

1           **Transcription-coupled repair and mismatch repair**  
2           **contribute towards preserving genome integrity at**  
3           **mononucleotide repeat tracts**

4  
5 Ilias Georgakopoulos-Soares<sup>1,2</sup>, Gene Koh<sup>1,3,4</sup>, Sophie E. Momen<sup>3,4</sup>, Josef Jiricny<sup>5</sup>, Martin  
6 Hemberg<sup>1,+</sup>, Serena Nik-Zainal<sup>3,4,+</sup>

7  
8 <sup>1</sup>Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton CB 10 1SA, UK.

9 <sup>2</sup>Department of Bioengineering and Therapeutic Sciences, Institute for Human Genetics, University of California San Francisco,  
10 San Francisco, California 94158, USA.

11 <sup>3</sup>Department of Medical Genetics, The Clinical School, University of Cambridge, Cambridge, CB2 0QQ, UK.

12 <sup>4</sup>MRC Cancer Unit, The Clinical School, University of Cambridge, Cambridge, CB2 0QQ, UK.

13 <sup>5</sup>Institute of Molecular Life Sciences, University of Zurich and Institute of Biochemistry of the ETH Zurich, Otto-Stern-Weg 3,  
14 CH-8093 Zurich, Switzerland.

15  
16 <sup>+</sup>Corresponding authors: Martin Hemberg ([mh26@sanger.ac.uk](mailto:mh26@sanger.ac.uk)) & Serena Nik-Zainal ([sn206@cam.ac.uk](mailto:sn206@cam.ac.uk))  
17  
18

19 **Abstract**

20  
21 The mechanisms that underpin how insertions or deletions (indels) become fixed in DNA have  
22 primarily been ascribed to replication-related and/or double-strand break (DSB)-related  
23 processes. Here we introduce a method to evaluate indels, orientating them relative to gene  
24 transcription. In so doing, we reveal a number of surprising findings: First, there is a  
25 transcriptional strand asymmetry in the distribution of mononucleotide repeat tracts in the  
26 reference human genome. Second, there is a strong transcriptional strand asymmetry of indels  
27 across 2,575 whole genome sequenced human cancers. We suggest that this is due to the activity  
28 of transcription-coupled nucleotide excision repair (TC-NER). Furthermore, TC-NER interacts  
29 with mismatch repair (MMR) under physiological conditions to produce strand bias. Finally, we  
30 show how insertions and deletions differ in their dependencies on these repair pathways. Our  
31 analytical approach reveals insights into the contribution of DNA repair towards indel  
32 mutagenesis in human cells.

33

## 34 **Introduction**

35

36 Mutations are not randomly distributed across the cancer genome. Their distribution is  
37 influenced by genomic, epigenomic and cellular physiological factors such as replication and  
38 transcription<sup>1-4</sup>. Transcription has been implicated in contributing to mutational strand  
39 asymmetries reflecting biases in DNA damage (transcription-associated damage) and DNA  
40 repair mechanisms (transcription-coupled repair) between the two strands<sup>3-5</sup>.

41

42 In this area, while substitutions in human cancers have been extensively studied,  
43 insertions/deletions (indels) have remained comparatively under-explored. This was historically  
44 due to the relative difficulty in obtaining high-quality indel data, further restricted by a limited  
45 repertoire of approaches to analyze indels as extensively as substitutions. Nevertheless, indels are  
46 common in human cancers and their location and sequence composition are non-random. Thus,  
47 like substitutions, they provide important insights into the mutational processes that have shaped  
48 the landscape of cancer genomes.

49

50 Here, we demonstrate that there is transcriptional strand asymmetry in the distribution of  
51 mononucleotide repeat tracts within the reference genome. We also observe transcriptional  
52 strand asymmetry in insertions and deletions at mononucleotide repeat tracts across cancer types,  
53 and are able to attribute the relative contributions of transcription-coupled nucleotide excision  
54 repair (TC-NER) and mismatch repair (MMR) pathways to indel patterns in human somatic  
55 cells.

56

## 57 **Results**

58

### 59 **Landscape of insertions and deletions across human cancers**

60 We utilized 2,416,257 indels from a highly curated set of 2,575 whole-genome sequenced  
61 (WGS) cancers of 21 different cancer-types. Median indel number per tumour was 386,  
62 corresponding to a conservative indel density of 0.127 per Mb per cancer genome. Deletions  
63 (median 222) were more prevalent than insertions (median 124) in the majority of cancers

64 (Mann-Whitney U p-value<0.05, Figure 1a, Supplementary Figure 1a). Moreover, deletion size  
65 showed greater variability than insertion size across and within tumour-types (Figure 1b,  
66 Supplementary Figure 1b-c, Levene's test, p-value<0.05).

67  
68 This first observation can be broadly explained by already-known mechanisms that generate  
69 indel lesions. Replication-related DNA polymerase slippage errors running through  
70 microsatellites tend to cause deletions<sup>6</sup>, because single-stranded DNA ahead of a polymerase can  
71 twist, causing a single repeat unit of a run of mono- or dinucleotides in the template strand to  
72 loop out. A polymerase passing over such a loop would generate a deletion<sup>6-9</sup>. Because these  
73 small insertion/deletion loops (IDLs) are efficiently repaired by MMR, the density of deletions in  
74 microsatellites is higher in cells lacking MMR. This phenomenon is referred to as microsatellite  
75 instability (MSI). The likelihood of formation of such loops increases with the length of the  
76 repeat and we confirm this by showing that indel frequency is augmented with increasing lengths  
77 of polynucleotide repeat tracts (Supplementary Figure 1d-e) and is more pronounced in MMR-  
78 deficient samples. By contrast, double-strand breaks (DSB) can give rise to deletions if repaired  
79 by non-homologous end-joining (NHEJ), or if addressed by homology-directed sub-pathways  
80 such as single-strand annealing (SSA) or micro-homology-directed end-joining (MMEJ)<sup>1,10-15</sup>.  
81 The latter result in larger deletions (3bp in size or more), thus explaining the broader spectrum of  
82 observed deletion sizes (Figure 1b, Levene's test, p-value<0.05).

83  
84 In general, it is more difficult to create an insertion. Transient dissociation of the primer and  
85 template strands and reannealing of the primer in a wrong register within the microsatellite could  
86 cause both an insertion or a deletion. These are less likely to arise during normal replication<sup>16</sup>  
87 because the end of the primer strand is tightly bound by the replisome. Thus, the occurrence and  
88 the relative frequencies of indels and their size spectra can be explained by known mechanisms.

