

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

# 5-hydroxymethylcytosine marks regions with reduced mutation frequency in human DNA

Marketa Tomkova<sup>a</sup>, Michael McClellan<sup>a</sup>, Skirmantas Kriaucionis<sup>a1</sup> and Benjamin Schuster-Böckler<sup>a1</sup>

## AFFILIATION

- a. Ludwig Cancer Research Oxford  
University of Oxford  
Old Road Campus Research Building  
Oxford OX3 7DQ  
United Kingdom

## KEYWORDS

Somatic mutation;DNA methylation; Cancer genomics;Hydroxymethylcytosine

## CORRESPONDING AUTHORS

1. Co-corresponding authors

- Benjamin Schuster-Böckler  
[benjamin.schuster-boeckler@ludwig.ox.ac.uk](mailto:benjamin.schuster-boeckler@ludwig.ox.ac.uk)
- Skirmantas Kriaucionis  
[skirmantas.kriaucionis@ludwig.ox.ac.uk](mailto:skirmantas.kriaucionis@ludwig.ox.ac.uk)

Ludwig Cancer Research Oxford  
University of Oxford  
Old Road Campus Research Building  
Oxford OX3 7DQ  
United Kingdom

26 **ABSTRACT**

27 CpG dinucleotides are the main mutational hot-spot in most cancers. The characteristic elevated C>T  
28 mutation rate in CpG sites has been related to 5-methylcytosine (5mC), an epigenetically modified  
29 base which resides in CpGs and plays a role in transcription silencing. In brain nearly a third of 5mCs  
30 have recently been found to exist in the form of 5-hydroxymethylcytosine (5hmC), yet the effect of  
31 5hmC on mutational processes is still poorly understood. Here we show that 5hmC is associated with  
32 an up to 53% decrease in the frequency of C>T mutations in a CpG context compared to 5mC. Tissue  
33 specific 5hmC patterns in brain, kidney and blood correlate with lower regional CpG>T mutation  
34 frequency in cancers originating in the respective tissues. Together our data reveal global and  
35 opposing effects of the two most common cytosine modifications on the frequency of cancer  
36 causing somatic mutations in different cell types.

37 **BACKGROUND**

38 Cancer genomics projects have revealed that the distribution of somatic mutations across the  
39 genome is not uniform (Lawrence et al., 2013). Apart from positive and negative selective pressure,  
40 a number of factors can influence regional mutation frequencies, such as chromatin organisation  
41 (Schuster-Böckler & Lehner, 2012), replication timing (Koren et al., 2012), metabolic load (Ames et  
42 al., 1993) and exposure to different mutagens (Poon et al., 2013). Furthermore, highly transcribed  
43 regions generally exhibit lower mutation frequencies due to the influence of transcription-coupled  
44 repair (Lawrence et al., 2013). In addition to the regional distribution of mutations, the local  
45 nucleotide contexts and mutation types (referred to as *mutational signatures*) have been  
46 investigated extensively since they provide critical clues about the mechanism of mutagenesis. For  
47 example, consensus motifs for cytidine deaminases (such as APOBEC and AID) were found enriched  
48 at mutational hot-spots, suggesting that activity of these enzymes could be the potential cause of  
49 those mutations (Nik-Zainal et al., 2012; Taylor et al., 2013). The most frequent mutational signature  
50 found in the majority of cancers is C to T transition in a CpG dinucleotide context (CpG>T)

51 (Alexandrov et al., 2013; Lawrence et al., 2013). This relates to the fact that cytosines in CpG  
52 dinucleotides are frequently methylated to form 5-methylcytosine (5mC). The rate of spontaneous  
53 deamination of 5mC into T is four fold higher than the rate of deamination of C into U (Lindahl &  
54 Nyberg, 1974). In the germline of vertebrates, this likely facilitated a general depletion of CpGs  
55 outside of CpG islands which are largely unmethylated.

56 The genomes of all examined vertebrate species feature DNA methylation, and loss of DNA  
57 methylation is incompatible with normal development in mice (E. Li et al., 1992; Okano et al., 1999).  
58 DNA methylation plays a role in gene expression, most notably by repressing one allele of imprinted  
59 genes. Moreover, it is involved in maintenance of genome stability, alternative splicing, X  
60 chromosome inactivation and suppression of retrotransposons (Klose & Bird, 2006; Jones, 2012). In  
61 2009, 5-hydroxymethylcytosine (5hmC) was indisputably shown to exist in DNA of brain and other  
62 tissues (Kriaucionis & Heintz, 2009). It was concurrently shown that ten-eleven translocation (TET)  
63 enzymes are able to convert 5mC into 5hmC (Tahiliani et al., 2009). Unlike 5mC, which is observed at  
64 similar levels in many cell types, the abundance of 5hmC varies widely. 5hmC was observed to be  
65 particularly enriched in brain cells (Kriaucionis & Heintz, 2009; Lister et al., 2013) and detectable in  
66 embryonic stem cells and all examined tissues (Tahiliani et al., 2009; Globisch et al., 2010;  
67 Szwagierczak et al., 2010; H. Wu & Zhang, 2011). 5hmC and higher oxidised states of methyl-cytosine  
68 have been proposed to play a role in de-methylation via ineffective re-methylation after replication  
69 or directly by thymine DNA glycosylase (TDG) (Tahiliani et al., 2009; He et al., 2011; Maiti & Drohat,  
70 2011; Shen et al., 2013; Hu et al., 2014). In addition to demethylation, 5hmC has been implicated in  
71 transcriptional regulation, and a number of DNA binding proteins recognising 5hmC have been  
72 identified (Mellén et al., 2012; Spruijt et al., 2013; Takai et al., 2014). 5hmC is found depleted in  
73 primary tumours and TET2 is frequently mutated in myelodysplastic syndrome, acute myelogenous  
74 leukaemia and T-cell lymphoma, indicating that 5hmC plays a role in carcinogenesis. However, the

75 molecular mechanism by which 5hmC affects carcinogenesis is poorly understood (Rasmussen &  
76 Helin, 2016).

77 5hmC is an important intermediate during demethylation in zygotes and ES cells (Tahiliani et al.,  
78 2009; Inoue & Zhang, 2011; Wossidlo et al., 2011), but the vast majority of 5hmC is found as a  
79 stable, long-lived modification in adult mouse tissue that undergoes little cell division (Bachman et  
80 al., 2014; Brazauskas & Kriaucionis, 2014). Thus, we hypothesised that – similar to 5mC – long-lived  
81 5hmC could have a substantial influence on the mutability of DNA. Little is known about the  
82 mutational properties of 5hmC, in part because until recently there has been a lack of information  
83 on the precise location of 5hmC in the genome. With the development of techniques for single-  
84 nucleotide resolution mapping of 5mC and 5hmC (Yu et al., 2012; Booth et al., 2014), it is now  
85 possible to differentiate mutation rates at 5mC and 5hmC sites. Recently, Supek et al. (Supek et al.,  
86 2014) reported elevated C>G transversion rates at 5hmC sites, using 5hmC maps from human and  
87 mouse embryonic stem cells. However, these findings are limited by the fact that embryonic stem  
88 cells differ substantially from the somatic tissues in which mutations were observed (Schultz et al.,  
89 2015).

90 A large proportion of mutations observed in any cancer genome originate in its pre-cancerous cell of  
91 origin (Nik-Zainal et al., 2012; Stephens et al., 2012; Tomasetti et al., 2013; S. Wu et al., 2015) and  
92 will have been influenced by its epigenetic landscape. The publication of single-base resolution maps  
93 of 5mC and 5hmC occupancy in samples of human brain, kidney and blood (Wen et al., 2014; Chen  
94 et al., 2015; Pacis et al., 2015) now enables us to interrogate the tissue-specific effect of cytosine  
95 modifications on somatic mutation rates in unprecedented detail.

96 Since 5hmC has been shown to be most abundant in human brain (W. Li & Liu, 2011; Nestor et al.,  
97 2012), we have initially focussed on assessing the relationship between mutability and DNA  
98 modifications in brain cancers. Based on a DNA sequencing data from five brain cancer types

99 encompassing 665 patients, we show that the dominant mutational signature in brain cancers is  
100 CpG>T, which is modulated by the modification state of cytosine. Strikingly, the CpG>T mutation  
101 frequency of 5-hydroxymethylcytosine is reduced nearly two-fold compared to the methylated state.  
102 We find that the ratio of 5hmC to 5mC in 100kb genomic intervals correlates with CpG>T mutation  
103 frequency even after accounting for confounding factors like gene density or CpG islands. When we  
104 expand our analysis to include mutations and 5hmC maps from kidney and myeloid lineage of blood  
105 cells, we observe a clear tissue-specific effect of 5hmC on mutagenicity. Finally, we measured 5mC  
106 and 5hmC levels using methodology of high accuracy in eight different human tissue types and show  
107 that reduced 5hmC levels associate with an increased proportion of CpG>T mutations in cancers of  
108 the corresponding tissue. Together, our findings suggest that hydroxymethylation has a significant  
109 influence on the likelihood of mutations at CpG sites across many human tissue types.

## 110 **RESULTS**

111 We compiled base-resolution maps of 5mC and 5hmC frequency in brain, kidney and myeloid cells  
112 from publicly available sources (Wen et al., 2014; Chen et al., 2015; Pacis et al., 2015). All three data  
113 sets are based on bisulfite (BS) and “Tet-assisted bisulfite” (TAB) sequencing, respectively. BS-Seq  
114 detects any modified cytosine (i.e., does not distinguish 5mC and 5hmC) whereas TAB-Seq  
115 specifically detects 5hmC. The combination of the two methods allows an estimation of the levels of  
116 both 5mC and 5hmC for all sufficiently covered cytosines. As 5hmC predominantly occurs in a CpG  
117 context, we focussed the analysis on CpG sites. Sequencing reads come from heterogeneous  
118 populations of cells. Hence, a single locus usually cannot be assigned a single state (C, 5mC or 5hmC).  
119 Instead, we estimated the frequency of modification, hydroxymethylation and methylation per site  
120 using the percentage of BS-Seq reads that showed a modification (referred to as *mod level*), the  
121 percentage of TAB-Seq reads that showed hydroxymethylation (referred to as *5hmC level*) and their  
122 difference ( $5mC\ level = mod\ level - 5hmC\ level$ ), respectively.

123 **5hmC sites in brain exhibit lower frequency of CpG>T mutations than 5mC sites**

124 Since brain tissue has been shown to exhibit particularly high levels of 5hmC (Fig. 1A), we first  
125 investigated the relationship between the regional distribution of 5hmC, 5mC and mutagenesis in  
126 brain tumours. We reasoned that this approach would provide the highest sensitivity to detect any  
127 correlation between 5hmC and mutation frequency.

128 We analysed 344370 somatic single nucleotide variants (SNVs) from 665 samples derived from  
129 exome and whole genome sequencing of the following cancer types: Glioblastoma (GBM), Glioma  
130 low grade (GLG), Neuroblastoma (NRB), Medulloblastoma (MDB) and Pilocytic astrocytoma (PA)  
131 (Alexandrov et al., 2013). The dominant point mutation type in these cancers were C>T transitions in  
132 a CpG context (Fig. 1B, 1C), similar to what was observed in combined analyses of all cancer types  
133 (Alexandrov et al., 2013).

134 Mutations and DNA modifications are not distributed uniformly along the chromosomes. Strikingly,  
135 5hmC levels were visibly and significantly anti-correlated with the frequency of CpG>T ( $r=-0.71$ ,  
136 chr3), while 5mC levels displayed a positive correlation ( $r=0.66$ , chr3, Fig. 1D, Fig. 1-figure  
137 supplements 1–2). This is not a simple consequence of the uneven distribution of genes, exons, CpG  
138 islands or levels of gene expression (Fig. 1-figure supplement 3 and additional analyses below).

139 Averaging over the entire genome, the frequency of C>T mutations differed substantially between  
140  $5mC_{high}$  ( $5mC_{high}$ : *mod level* > 10% and *5hmC level / mod level* <= 0.3) and  $5hmC_{high}$  ( $5hmC_{high}$ : *mod*  
141 *level* > 10% and *5hmC level / mod level* >= 0.5) sites. The fraction of mutated  $5hmC_{high}$  sites was  
142 significantly lower than the fraction of mutated  $5mC_{high}$  sites (Fig. 1E). The lower mutation frequency  
143 was consistently observed in data derived from both exome and whole genome sequencing projects  
144 ( $P<0.001$ , Wilcoxon signed-rank test). Moreover, all brain cancer types individually displayed a  
145 significant (28–53%,  $P<0.05$  in all types) reduction of C>T mutations in  $5hmC_{high}$  sites (Fig. 2A, 2B).

146 It has been shown that CpG>T mutations are one of the two mutational signatures correlating with  
147 age (Alexandrov et al., 2015), supporting the fact that these mutations were gathered during the  
148 entire lives of the patients, not only after the origin of cancer. Here we observed that this correlation  
149 is present in both methylated and hydroxymethylated sites (Fig. 2C). Moreover, the slope for 5mC  
150 was steeper than for 5hmC, suggesting that even the mechanisms causing the difference of CpG>T  
151 mutability between 5mC and 5hmC were present in the pre-cancerous cell of origin.

152 We also compared the fraction of mutated 5mC<sub>high</sub> and 5hmC<sub>high</sub> sites for the other two possible  
153 types of mutations: C>A and C>G. As shown in Fig. 1E, C>A or C>G transversions are an order of  
154 magnitude less frequent than C>T transitions in both 5mC and 5hmC sites. The relationship between  
155 C>A and C>G mutations and 5hmC varied between cancer types. In GBM and LGG the frequency of  
156 C>A mutations was significantly higher in 5mC<sub>high</sub> compared to 5hmC<sub>high</sub> sites, but in NRB, MDB and  
157 PA we detected no significant difference. The frequency of C>G mutations in 5mC<sub>high</sub> sites differed  
158 significantly from 5hmC<sub>high</sub> sites only in MDB, PA and GBM. In MDB and PA, 5hmC<sub>high</sub> sites were  
159 slightly enriched for C>G mutations, whereas in GBM an enrichment was observed at 5mC<sub>high</sub> sites.  
160 Since C>T transitions are the most common somatic mutation type in brain and the difference in C>T  
161 mutations between 5mC<sub>high</sub> and 5hmC<sub>high</sub> sites is more consistent among cancer types than in the  
162 C>A and C>G transversions, we focus mainly on C>T mutations in the remainder of this report.

163 We confirmed that C>T mutations are significantly depleted at 5hmC sites across a wide range of  
164 thresholds in definitions of 5mC<sub>high</sub> and 5hmC<sub>high</sub> (Fig. 2-figure supplement 1A–F). In fact, more  
165 stringent definitions of 5hmC (e.g., 5hmC<sub>high</sub>: *5hmC level / mod level*  $\geq$  0.7) result in even greater  
166 differences (42–59%) in C>T mutation frequencies between 5mC<sub>high</sub> and 5hmC<sub>high</sub> sites (Fig. 2-figure  
167 supplement 1G–I, Fig. 3-figure supplement 1A–D), but these definitions would reduce the number of  
168 sites too much for our subsequent statistical analyses.

