *S. cerevisiae*, where an analog of the human P286R variant (but not other
273 variants) in the yeast Pol ϵ produced a strong mutator phenotype, increasing the mutation rate
274 beyond that of the proofreading-null allele. However, we observed a marked increase of C>T
275 mutation frequency in modified CpG sites in a wide range of Pol ϵ variants (Fig. 3A). Furthermore, a
276 strong correlation of GCG>GTG mutations with DNA modification levels was observed across *POLE*-
277 MUT and MSI samples from multiple cancer types. It therefore seems unlikely that multiple different
278 Pol ϵ and MMR mutations all result in the same mutator phenotype.

279 The third model posits that wildtype Pol ϵ has a slightly decreased fidelity when encountering 5mC,
280 particularly in a GCG context, on the template strand and incorrectly pairs it with A, leading to
281 5mC:A mismatches. This could potentially be a consequence of the high structural similarity
282 between 5mC and T, both of which present a methyl group at the same position of pyrimidine ring. If
283 the resulting 5mC:A mismatches were not repaired before the next round of replication, for example

284 because of a lack of mismatch repair in MSI tumours, one would expect an enrichment of GCG>GTG
285 mutations on the leading strand, as we observe in our data. Similarly, a lack of proofreading by Pol ϵ
286 itself might overwhelm the capacity of downstream repair pathways and thus, too, lead to an
287 increased CpG>TpG mutations rate. The fact that we also detected a leading strand bias for
288 GCG>GTG mutations in a majority of Pol ϵ and MMR proficient tumours hints at the possibility that
289 the mechanism described above does contribute to the overall CpG>TpG mutation burden. This
290 model is also consistent with observations from samples with a mutation in the proofreading
291 domain of *POLD1*, a gene encoding the catalytic subunit of Pol δ . *POLD1*-MUT samples are also
292 highly mutated, but, unlike in *POLE*-MUT samples, CpG>TpG mutations form only a small percentage
293 of the mutation burden [12]. This observation supports the notion that the CpG>TpG mutagenesis is
294 specifically linked to the leading strand synthesis.

295 **5. CONCLUSIONS**

296 To conclude, we have presented evidence suggesting that replication of methylated cytosines is
297 likely to contribute to the higher mutation rate of CpGs in the genome. This unanticipated finding
298 changes the commonly accepted paradigm in the field, where spontaneous deamination has been
299 proposed as the only reason for the mutagenicity of methylated CpG sites. While replication-linked
300 CpG>TpG mutations dominate in Pol ϵ mutated or MMR deficient cells, the relative contribution of
301 replication-linked mutations compared to deamination-induced mutations in repair-proficient cells is
302 less clear. Pol ϵ proofreading and MMR both repair mutations originating during replication, while
303 *MBD4* and *TDG* are glycosylases repairing lesions caused by spontaneous deamination of 5mC. Pol ϵ
304 mutations increase CpG mutation rate by 210-fold in human cancers, while *Mbd4* deficient mice
305 exhibit an increase in mutation frequency by 3-fold [30], suggesting that replication might be more
306 mutagenic at methylated CpGs than deamination, unless *TDG* plays a dominant role in repair of
307 deamination lesions. Thus, Pol ϵ might even be the primary source of C>T mutations in methylated
308 CpGs, which could also explain that cancers from tissues with higher turnover rates exhibit an

309 increased rate of CpG>TpG mutations [31]. Further experimental work will be required to fully
310 elucidate the fidelity of Pol ϵ when replicating 5mC.

311 **SUPPORTING INFORMATION**

312 Fig. 1-supplement 1: Frequency of C to T mutations in a CpG context is unexpectedly high in *POLE*-
313 MUT and MSI samples.

314 Fig. 1-supplement 2: Frequency of C to T mutations in a CpG context in *POLE*-MUT and MSI samples
315 correlates with DNA modification levels: comparison of linear models.

316 Fig. 2-supplement 1: Frequency of C to T mutations in a CpG context in *POLE*-MUT and MSI samples
317 is higher on the leading strand than on the lagging strand, especially in modified CpG sites.

318 Fig. 3-supplement 1: CpG>TpG mutation frequency in different sequence contexts.

319 Fig. 3-supplement 2: Increase of C to T mutations in modified cytosine on the leading strand is most
320 consistent in a GCG sequence context in *POLE*-MUT and MSI samples.