89  
90 In addition to classical contributions of replication and DSB-repair pathways to indel formation,  
91 we introduce another dimension to exploratory analyses of indel mutagenesis: the contribution of  
92 transcription. Transcription has been implicated in asymmetric distribution of substitutions  
93 between strands for decades<sup>4,17-18</sup>. In particular, transcription-coupled nucleotide excision repair  
94 (TC-NER) is believed to preferentially repair DNA damage on the template (transcribed or non-

95 coding) strand. TC-NER activity is thus inferred from the excess of mutagenesis on the non-  
96 template (coding) as compared to the template (non-coding) strand, particularly for those  
97 environmental mutagens where the target of primary DNA adduct formation is known. For  
98 example, guanines adducted by tobacco carcinogens result in an excess of G>T mutations on the  
99 non-template strand<sup>19-21</sup>. Likewise, primary covalent modifications of cytosines forming 6,4  
100 pyrimidine-pyrimidone dimers (6,4-PPs) and cyclobutane pyrimidine dimers (CPDs) by  
101 ultraviolet light are preferentially repaired on the template strand resulting in an excess of C>T  
102 transitions on the non-template strand<sup>22</sup>. However, to the best of our knowledge, transcriptional  
103 strand asymmetry in indels has not been investigated, primarily due to the technical challenge of  
104 being unable to orientate each indel with respect to transcriptional strand.

105

### 106 **Asymmetries of repetitive tracts in the reference genome**

107 We set out to determine transcriptional strand asymmetry of indels by focusing on  
108 mononucleotide repeats of up to ten base pairs in length. We first analyzed the distribution of  
109 mononucleotide repeats across the gene body in the reference human genome. Each gene was  
110 divided into ten equal-sized bins to correct for differences in gene length. Two additional bins  
111 were added upstream of the transcription start site (TSS) and two downstream of the  
112 transcription end site (TES), each 10kB in length, resulting in a total of 14 bins.

113

114 We observed a strong enrichment of polyG/polyC motifs directly upstream and downstream of  
115 the TSS and downstream of the TES; this contrasted with the distribution of polyT/polyA motifs,  
116 which were found to be enriched throughout gene body (Figure 2a, Supplementary Figure 3a).

117

118 We calculated the frequencies of each polyN motif (where N is any nucleotide) on the template  
119 and non-template strands in the reference genome. Because the direction of transcription for each  
120 gene in the genome is known, each polyN motif can be orientated (Figure 2b). For example, for a  
121 gene on the (+) strand, the template strand is the (-) strand. A polyT motif that is on the (-) strand  
122 of this gene is therefore described as being on the template strand. It can also be described as a  
123 poly-A motif on the non-template strand (Figure 2b). Using this reasoning, we assigned each  
124 polyN motif to either the template or non-template strand of the reference genome. If there were  
125 no asymmetries, polyN tracts would occur with equal probabilities on either strand.

126

127 Intriguingly, we found that polyT motifs displayed a bias towards the non-template strand in the  
128 reference genome, with a non-template to template asymmetry enrichment for short polyT motifs  
129 of ~1.15-fold (Figure 2c-d). This was tract-length-dependent, where longer repetitive tracts were  
130 associated with greater strand bias of up to ~1.4-fold at >5nt polyT motifs (weighted average  
131 asymmetry of 1.14-fold, Figure 2b, Supplementary Figure 3b-g). In contrast, we did not observe  
132 a similarly pronounced asymmetry of polyG motifs across gene bodies, although a skew in  
133 polyG motifs was noted at the boundaries of gene bodies, in-keeping with previous reports of  
134 GC-skewing at either end of genes<sup>23</sup> (Figure 2c-d, Supplementary Figure 2a-b, Supplementary  
135 Figure 3d-g). The marked variation in strand distributions of the polyN motifs in the reference  
136 genome is appreciated particularly around the TSS and TES (Figure 2c-d, Supplementary Figure  
137 3d-g).

138

### 139 **Transcriptional strand asymmetries of small indels occur at polynucleotide repeat tracts**

140

141 We next investigated whether there was strand asymmetry for indel occurrences at polyN tracts.  
142 All analyses henceforth, correct for the skewed background distributions of polyT and polyG  
143 motifs. Across cancers, polyT motifs of 2-10bp in length were consistently more mutable on the  
144 non-template strand (binomial test,  $p\text{-value} < e^{-5}$ ). Strand asymmetry was more pronounced for  
145 longer polyT tracts in all cancers (Kruskal-Wallis H-test with Bonferroni correction,  $p\text{-value} < e^{-9}$ ),  
146 (Supplementary Figure 4a). The levels of asymmetry varied by cancer type, with increased  
147 indel mutagenesis on the non-template over the template strand ranging from 2.1% in ovarian  
148 cancers to 16.5% in uterine cancers (Figure 3a, e). This was surprising, given that the prevailing  
149 dogma on indel formation, particularly at poly-nucleotide repeat tracts, involves the formation of  
150 small IDLs that are substrates for MMR<sup>24-27</sup>. Rather, our analysis showing marked transcriptional  
151 strand asymmetry implicates either the activity of TC-NER at these motifs or the activity of  
152 transcription-associated damage.

153

154 We noted that uterine, colorectal, biliary and stomach cancers showed the highest levels of  
155 transcriptional strand asymmetry (binomial test with Bonferroni correction,  $p\text{-value} < 0.001$  for all  
156 four cancer types), with 16.5%, 15.5%, 16.3% and 15.5% more indels occurring on non-template  
157 than template polyT motifs (Figure 3a). Notably, these cancer types are often associated with

158 incidences of MMR deficiency<sup>28</sup>. To explore the contribution of MMR to transcriptional strand  
159 asymmetries of indels at polyT tracts, we compared samples with MMR deficiency (or MSI) to  
160 microsatellite stable (MSS) samples. Surprisingly, in MSI samples, transcriptional strand bias  
161 towards the non-template strand for polyT motifs was more pronounced than in MSS samples,  
162 with a 7.9-12.8% increased indel occurrence (Figure 3b), (Mann-Whitney U p-value<0.001 in all  
163 cases, Bonferroni corrected). This suggests that not only is TC-NER implicated in the repair at  
164 polyN motifs, it is also dependent on the normal physiological functioning of MMR. In the  
165 absence of MMR, damage at these sites relies more heavily on TC-NER alone, resulting in an  
166 increase in strand bias.