169 **Reduced 5hmC mutability in brain is not accounted for by genomic regions or gene expression**

170 We next examined whether the decreased frequency of C>T transitions at 5hmC vs. 5mC sites might  
171 be an indirect effect of 5hmC being associated with genomic regions of lower mutability. Levels of  
172 5mC and 5hmC vary across genomic regions. For example, 5hmC density is elevated in highly  
173 expressed genes in brain (Mellén et al., 2012; Yu et al., 2012; Lister et al., 2013; Wen et al., 2014).  
174 The observed decrease in C>T mutation frequencies might therefore be attributable to higher gene  
175 expression, which would correlate with higher transcription coupled repair. We therefore performed  
176 the analysis described above separately for regions with high vs. low gene expression in human brain  
177 (see Methods). There was a lower overall mutation frequency in highly expressed genes (Fig. 3A–B),  
178 but both highly and lowly expressed genes exhibited significantly lower C>T transition rates at 5hmC  
179 sites compared to 5mC sites (Fig. 3A–D). This suggests that the observed effect is independent of  
180 transcription and thus not a result of transcription coupled repair.

181 Gene expression is only one of many possible region-related confounding factors. Hence, to correct  
182 for any regional variation, we performed the analysis on groups of sites generated by pairing the  
183 modified CpGs: each 5hmC site was paired with the nearest yet unpaired 5mC site from an  
184 equivalent genomic and sequence context (an approach adapted from (Supek et al., 2014), see  
185 Methods). Thereby we compared the mutation frequencies of two groups (one group comprising  
186 5mC sites and one group comprising 5hmC sites) containing the same number of loci (6801374  
187 cytosines in each group). As a result of this experimental setup, a substantial fraction of mutated  
188 5mC sites were excluded, greatly reducing the statistical power of this “paired” analysis.  
189 Nevertheless, the frequency of C>T mutations in 5hmC remained significantly lower than in 5mC in  
190 both exomes and genomes (Fig. 3E–F), supporting that the difference between 5mC and 5hmC  
191 mutation frequency is not caused by regional differences.

192 To ensure that there is no confounding bias in the spatial distribution of mutations around 5mC or  
193 5hmC sites, respectively, we plotted mutation frequencies in a 2kb radius up and downstream of  
8

194 modified loci (Fig. 3-figure supplement 1G, Methods). The mutation frequency differed substantially  
195 in the aligned positions of DNA modifications but was indistinguishable in the surrounding area. In  
196 conclusion, regional mutation rate variability is unlikely to account for the difference in C>T  
197 mutational load in 5mC and 5hmC sites.

### 198 **Relative 5hmC correlates with CpG>T mutation frequency**

199 The 5mC and 5hmC frequency at each base reflect the prevalence of each modification within the  
200 sequenced cell population. We hypothesised that if 5hmC had a direct effect on C>T mutation  
201 likelihood, we would observe an increase in mutation frequency with decreasing 5hmC occupancy.  
202 To test this, we formally defined  $5hmC_{rel}$  as the ratio of the proportion of reads supporting 5hmC,  
203 relative to the proportion of reads supporting any modification at a particular cytosine ( $5hmC_{rel} =$   
204 *5hmC level / mod level*). We then divided brain CpG sites into bins according to their  $5hmC_{rel}$  level  
205 and computed the fraction of mutated sites in each bin (Fig. 4A). We observed a clear linear  
206 relationship between  $5hmC_{rel}$  values and C>T mutation frequencies. Notably, the correlation was  
207 present in all the tested brain cancer types in exome- and whole genome-sequenced samples. A  
208 regression slope test confirmed significance of this relationship in all the cancer types. To confirm  
209 that the results are not influenced by an uneven distribution of information across bins, we  
210 performed quantile binning so that each bin contains an approximately equal number of positions  
211 (see Methods). The results of quantile bins were equivalent to evenly spaced bins (Fig. 4-figure  
212 supplement 1H).

213 For comparison, we also evaluated the relationship between  $5hmC_{rel}$  and the frequency of C>A and  
214 C>G mutations (Fig. 4A). Consistent with our previous results, an increase in  $5hmC_{rel}$  is associated  
215 with an increase in C>G mutations in whole genomes (from MDB and PA samples), but the  
216 relationship in other cancer types shows no significant trend. C>A mutations decrease with  $5hmC_{rel}$   
217 levels in GBM but exhibit no significant signal in the remaining tumour types.

218 This result supports the conclusion that the decrease in C>T mutation frequency at 5hmC sites is not  
219 an artefact of our chosen definition of 5mC or 5hmC. Even more importantly, it supports the notion  
220 that this decrease is directly caused by the properties of these DNA modifications.

#### 221 **Mutation load of 5hmC sites is similar to unmodified cytosines**

222 The findings reported so far could be attributed to an elevated mutation rate in 5mC, to a lowered  
223 mutagenicity of 5hmC or a combination of the two. To investigate this question, we compared  
224 mutation frequencies at 5mC and 5hmC sites to that of unmodified cytosines. We divided all the  
225 sequenced CpG sites into 9x9 bins according to their levels of 5mC and 5hmC. We observed that the  
226 mutation frequency of unmodified cytosine is similar to 5hmC, whereas 5mC exhibited much higher  
227 mutation frequency (Fig. 4B). Further, we calculated the mutation frequency distribution in sites that  
228 exhibited almost no methylation or almost no hydroxymethylation, respectively. When methylated  
229 sites are excluded, the mutation frequency does not show any significant trend with increasing levels  
230 of 5hmC (Fig. 4C). Conversely, excluding hydroxymethylated sites leads to a significant gradient in  
231 mutation frequency with increasing levels of 5mC (Fig. 4D). When only fully modified sites (*mod level*  
232  $\geq 90\%$ ) are taken into account, increasing levels of 5hmC (i.e., decreasing levels of 5mC) are  
233 associated with a significant decrease in C>T mutation frequency (Fig. 4E).

#### 234 **5hmC is a predictor of CpG>T mutation frequency across the genome**

235 To examine the exclusive impact of DNA modifications on regional frequencies of mutations, we split  
236 the genome into 100kb windows and fitted a generalised linear model to explain the observed per-  
237 window CpG>T mutation frequency from a combination of features including average 5mC and  
238 5hmC levels, 5hmC<sub>rel</sub>, gene density, CpG island density amongst others. Only whole genome  
239 sequencing data were used for this analysis. To compare the resulting models, we calculated their  
240 respective “explained deviance”  $D^2$ , a generalisation of explained variance that is more appropriate  
241 for comparing generalised linear models (see Methods).

242 The best individual predictor of CpG>T mutation frequency was 5hmC<sub>rel</sub> ( $D^2 = 0.11$ ), outperforming all  
243 other features (Fig. 5A). Interestingly, the sum of 5mC and 5hmC levels (“mod”) performed worst,  
244 suggesting that bisulfite sequencing measurements alone are a poor predictor of mutagenicity.  
245 When combining all 11 features into one model, the total explained deviance for 100kb windows  
246 was 16%.

247 Varying the chosen window size (3kbp–3Mbp; Fig. 5B, Fig. 5-figure supplement 1A–C) does not  
248 substantially change the comparison of the predictive power of the respective features. In all cases,  
249 5mC and 5hmC<sub>rel</sub> were the two best predictors, with 5hmC<sub>rel</sub> performing slightly better with smaller  
250 windows. However, the total explained deviance increased with window size, reaching values as high  
251 as 45% for univariate models and 60% for models with all predictors. This led us to question whether  
252 the increasing predictive power of larger windows has a biological reason, or whether it is a  
253 consequence of the lower data density in small windows.

254 Since many smaller windows contain no observed mutations, low  $D^2$  values could simply reflect a  
255 lack of data. To test this, we generated simulated mutations so that a “perfect” predictor was  
256 linearly related to the mutation likelihood per window (see Methods). We then measured the effect  
257 of window and sample size (number of patients) on the observed  $D^2$ , repeating the simulations 10  
258 times. The resulting curves of the explained deviance resemble those measured in the real data (Fig.  
259 5-figure supplement 1D). Moreover, in the simulated data, higher numbers of patients lead to higher  
260  $D^2$  even for smaller window sizes, suggesting that lower  $D^2$  values in smaller windows indeed are a  
261 consequence of lower data density.

## 262 **Level of genic 5hmC correlates with decrease of CpG>T**

263 It has been reported that 5hmC is enriched in gene bodies, and several brain cancer sequencing data  
264 sets. We therefore tested whether the relationship between 5hmC and mutations, which we  
265 observed across the whole genome, is also detectable in exonic regions alone.

266 In line with our earlier results, we found that 5hmC<sub>rel</sub> significantly contributes to the deviance  
267 explained by the model, beyond covariation with gene expression (Fig. 5C–D; F-test  $p < 2e-100$ ). We  
268 hypothesised that this effect should be most pronounced when using modC>T and CpG>T as the  
269 response variable, whereas it should decrease when using definitions of mutations that include a  
270 progressively wider range of loci (C>T, C>N, N>N). Indeed, the unique contribution of 5hmC<sub>rel</sub> to the  
271 explained gene mutation frequency decreased as the mutation sets became larger and more distant  
272 from modC>T (Fig. 5C–D, Fig. 5-figure supplements 2–3). Nevertheless, in all of the cases, 5hmC<sub>rel</sub>  
273 significantly improved the fit of the model. Conversely, we confirmed that 5hmC<sub>rel</sub> had no significant  
274 predictive power for the frequency of T>N mutations (Fig. 5C–D; column T>N), supporting the  
275 hypothesis that 5hmC<sub>rel</sub> selectively affects mutations in modified cytosines.

#### 276 **Decreased CpG>T mutation frequency in 5hmC is not limited to brain tissue**

277 Two recently published datasets allowed us to address the question of mutational properties of 5mC  
278 and 5hmC also in two other tissues: kidney (Chen et al., 2015) and blood (Pacis et al., 2015). For  
279 blood we used 174 sequencing samples from Acute Myeloid Leukaemia (AML) as the cancer type  
280 closest to the blood dendritic cells in which the BS-Seq and TAB-Seq measurements were performed.  
281 For kidney we combined 585 samples from four distinct sequencing projects, covering Kidney Clear  
282 Cell, Kidney Papillary and Kidney Chromophobe carcinomas.

283 Matching our findings in brain, 5hmC sites were mutated significantly less frequently than 5mC sites  
284 in both tissue types (Fig. 6B), irrespective of whether genome or exome sequencing data were used.  
285 Moreover, a similar difference was present in all available replicates of the BS-Seq and TAB-Seq  
286 measurements (6 for blood, 2 for kidney, Fig. 6-figure supplement 1A).

287 Genomic distribution of 5hmC differs substantially between the three tissue types (Fig. 6-figure  
288 supplement 2). Consequently, we expected the association between mutation frequency and 5hmC  
289 to be highest when mutation and modification data are derived from matching tissue types. To test

290 this hypothesis, we used a GLM on genomic windows of 100kbp to predict CpG>T mutation rate  
291 from a combination of 5hmC<sub>rel</sub> maps of all three tissues. Strikingly, for each cancer type, the best  
292 predictor came from the same tissue type (Fig. 6A), suggesting that tissue differences in 5hmC are  
293 reflected in the CpG>T mutation landscape. The same results were obtained in all available  
294 replicates of the 5hmC<sub>rel</sub> maps (Fig. 6-figure supplement 1B). Finally, we added a 5hmC<sub>rel</sub> map  
295 derived from embryonic stem cells (ESC) as an additional predictor, to compare our findings to  
296 previously reported results (Supek et al., 2014). As we anticipated, the ESC-derived 5hmC levels have  
297 substantially lower predictive power on CpG>T mutation rate than any of the tissue-derived maps.

298 While base-resolution maps of 5hmC for human tissue are still rare, there is a wide range of BS-Seq  
299 data sets available in public databases. Given our findings thus far, we predicted that tissues with  
300 high levels of 5hmC relative to 5mC would exhibit fewer CpG>T mutations in modified sites than  
301 tissues with low total 5hmC. To test this hypothesis, we measured total levels of 5mC and 5hmC  
302 using High Pressure Liquid Chromatography (HPLC-UV) in DNA of eight human tissue types for which  
303 BS-Seq maps are publicly available (Fig. 6-figure supplement 3). In order to account for unrelated  
304 cancer-type specific differences in CpG>T mutability, we normalised the mutation frequency in  
305 modified sites by the mutation frequency in unmodified sites. The analysis of association between  
306 genomic 5hmC and enrichment of CpG>T mutations revealed a strong negative correlation (Fig. 6C)  
307 in all tissue types except lung. Interestingly, this difference seems to stem from smoking-related  
308 effects. Lung cancer mutation data from heavy smokers revealed a markedly lower frequency of  
309 CpG>T mutations in modified sites, relative to other mutations. It has been reported that the typical  
310 C>A mutational signature associated with smoking was found significantly enriched in CpGs outside  
311 CpG islands, suggesting that it preferentially occurs at modified CpG sites (Pleasance et al., 2010).  
312 Accordingly, our data indicate that CpG>T mutations might also be differentially affected by  
313 smoking-related mutagens.

314 **DISCUSSION**

315 Here we have established a link between the landscape of DNA modifications and the mutational  
316 profile of somatic human cells. Our measurements indicate that 5hmCs carry between 28 and 53%  
317 fewer mutations than methylated cytosines in brain. This results in a mutational load at 5hmC sites  
318 that is comparable to that of unmodified cytosines in CpG dinucleotides. This effect is not only  
319 observable in brain, but also in kidney cancers and myeloid leukaemias. The relationship between  
320 5hmC and CpG>T mutation rate can be detected at the scale of the exome as well as genome-wide  
321 and is independent of other region-specific influences on mutation frequency. We show that the  
322 relative impact of hydroxymethylation on mutagenesis decreases proportionally to the level of 5hmC  
323 in the tissue, suggesting that it represents a general property of this DNA modification.

324 Two possible scenarios could explain the striking difference in mutability between 5mC and 5hmC.  
325 Firstly, spontaneous and enzymatic deamination reactions of 5hmC could be less favourable than  
326 5mC. As a consequence, fewer new mutation events would be expected at 5hmC sites. Indeed,  
327 cytosine deaminases (namely, AID and APOBEC1-3) have 4.4–38x lower activity on sites with 5hmC  
328 compared to 5mC, supporting this possibility (Nabel et al., 2012; Rangam et al., 2012). Secondly,  
329 deamination of 5mC produces thymine whereas 5hmC deaminates to 5-hydroxymethyluracil  
330 (5hmU). This atypical base in DNA could be more efficiently recognised and replaced by the DNA  
331 glycosylases initiating base-excision repair (BER). Determining the relative contribution of DNA  
332 glycosylases to the lower mutation rate would be challenging, since some of these enzymes  
333 recognise several types of mismatches. TDG and MBD4 excise both T and 5hmU when mis-paired  
334 with G (Hardeland et al., 2003; Cortellino et al., 2011; Guo et al., 2011; Hashimoto et al., 2012;  
335 Moréra et al., 2012), whereas SMUG1 does not repair T:G but has a robust activity for 5hmU:G  
336 (Nilsen et al., 2001; Kemmerich et al., 2012). Therefore, there might be more efficient repair of  
337 5hmU in the genome. Further genome sequencing efforts might identify patients with rare

338 inactivating mutations in BER and/or mismatch-repair pathways that could be valuable for future  
339 investigations of the relationship between DNA repair and cytosine mutability.