321 Fig. 4-supplement 1: GCG>GTG mutations are more frequent on the leading strand than on the
322 lagging strand, even in Pol ϵ and MMR proficient samples.

323 Supplementary Table 1: Overview of BS-Seq and TAB-Seq data used to generate modification maps.

324 Supplementary Table 2: Overview of whole genome sequencing data used for mutation information.

325 **ACKNOWLEDGMENTS**

326 We thank Dr. Mary Muers and Jakub Tomek for comments on the manuscript. S.K. and B.S.-B. are
327 funded by Ludwig Cancer Research. S.K. received funding from BBSRC grant BB/M001873/1. M.T. is
328 funded by EPSRC grant EP/F500394/1 and the Bakala Foundation.

329 **CONFLICT OF INTEREST STATEMENT**

330 The authors declare that there are no conflicts of interest

331 **REFERENCES**

332 [1] L.B. Alexandrov, S. Nik-Zainal, D.C. Wedge, S.A.J.R. Aparicio, S. Behjati, A. V Biankin, G.R.
333 Bignell, N. Bolli, A. Borg, A.-L. Børresen-Dale, S. Boyault, B. Burkhardt, A.P. Butler, C. Caldas,
334 H.R. Davies, C. Desmedt, R. Eils, J.E. Eyfjörd, J.A. Foekens, M. Greaves, F. Hosoda, B. Hutter, T.

335 Ilicic, S. Imbeaud, M. Imielinski, M. Imielinsk, N. Jäger, D.T.W. Jones, D. Jones, S. Knappskog,
336 M. Kool, S.R. Lakhani, C. López-Otín, S. Martin, N.C. Munshi, H. Nakamura, P.A. Northcott, M.
337 Pajic, E. Papaemmanuil, A. Paradiso, J. V Pearson, X.S. Puente, K. Raine, M. Ramakrishna, A.L.
338 Richardson, J. Richter, P. Rosenstiel, M. Schlesner, T.N. Schumacher, P.N. Span, J.W. Teague,
339 Y. Totoki, A.N.J. Tutt, R. Valdés-Mas, M.M. van Buuren, L. van 't Veer, A. Vincent-Salomon, N.
340 Waddell, L.R. Yates, J. Zucman-Rossi, P.A. Futreal, U. McDermott, P. Lichter, M. Meyerson,
341 S.M. Grimmond, R. Siebert, E. Campo, T. Shibata, S.M. Pfister, P.J. Campbell, M.R. Stratton,
342 Signatures of mutational processes in human cancer., *Nature*. 500 (2013) 415–21.
343 doi:10.1038/nature12477.

344 [2] M.S. Lawrence, P. Stojanov, P. Polak, G. V Kryukov, K. Cibulskis, A. Sivachenko, S.L. Carter, C.
345 Stewart, C.H. Mermel, S.A. Roberts, A. Kiezun, P.S. Hammerman, A. McKenna, Y. Drier, L. Zou,
346 A.H. Ramos, T.J. Pugh, N. Stransky, E. Helman, J. Kim, C. Sougnez, L. Ambrogio, E. Nickerson,
347 E. Shefler, M.L. Cortés, D. Auclair, G. Saksena, D. Voet, M. Noble, D. DiCara, P. Lin, L.
348 Lichtenstein, D.I. Heiman, T. Fennell, M. Imielinski, B. Hernandez, E. Hodis, S. Baca, A.M.
349 Dulak, J. Lohr, D.-A. Landau, C.J. Wu, J. Melendez-Zajgla, A. Hidalgo-Miranda, A. Koren, S.A.
350 McCarroll, J. Mora, R.S. Lee, B. Crompton, R. Onofrio, M. Parkin, W. Winckler, K. Ardlie, S.B.
351 Gabriel, C.W.M. Roberts, J.A. Biegel, K. Stegmaier, A.J. Bass, L.A. Garraway, M. Meyerson, T.R.
352 Golub, D.A. Gordenin, S. Sunyaev, E.S. Lander, G. Getz, Mutational heterogeneity in cancer
353 and the search for new cancer-associated genes., *Nature*. 499 (2013) 214–8.
354 doi:10.1038/nature12213.

355 [3] A.P. Bird, M.H. Taggart, Variable patterns of total DNA and rDNA methylation in animals,
356 *Nucleic Acids Research*. 8 (1980) 1485–1497. doi:10.1093/nar/8.7.1485.