167

### 168 **Nucleotide excision repair and mismatch repair influence the indel landscape**

169 To validate this hypothesis regarding the reliance of TC-NER on MMR, we examined  
170 experimentally-generated mutation patterns from CRISPR-Cas9 knockouts of a human cancer  
171 cell line, HAP1<sup>29</sup>. We would expect the presence of transcriptional strand bias under normal  
172 conditions but that the magnitude of the effect would be increased in MMR gene knockouts.  
173 Indeed, our analytical findings are recapitulated in the experimental setting. In a knock-out  
174 model of MutS homolog 6 (*MSH6*), a key MMR gene, 1,663 indels occurred on the non-template  
175 strand polyT tracts, whereas 1,165 indels occurred on template polyT tracts, corresponding to a  
176 16.9% corrected increase in frequency on the non-template strand (binomial test, p-value<e-6), a  
177 similar magnitude to that observed in cancers. However, when this is divided by polynucleotide  
178 tract lengths (T, TT, TTT, T<sub>n</sub>), the numbers are low and the experimental data is under-powered  
179 to demonstrate the effect at all repeat lengths, even though the effect is there in aggregate.

180

181 Interestingly, in contrast to polyT motifs, we did not observe transcriptional strand bias for indels  
182 at polyG motifs across cancers (binomial test, p-value>0.05), (Figure 3a, c, binomial test, p-  
183 value<e-5), (Supplementary Figure 4b), with lung cancers being the notable exception; they  
184 exhibited a large excess of G indels on the non-template strand (15.8% greater indel occurrence  
185 at polyGs on non-template compared to template strand (binomial test with Bonferroni  
186 correction, p-value<e-30), (Figure 3c). This pattern of asymmetry mirrors what was observed for  
187 G>T substitutions in lung cancers, which are attributed to the formation of bulky adducts on  
188 guanines from tobacco-related carcinogens. This type of helix-distorting damage is classically

189 repaired by TC-NER<sup>19-21</sup>. The observation of transcriptional strand asymmetry for G indels at  
190 polyG tracts in tobacco-associated lung cancers reinforces how TC-NER can be involved in  
191 maintenance of genome integrity at polyN motifs, and could therefore also be acting at polyT  
192 tracts as hypothesized earlier.

193

194 To validate this observation of indel transcriptional strand asymmetry with tobacco exposure, we  
195 analyzed indel mutational profiles of non-cancerous human cells exposed to various polycyclic  
196 aromatic hydrocarbons (PAHs) including benzo[a]pyrene [0.39  $\mu$ M and 2  $\mu$ M] and  
197 benzo[a]pyrene diol epoxide [0.125  $\mu$ M], believed to be the carcinogenic components of tobacco  
198 smoke. We observed that 77 indels occurred on non-template polyG tracts, in contrast to only 39  
199 indels on template polyG tracts. This corresponds to nearly double the number of indels on the  
200 non-template strand over the template strand (binomial test, p-value<0.001), supporting our  
201 analytical observations of *in vivo* patterns derived from studying human cancers (Figure 3c, f).

202

203 The activity of TC-NER is linked to gene expression levels<sup>30</sup> where higher levels of transcription  
204 are associated with increased TC-NER activity. To seek further support that TC-NER plays a  
205 role in the repair of polyG motifs in lung cancers, we explored the degree of asymmetry in  
206 relation to gene expression levels. We used gene expression data from a representative cell-of-  
207 origin (Supplementary Table 2). In keeping with our hypothesis that TC-NER plays a pivotal  
208 role in the repair at polyG tracts in lung cancers, there was minimal transcriptional strand  
209 asymmetry for polyG motifs at genes that were not expressed or lowly expressed, and strong  
210 asymmetry for medium- and highly-expressed genes (Mann-Whitney U p-value<0.001). This  
211 effect was also strongly-dependent on the length of the polyG motifs (Figure 3d), (Kruskal-  
212 Wallis H-test with Bonferroni correction, p-value<e-05 for medium and high expression genes,  
213 p-value>0.05 for low expression genes).

214

215 Replication has also been reported to induce asymmetric mutation distributions between leading /  
216 lagging strands<sup>3-4</sup>. To exclude the possibility of replication strand orientation confounding our  
217 observations, we investigated whether indel transcriptional strand asymmetries at selected polyN  
218 motifs were related to leading and lagging replicative orientation. We found that replication  
219 strand orientation had limited effect on the observed indel transcriptional strand asymmetry of

220 these tracts (Mann-Whitney U with Bonferroni correction,  $p\text{-value}>0.05$  in all cases),  
221 (Supplementary Figure 5a-b). This supports the role of transcription and excludes the influence  
222 of replication in the generation of indel transcriptional strand asymmetries.

223

### 224 **Insertions and deletions are differentially-dependent on DNA repair pathways**

225 Next, we distinguished insertions from deletions at polyT and polyG tracts to find that  
226 transcriptional strand asymmetry differed between these classes of indels (Figure 4a,  
227 Supplementary Figure 6a-b). Insertions showed aggravated asymmetries at polyT tracts across all  
228 cancer types and were independent of MMR status, suggesting that mutagenesis associated with  
229 polyT tracts may be largely dependent on TC-NER (Wilcoxon signed-rank,  $p\text{-value}<e-5$ ) (Figure  
230 4a). By contrast, non-template strand bias of deletions at poly-T tracts was restricted to tumour  
231 types that had a high incidence of MSI (biliary, colorectal, stomach and endometrial). Thus,  
232 mutagenesis that results in deletions is more heavily dependent on the MMR pathway.

233

234 To support this hypothesis, we investigated the relationship between transcriptional strand  
235 asymmetry for insertions and deletions at polyT motifs, and gene expression levels  
236 (Supplementary Table 2). Genes with higher expression levels displayed stronger transcriptional  
237 strand asymmetry of insertions at polyT tracts across all inspected cancer types (Figure 4b,  
238 Mann-Whitney U with Bonferroni correction,  $p\text{-value}<0.05$ ), implicating TC-NER, which is  
239 linked to expression levels<sup>30</sup>. However, a relationship between expression levels and  
240 transcriptional strand asymmetry of deletions at polyT tracts could only be observed for a subset  
241 of cancer types (Figure 4c, Mann-Whitney U with Bonferroni correction,  $p\text{-value}>0.05$ ) and the  
242 strand bias was less apparent. In contrast, at polyG motifs, we did not observe consistent  
243 associations between expression levels and transcription strand asymmetry of insertions or  
244 deletions across cancer types (Supplementary Figure 6c-d, Mann-Whitney U with Bonferroni  
245 correction  $p\text{-value}>0.05$ ), with the exception of lung cancers; this we expected because of the  
246 influence of bulky adducts from tobacco carcinogens (Figure 3d).