340 It has previously been suggested that 5hmC levels increase the frequency of C>G mutations (Supek  
341 et al., 2014). As part of this analysis, only a very small (albeit statistically significant) decrease of C>T  
342 mutations in 5hmC sites in both SNPs and cancer SNVs was observed. There are two factors that  
343 could explain why we observe very different effects sizes for C>T and C>G mutations in 5hmC sites.  
344 Firstly, Supek et al. consider all sites with as little as one 5hmC read to be hydroxymethylated,  
345 whereas we require the level of 5hmC to exceed 5mC. In fact, when examining the effect of variation  
346 in these thresholds (Fig. 2-figure supplement 1A–F), we noticed that the results for C>G fluctuate  
347 substantially across the range of tested cut-off values (see also Fig. 3-figure supplement 1E–F).  
348 Secondly, we present evidence that tissue-specific changes in 5hmC patterns have great influence on  
349 the extent of correlation between 5hmC and mutability (Fig. 6B). Specifically, 5hmC genomic  
350 localisation in embryonic stem cells was a poor predictor of CpG>T mutations in brain, kidney and  
351 blood, compared to the respective tissue-specific 5hmC patterns.

352 The best predictor of CpG>T mutations in any of the three tested tissues was the 5hmC<sub>rel</sub> map from  
353 the corresponding anatomical site. This provides evidence that the slow accumulation of CpG>T  
354 mutations in the pre-cancerous tissue was strongly influenced by the DNA modification landscape.  
355 However, any bulk tissue sample encompasses a mixture of different cell types. Mounting evidence  
356 suggests that solid tumours originate from a defined subset of cells within any one tissue type. For  
357 example, glioblastomas were proposed to originate from stem or progenitor cell types enriched in  
358 the subventricular zone, while medulloblastomas have mixed cells of origin (Visvader, 2011). Those  
359 cell types are of low abundance in normal tissue biopsies. The fact that we observe a clear inverse  
360 relationship between CpG>T mutations and the location of 5hmC in multiple tissue types suggests  
361 that the DNA modification landscape in cancer-progenitor cells is sufficiently similar to the tissue  
362 average to be informative about the mutation frequencies in cancer.

363 Under this assumption we predict that the impact of DNA modifications on the frequency of CpG>T  
364 mutations is likely to be bigger than measured here, since the terminally differentiated cells that  
365 make up the bulk of the tissue may have diverged further from cancer-progenitors cells.  
366 Advancements in the identification of cancer origins and isolation of single cells, combined with  
367 single-cell bisulfite sequencing, will enable an improved assessment of the impact of DNA  
368 modifications on mutability.

369 The strong correlation between relative 5hmC levels in a tissue and the mutability of modified  
370 cytosine also points towards a shared underlying mutagenic process. The notable deviation of  
371 smoking-induced lung-cancers supports this hypothesis. We speculate that a yet undefined smoking-  
372 induced mutagenic mechanism preferentially affects unmethylated CpG sites. More experimental  
373 work will be needed to elucidate the biochemical causes for this phenomenon. In the future, the  
374 linear relationship between 5hmC levels and CpG>T mutation rate could thus be used to identify  
375 other environmental mutagens with a differential effect on modified cytosines.

376

377

378 **METHODS**

379 **Code.** Most of the analyses were performed using Matlab. Code and other required files are available on  
380 Figshare under doi 10.6084/m9.figshare.c.3249394 (Tomkova et al., 2016).

381 **Mutation data.** Cancer somatic mutations (see Table 2) were obtained from a dataset compiled by Alexandrov  
382 et al. (Alexandrov et al., 2013), complemented with whole genome samples from ICGC, (Wang et al., 2014),  
383 and TCGA. Briefly, aligned reads for 49 AML tumour and normal samples were downloaded from the UCSC  
384 CGHub website under TCGA access request #10140. Somatic variants were called using Strelka (Saunders et al.,  
385 2012) with default parameters. All variants were classified by the pyrimidine of the mutated Watson-Crick  
386 base pair (C or T) and the immediate 5' and 3' sequence context into 96 possible mutation types as described  
387 by Alexandrov et al. (2013).

388 **Modification data.** DNA modification information (see Table 1) for brain was extracted from supplementary  
389 information provided by Wen et al. (2014). Only sites with more than 5 TAB-Seq reads were taken into  
390 account. 5hmC<sub>high</sub> and 5mC<sub>high</sub> sites were defined based on values of *mod level* (unconverted/total BS-Seq  
391 reads) and *5hmC level* (unconverted/total TAB-Seq reads) per site:

- 392
- 5mC<sub>high</sub>: *mod level* > 10% and *5hmC level* / *mod level* ≤ threshold<sub>5mC</sub>
  - 5hmC<sub>high</sub>: *mod level* > 10% and *5hmC level* / *mod level* ≥ threshold<sub>5hmC</sub>
- 393

394 Effects of the choice of both thresholds were explored and then the values of threshold<sub>5mC</sub> = 0.3 and  
395 threshold<sub>5mC</sub> = 0.5 were used. In blood, BS-Seq and TAB-Seq values in CpG sites were taken from  
396 supplementary files provided by (Pacis et al., 2015). For kidney and ESC maps, raw reads were processed with  
397 Trim galore, Bismark (Krueger et al., 2012) and Mark duplicates from Picard tools. Multiple replicates were  
398 processed both independently and together (adding the reads from the replicates together). Only sites with at  
399 least 10% *mod level* were taken into account to compute 5hmC<sub>rel</sub>.

400 To compute the number of modified sites inside the exome, the reference Illumina Truseq definition of exon  
401 loci was downloaded from the Illumina website. Overlapping exons were merged using bedtools so that each  
402 genomic site is covered by at most one exon. Two-sided paired Wilcoxon signed-rank test was used for testing  
403 significance between mutation frequency of 5mC<sub>high</sub> and 5hmC<sub>high</sub> sites. The same test was used for all the  
404 following statistics, if not stated otherwise.

405 **Gene expression data.** Gene expression (in FPKM) from RNAseq experiments on 630 brain tissue samples were  
406 downloaded from the GTEx human tissue expression project (<http://www.gtexportal.org/home/>).

407 **Visualisation on genome.** The following genomic features were computed in 100kbp windows: average 5hmC,  
408 5mC, 5hmCrel (all from the supplementary information provided by Wen et al. (2014), i.e.  $\text{mod} \geq 10\%$ ),  
409 average  $\log(1 + \text{gene expression})$ , gene density, exon density, CpG density, modCpG density, CpG island (CGI)  
410 density, and average modification level (from raw BS-Seq reads). These features and CpG>T mutation  
411 frequency (from MDB and PA whole-genome sequencing datasets) were z-score normalised and plotted per  
412 chromosome after Gaussian smoothing with parameters  $n = 50$ ,  $\text{sigma} = 2.5$ .

413 **Mutation frequency in highly and lowly expressed genes.** Genes were sorted according to their median  
414 expression values. The upper 50-percentile (9701 genes) were classified as highly expressed, the rest as lowly  
415 expressed. Introns were included only for whole genome samples.

416 **Pairing of 5mC and 5hmC sites.** For each 5hmC<sub>high</sub> site in random order, the nearest not previously selected  
417 5mC<sub>high</sub> site was selected such that the 5mC-5hmC pair fulfilled the following conditions: both 5hmC<sub>high</sub> and  
418 5mC<sub>high</sub> sites are inside an exon or both are outside exons, and both share the same context (CG, CHG, and  
419 CHH, where H is T, A or C). This resulted in 6801374 pairs with a median distance of 1 and 25<sup>th</sup> and 75<sup>th</sup>  
420 quantiles of -177 and +177, respectively.

421 **Mutation frequency around aligned 5mC and 5hmC.** Modified sites with no other modifications in a 2kb  
422 radius were selected (374 000 sites with 5mC and the same number of 5hmC sites), and the mutation  
423 frequency up to 2kbp upstream and downstream (in bins without other modifications) was plotted.

424 **Gradients analysis.** All modified cytosines (i.e.,  $\text{mod level} > 10\%$ ) in the CpG context were divided into 9 right-  
425 open intervals according to their ratio of 5hmC level to mod level. The leftmost bin contained cytosines where  
426 the major modification is 5mC, while the rightmost bin contained cytosines where the major modification is  
427 5hmC. In each bin, the frequency of mutations was computed and plotted. A linear regression model was  
428 fitted to the data (function `fitlm` in MatLab) and the significance of the linear coefficient was tested using F-  
429 test for the hypothesis that the regression coefficient is zero (function `coefTest` in MatLab). For gradients  
430 with equal binning each interval contained approximately the same number of sites (apart from the first bin,  
431 which included all values with 5hmC<sub>rel</sub>=0).

432 **Prediction of mutation frequency in genomic windows.** CpG>T mutation frequency (response variable) and  
 433 genomic features (predictors; same as above in *Visualisation on genome*) were computed in genomic windows  
 434 of sizes 3kbp–3Mbp. Then a generalised linear model (`fitglm`) assuming Poisson distribution of the response  
 435 variable was fitted with a linear model specification (i.e., intercept + linear term for each predictor) and  
 436 `DispersionFlag` set to true. Model fits were compared in terms of  $D^2$  and p-value  
 437 (`model.devianceTest`), as recommended, e.g., in (Guisan & Zimmermann, 2000; Mittlböck & Heinzl,  
 438 2004).

439 **Simulation of effects of number of patients on GLM.** Each chromosome was split into windows of a given  
 440 window size. For each window, all CpG sites were counted. A random predictor was generated in each window  
 441 with a beta distribution ( $\text{Beta}(3,4)$ ). For each patient, a random number of mutations in each window was  
 442 generated as

$$\text{Binomial}(n = \text{windowSize}(i\text{Window}), p = \frac{\text{predictor}(i\text{Window})}{\text{coefficient}})$$

443 where:

$$\text{coefficient} = \frac{\sum_{i\text{Window}} \text{windowSize}(i\text{Window}) * \text{predictor}(i\text{Window})}{174}$$

444 The coefficient was set so that the expected total number of mutations per patient summed to 174, the  
 445 observed average number of CpG>T mutations in brain WGS data. The response variable was set as the  
 446 average CpG>T mutation frequency over all patients. A GLM was fit on the given predictor and response  
 447 variable and  $D^2$  was measured. The process was repeated 10 times for each combination of window size and  
 448 number of patients.

449 **Gene-wise prediction of mutation frequency.** Mutation frequency was modelled with two predictor variables:  
 450 average 5hmC<sub>rel</sub> per gene and log<sub>e</sub>-transformed gene expression. The following response variables computed  
 451 in exons of each gene were compared:

- 452 • modC>T: number of C>T mutations in modified C sites / number of modified C sites
- 453 • CpG>T: number of C>T mutations in CpG sites / number of CpG sites
- 454 • C>T: number of C>T mutations / number of C sites
- 455 • C>N: number of mutations from C / number of C sites
- 456 • N>N: number of mutations / number of sites

457 • T>N: number of mutations from T / number of T sites

458 Genes with missing values in at least one of the predictors and genes classified as outliers in at least one  
459 response variable were excluded from the analysis. Outliers were classified in the following way:  $y \geq$   
460  $quantile(y, 0.999) + 2.5 * (quantile(y, 0.999) - quantile(y, 0.001))$ . Out of 17,605 genes, 10 were classified and  
461 removed as outliers: ASPN, BBOX1, CCL4, ESPN, FOLH1, HLA-DPB1, IDH1, NLRP6, S100P, and TP53. The same  
462 GLM model as above was used. To calculate the relative contribution of one predictor variable over the other,  
463 two models were fitted with either one or both predictor variables and the difference in  $D^2$  was used.

464 **HPLC measurements of total 5hmC and 5mC in eight tissues.** 10ug of genomic DNA (amsbio; D1234003,  
465 D1234004, D1234035, D1234086, D1234090, D1234122, D1234142, D1234148, D1234149, D1234152,  
466 D1234171, D1234188, D1234200, D1234206, D1234226, D1234227, D1234246, D1234248, D1234260,  
467 D1234274, HG-101) was treated with 1U RNase A (Thermo Scientific) , purified by phenol chloroform ethanol  
468 precipitation and incubated overnight in hydrolysis solution (45 mM NaCl, 9 mM MgCl<sub>2</sub>, 9 mM Tris pH 7.9,  
469  $\geq 250$  U/ml Benzonase (sigma), 50 mU/ml Phosphodiesterase I,  $\geq 20$  U/ml Alkaline phosphatase, 46.8 ng/ml  
470 EHNA hydrochloride, 8.64  $\mu$ M deferoxamine). Protein components were removed by centrifugation through  
471 Amicon centrifugal filter unit (3 kDa cut-off, Millipore) before samples were lyophilised and resuspended in  
472 buffer A. Nucleosides were resolved with an Agilent UHPLC 1290 instrument fitted with Eclipse Plus C18 RRHD  
473 1.8  $\mu$ m (2.1  $\times$  150 mm column) and detected and quantified with Agilent 1290 DAD fitted with a Max-Light 60  
474 mm cell. Buffer A was 100 mM ammonium acetate, pH 6.5; buffer B was 40% acetonitrile, and the flow rate  
475 0.4 ml min<sup>-1</sup>. The gradient was between 1.8–100% of 40% acetonitrile with the following steps: 1–2 min,  
476 100% A; 2–16 min 98.2% A, 1.8% B; 16–18 min 70% A, 30% B; 18–20 min 50% A, 50% B; 20–21.5 min 25% A,  
477 75% B; 21.5–22.5 min 100% B; 22.5-24.5 min 100% A. Relative abundance of 5mC and 5hmC were established  
478 by detection of adenosine at 280nm allowing determination of total cytosine by extinction coefficient  
479 calculation using standards.

480

481

482 **AUTHOR CONTRIBUTIONS**

483 Conceptualization: S.K. and B.S.-B.; Investigation: M.T.; HPLC measurements: M.McM.; Writing –  
484 Original Draft: B.S.-B., S.K. and M.T.

485 **ACKNOWLEDGEMENTS**

486 We would like to thank Jakub Tomek, Pijus Brazauskas and David Severson for helpful discussions.

487 S.K. and B.S.-B. are funded by Ludwig Cancer Research. S.K. received funding from BBSRC grant  
488 BB/M001873/1. M.T. is funded by EPSRC (EP/F500394/1) and Bakala Foundation.