357 [4] M. Tomkova, M. McClellan, S. Kriaucionis, B. Schuster-Boeckler, 5-hydroxymethylcytosine
358 marks regions with reduced mutation frequency in human DNA, *eLife*. 5 (2016) 1–23.
359 doi:10.7554/eLife.17082.

- 360 [5] T. Lindahl, B. Nyberg, Heat-induced deamination of cytosine residues in deoxyribonucleic
361 acid, *Biochemistry*. 13 (1974) 3405–3410. doi:10.1021/bi00713a035.
- 362 [6] K. Wiebauer, J. Jiricny, In vitro correction of GT mispairs to GC pairs in nuclear extracts from
363 human cells, *Nature*. 339 (1989) 234–236. doi:10.1038/339234a0.
- 364 [7] B. Hendrich, U. Hardeland, H. Ng, J. Jiricny, A. Bird, The thymine glycosylase MBD4 can bind to
365 the product of deamination at methylated CpG sites, *Nature*. 401 (1999) 525–525.
366 doi:10.1038/35006691.
- 367 [8] A. Bellacosa, A.C. Drohat, Role of base excision repair in maintaining the genetic and
368 epigenetic integrity of CpG sites, *DNA Repair*. 32 (2015) 33–42.
369 doi:10.1016/j.dnarep.2015.04.011.
- 370 [9] T. Pfaffeneder, F. Spada, M. Wagner, C. Brandmayr, S.K. Laube, D. Eisen, M. Truss, J.
371 Steinbacher, B. Hackner, O. Kotljarova, D. Schuermann, S. Michalakis, O. Kosmatchev, S.
372 Schiesser, B. Steigenberger, N. Raddaoui, G. Kashiwazaki, U. Müller, C.G. Spruijt, M.
373 Vermeulen, H. Leonhardt, P. Schär, M. Müller, T. Carell, Tet oxidizes thymine to 5-
374 hydroxymethyluracil in mouse embryonic stem cell DNA., *Nature Chemical Biology*. 10 (2014)
375 574–81. doi:10.1038/nchembio.1532.
- 376 [10] E. Rayner, I.C. van Gool, C. Palles, S.E. Kearsey, T. Bosse, I. Tomlinson, D.N. Church, A panoply
377 of errors: polymerase proofreading domain mutations in cancer, *Nature Reviews Cancer*. 16
378 (2016) 71–81. doi:10.1038/nrc.2015.12.
- 379 [11] E. Shinbrot, E.E. Henninger, N. Weinhold, K.R. Covington, A.Y. Göksenin, N. Schultz, H. Chao,
380 H. Doddapaneni, D.M. Muzny, R.A. Gibbs, C. Sander, Z.F. Pursell, D.A. Wheeler, Exonuclease
381 mutations in DNA Polymerase epsilon reveal replication strand specific mutation patterns
382 and human origins of replication., *Genome Research*. (2014) 1740–1750.
383 doi:10.1101/gr.174789.114.