247

248 To provide further evidence for the role of TC-NER in the observed transcriptional strand  
249 asymmetry at polyT tracts for insertions relative to deletions, we reasoned that patients with  
250 defects in the TC-NER pathway would have indel patterns that should not demonstrate

251 transcriptional strand bias because of defective NER. By contrast, patients with defects in global  
252 genome NER (GG-NER), may not manifest any changes in transcriptional strand bias. Tumor  
253 samples from these rare syndromes are however extremely difficult to obtain and systematic  
254 WGS data are not widely available to perform such analyses. We sequenced a cutaneous  
255 malignancy derived from a patient with an autosomal recessive DNA repair defect called  
256 Xeroderma Pigmentosum (XP). The patient was a compound heterozygote for the XPC gene,  
257 involved in GG-NER. Intriguingly, non-template strand bias was not observed for insertions at  
258 polyT tracts in this tumour, in contrast to what we had observed across cancer types (Figure 4a,  
259 binomial test  $p$ -value $<0.001$ ). The numbers of insertions however were small in this single  
260 sample. We thus performed a down-sampling of the numbers of insertions for all cancer types to  
261 similar levels as the XP-mutant tumour to examine whether the difference in transcriptional  
262 strand asymmetry remained significant. The XP-deficient tumour consistently displayed  
263 decreased levels of non-template strand bias at insertions in comparison to all cancer types  
264 (Supplementary Figure 7a-b), whereas for deletions there were no significant differences relative  
265 to other cancers (Supplementary Figure 7c-d). A more robust assessment will be required in due  
266 course following collection of more XP-deficient tumours of the various XP proteins in the NER  
267 pathway. These tumours are extremely rare however and is beyond the scope of this paper for  
268 comprehensive collection and analysis.

269

270

### 271 **Transcriptional strand asymmetries at indels may be a general feature**

272 Finally, to understand whether our observations were restricted to mononucleotide tracts or if  
273 they could be a more generic mechanistic feature of indel mutagenesis, we attempted to explore  
274 other types of indels. The limitation is the difficulty in assigning other types of indels to specific  
275 strands. It was, however, possible to ascribe strandedness to dinucleotide repeat tracts. There  
276 were some caveats: palindromic GC/CG and AT/TA dinucleotides could not be oriented, and  
277 AA/TT/GG/CC dinucleotides were excluded because these are similar to mononucleotide  
278 polyA/T/G/C's respectively. This left us with eight types of poly-dinucleotide repeat tracts that  
279 we could analyse (GT/TG/AC/CA/CT/TC/AG/GA), (Supplementary Figure 8). Indeed,  
280 correcting for background asymmetries in the genome, we observed transcriptional strand  
281 asymmetries for several poly-dinucleotide repeat tracts (Figure 4d). This was most marked

282 amongst tumour types where MSI was prevalent (Figure 4d, Supplementary Figure 9, Mann-  
283 Whitney U tests with Bonferroni correction). Furthermore, strand asymmetry in insertions was  
284 stronger than in deletions (Wilcoxon signed-rank tests with Bonferroni correction,  
285 Supplementary Figure 10a-b), with the exception of MSI tumours (Supplementary Figure 10a-b).  
286 Thus, our findings appear to be applicable to motifs other than mononucleotide repeat tracts.

287

288

## 289 **Discussion**

290

291 In this work, we have described a method to investigate transcriptional strand asymmetries for  
292 indels. Unexpectedly, we found biases in the distribution of mononucleotide repeat tracts in the  
293 reference genome at transcribed regions, and the bias is more pronounced for longer tracts. This  
294 bias needs to be considered when exploring transcriptional strand asymmetries for indels  
295 overlapping mononucleotide repeat tracts.

296

297 Our analysis demonstrates strong and previously-undescribed transcriptional strand asymmetries  
298 of indels. Our results implicate particular DNA repair pathways, namely TC-NER and MMR as  
299 contributing factors to the observed strand biases at indels (Figures 3-4). We further reveal that  
300 the formation of insertions is largely TC-NER- dependent, while the formation of deletions is  
301 additionally reliant on MMR, thus reinforcing how distinct mechanisms underpin the formation  
302 of different classes of indel.

303

304

305

306

307

308

309

310

311

312

313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358

359 **Methods**

360  
361 **Mutation calling**  
362

363 Data were obtained from whole genome sequenced (WGS) cancers from ICGC under the project  
364 PanCancer Analysis of Whole Genomes (PCAWG)<sup>31</sup>. They included 46 cancer projects from 21  
365 organs. In total, 2,575 whole genome sequenced patients were analysed using the GRCh37  
366 (hg19) reference assembly of the human genome.

367  
368 Somatic indel calls were performed using three pipelines from four somatic variant callers. These  
369 were the Wellcome Sanger Institute pipeline, the DKFZ/ EMBL pipeline and the Broad Institute  
370 pipeline<sup>31</sup>, with somatic variant false discovery rate of 2.5%. Indel calling was performed by  
371 those algorithms and only indels called by at least two of the callers were analysed<sup>31</sup>, therefore  
372 generating a conservative dataset (Supplementary Table 1). As a result, the false negative rate of  
373 indel detection could be higher than that of other methods, and of each pipeline separately, which  
374 implies that many indels present in the samples were not identified successfully. However,  
375 because of the large number of WGS tumour samples available, a sufficient number of indels  
376 remained (Supplementary Table 1). Finally, for a small subset of indels, the indel calls were  
377 visually examined using JBrowse Genome Browser<sup>32</sup>, to inspect the number of reads reporting  
378 the indel, if the indel calls were biased towards the end of the sequencing reads or if there were  
379 other systematic biases between the normal and tumour sequencing reads; such biases could not  
380 be identified.

381  
382 Bedtools intersect utility was used to measure overlap between indels and polyN tracts. The term  
383 overlap in this context refers to deleted bases occurring at any position across the entire length of  
384 the repeat or inserted bases occurring at any position across the length of the repeat and  
385 immediately before or after the repeat. Indel density was defined as the number of indel  
386 mutations for a given number of bases.