489 **REFERENCES**

- 490 Alexandrov, L. B., Jones, P. H., Wedge, D. C., Sale, J. E., & Peter, J. (2015). Clock-like mutational  
491 processes in human somatic cells. *Nature Publishing Group*, 47(12), 1402–1407.  
492 doi:10.1038/ng.3441
- 493 Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Aparicio, S. a J. R., Behjati, S., Biankin, A. V, Bignell, G.  
494 R., Bolli, N., Borg, A., Børresen-Dale, A.-L., Boyault, S., Burkhardt, B., Butler, A. P., Caldas, C.,  
495 Davies, H. R., Desmedt, C., Eils, R., ... Stratton, M. R. (2013). Signatures of mutational processes  
496 in human cancer. *Nature*, 500(7463), 415–21. doi:10.1038/nature12477
- 497 Ames, B. N., Shigenaga, M. K., & Hagen, T. M. (1993). Oxidants, antioxidants, and the degenerative  
498 diseases of aging. *Proceedings of the National Academy of Sciences of the United States of*  
499 *America*, 90(17), 7915–7922. doi:10.1073/pnas.90.17.7915
- 500 Bachman, M., Uribe-lewis, S., Yang, X., Williams, M., & Murrell, A. (2014). 5-Hydroxymethylcytosine  
501 is a predominantly stable DNA modification. *Nature Chemistry*, 6(12), 1049–1055.  
502 doi:10.1038/nchem.2064
- 503 Booth, M. J., Marsico, G., Bachman, M., Beraldi, D., & Balasubramanian, S. (2014). Quantitative  
504 sequencing of 5-formylcytosine in DNA at single-base resolution. *Nature Chemistry*, 6(5), 435–  
505 40. doi:10.1038/nchem.1893
- 506 Brazauskas, P., & Kriaucionis, S. (2014). DNA modifications: Another stable base in DNA. *Nature*  
507 *Chemistry*, 6(12), 1031–1033. doi:10.1038/nchem.2115
- 508 Chen, K., Zhang, J., Guo, Z., Ma, Q., Xu, Z., Zhou, Y., Xu, Z., Li, Z., Liu, Y., Ye, X., Li, X., Yuan, B., Ke, Y.,  
509 He, C., Zhou, L., Liu, J., & Ci, W. (2015). Loss of 5-hydroxymethylcytosine is linked to gene body  
510 hypermethylation in kidney cancer. *Cell Research*, 103–118. doi:10.1038/cr.2015.150
- 511 Cortellino, S., Xu, J., Sannai, M., Moore, R., Caretti, E., Cigliano, A., Le Coz, M., Devarajan, K., Wessels,  
512 A., Soprano, D., Abramowitz, L. K., Bartolomei, M. S., Rambow, F., Bassi, M. R., Bruno, T.,  
513 Fanciulli, M., Renner, C., ... Bellacosa, A. (2011). Thymine DNA glycosylase is essential for active  
514 DNA demethylation by linked deamination-base excision repair. *Cell*, 146(1), 67–79.  
515 doi:10.1016/j.cell.2011.06.020
- 516 Globisch, D., Münzel, M., Müller, M., Michalakakis, S., Wagner, M., Koch, S., Brückl, T., Biel, M., &

- 517 Carell, T. (2010). Tissue distribution of 5-hydroxymethylcytosine and search for active  
518 demethylation intermediates. *PLoS One*, 5(12), e15367. doi:10.1371/journal.pone.0015367
- 519 Guisan, A., & Zimmermann, N. E. (2000). Predictive habitat distribution models in ecology. *Ecological*  
520 *Modelling*, 135, 147–186. doi:10.1016/S0304-3800(00)00354-9
- 521 Guo, J. U., Su, Y., Zhong, C., Ming, G. L., & Song, H. (2011). Hydroxylation of 5-methylcytosine by  
522 TET1 promotes active DNA demethylation in the adult brain. *Cell*, 145(3), 423–434.  
523 doi:10.1016/j.cell.2011.03.022
- 524 Hardeland, U., Bentele, M., Jiricny, J., & Schär, P. (2003). The versatile thymine DNA-glycosylase: A  
525 comparative characterization of the human, Drosophila and fission yeast orthologs. *Nucleic*  
526 *Acids Research*, 31(9), 2261–2271. doi:10.1093/nar/gkg344
- 527 Hashimoto, H., Zhang, X., & Cheng, X. (2012). Excision of thymine and 5-hydroxymethyluracil by the  
528 MBD4 DNA glycosylase domain: Structural basis and implications for active DNA  
529 demethylation. *Nucleic Acids Research*, 40(17), 8276–8284. doi:10.1093/nar/gks628
- 530 He, Y.-F., Li, B.-Z., Li, Z., Liu, P., Wang, Y., Tang, Q., Ding, J., Jia, Y., Chen, Z., Li, L., Sun, Y., Li, X., Dai, Q.,  
531 Song, C.-X., Zhang, K., He, C., & Xu, G.-L. (2011). Tet-mediated formation of 5-carboxylcytosine  
532 and its excision by TDG in mammalian DNA. *Science*, 333(6047), 1303–7.  
533 doi:10.1126/science.1210944
- 534 Hu, X., Zhang, L., Mao, S. Q., Li, Z., Chen, J., Zhang, R. R., Wu, H. P., Gao, J., Guo, F., Liu, W., Xu, G. F.,  
535 Dai, H. Q., Shi, Y. G., Li, X., Hu, B., Tang, F., Pei, D., & Xu, G. L. (2014). Tet and TDG mediate DNA  
536 demethylation essential for mesenchymal-to-epithelial transition in somatic cell  
537 reprogramming. *Cell Stem Cell*, 14(4), 512–522. doi:10.1016/j.stem.2014.01.001
- 538 Inoue, A., & Zhang, Y. (2011). Replication-Dependent Loss of 5-Hydroxymethylcytosine in Mouse  
539 Preimplantation Embryos. *Science*, 334(6053), 194–194. doi:10.1126/science.1212483
- 540 Jones, P. A. (2012). Functions of DNA methylation: islands, start sites, gene bodies and beyond.  
541 *Nature Reviews Genetics*, 13(7), 484–492. doi:10.1038/nrg3230
- 542 Kemmerich, K., Dingler, F. a., Rada, C., & Neuberger, M. S. (2012). Germline ablation of SMUG1 DNA  
543 glycosylase causes loss of 5-hydroxymethyluracil-and UNG-backup uracil-excision activities and  
544 increases cancer predisposition of Ung<sup>-/-</sup>Msh2<sup>-/-</sup> mice. *Nucleic Acids Research*, 40(13), 6016–  
545 6025. doi:10.1093/nar/gks259
- 546 Klose, R. J., & Bird, A. P. (2006). Genomic DNA methylation: The mark and its mediators. *Trends in*  
547 *Biochemical Sciences*, 31(2), 89–97. doi:10.1016/j.tibs.2005.12.008
- 548 Koren, A., Polak, P., Nemesh, J., Michaelson, J. J., Sebat, J., Sunyaev, S. R., & McCarroll, S. a. (2012).  
549 Differential relationship of DNA replication timing to different forms of human mutation and  
550 variation. *American Journal of Human Genetics*, 91(6), 1033–1040.  
551 doi:10.1016/j.ajhg.2012.10.018
- 552 Kriaucionis, S., & Heintz, N. (2009). The nuclear DNA base 5-hydroxymethylcytosine is present in  
553 Purkinje neurons and the brain. *Science*, 324(5929), 929–30. doi:10.1126/science.1169786
- 554 Krueger, F., Kreck, B., Franke, A., & Andrews, S. R. (2012). DNA methylome analysis using short  
555 bisulfite sequencing data. *Nature Methods*, 9(2), 145–151. doi:10.1038/nmeth.1828
- 556 Lawrence, M. S., Stojanov, P., Polak, P., Kryukov, G. V., Cibulskis, K., Sivachenko, A., Carter, S. L.,  
557 Stewart, C., Mermel, C. H., Roberts, S. A., Kiezun, A., Hammerman, P. S., McKenna, A., Drier, Y.,  
558 Zou, L., Ramos, A. H., Pugh, T. J., ... Getz, G. (2013). Mutational heterogeneity in cancer and the  
559 search for new cancer-associated genes. *Nature*, 499(7457), 214–8. doi:10.1038/nature12213

- 560 Li, E., Bestor, T. H., & Jaenisch, R. (1992). Targeted mutation of the DNA methyltransferase gene  
561 results in embryonic lethality. *Cell*, *69*(6), 915–926. doi:10.1016/0092-8674(92)90611-F
- 562 Li, W., & Liu, M. (2011). Distribution of 5-hydroxymethylcytosine in different human tissues. *Journal*  
563 *of Nucleic Acids*, *2011*, 870726. doi:10.4061/2011/870726
- 564 Lindahl, T., & Nyberg, B. (1974). Heat-induced deamination of cytosine residues in deoxyribonucleic  
565 acid. *Biochemistry*, *13*(16), 3405–3410. doi:10.1021/bi00713a035
- 566 Lister, R., Mukamel, E. a, Nery, J. R., Urich, M., Puddifoot, C. a, Johnson, N. D., Lucero, J., Huang, Y.,  
567 Dwork, A. J., Schultz, M. D., Yu, M., Tonti-Filippini, J., Heyn, H., Hu, S., Wu, J. C., Rao, A., Esteller,  
568 M., ... Ecker, J. R. (2013). Global epigenomic reconfiguration during mammalian brain  
569 development. *Science*, *341*(6146), 1237905. doi:10.1126/science.1237905
- 570 Maiti, A., & Drohat, A. C. (2011). Thymine DNA glycosylase can rapidly excise 5-formylcytosine and 5-  
571 carboxylcytosine: Potential implications for active demethylation of CpG sites. *Journal of*  
572 *Biological Chemistry*, *286*(41), 35334–35338. doi:10.1074/jbc.C111.284620
- 573 Mellén, M., Ayata, P., Dewell, S., Kriaucionis, S., & Heintz, N. (2012). MeCP2 binds to 5hmC enriched  
574 within active genes and accessible chromatin in the nervous system. *Cell*, *151*(7), 1417–1430.  
575 doi:10.1016/j.cell.2012.11.022
- 576 Mittlböck, M., & Heinzl, H. (2004). Pseudo R-squared measures for generalized linear models. In  
577 *Proceedings of the 1st European Workshop on the Assessment of Diagnostic Performance* (pp.  
578 71–80). Milan, Italy.
- 579 Moréra, S., Grin, I., Vigouroux, A., Couvé, S., Henriot, V., Saparbaev, M., & Ishchenko, A. a. (2012).  
580 Biochemical and structural characterization of the glycosylase domain of MBD4 bound to  
581 thymine and 5-hydroxymethyluracil-containing DNA. *Nucleic Acids Research*, *40*(19), 9917–  
582 9926. doi:10.1093/nar/gks714
- 583 Nabel, C. S., Jia, H., Ye, Y., Shen, L., Goldschmidt, H. L., Stivers, J. T., Zhang, Y., & Kohli, R. M. (2012).  
584 AID/APOBEC deaminases disfavor modified cytosines implicated in DNA demethylation. *Nature*  
585 *Chemical Biology*, *8*(9), 751–758. doi:10.1038/nchembio.1042
- 586 Nestor, C. E., Ottaviano, R., Reddington, J., Sproul, D., Reinhardt, D., Dunican, D., Katz, E., Dixon, J.  
587 M., Harrison, D. J., & Meehan, R. R. (2012). Tissue type is a major modifier of the 5-  
588 hydroxymethylcytosine content of human genes. *Genome Research*, *467*–477.  
589 doi:10.1101/gr.126417.111
- 590 Nik-Zainal, S., Alexandrov, L. B., Wedge, D. C., Van Loo, P., Greenman, C. D., Raine, K., Jones, D.,  
591 Hinton, J., Marshall, J., Stebbings, L. a, Menzies, A., Martin, S., Leung, K., Chen, L., Leroy, C.,  
592 Ramakrishna, M., Rance, R., ... Stratton, M. R. (2012). Mutational processes molding the  
593 genomes of 21 breast cancers. *Cell*, *149*(5), 979–93. doi:10.1016/j.cell.2012.04.024
- 594 Nilsen, H., Haushalter, K. a., Robins, P., Barnes, D. E., Verdine, G. L., & Lindahl, T. (2001). Excision of  
595 deaminated cytosine from the vertebrate genome: Role of the SMUG1 uracil-DNA glycosylase.  
596 *EMBO Journal*, *20*(15), 4278–4286. doi:10.1093/emboj/20.15.4278
- 597 Okano, M., Bell, D. W., Haber, D. a., & Li, E. (1999). DNA methyltransferases Dnmt3a and Dnmt3b are  
598 essential for de novo methylation and mammalian development. *Cell*, *99*(3), 247–257.  
599 doi:10.1016/S0092-8674(00)81656-6
- 600 Pacis, A., Tailleux, L., Morin, A. M., Lambourne, J., Maclsaac, J. L., Yotova, V., Dumaine, A., Danckaert,  
601 A., Luca, F., Grenier, J., Hansen, K. D., Gicquel, B., Yu, M., Pai, A., He, C., Tung, J., Pastinen, T., ...  
602 Barreiro, L. B. (2015). Bacterial infection remodels the DNA methylation landscape of human

603 dendritic cells. *Genome Research*, 25(12), 1801–11. doi:10.1101/gr.192005.115

604 Pleasance, E. D., Stephens, P. J., O’Meara, S., McBride, D. J., Meynert, A., Jones, D., Lin, M. L., Beare,  
605 D., Lau, K. W., Greenman, C., Varella, I., Nik-Zainal, S., Davies, H. R., Ordonez, G. R., Mudie, L. J.,  
606 Latimer, C., Edkins, S., ... Campbell, P. J. (2010). A small-cell lung cancer genome with complex  
607 signatures of tobacco exposure. *Nature*, 463(7278), 184–190. doi:10.1038/nature08629

608 Poon, S. L., Pang, S.-T., McPherson, J. R., Yu, W., Huang, K. K., Guan, P., Weng, W.-H., Siew, E. Y., Liu,  
609 Y., Heng, H. L., Chong, S. C., Gan, A., Tay, S. T., Lim, W. K., Cutcutache, I., Huang, D., Ler, L. D., ...  
610 Teh, B. T. (2013). Genome-wide mutational signatures of aristolochic acid and its application as  
611 a screening tool. *Science Translational Medicine*, 5(197), 197ra101.  
612 doi:10.1126/scitranslmed.3006086

613 Rangam, G., Schmitz, K. M., Cobb, A. J. a, & Petersen-Mahrt, S. K. (2012). AID enzymatic activity is  
614 inversely proportional to the size of cytosine c5 orbital cloud. *PLoS ONE*, 7(8), 3–8.  
615 doi:10.1371/journal.pone.0043279

616 Rasmussen, K. D., & Helin, K. (2016). Role of TET enzymes in DNA methylation , development , and  
617 cancer. *Genes & Development*, 30, 733–750. doi:10.1101/gad.276568.115.maintain