- 384 [12] A. Shlien, B.B. Campbell, R. de Borja, L.B. Alexandrov, D. Merico, D. Wedge, P. Van Loo, P.S.
385 Tarpey, P. Coupland, S. Behjati, A. Pollett, T. Lipman, A. Heidari, S. Deshmukh, N. Avitzur, B.
386 Meier, M. Gerstung, Y. Hong, D.M. Merino, M. Ramakrishna, M. Remke, R. Arnold, G.B.
387 Panigrahi, N.P. Thakkar, K.P. Hodel, E.E. Henninger, A.Y. Göksenin, D. Bakry, G.S. Charames, H.
388 Druker, J. Lerner-Ellis, M. Mistry, R. Dvir, R. Grant, R. Elhasid, R. Farah, G.P. Taylor, P.C.
389 Nathan, S. Alexander, S. Ben-Shachar, S.C. Ling, S. Gallinger, S. Constantini, P. Dirks, A. Huang,
390 S.W. Scherer, R.G. Grundy, C. Durno, M. Aronson, A. Gartner, M.S. Meyn, M.D. Taylor, Z.F.
391 Pursell, C.E. Pearson, D. Malkin, P.A. Futreal, M.R. Stratton, E. Bouffet, C. Hawkins, P.J.
392 Campbell, U. Tabori, Combined hereditary and somatic mutations of replication error repair
393 genes result in rapid onset of ultra-hypermuted cancers, *Nature Genetics*. 47 (2015) 257–
394 262. doi:10.1038/ng.3202.
- 395 [13] T.M. Albertson, M. Ogawa, J.M. Bugni, L.E. Hays, Y. Chen, Y. Wang, P.M. Treuting, J.A. Heddle,
396 R.E. Goldsby, B.D. Preston, DNA polymerase epsilon and delta proofreading suppress discrete
397 mutator and cancer phenotypes in mice., *Proceedings of the National Academy of Sciences of*
398 *the United States of America*. 106 (2009) 17101–4. doi:10.1073/pnas.0907147106.
- 399 [14] A.J. Bass, M.S. Lawrence, L.E. Brace, A.H. Ramos, Y. Drier, K. Cibulskis, C. Sougnez, D. Voet, G.
400 Saksena, A. Sivachenko, R. Jing, M. Parkin, T. Pugh, R.G. Verhaak, N. Stransky, A.T. Boutin, J.
401 Barretina, D.B. Solit, E. Vakiani, W. Shao, Y. Mishina, M. Warmuth, J. Jimenez, D.Y. Chiang, S.
402 Signoretti, W.G. Kaelin, N. Spardy, W.C. Hahn, Y. Hoshida, S. Ogino, R.A. Depinho, L. Chin, L.A.
403 Garraway, C.S. Fuchs, J. Baselga, J. Tabernero, S. Gabriel, E.S. Lander, G. Getz, M. Meyerson,
404 Genomic sequencing of colorectal adenocarcinomas identifies a recurrent VTI1A-TCF7L2
405 fusion., *Nature Genetics*. 43 (2011) 964–8. doi:10.1038/ng.936.
- 406 [15] A.M. Dulak, P. Stojanov, S. Peng, M.S. Lawrence, C. Fox, C. Stewart, S. Bandla, Y. Imamura,
407 S.E. Schumacher, E. Shefler, A. McKenna, S.L. Carter, K. Cibulskis, A. Sivachenko, G. Saksena,
408 D. Voet, A.H. Ramos, D. Auclair, K. Thompson, C. Sougnez, R.C. Onofrio, C. Guiducci, R.
409 Beroukhim, Z. Zhou, L. Lin, J. Lin, R. Reddy, A. Chang, R. Landrenau, A. Pennathur, S. Ogino,

410 J.D. Luketich, T.R. Golub, S.B. Gabriel, E.S. Lander, D.G. Beer, T.E. Godfrey, G. Getz, A.J. Bass,
411 Exome and whole-genome sequencing of esophageal adenocarcinoma identifies recurrent
412 driver events and mutational complexity., *Nature Genetics*. 45 (2013) 478–86.
413 doi:10.1038/ng.2591.

414 [16] K. Wang, S.T. Yuen, J. Xu, S.P. Lee, H.H.N. Yan, S.T. Shi, H.C. Siu, S. Deng, K.M. Chu, S. Law,
415 K.H. Chan, A.S.Y. Chan, W.Y. Tsui, S.L. Ho, A.K.W. Chan, J.L.K. Man, V. Foglizzo, M.K. Ng, A.S.
416 Chan, Y.P. Ching, G.H.W. Cheng, T. Xie, J. Fernandez, V.S.W. Li, H. Clevers, P.A. Rejto, M. Mao,
417 S.Y. Leung, Whole-genome sequencing and comprehensive molecular profiling identify new
418 driver mutations in gastric cancer., *Nature Genetics*. 46 (2014) 573–82. doi:10.1038/ng.2983.

419 [17] N.J. Haradhvala, P. Polak, P. Stojanov, K.R. Covington, E. Shinbrot, J.M. Hess, E. Rheinbay, J.
420 Kim, Y.E. Maruvka, L.Z. Braunstein, A. Kamburov, P.C. Hanawalt, D.A. Wheeler, A. Koren, M.S.
421 Lawrence, G. Getz, Mutational Strand Asymmetries in Cancer Genomes Reveal Mechanisms
422 of DNA Damage and Repair, *Cell*. 164 (2016) 538–549. doi:10.1016/j.cell.2015.12.050.

423 [18] C.T. Saunders, W.S.W. Wong, S. Swamy, J. Becq, L.J. Murray, R.K. Cheetham, Strelka: Accurate
424 somatic small-variant calling from sequenced tumor-normal sample pairs, *Bioinformatics*. 28
425 (2012) 1811–1817. doi:10.1093/bioinformatics/bts271.