387  
388 The distance between each pair of consecutive indels was calculated per patient. Indels in  
389 different chromosomes were excluded since we could not define their pairwise distance. The  
390 same analysis was also performed separately for insertions and deletions to generate  
391 Supplementary Figure 1a.

392  
393 Indels from HAP1 cells with MutS homolog 6 (*MSH6*) knock-out were obtained from<sup>29</sup>. Indels  
394 from cells exposed to various polycyclic aromatic hydrocarbons (PAHs), namely benzo[a]pyrene  
395 [0.39  $\mu$ M and 2  $\mu$ M] and benzo[a]pyrene diolepoxide [0.125  $\mu$ M], were obtained from<sup>33</sup> to  
396 examine transcriptional strand asymmetry at indels overlapping polyN motifs in experimental  
397 settings.

398  
399 Substitution calling was performed using four somatic mutation-calling algorithms, with  
400 mutation calls being shared by at least two algorithms<sup>31</sup>. For lung cancers, C>A substitutions  
401 were examined with respect to transcriptional strand asymmetries at polyG tracts and replication  
402 timing (Supplementary Figure 6e).

403  
404 Mutational enrichment at MSI over MSS samples was defined as:

405 Ratio: (Proportion of indels overlapping polyA/T at MSI samples) / (Proportion of indels  
406 overlapping polyA/T at MSS samples)

407

408

### 409 **Transcriptional strand asymmetries at the human genome.**

410

411 Gene annotation from Ensembl was followed<sup>34</sup> and genes were downloaded from Biomart  
412 (<http://grch37.ensembl.org/biomart/martview/c1d06f3affb6260c0cd7147bb4c3b6a8>) using Gene  
413 start and Gene end to define genes and filtering by only including protein-coding genes and we  
414 also selected the attributes Strand and Gene Name. BEDTools utilities v2.21.0 were used to  
415 manipulate genomic files and intervals<sup>35</sup>. GC-skew is a measure of bias in the number of Gs or  
416 Cs between the template and non-template strands. GC-skew was calculated as  $(G-C) / (G+C)$   
417 for windows of 100 bp around the TSS and TES. Similarly, AT-skew was calculated as  $(A-T) /$   
418  $(A+T)$  for windows of 100 bp around the TSS and TES (Supplementary Figure 2a-b).

419

420 Genes in the positive and negative orientations were separated to determine the direction of gene  
421 transcription. Scripts were written in python to identify non-overlapping polyN motifs of size 1-  
422 10bp as well as dinucleotide motifs of length 2-10bp genome-wide and orient them in terms of  
423 transcription direction at genic regions.

424

425 Template motifs were the motifs in: i) positive gene orientation and negative genome strand, ii)  
426 negative gene orientation and positive (reference) genome strand. Non-template motifs were the  
427 motifs in: i) positive gene orientation and positive genome strand (reference), ii) negative gene  
428 orientation and negative genome strand. Bedtools *intersect* utility was used to calculate motif  
429 occurrences in template and non-template strands across genic regions.

430

431 To investigate the effect of the distance from the TSS and the TES across the gene length, for  
432 genes with unequal gene length, we divided each gene into ten genomic bins of equal size. Also,  
433 two additional bins upstream from the TSS and two bins downstream of the TES, each 10kB  
434 in size, were added. Then, we calculated the frequency and the strand asymmetry bias of polyN  
435 motifs in each genic bin (Figure 2a, Supplementary Figure 3). In particular, we calculated the  
436 density of polyNs at a bin as the number of polyNs over the total number of bases at that bin.  
437 However, we derived the enrichment of polyNs at the bin by comparing the ratio of the density at  
438 the bin against the density across all bins.

439

440 Relative enrichment of a polyN tract at a bin was calculated as:

441  $\text{Enrichment} = (\text{Density of polyN motif at bin}) / (\text{Density of polyN motif across all bins})$

442

443 Strand asymmetry bias was calculated as:

444  $(\text{motif occurrences at non-template strand}) / (\text{motif occurrences at template strand})$

445

446 The distribution of polyN motifs at the template and non-template strands relative to the TSS and  
447 the TES were calculated with bedtools *intersect* command using the gene orientation approach  
448 described earlier to generate (Figure 2c-d), (Supplementary Figure 3a-g). Bootstrapping using  
449 random sampling of genes with replacement was performed from which the standard deviation

450 of the strand asymmetry bias was calculated. For Figures 2c-d the interval used was 100 bp and  
451 error bars represent standard error from bootstrapping with replacement (1,000 fold).

452

453

#### 454 **Template / Non-template strand asymmetries in cancer.**

455

456 The numbers of indels overlapping motifs found in the template or non-template strands were  
457 obtained using the bedtools *intersect* command. Strand bias was calculated for the vector of  
458 genes, reporting the number of polyN motif occurrences and the number of overlapping motifs  
459 as:

460

461  $A = (\text{indels overlapping motif at non-template}) / (\text{motif occurrences at non-template})$

462  $B = (\text{indels overlapping motif at template}) / (\text{motif occurrences at template})$

463  $\text{Strand bias} = A / (A+B)$

464 with motifs representing polyN repeat tracts of size 2-10bp and dinucleotide repeat tracts of 1-5  
465 repeated units, at genic regions (Figure 3a-d, Figure 4a-d).

466

467 We performed bootstrapping with replacement, randomly selecting the indels overlapping motifs  
468 at template and non-template strands from each randomly-selected gene, for equal number of  
469 genes in multiple iterations, from which we calculated the standard deviation for the strand bias.

470

471 MMR-deficient samples were identified using genome plots and mutational signature profiles of  
472 each patient for stomach, uterus and colorectal tumours. Subsequently, transcriptional strand  
473 asymmetry levels at indels overlapping polyT tracts were compared between MSS and MSI  
474 samples to investigate the role of mismatch repair in transcriptional strand asymmetries (Figure  
475 3b, Supplementary Figure 9).

476

477

#### 478 **RNA-seq analysis.**

479

480 For the comparative analysis between expression levels and transcriptional strand asymmetry,  
481 cell of origin cell lines, where available, were used from Roadmap Epigenomics project<sup>36</sup>  
482 (Supplementary Table 2). For each cell line, genes were grouped in expression level quantiles,  
483 namely “low”, “medium” and “high” based on the associated RPKM gene expression values.  
484 The groups were defined using the 33<sup>rd</sup> and 66<sup>th</sup> percentiles from the RPKM gene expression  
485 values for protein-coding genes.