618 Saunders, C. T., Wong, W. S. W., Swamy, S., Becq, J., Murray, L. J., & Cheetham, R. K. (2012). Strelka:  
619 Accurate somatic small-variant calling from sequenced tumor-normal sample pairs.  
620 *Bioinformatics*, 28(14), 1811–1817. doi:10.1093/bioinformatics/bts271

621 Schultz, M. D., He, Y., Whitaker, J. W., Hariharan, M., Mukamel, E. a., Leung, D., Rajagopal, N., Nery,  
622 J. R., Urich, M. a., Chen, H., Lin, S., Lin, Y., Jung, I., Schmitt, A. D., Selvaraj, S., Ren, B., Sejnowski,  
623 T. J., ... Ecker, J. R. (2015). Human body epigenome maps reveal noncanonical DNA methylation  
624 variation. *Nature*, 523(7559), 212–216. doi:10.1038/nature14465

625 Schuster-Böckler, B., & Lehner, B. (2012). Chromatin organization is a major influence on regional  
626 mutation rates in human cancer cells. *Nature*, 488(7412), 504–7. doi:10.1038/nature11273

627 Shen, L., Wu, H., Diep, D., Yamaguchi, S., D’Alessio, A. C., Fung, H. L., Zhang, K., & Zhang, Y. (2013).  
628 Genome-wide analysis reveals TET- and TDG-dependent 5-methylcytosine oxidation dynamics.  
629 *Cell*, 153(3), 692–706. doi:10.1016/j.cell.2013.04.002

630 Spruijt, C. G., Gnerlich, F., Smits, A. H., Pfaffeneder, T., Jansen, P. W. T. C., Bauer, C., Münzel, M.,  
631 Wagner, M., Müller, M., Khan, F., Eberl, H. C., Mensinga, A., Brinkman, A. B., Lephikov, K.,  
632 Müller, U., Walter, J., Boelens, R., ... Vermeulen, M. (2013). Dynamic readers for 5-  
633 (hydroxy)methylcytosine and its oxidized derivatives. *Cell*, 152(5), 1146–59.  
634 doi:10.1016/j.cell.2013.02.004

635 Stephens, P. J., Tarpey, P. S., Davies, H., Van Loo, P., Greenman, C., Wedge, D. C., Zainal, S. N.,  
636 Martin, S., Varella, I., Bignell, G. R., Yates, L. R., Papaemmanuil, E., Beare, D., Butler, A.,  
637 Cheverton, A., Gamble, J., Hinton, J., ... Ottestad, L. (2012). The landscape of cancer genes and  
638 mutational processes in breast cancer. *Nature*, 486(7403), 400–404. doi:10.1038/nature11017

639 Supek, F., Lehner, B., Hajkova, P., & Warnecke, T. (2014). Hydroxymethylated cytosines are  
640 associated with elevated C to G transversion rates. *PLoS Genetics*, 10(9), e1004585.  
641 doi:10.1371/journal.pgen.1004585

642 Szwagierczak, A., Bultmann, S., Schmidt, C. S., Spada, F., & Leonhardt, H. (2010). Sensitive enzymatic  
643 quantification of 5-hydroxymethylcytosine in genomic DNA. *Nucleic Acids Research*, 38(19),  
644 e181. doi:10.1093/nar/gkq684

645 Tahiliani, M., Koh, K. P., Shen, Y., Pastor, W. A., Bandukwala, H., Brudno, Y., Agarwal, S., Iyer, L. M.,

646 Liu, D. R., Aravind, L., & Rao, A. (2009). Conversion of 5-methylcytosine to 5-  
647 hydroxymethylcytosine in mammalian DNA by MLL partner TET1. *Science*, *324*(5929), 930–5.  
648 doi:10.1126/science.1170116

649 Takai, H., Masuda, K., Sato, T., Sakaguchi, Y., Suzuki, T., Suzuki, T., Koyama-Nasu, R., Nasu-Nishimura,  
650 Y., Katou, Y., Ogawa, H., Morishita, Y., Kozuka-Hata, H., Oyama, M., Todo, T., Ino, Y., Mukasa,  
651 A., Saito, N., ... Akiyama, T. (2014). 5-Hydroxymethylcytosine Plays a Critical Role in  
652 Glioblastomagenesis by Recruiting the CHTOP-Methylosome Complex. *Cell Reports*, *9*(1), 48–  
653 60. doi:10.1016/j.celrep.2014.08.071

654 Taylor, B. J. M., Nik-Zainal, S., Wu, Y. L., Stebbings, L. a., Raine, K., Campbell, P. J., Rada, C., Stratton,  
655 M. R., & Neuberger, M. S. (2013). DNA deaminases induce break-associated mutation showers  
656 with implication of APOBEC3B and 3A in breast cancer kataegis. *eLife*, *2013*(2), 1–14.  
657 doi:10.7554/eLife.00534

658 Tomasetti, C., Vogelstein, B., & Parmigiani, G. (2013). Half or more of the somatic mutations in  
659 cancers of self-renewing tissues originate prior to tumor initiation. *Proceedings of the National  
660 Academy of Sciences*, *110*(6), 1999–2004. doi:10.1073/pnas.1221068110

661 Tomkova, M., McClellan, M., Kriaucionis, S., & Schuster-Böckler, B. (2016). Code and intermediate  
662 files for “5-hydroxymethylcytosine marks regions with reduced mutation frequency.”  
663 doi:10.6084/m9.figshare.c.3249394

664 Visvader, J. E. (2011). Cells of origin in cancer. *Nature*, *469*(7330), 314–322.  
665 doi:10.1038/nature09781

666 Wang, K., Yuen, S. T., Xu, J., Lee, S. P., Yan, H. H. N., Shi, S. T., Siu, H. C., Deng, S., Chu, K. M., Law, S.,  
667 Chan, K. H., Chan, A. S. Y., Tsui, W. Y., Ho, S. L., Chan, A. K. W., Man, J. L. K., Foglizzo, V., ...  
668 Leung, S. Y. (2014). Whole-genome sequencing and comprehensive molecular profiling identify  
669 new driver mutations in gastric cancer. *Nature Genetics*, *46*(6), 573–82. doi:10.1038/ng.2983

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

674 Wossidlo, M., Nakamura, T., Lepikhov, K., Marques, C. J., Zakhartchenko, V., Boiani, M., Arand, J.,  
675 Nakano, T., Reik, W., & Walter, J. (2011). 5-Hydroxymethylcytosine in the mammalian zygote is  
676 linked with epigenetic reprogramming. *Nature Communications*, *2*, 241.  
677 doi:10.1038/ncomms1240

678 Wu, H., & Zhang, Y. (2011). Tet1 and 5-hydroxymethylation: A genome-wide view in mouse  
679 embryonic stem cells. *Cell Cycle*, *10*(15), 2428–2436. doi:10.4161/cc.10.15.16930

680 Wu, S., Powers, S., Zhu, W., & Hannun, Y. A. (2015). Substantial contribution of extrinsic risk factors  
681 to cancer development. *Nature*, *529*(7584), 43–47. doi:10.1038/nature16166

682 Yu, M., Hon, G. C., Szulwach, K. E., Song, C. X., Zhang, L., Kim, A., Li, X., Dai, Q., Shen, Y., Park, B., Min,  
683 J. H., Jin, P., Ren, B., & He, C. (2012). Base-resolution analysis of 5-hydroxymethylcytosine in the  
684 mammalian genome. *Cell*, *149*(6), 1368–1380. doi:10.1016/j.cell.2012.04.027

685

686 **FIGURE LEGENDS**

687 **Figure 1**

688 **C>T mutations are common in the genome but depleted in 5hmC sites compared to 5mC sites. A:**  
689 Distribution of 5hmC in a CpG context in brain compared to kidney and blood. **B:** Frequency of SNVs  
690 in brain cancer exomes, stratified by sequence context, normalised by frequency of trinucleotides. **C:**  
691 Distribution of single-nucleotide variants (whole genomes) in brain cancer according to type, context  
692 and modification state. **D:** CpG>T mutation frequency (black), 5hmC (blue) and 5mC (orange) density  
693 in 100kbp windows of chromosome 3, smoothed with a Gaussian filter (n=50, sigma=2.5). **E:** Average  
694 fraction of mutated sites for 5mC<sub>high</sub> vs. 5hmC<sub>high</sub> over all patient samples (CpG sites only; \*\*\*P <  
695 0.001; \*\*P < 0.01; \*P < 0.05, see Methods).

696 **Figure 2**

697 **Differential mutation frequency between 5mC and 5hmC is present in all 5 brain cancer types and**  
698 **correlates with age at diagnosis. A:** Average fraction of mutated CpG sites for 5mC<sub>high</sub> vs. 5hmC<sub>high</sub>  
699 computed separately for each cancer type. **B:** Box plot of C>T mutation frequency, as shown in A. **C.:**  
700 Correlation of whole genome CpG>T mutation frequency with age at the time of diagnosis in  
701 patients with Medulloblastoma and Pilocytic Astrocytoma.

702 **Figure 3**

703 **Depletion of C>T mutations in 5hmC sites is not explained by gene expression or regional mutation**  
704 **rate variation. A–B:** Frequency of mutations in 5mC<sub>high</sub> vs 5hmC<sub>high</sub> sites within highly expressed (A)  
705 or lowly expressed (B) genes (see Methods). **C–D:** Boxplot visualisation of C>T mutation frequency  
706 for each cancer type. **E:** For each patient sample, the overall difference in mutations in paired sites  
707 was calculated and compared using a Wilcoxon signed-rank test. Shown here is a histogram of  
708 samples by the difference in mutations for paired 5mC and 5hmC sites (negative values shown blue,  
709 positive in orange; see Methods for details). Mutations in 5mC sites exceed paired 5hmC sites,

710 causing a shift to the right. **F:** Same as E but using a more stringent definition of 5mC (only sites with  
711  $\text{threshold}_{5\text{mC}} \leq 0.2$ ).

712 **Figure 4**

713 **Mutation frequency negatively correlates with 5hmC<sub>rel</sub> level per base.** **A:** Fraction of mutated CpG  
714 sites as a function of 5hmC<sub>rel</sub> levels by mutation and cancer type. Bins to the left represent sites  
715 predominantly methylated, while bins to the right contain increasingly hydroxymethylated sites.  
716 Black line denotes linear regression fit (F-test for coefficient deviation from 0, see Methods). **B:**  
717 Distribution of CpG>T mutation frequency by modification type. The top left bin contains cytosines  
718 that are mostly unmodified, the bottom left bin contains exclusively methylated cytosines and the  
719 top right bin contains cytosines that are mostly hydroxymethylated. **C:** Top row of B, *i.e.* distribution  
720 of mutations in unmethylated sites. **D:** First column of B, *i.e.* distribution of mutations in sites  
721 without 5hmC. **E:** Diagonal of B, *i.e.* distribution of mutations in highly modified sites.

722 **Figure 5**

723 **Predictors of mutations: 5hmC<sub>rel</sub> compared to other genomic features.** **A:** Prediction of CpG>T  
724 mutation frequency (using whole genome sequencing only) in 100kbp genomic windows. Predictors  
725 are sorted according to the  $D^2$  in a univariate model. The height of the  $k^{\text{th}}$  bar denotes the  $D^2$  of a  
726 model with the first  $k$  predictors. **B:** Comparison of the nine predictors of CpG>T mutation features  
727 by  $D^2$  in a univariate models, in a range of window sizes. **C:** Prediction of different types of mutation  
728 frequency in genes. Increase in  $D^2$  of a generalised linear model including 5hmC<sub>rel</sub> over gene  
729 expression (purple) or gene expression over 5hmC<sub>rel</sub> (green) (see Methods). **D:** Significance of  
730 observations in C (see Methods).

731 **Figure 6**

732 **Decreased CpG>T mutation frequency in 5hmC is not limited to brain tissue.** **A:** Predictions of  
733 CpG>T mutation frequency in whole genome cancers in blood (AML), kidney and brain using 5hmC<sub>rel</sub>

734 maps from blood, kidney, brain and embryonic stem cells (ESC) in 100kbp genomic windows. The  
735 values are z-score normalised per rows in order to normalise for different number of patients and  
736 mutations in each cancer type (the original  $D^2$  values are in parentheses); the higher values of  $D^2$   
737 (green), the better predictions. B: CpG>T mutation frequency in 5mC vs. 5hmC in kidney and blood.  
738 **C:** Correlation of total 5hmC<sub>rel</sub> levels (measured with HPLC) with frequency of CpG>T mutations in  
739 modified cytosines normalised by the frequency in unmodified cytosines in different tissues (see  
740 Methods).

#### 741 **TABLE LEGENDS**

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

743 Table 2: Overview of whole genome and exome sequencing data used for mutation information.

#### 744 **SUPPLEMENTARY FIGURE LEGENDS**

745 Figure 1-figure supplements 1–2: **Distribution of CpG>T mutations vs modifications across all**  
746 **chromosomes:** CpG>T mutation frequency (black), 5hmC (blue) and 5mC (orange) density in 100kbp  
747 windows, smoothed with a Gaussian filter (n=50, sigma=2.5).

748 Figure 1-figure supplement 3: **Distribution of CpG>T mutations vs other genomic features:** CpG>T  
749 mutation frequency (black) and several genomic features in 100kbp windows on chromosome 3,  
750 smoothed with a Gaussian filter (n=50, sigma=2.5). CGIs: density of CpG islands, EXONS: density of  
751 exons, GENEs: density of genes, CpG: density of CpGs, modCpG: density of CpGs with *mod level*  $\geq$   
752 10%; and average modification levels: mod, 5hmC, 5mC, and 5hmC<sub>rel</sub>.

753 Figure 2-figure supplement 1: **Depletion of C>T mutations in 5hmC<sub>high</sub> is relatively insensitive to**  
754 **varying definitions of 5mC<sub>high</sub> and 5hmC<sub>high</sub>.** A–F: Significance of a difference in mutation frequency  
755 in 5mC<sub>high</sub> and 5hmC<sub>high</sub>, for a range of values of threshold<sub>5mC</sub> and threshold<sub>5hmC</sub> (5mC<sub>high</sub> is defined as  
756 sufficiently modified sites with 5hmC<sub>rel</sub>  $\leq$  threshold<sub>5mC</sub>; 5hmC<sub>high</sub> is defined as sufficiently modified

757 sites with  $5\text{hmC}_{\text{rel}} \geq \text{threshold}_{5\text{hmC}}$ ). One-sided paired Wilcoxon sign-rank test was used. Red colour  
758 represents a significant increase of mutation frequency in  $5\text{mC}_{\text{high}}$  (right tail test) whereas blue  
759 colour represents elevated mutations in  $5\text{hmC}_{\text{high}}$  (left tail test). **G–I**: C>T mutation frequency for  
760  $5\text{mC}_{\text{high}}$  vs.  $5\text{hmC}_{\text{high}}$  with  $\text{threshold}_{5\text{mC}} = 0.3$  and  $\text{threshold}_{5\text{hmC}} = 0.7$ .