426 [19] L. Wen, X. Li, L. Yan, Y. Tan, R. Li, Y. Zhao, Y. Wang, J. Xie, Y. Zhang, C. Song, M. Yu, X. Liu, P.
427 Zhu, X. Li, Y. Hou, H. Guo, X. Wu, C. He, R. Li, F. Tang, J. Qiao, Whole-genome analysis of 5-
428 hydroxymethylcytosine and 5-methylcytosine at base resolution in the human brain.,
429 *Genome Biology*. 15 (2014) R49. doi:10.1186/gb-2014-15-3-r49.

430 [20] K. Chen, J. Zhang, Z. Guo, Q. Ma, Z. Xu, Y. Zhou, Z. Xu, Z. Li, Y. Liu, X. Ye, X. Li, B. Yuan, Y. Ke, C.
431 He, L. Zhou, J. Liu, W. Ci, Loss of 5-hydroxymethylcytosine is linked to gene body
432 hypermethylation in kidney cancer, *Cell Research*. (2015) 103–118. doi:10.1038/cr.2015.150.

433 [21] R. Pidsley, E. Zotenko, T.J. Peters, M.G. Lawrence, G.P. Risbridger, P. Molloy, S. Van Djik, B.
434 Muhlhausler, C. Stirzaker, S.J. Clark, P. Jones, S. Baylin, Y. Ko, D. Mohtat, M. Suzuki, A. Park,

435 M. Izquierdo, S. Han, T. Dayeh, P. Volkov, S. Salo, E. Hall, E. Nilsson, A. Olsson, R. Pidsley, J.
436 Viana, E. Hannon, H. Spiers, C. Troakes, S. Al-Saraj, C. Stirzaker, P. Taberlay, A. Statham, S.
437 Clark, S. Clark, J. Harrison, C. Paul, M. Frommer, R. Lister, M. Pelizzola, R. Downen, R. Hawkins,
438 G. Hon, J. Tonti-Filippini, M. Bibikova, J. Le, B. Barnes, S. Saedinia-Melnyk, L. Zhou, R. Shen, T.
439 Hinoue, D. Weisenberger, C. Lange, H. Shen, H. Byun, D. Berg, L. Breitling, R. Yang, B. Korn, B.
440 Burwinkel, H. Brenner, V. Rakyan, T. Down, S. Maslau, T. Andrew, T. Yang, H. Beyan, M.
441 Bibikova, B. Barnes, C. Tsan, V. Ho, B. Klotzle, J. Le, T. Morris, S. Beck, Y. Chen, S. Choufani, D.
442 Grafodatskaya, D. Butcher, J. Ferreira, R. Weksberg, Y. Chen, M. Lemire, S. Choufani, D.
443 Butcher, D. Grafodatskaya, B. Zanke, H. Naeem, N. Wong, Z. Chatterton, M. Hong, J.
444 Pedersen, N. Corcoran, T. Peters, M. Buckley, A. Statham, R. Pidsley, K. Samaras, R. V Lord, D.
445 Wang, L. Yan, Q. Hu, L. Sucheston, M. Higgins, C. Ambrosone, C. Warden, H. Lee, J. Tompkins,
446 X. Li, C. Wang, A. Riggs, M. Lizio, J. Harshbarger, H. Shimoji, J. Severin, T. Kasukawa, S. Sahin,
447 L. Siggins, K. Ekwall, S. Dedeurwaerder, M. Defrance, E. Calonne, H. Denis, C. Sotiriou, F.
448 Fuks, R. Pidsley, Y.W. CC, M. Volta, K. Lunnon, J. Mill, L. Schalkwyk, A. Teschendorff, F.
449 Marabita, M. Lechner, T. Bartlett, J. Tegner, D. Gomez-Cabrero, N. Touleimat, J. Tost, R.
450 Thurman, E. Rynes, R. Humbert, J. Vierstra, M. Maurano, E. Haugen, R. Andersson, C.
451 Gebhard, I. Miguel-Escalada, I. Hoof, J. Bornholdt, M. Boyd, A. Kundaje, W. Meuleman, J.
452 Ernst, M. Bilenky, A. Yen, M. Ritchie, B. Phipson, D. Wu, Y. Hu, C. Law, W. Shi, M. Stadler, R.
453 Murr, L. Burger, R. Ivanek, F. Lienert, A. Schöler, M. Ziller, H. Gu, F. Müller, J. Donaghey, L.-Y.
454 Tsai, O. Kohlbacher, S. Huang, B. Bao, T. Hour, C. Huang, C. Yu, C. Liu, S. Neuhausen, M.
455 Slattery, C. Garner, Y. Ding, M. Hoffman, A. Brothman, R. Reams, K. Kalari, H. Wang, F.
456 Odedina, K. Soliman, C. Yates, J. Song, C. Stirzaker, J. Harrison, J. Melki, S. Clark, M. Coolen, C.
457 Stirzaker, J. Song, A. Statham, Z. Kassir, C. Moreno, M. Makrides, R. Gibson, A. McPhee, L.
458 Yelland, J. Quinlivan, P. Ryan, M. Lawrence, R. Taylor, R. Toivanen, J. Pedersen, S. Norden, D.
459 Pook, S. Clark, A. Statham, C. Stirzaker, P. Molloy, M. Frommer, A. Auton, L. Brooks, R. Durbin,
460 E. Garrison, H. Kang, W. Kent, Critical evaluation of the Illumina MethylationEPIC BeadChip
461 microarray for whole-genome DNA methylation profiling, *Genome Biology*. 17 (2016) 208.