486

487 Transcriptional strand asymmetry at indels overlapping polyN motifs was investigated in relation  
488 to gene expression levels to generate Figure 3d, Figure 4b-c, Supplementary Figure 6c-d.

489

490 For lung cancer, using cell of origin RNA-seq data (IMR-90) from Roadmap Epigenomics  
491 project<sup>36</sup>, polyG tracts were grouped according to their length to investigate if the length of  
492 polyG tracts was associated with transcriptional strand asymmetry at indels across the gene  
493 expression quantiles (Figure 3d).

494

495

496 **XPC dataset.**

497 A cutaneous malignancy derived from a patient with an autosomal recessive DNA repair defect  
498 called Xeroderma Pigmentosum (XP) mutation was obtained from <sup>37</sup>. The patient was a  
499 compound heterozygote for the XPC gene. We performed the non-template strand asymmetry  
500 analysis for insertions and deletions overlapping polyT tracts (Supplementary Figure 7a-d). To  
501 control for the lower number of indels in the patient we randomly selected equal number of  
502 insertions and deletions in each cancer type, weighting for the observed transcriptional strand  
503 asymmetry at polyT tracts in each cancer type and we compared the transcriptional strand  
504 asymmetry profile of the XP sample to that of each cancer type and calculated the associated z-  
505 score and p-value from 10,000-fold bootstrapping this process for each cancer type  
506 (Supplementary Figure 7a-d).

507  
508

509 **Replication timing analysis.**

510

511 Repli-Seq data for IMR-90 cell line were obtained from The ENCODE Project Consortium<sup>38</sup> and  
512 replication domains were generated using the observed Repli-seq signal<sup>4</sup>. Genes were grouped  
513 across five replication timing quantiles and transcriptional strand asymmetry at indels  
514 overlapping polyG tracts within transcribed regions was calculated for each quantile  
515 (Supplementary Figure 6e). The same type of analysis was performed for lung cancer C>A  
516 (G>T) substitutions to investigate the contribution of replication timing to the levels of  
517 transcriptional strand asymmetries at substitutions and indels overlapping polyG motifs of 2-  
518 10bp length (Supplementary Figure 6e).

519

520 Leading and lagging orientation of the replication machinery across the human genome was  
521 inferred for MCF-7 cell line with Repli-seq data by using the finite difference approximations of  
522 second and first derivatives<sup>4</sup>. Subsequently, polyN motifs were separated into those occurring in  
523 the leading orientation and those in the lagging orientation. The indel transcriptional strand  
524 asymmetry analysis was performed separately for polyT and polyG motifs occurring at the  
525 leading and lagging orientations therefore controlling for the effect of replication orientation  
526 (Supplementary Figure 5).

527

528

529

530 Statistical analyses across the manuscript were performed in python with packages “math”,  
531 “scipy”, “pandas”, “scikit-learn” and “numpy”. Figures across the manuscript were generated in  
532 python using packages “matplotlib”, “seaborn” and “pandas”.

533

534

535 **Data Availability.**

536

537 Relevant files including mutation data count tables can be found here:

538 <https://data.mendeley.com/datasets/kdywxnn729/3>

539

540 Primary mutation data were obtained from ICGC under the project PanCancer Analysis of  
541 Whole Genomes (PCAWG)<sup>31</sup>. A cutaneous malignancy derived from a patient with an autosomal

542 recessive DNA repair defect called Xeroderma Pigmentosum (XP) mutation was obtained from  
543 <sup>37</sup>. Indel mutational profiles of non-cancerous human cells exposed to various polycyclic  
544 aromatic hydrocarbons (PAHs) including benzo[a]pyrene [0.39  $\mu$ M and 2  $\mu$ M] and  
545 benzo[a]pyrene diol epoxide [0.125  $\mu$ M] were derived from  
546 [ftp://ftp.sanger.ac.uk/pub/cancer/Zou et al 2017](ftp://ftp.sanger.ac.uk/pub/cancer/Zou_et_al_2017) and experimentally-generated mutation  
547 patterns from CRISPR-Cas9 knockouts of a human cancer cell line for *MSH6* were derived from  
548 <https://data.mendeley.com/datasets/m7r4msjb4c/2> . All data is available from the authors upon  
549 reasonable request.

550  
551  
552

### 553 **Code Availability.**

554  
555

All associated code has been deposited in <https://data.mendeley.com/datasets/kdywxnm729/3> and  
556 are available from the authors upon reasonable request.

557  
558  
559  
560  
561

### 562 **Author contributions**

563

IGS, MH and SNZ conceived the concepts and analytical framework and drove the intellectual  
564 exercise. IGS wrote the code for analysing and presenting the data. SEM generated the XPC-  
565 deficient tumour data under the supervision of SNZ. IGS, MH and SNZ wrote the manuscript  
566 with the help of GK, SEM and JJ.

567  
568  
569

### 569 **Acknowledgements**

570

MH is supported by the Wellcome Trust Sanger Institute core grant. SNZ is funded by a CRUK  
571 Advanced Clinician Scientist Award (C60100/A23916) and a CRUK Grand Challenge Award  
572 (C60100/A25274). JJ is funded by the Swiss National Science Foundation (31003B-170267).

573  
574  
575

### 575 **Competing Interests**

576

SNZ has patent applications with the UK IPO. SNZ is also a consultant for Artios Pharma Ltd,  
577 Astra Zeneca and the Scottish Genomes Partnership. The remaining authors declare no  
578 competing interests.