761 Figure 3-figure supplement 1: **Depletion of C>T mutations in  $5\text{hmC}_{\text{high}}$  is relatively insensitive to**  
762 **varying definitions of  $5\text{mC}_{\text{high}}$  and  $5\text{hmC}_{\text{high}}$ .** **A–D**: C>T mutation frequency for  $5\text{mC}_{\text{high}}$  vs.  $5\text{hmC}_{\text{high}}$  in  
763 highly vs. lowly expressed genes with  $\text{threshold}_{5\text{mC}} = 0.3$  and  $\text{threshold}_{5\text{hmC}} = 0.7$ . **E–F**: C>G mutation  
764 frequency with  $\text{threshold}_{5\text{mC}} = 0.0$  and  $\text{threshold}_{5\text{hmC}} = 0.5$ . **G**: Mutation frequency around aligned 5mC  
765 and 5hmC sites.

766 Figure 4-figure supplement 1: **CpG>T mutation frequency as a function of  $5\text{hmC}_{\text{rel}}$  levels with equal**  
767 **binning** (each bin contains approximately the same number of sites).

768 Figure 5-figure supplement 1: **Genome-wide prediction of CpG>T mutation frequency:  $5\text{hmC}_{\text{rel}}$**   
769 **compared to other genomic features.** **A–C**: Comparison of nine predictors of CpG>T mutation  
770 frequency in a range of window sizes by p-value of univariate generalised linear models (**A**),  
771 Spearman correlation (**B**), and Pearson correlation (**C**). **D**: Effects of window size and patient  
772 numbers on  $D^2$  of GLM with one response variable (simulated mutation frequency) generated  
773 proportionally from a single ideal predictor (see methods for details).

774 Figure 5-figure supplement 2: **Effects of  $5\text{hmC}_{\text{rel}}$  levels on gene mutability.** Data for GLM with  
775 Poisson distribution (the fitted curve is in green). Genes defined as outliers in at least one definition  
776 of mutation frequency (above the red line) are plotted in red. For convenience, the mutation  
777 frequency is plotted on log-scale.

778 Figure 5-figure supplement 3: **Effects of  $5\text{hmC}_{\text{rel}}$  levels on gene mutability.** **A–B**: GLM results fitted  
779 separately for  $5\text{hmC}_{\text{rel}}$  (purple) and gene expression (green) and both of them together (yellow). **C–**  
780 **D**: Frequency of modC>T mutations of all genes (**C**) and gene density (**D**) in the space of  $5\text{hmC}_{\text{rel}}$  and

781 gene expression. Briefly, for figures **C** and **D** the space was limited to [quantile(x, 0.05), quantile(x,  
782 0.95)] on both axes and then binned into 100x100 bins. In each bin, the average mutation frequency  
783 (in the form of  $\log(\text{mutFreq} + \min(\text{mutFreq}(\text{mutFreq} > 0)))$ ) was computed. The resulting matrix was  
784 smoothed by applying a Gaussian filter (radius 5 bins, sigma 2) weighted by the number of genes in  
785 each bin (bins with  $\geq 2/3$  missing values in their neighbourhood were set to NaN) and plotted with  
786 `pcolor` (NaN bins are shown in black).

787 Figure 6-figure supplement 1: **Decreased CpG>T mutation frequency in 5hmC is present in three**  
788 **tissues consistently for different replicates of modification maps. A:** CpG>T mutation frequency in  
789 5mC compared to 5hmC in blood and kidney using modification maps from different replicates  
790 merged together (A) and used separately (B). **C:** Predictions of CpG>T mutation frequency in whole  
791 genome cancers in blood (AML), kidney and brain using different replicates of 5hmC<sub>rel</sub> maps from  
792 blood, kidney, brain and embryonic stem cells (ESC) in 100kbp genomic windows. The values are z-  
793 score normalised per rows in order to normalise for different number of patients and mutations in  
794 each cancer type (the original  $D^2$  values are in parentheses); the higher values of  $D^2$  (green colour),  
795 the better predictions.

796 Figure 6-figure supplement 2: **Comparison of 5hmC in 10kbp windows in blood, kidney (2**  
797 **replicates), and brain:** distribution of 5hmC values in each map and Pearson correlation of pairs of  
798 maps.

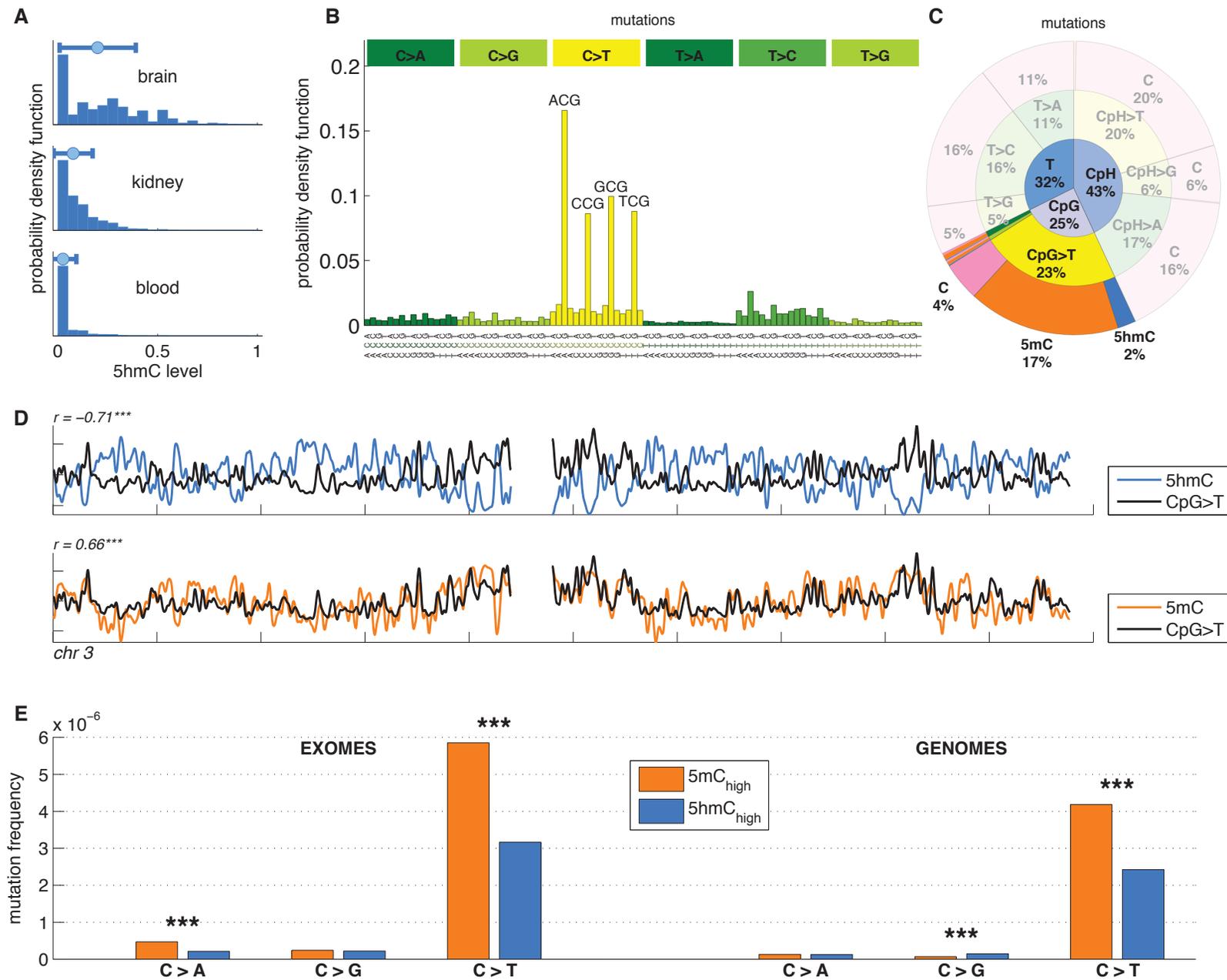
799 Figure 6-figure supplement 3: **HPLC measurements of total 5hmC and 5mC in eight tissues:** average  
800 values with standard deviation of 5mC and 5hmC (as a percentage of total cytosine).

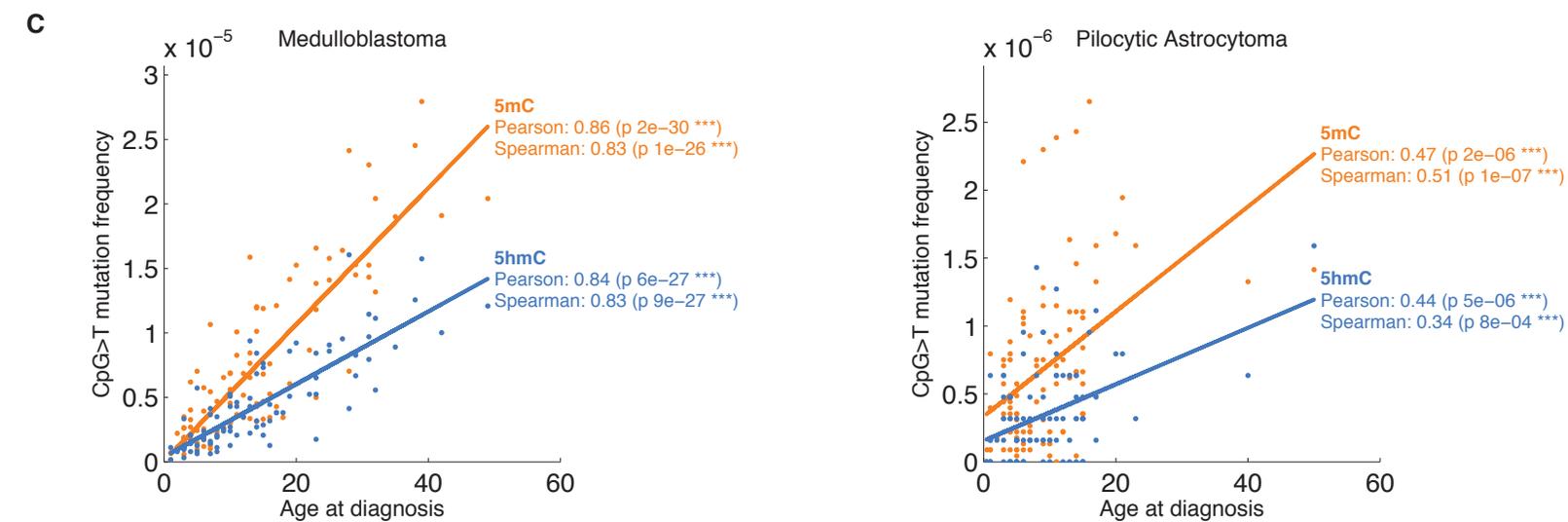
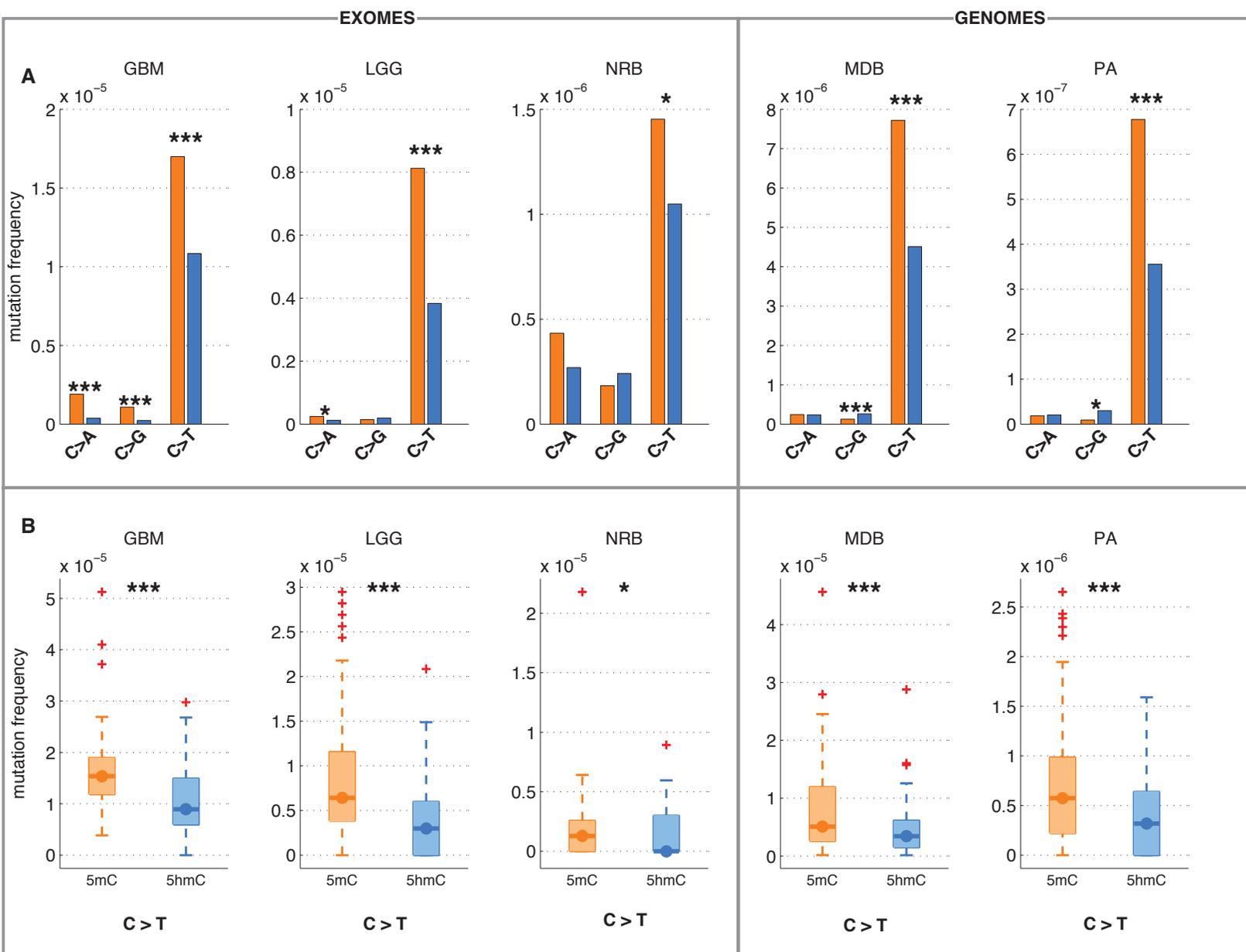
801