462 doi:10.1186/s13059-016-1066-1.

463 [22] A.R. Vandiver, R.A. Irizarry, K.D. Hansen, L.A. Garza, A. Runarsson, X. Li, A.L. Chien, T.S. Wang,
464 S.G. Leung, S. Kang, A.P. Feinberg, Age and sun exposure-related widespread genomic blocks
465 of hypomethylation in nonmalignant skin., *Genome Biology*. 16 (2015) 80.
466 doi:10.1186/s13059-015-0644-y.

467 [23] F. Krueger, B. Kreck, A. Franke, S.R. Andrews, DNA methylome analysis using short bisulfite
468 sequencing data, *Nature Methods*. 9 (2012) 145–151. doi:10.1038/nmeth.1828.

469 [24] J.-C. Shen, W.M. Rideout, P.A. Jones, The rate of hydrolytic deamination of 5-methylcytosine
470 in double-stranded DNA, *Nucleic Acids Research*. 22 (1994) 972–976.
471 doi:10.1093/nar/22.6.972.

472 [25] B. Stillman, DNA Polymerases at the Replication Fork in Eukaryotes, *Molecular Cell*. 30 (2008)
473 259–260. doi:10.1016/j.molcel.2008.04.011.

474 [26] R.E. Georgescu, G.D. Schauer, N.Y. Yao, L.D. Langston, O. Yurieva, D. Zhang, J. Finkelstein,
475 M.E. O'Donnell, Reconstitution of a eukaryotic replisome reveals suppression mechanisms
476 that define leading/lagging strand operation, *eLife*. 2015 (2015) 1–20.
477 doi:10.7554/eLife.04988.

478 [27] S.A. Lujan, J.S. Williams, Z.F. Pursell, A.A. Abdulovic-Cui, A.B. Clark, S.A. Nick McElhinny, T.A.
479 Kunkel, Mismatch Repair Balances Leading and Lagging Strand DNA Replication Fidelity, *PLoS*
480 *Genetics*. 8 (2012) e1003016. doi:10.1371/journal.pgen.1003016.

481 [28] R.C. Poulos, J. Olivier, J.W.H. Wong, The interaction between cytosine methylation and
482 processes of DNA replication and repair shape the mutational landscape of cancer genomes,
483 *Nucleic Acids Research*. (2017) 1–10. doi:10.1093/nar/gkx463.

484 [29] D.P. Kane, P. V. Shcherbakova, A common cancer-associated DNA polymerase ϵ mutation
485 causes an exceptionally strong mutator phenotype, indicating fidelity defects distinct from

486 loss of proofreading, *Cancer Research*. 74 (2014) 1895–1901. doi:10.1158/0008-5472.CAN-
487 13-2892.

488 [30] C.B. Millar, Enhanced CpG Mutability and Tumorigenesis in MBD4-Deficient Mice, *Science*.
489 297 (2002) 403–405. doi:10.1126/science.1073354.

490 [31] L.B. Alexandrov, P.H. Jones, D.C. Wedge, J.E. Sale, J. Peter, Clock-like mutational processes in
491 human somatic cells, *Nature*. 47 (2015) 1402–1407. doi:10.1038/ng.3441.

492