579  
580  
581  
582  
583  
584  
585  
586  
587

588  
589  
590  
591

## References.

592. Nik-Zainal, S. et al. Mutational Processes Molding the Genomes of 21 Breast Cancers. *Cell* vol.  
593 149 979–993 (2012).
594. Polak, P. et al. Cell-of-origin chromatin organization shapes the mutational landscape of cancer.  
595 *Nature* **518**, 360–364 (2015).
596. Haradhvala, N. J. et al. Mutational Strand Asymmetries in Cancer Genomes Reveal  
597 Mechanisms of DNA Damage and Repair. *Cell* **164**, 538–549 (2016).
598. Morganella, S. et al. The topography of mutational processes in breast cancer genomes. *Nat.*  
599 *Commun.* **7**, 11383 (2016).
600. Polak, P. & Arndt, P. F. Transcription induces strand-specific mutations at the 5' end of human  
601 genes. *Genome Res.* **18**, 1216–1223 (2008).
602. Strand, M., Earley, M. C., Crouse, G. F. & Petes, T. D. Mutations in the MSH3 gene  
603 preferentially lead to deletions within tracts of simple repetitive DNA in *Saccharomyces*  
604 *cerevisiae*. *Proceedings of the National Academy of Sciences* vol. 92 10418–10421 (1995).
605. Tran, H. T., Gordenin, D. A. & Resnick, M. A. The prevention of repeat-associated deletions in  
606 *Saccharomyces cerevisiae* by mismatch repair depends on size and origin of deletions.  
607 *Genetics* **143**, 1579–1587 (1996).
608. Tran, H. T., Keen, J. D., Krickler, M., Resnick, M. A. & Gordenin, D. A. Hypermutability of  
609 homonucleotide runs in mismatch repair and DNA polymerase proofreading yeast mutants.  
610 *Molecular and Cellular Biology* vol. 17 2859–2865 (1997).
611. Sia, E. A., Kokoska, R. J., Dominska, M., Greenwell, P. & Petes, T. D. Microsatellite instability in  
612 yeast: dependence on repeat unit size and DNA mismatch repair genes. *Molecular and Cellular*  
613 *Biology* vol. 17 2851–2858 (1997).

6140. Delacote, F. An *xrcc4* defect or Wortmannin stimulates homologous recombination specifically  
615 induced by double-strand breaks in mammalian cells. *Nucleic Acids Research* vol. 30 3454–  
616 3463 (2002).
61711. Simsek, D. & Jasin, M. Alternative end-joining is suppressed by the canonical NHEJ component  
618 *Xrcc4*–ligase IV during chromosomal translocation formation. *Nature Structural & Molecular*  
619 *Biology* vol. 17 410–416 (2010).
62012. Jasin, M. & Rothstein, R. Repair of strand breaks by homologous recombination. *Cold Spring*  
621 *Harb. Perspect. Biol.* **5**, a012740 (2013).
62213. Ghezraoui, H. *et al.* Chromosomal translocations in human cells are generated by canonical  
623 nonhomologous end-joining. *Mol. Cell* **55**, 829–842 (2014).
62414. Kass, E. M., Lim, P. X., Helgadottir, H. R., Moynahan, M. E. & Jasin, M. Robust homology-  
625 directed repair within mouse mammary tissue is not specifically affected by *Brca2* mutation.  
626 *Nature Communications* vol. 7 (2016).
62715. Nik-Zainal, S. *et al.* Landscape of somatic mutations in 560 breast cancer whole-genome  
628 sequences. *Nature* **534**, 47–54 (2016).
62916. Petrov, D. A. Mutational Equilibrium Model of Genome Size Evolution. *Theoretical Population*  
630 *Biology* vol. 61 531–544 (2002).
63117. Francino, M. P., Chao, L., Riley, M. A. & Ochman, H. Asymmetries generated by transcription-  
632 coupled repair in enterobacterial genes. *Science* **272**, 107–109 (1996).
63318. Majewski, J. Dependence of mutational asymmetry on gene-expression levels in the human  
634 genome. *Am. J. Hum. Genet.* **73**, 688–692 (2003).
63519. Denissenko, M. F. *et al.* The p53 codon 249 mutational hotspot in hepatocellular carcinoma is  
636 not related to selective formation or persistence of aflatoxin B1 adducts. *Oncogene* **17**, 3007–  
637 3014 (1998).
63820. Rodin, S. N. & Rodin, A. S. Origins and selection of p53 mutations in lung carcinogenesis.  
639 *Semin. Cancer Biol.* **15**, 103–112 (2005).

6401. Pleasance, E. D. et al. A comprehensive catalogue of somatic mutations from a human cancer  
641 genome. *Nature* **463**, 191–196 (2010).
6422. Fousteri, M. & Mullenders, L. H. F. Transcription-coupled nucleotide excision repair in  
643 mammalian cells: molecular mechanisms and biological effects. *Cell Res.* **18**, 73–84 (2008).
64423. GINNO, P. A., LIM, Y. W., LOTT, P. L., KORF, I. & CHÉDIN, F. GC skew at the 5' and 3' ends of  
645 human genes links R-loop formation to epigenetic regulation and transcription termination.  
646 *Genome Res.* **23**, 1590–1600 (2013).
64724. Ellegren, H. Microsatellites: simple sequences with complex evolution. *Nature Reviews*  
648 *Genetics* vol. 5 435–445 (2004).
64925. Garcia-Diaz, M. & Kunkel, T. A. Mechanism of a genetic glissando: structural biology of indel  
650 mutations. *Trends Biochem. Sci.* **31**, 206–214 (2006).
6526. Romanova, N. V. & Crouse, G. F. Different roles of eukaryotic MutS and MutL complexes in  
652 repair of small insertion and deletion loops in yeast. *PLoS Genet.* **9**, e1003920 (2013).
65327. Kunkel, T. A. & Erie, D. A. Eukaryotic Mismatch Repair in Relation to DNA Replication. *Annu.*  
654 *Rev. Genet.* **49**, 291–313 (2015).
65528. Baretti, M. & Le, D. T. DNA mismatch repair in cancer. *Pharmacol. Ther.* **189**, 45–62 (2018).
65629. Zou, X. et al. Validating the concept of mutational signatures with isogenic cell models. *Nat.*  
657 *Commun.* **9**, 1744 (2018).
65830. Hanawalt, P. C. & Spivak, G. Transcription-coupled DNA repair: two decades of progress and  
659 surprises. *Nat. Rev. Mol. Cell Biol.* **9**, 958–970 (2008).
66031. ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analysis of  
661 whole genomes. *Nature* **578**, 82–93 (2020).
66232. Buels, R. et al. JBrowse: a dynamic web platform for genome visualization and analysis.  
663 *Genome Biol.* **17**, 66 (2016).
66433. Kucab, J. E. et al. A Compendium of Mutational Signatures of Environmental Agents. *Cell* **177**,  
665 821–836.e16 (2019).

6664. Yates, A. et al. Ensembl 2016. *Nucleic Acids Res.* **44**, D710–6 (2016).

6675. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.  
668 *Bioinformatics* vol. 26 841–842 (2010).

6696. Roadmap Epigenomics Consortium et al. Integrative analysis of 111 reference human  
670 epigenomes. *Nature* **518**, 317–330 (2015).

6717. Momen, S. et al. Dramatic response of metastatic cutaneous angiosarcoma to an immune  
672 checkpoint inhibitor in a patient with xeroderma pigmentosum: whole-genome sequencing aids  
673 treatment decision in end-stage disease. *Molecular Case Studies* vol. 5 a004408 (2019).