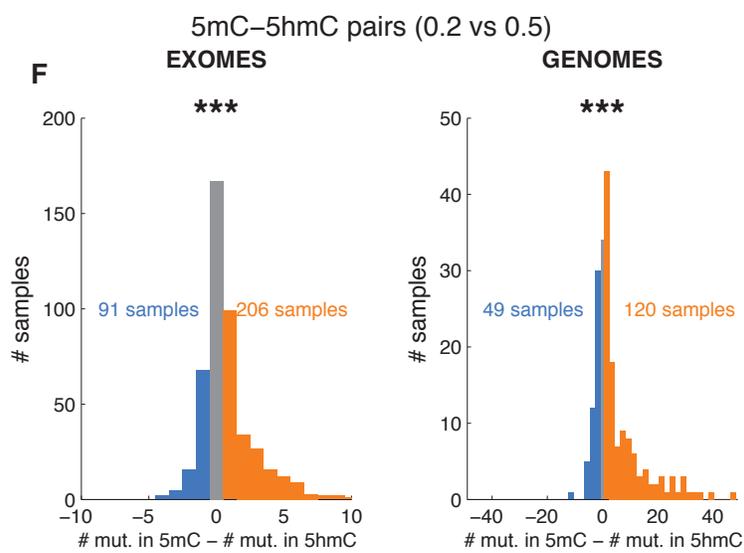
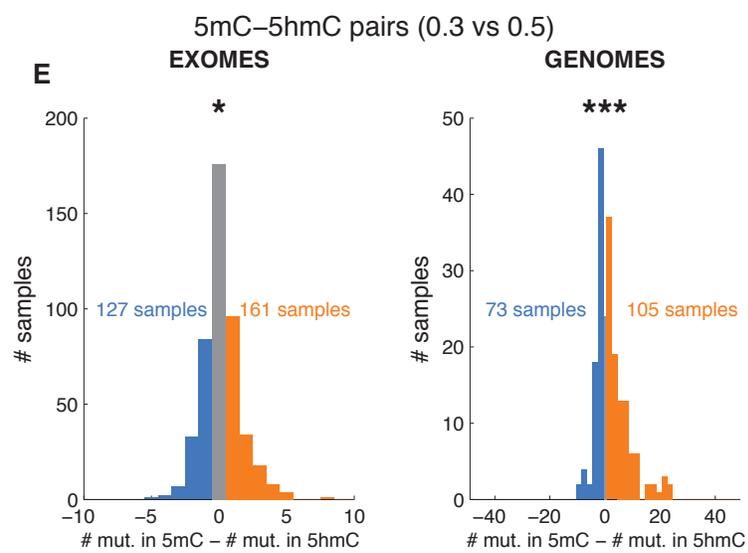
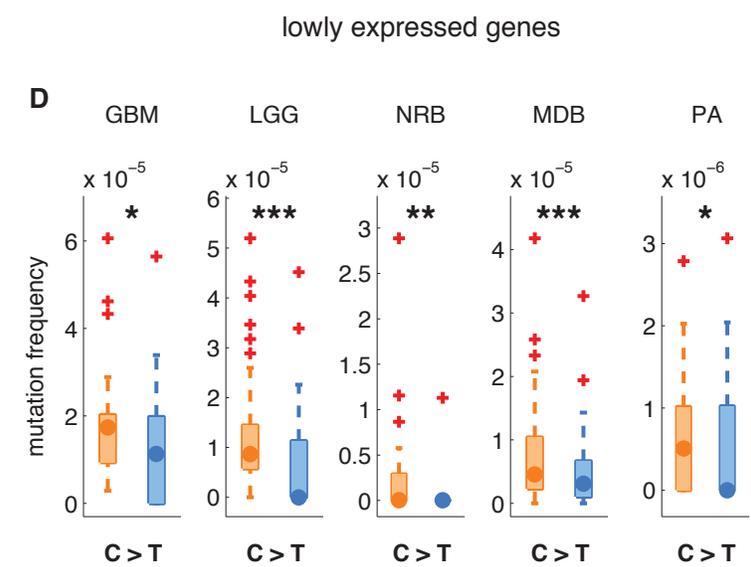
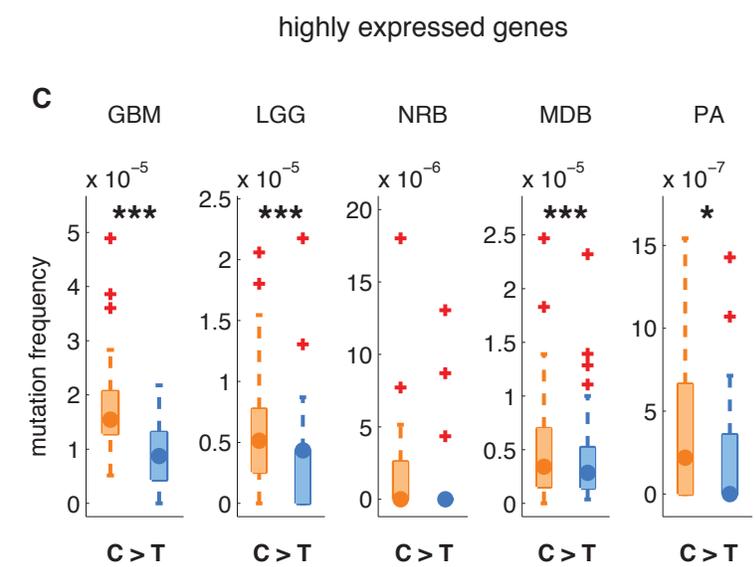
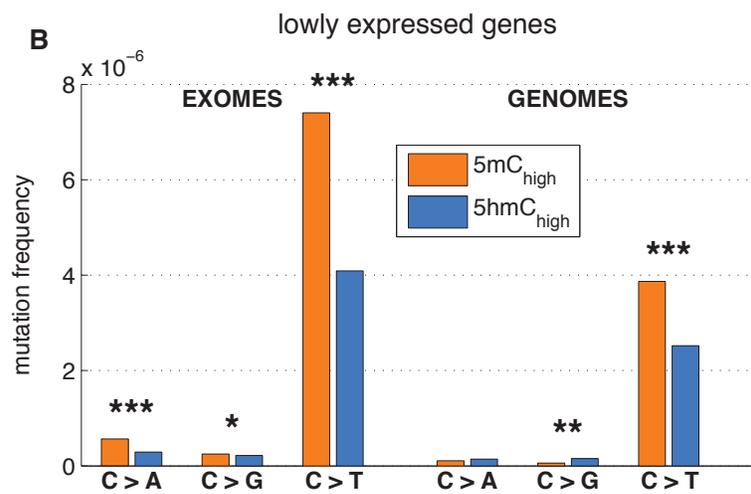
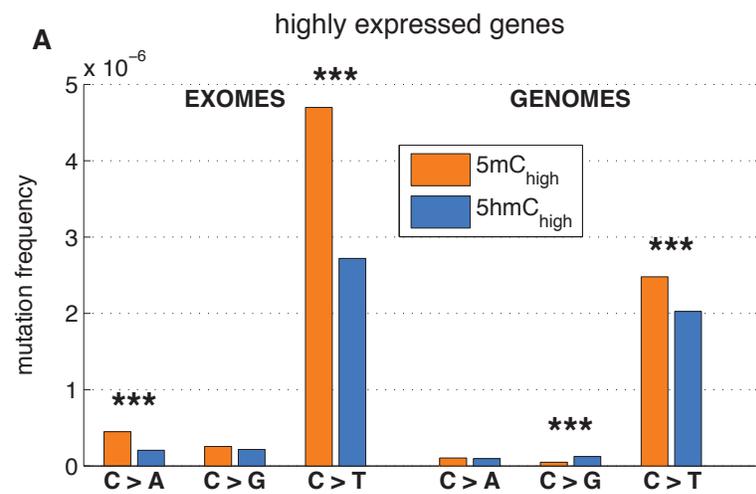
Cancer type	Exome		Whole genome		Tissue	Source
	Patients	SNVs	Patients	SNVs		
<b>Glioblastoma</b>	39	22954	-	-	brain	(Alexandrov et al., 2013)
<b>Glioma Low Grade</b>	215	9678	-	-	brain	(Alexandrov et al., 2013)
<b>Medulloblastoma</b>	-	-	100	139553	brain	(Alexandrov et al., 2013)
<b>Neuroblastoma</b>	210	5027	-	-	brain	(Alexandrov et al., 2013)
<b>Pilocytic Astrocytoma</b>	-	-	101	12989	brain	(Alexandrov et al., 2013)
<b>Kidney Chromophobe</b>	65	1646	-	-	kidney	(Alexandrov et al., 2013)
<b>Kidney Clear Cell</b>	325	30273	-	-	kidney	(Alexandrov et al., 2013)
<b>Kidney Papillary</b>	100	6479	-	-	kidney	(Alexandrov et al., 2013)
<b>RECA-EU</b>	-	-	95	488922	kidney	ICGC
<b>Acute Myeloid Leukaemia</b>	147	2214	7	3659	blood myeloid	(Alexandrov et al., 2013)
<b>Acute Myeloid Leukaemia</b>	-	-	49	176164	blood myeloid	TCGA
<b>Acute Lymphoblastic Leukaemia</b>	140	1869	1	7881	blood other	(Alexandrov et al., 2013)
<b>Chronic Lymphoid Leukaemia</b>	103	3998	28	54746	blood other	(Alexandrov et al., 2013)
<b>Lymphoma B-cell</b>	24	824	24	142753	blood other	(Alexandrov et al., 2013)
<b>Myeloma</b>	69	3973	-	-	blood other	(Alexandrov et al., 2013)
<b>Breast</b>	844	55731	119	647692	breast	(Alexandrov et al., 2013)
<b>Liver</b>	-	-	88	899445	liver	(Alexandrov et al., 2013)
<b>Lung Adeno</b>	636	248519	24	1505512	lung	(Alexandrov et al., 2013)
<b>Lung Small Cell</b>	70	17639	-	-	lung	(Alexandrov et al., 2013)
<b>Lung Squamous</b>	176	70485	-	-	lung	(Alexandrov et al., 2013)
<b>Pancreas</b>	98	5093	15	122787	pancreas	(Alexandrov et al., 2013)
<b>Stomach</b>	212	102110	0	0	stomach	(Alexandrov et al., 2013)
<b>Stomach</b>	0	0	100	1995618	stomach	(Wang et al., 2014)

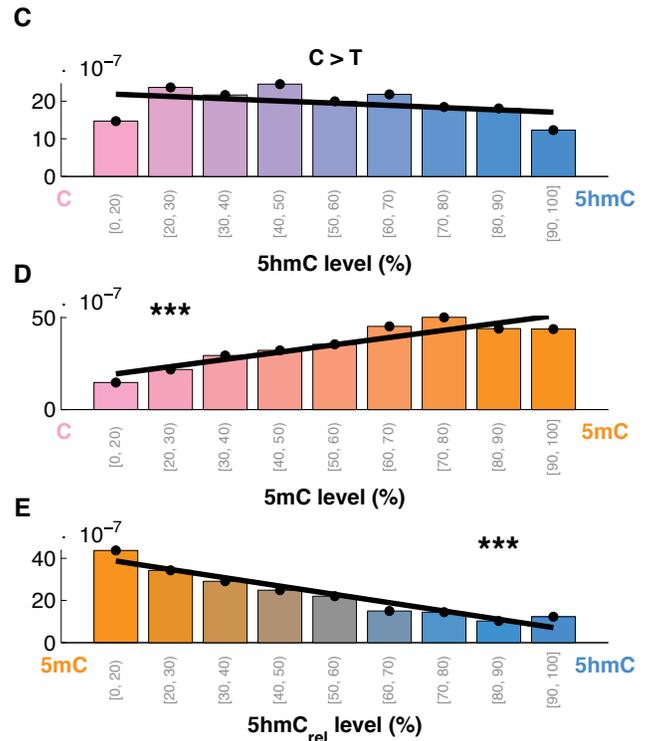
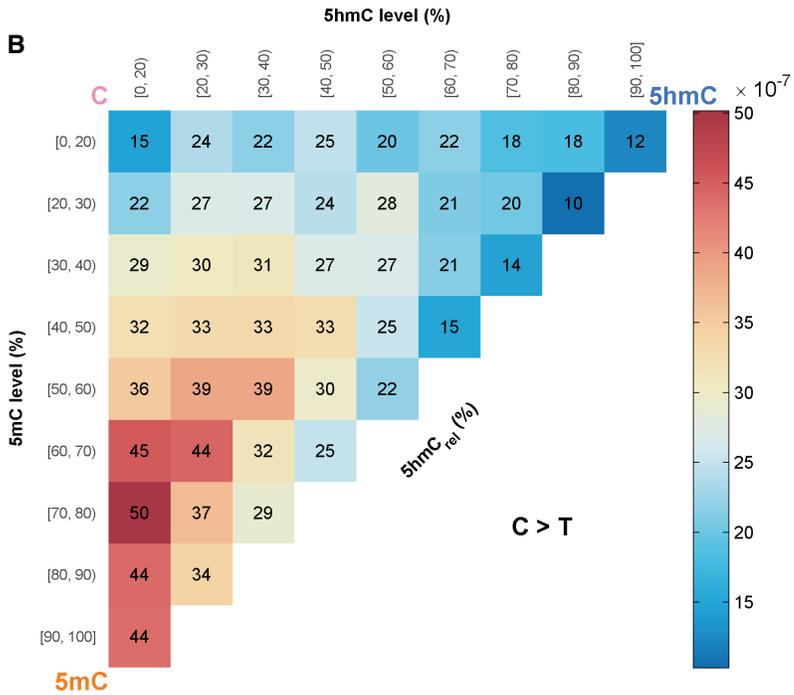
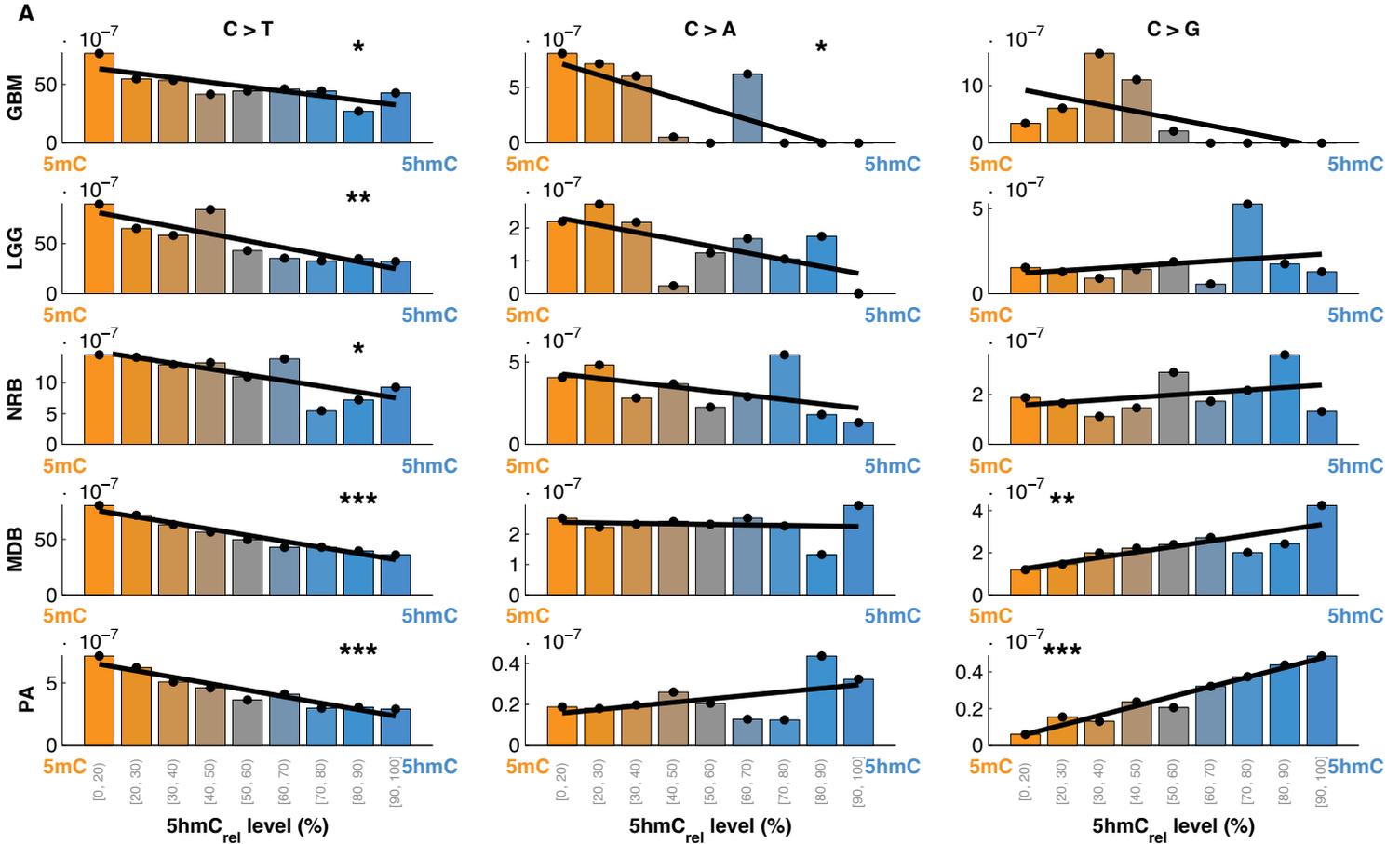
Tissue	BS-Seq		TAB-Seq		Source	Link
	CpGs	Average	CpGs	Average		
<b>Brain</b>	53388534	0.791200	53847986	0.221845	(Wen et al., 2014)	<a href="#">BS</a> , <a href="#">TAB</a>
<b>Kidney r1</b>	54117976	0.759816	46303252	0.088314	(Chen et al., 2015)	<a href="#">BS</a> , <a href="#">TAB</a>
<b>Kidney r2</b>	53861360	0.756454	54585341	0.094624	(Chen et al., 2015)	<a href="#">BS</a> , <a href="#">TAB</a>
<b>Kidney total</b>	54857866	0.757601	54928295	0.092868	(Chen et al., 2015)	
<b>Blood dendritic r1</b>	24586388	0.791371	-	-	(Pacis et al., 2015)	<a href="#">BS</a>
<b>Blood dendritic r2</b>	23524704	0.786322	-	-	(Pacis et al., 2015)	<a href="#">BS</a>
<b>Blood dendritic r3</b>	24419613	0.782467	-	-	(Pacis et al., 2015)	<a href="#">BS</a>
<b>Blood dendritic r4</b>	24452547	0.787728	-	-	(Pacis et al., 2015)	<a href="#">BS</a>
<b>Blood dendritic r5</b>	24399654	0.774058	-	-	(Pacis et al., 2015)	<a href="#">BS</a>
<b>Blood dendritic r6</b>	24533745	0.790145	-	-	(Pacis et al., 2015)	<a href="#">BS</a>
<b>Blood dendritic total</b>	25586845	0.789083	27754454	0.029103	(Pacis et al., 2015)	<a href="#">TAB</a>
<b>Breast</b>	53222114	0.735273	-	-	Epigenome Roadmap	<a href="#">BS</a>
<b>Pancreas</b>	54341922	0.697484	-	-	Epigenome Roadmap	<a href="#">BS</a>
<b>Lung</b>	54236520	0.774050	-	-	Epigenome Roadmap	<a href="#">BS</a>
<b>Liver</b>	51884076	0.757123	-	-	Epigenome Roadmap	<a href="#">BS</a>
<b>Stomach</b>	54054176	0.762082	-	-	Epigenome Roadmap	<a href="#">BS</a>
<b>Blood (hematopoietic multipotent progenitor cell)</b>	51822931	0.852170	-	-	Blueprint Epigenome	<a href="#">BS</a>