6748. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human  
675 genome. *Nature* **489**, 57–74 (2012).

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692 **Figure 1: Indel characteristics across cancer types. A)** The ratio of deletions to insertions for  
693 each tumour-type. (Mann-Whitney U test, p-value<0.05 per cancer type). **B)** Distribution of size  
694 of insertions and of deletions for each tumour-type. Deletions displayed greater size variance in  
695 comparison to insertions across cancer types (Levene's test, p-value<0.05) and for individual  
696 cancers types (p-value<0.001 in breast, pancreas, liver, ovary, skin, lung, cervix, bone, head /  
697 neck, colorectal, p-value<0.05 in biliary and lymphoid cancers).

698

699

700

701 **Figure 2: Strand asymmetries of polynucleotide (polyN) repeat tracts within transcribed**  
702 **regions. A)** Enrichment of various polyN motifs across genes. Each gene is divided into ten  
703 bins, and two additional bins are added at either end of each gene. For any given bin, blue  
704 indicates relative enrichment in comparison to all other bins for that polyN, whereas red  
705 indicates relative depletion. **B)** Scheme depicting the identification of polyN motifs on the  
706 template (blue) or non-template (orange) strands, dependent on the direction of the gene. RNA-  
707 polymerase II (RNAPII) binds to the template strand and mediates transcription. Thus, in the  
708 panel above, where the gene is on the (+) strand, the polyA tracts are on the non-template  
709 strand. In the panel below, where the gene is on the (-) strand, the polyA tracts on the template  
710 strand. **C)** Density of polyT and polyG motifs around the transcription start site (TSS). The  
711 gradient of pink to purple represents polyG tracts of 1-5bp length, whereas the gradient of light  
712 blue to dark blue represents polyT tracts of 1 to 5bp length. Error bars represent standard error  
713 from 1,000-fold bootstrapping. **D)** Density of polyT and polyG motifs around the transcription  
714 end site (TES). Error bars represent standard error from 1,000-fold bootstrapping.

715

716

717

718 **Figure 3: Transcriptional strand asymmetry of indels that occur at polyN tracts across**  
719 **multiple cancer types. A)** Transcriptional strand asymmetries of indels occurring at polyT  
720 motifs. Average bias is shown with error bars showing standard deviation after 1,000  
721 bootstraps. Myeloid, cervix and thyroid cancers were excluded due to low numbers of total  
722 indels (Supplementary table 1). T=template, NT=non-template. Strand bias was calculated as  
723 mutational density of non-template strand over total mutational density (of non-template and  
724 template strands). **B)** Strand bias of MSI and MSS samples in stomach, biliary, uterus and  
725 colorectal tumours (Mann-Whitney U p-value<0.001 in all cases, Bonferroni corrected). **C)**  
726 Transcriptional strand asymmetries of indels occurring at polyG motifs. Average bias is shown,  
727 with error bars showing standard deviation from bootstrapping. **D)** Relationship between indel  
728 strand bias and gene expression levels in lung cancer (Mann-Whitney U p-value<0.001 for  
729 comparisons between low and medium expressed genes and between medium and highly  
730 expressed genes) according to length of polyG tracts (Kruskal-Wallis H-test with Bonferroni  
731 correction, p-value<0.001 for medium and high expression genes, p-value>0.05 for low  
732 expression genes). **E)** Scheme depicting mechanism of indel mutagenesis at poly-T tracts. DNA  
733 damage, shown as asterisks (\*) that arise at T nucleotides of poly-T tracts can occur on both  
734 template and non-template strands. The subsequent DNA repair, postulated to be TC-NER,  
735 results in preferential correction of DNA damage on the template strand, leaving T insertions  
736 (highlighted in as red T's) and T deletions (shown as red -) on the non-template strand. **F)**  
737 Schematic depicting mechanism of indel mutagenesis at poly-G tracts in lung cancers from  
738 smokers. DNA damage in the form of adducted guanines (\*) are asymmetrically repaired by TC-  
739 NER, with preferential repair of the template strand, thus accumulating more G indels on the  
740 non-template strand.

741  
742  
743

744

745 **Figure 4: Transcriptional strand asymmetry at insertions and deletions. A)** Transcriptional  
746 strand asymmetry of insertions and deletions at polyT tracts. Error bars represent standard  
747 deviation from bootstrapping with replacement. Both insertions and deletions displayed a strand  
748 asymmetry bias towards the non-template strand for polyT tracts across cancer types (Binomial  
749 test with Bonferroni correction,  $p$ -value $<0.001$  for insertions and  $p$ -value $<0.05$  for deletions). **B)**  
750 Transcriptional strand asymmetry occurring at polyT tracts according to level of gene expression  
751 for insertions. Mann-Whitney U with Bonferroni correction,  $p$ -value $<0.001$  when comparing low  
752 and high expression gene sets across all cancer types except skin, ovarian and lymphoid  
753 cancers ( $p$ -value $<0.05$ ) and CNS ( $p$ -value $>0.05$ ). **C)** Transcriptional strand asymmetry occurring  
754 at polyT tracts according to level of gene expression for deletions. Mann-Whitney U with  
755 Bonferroni correction,  $p$ -value $<0.001$  when comparing low and high expression gene sets for  
756 skin and  $p$ -value $<0.05$  for stomach and pancreatic cancers. **D)** Hierarchical clustering displaying  
757 transcriptional strand asymmetries for indels overlapping dinucleotide motifs. Dinucleotide  
758 repeat tracts of up to five repeated units are displayed. Purple represents asymmetry towards  
759 the non-template strand, whereas orange represents asymmetry towards the template strand. In  
760 the dendrogram of cancers, biliary, uterus, colorectal and stomach cancers are more distant  
761 from the other cancers, and contain MSI samples, while lung cancers are also separable from  
762 other cancer types, further reinforcing our observations regarding the DNA damage and repair  
763 processes that contribute to the observed asymmetries. Across cancer types a non-template  
764 strand asymmetry preference was observed for TG, TC and CT motifs (Binomial test with  
765 Bonferroni correction,  $p$ -value $<0.001$ ) and for GT motifs (Binomial test with Bonferroni  
766 correction,  $p$ -value $<0.05$ ) and a template strand asymmetry for CA, GA and AG motifs (Binomial  
767 test with Bonferroni correction,  $p$ -value $<0.001$ ) and for AC motifs (Binomial test with Bonferroni  
768 correction,  $p$ -value $<0.05$ ).