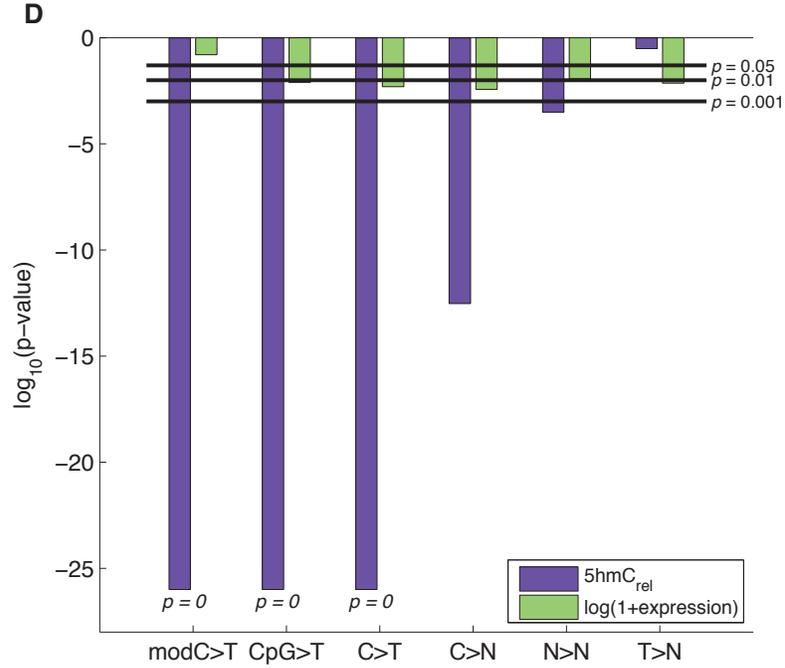
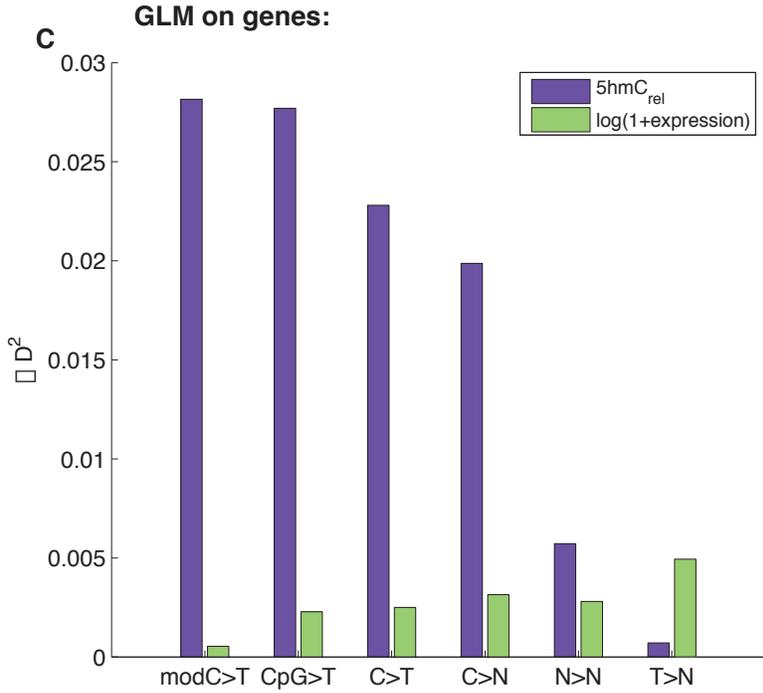
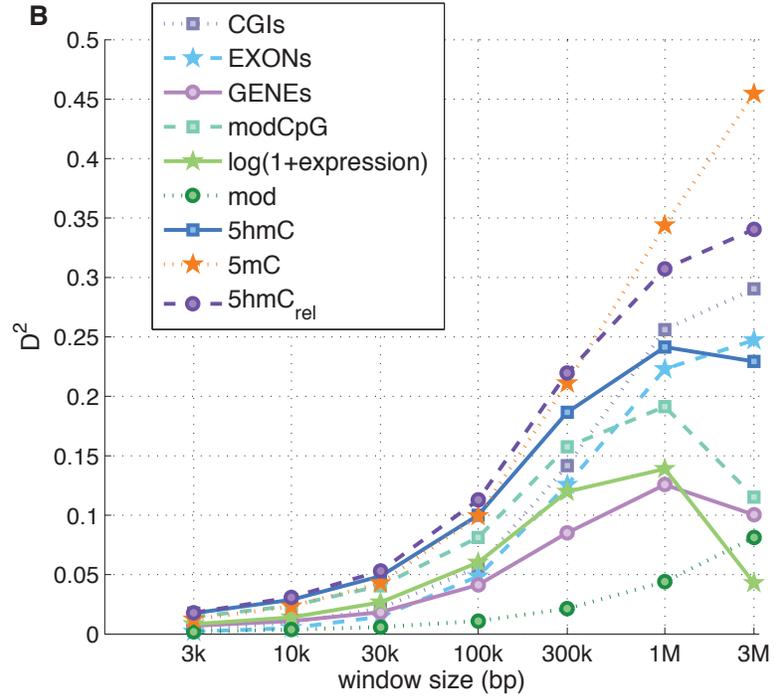
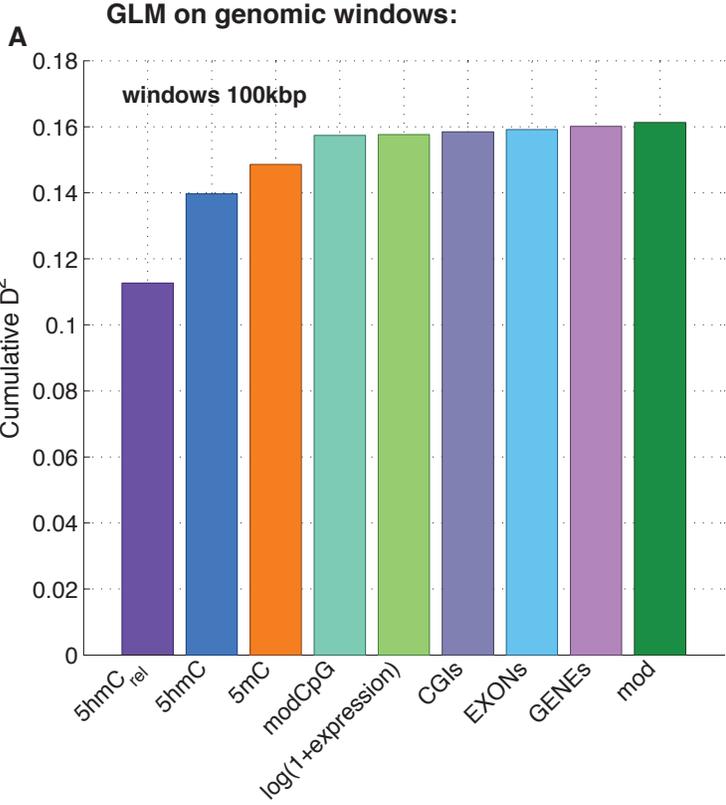
803



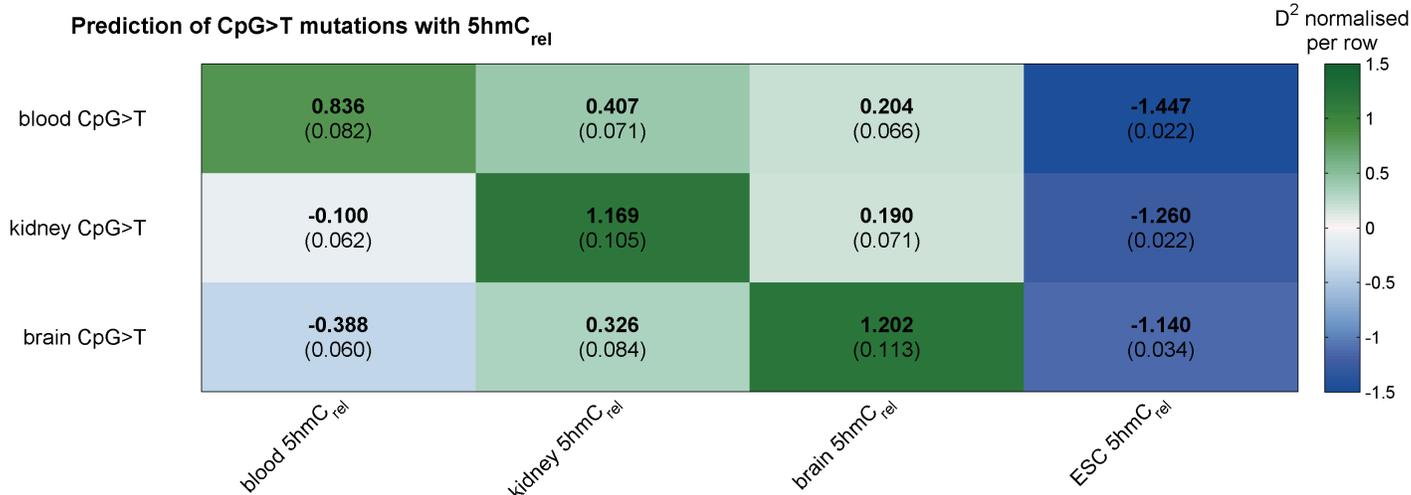




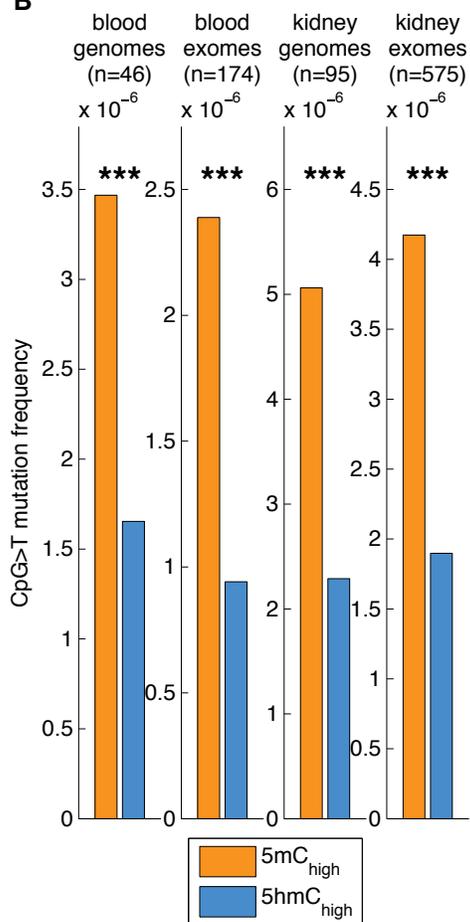




**A**



**B**



**C**

