

1 **Comprehensive longitudinal study of epigenetic mutations in**
2 **aging**

3 Yunzhang Wang¹, Robert Karlsson¹, Juulia Jylhävä¹, Åsa K. Hedman^{2,3}, Catarina Almqvist^{1,4}, Ida K.
4 Karlsson^{1,5}, Nancy L. Pedersen¹, Malin Almgren⁶, Sara Hägg¹

5

6 1. Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm,
7 Sweden

8 2. Rheumatology Unit, Department of Medicine Solna, Karolinska Institutet, Stockholm, Sweden

9 3. Pfizer Worldwide Research and Development, Stockholm, Sweden

10 4. Astrid Lindgren Children's Hospital, Karolinska University Hospital, Stockholm, Sweden

11 5. Institute of Gerontology and Aging Research Network – Jönköping (ARN-J), School of Health
12 and Welfare, Jönköping University, Jönköping, Sweden

13 6. Department of Clinical Neuroscience, Centrum for Molecular Medicine, Karolinska Institutet,
14 Stockholm, Sweden

15

16 Email: yunzhang.wang@ki.se; robert.karlsson@ki.se; juulia.jylhava@ki.se; asa.hedman@ki.se;
17 catarina.almqvist@ki.se; ida.karlsson@ki.se; nancy.pedersen@ki.se; malin.almgren@ki.se;
18 sara.hagg@ki.se

19

20 *Corresponding author:

21 Dr. Sara Hägg, Associate Professor, Department of Medical Epidemiology and Biostatistics,
22 Karolinska Institutet, Nobels väg 12A, Stockholm 17177, Sweden,

23 Phone: +46-8-524 82236

24

25 **Abstract**

26 **Background:** The role of DNA methylation in aging has been widely studied. However, epigenetic
27 mutations, here defined as aberrant methylation levels compared to the distribution in a
28 population, are less understood. Hence, we investigated longitudinal accumulation of epigenetic
29 mutations, using 994 blood samples collected at up to five time points from 375 individuals in old
30 ages.

31 **Results:** We verified earlier cross-sectional evidence on the increase of epigenetic mutations with
32 age, and identified important contributing factors including sex, CD19+ B cells, genetic
33 background, cancer diagnosis and technical artifacts. We further classified epigenetic mutations
34 into High/Low Methylation Outliers (HMO/LMO) according to their changes in methylation, and
35 specifically studied methylation sites (CpGs) that were prone to mutate (frequently mutated
36 CpGs). We validated four epigenetically mutated CpGs using pyrosequencing in 93 samples.
37 Furthermore, by using twins, we concluded that the age-related accumulation of epigenetic
38 mutations was not related to genetic factors, hence driven by stochastic or environmental effects.

39 **Conclusions:** Here we conducted a comprehensive study of epigenetic mutation and highlighted
40 its important role in aging process and cancer development.

41

42 **Key words:** Epigenetic mutation, aging, cancer, twin study

43

44 Introduction

45 Epigenetic processes, among which DNA methylation is one of the most well studied, are
46 fundamental in human aging [1]. Studies on DNA methylation have identified age-associated
47 changes in methylation levels shared by individuals [2,3], and have also reported an increasing
48 divergence of methylation levels between individuals with age [4,5].

49 Epigenetic mutations, defined as aberrant methylation levels that can lead to unusual gene
50 expression, may be involved in cancer development and important for human aging [6,7]. Unlike
51 age-associated changes in methylation levels that are shared among individuals, the incidences of
52 epigenetic mutations are rare, stochastic and inconsistent between individuals. Epigenetic
53 mutations can partly explain the increasing variability of methylation levels between individuals
54 over time, but the extreme methylation levels may incur stronger biological consequences, such
55 as cancer. Epigenetic mutations could contribute to the aging process through the accumulation
56 of abnormally methylated CpGs (cytosine-phosphate-guanine sites), which could further cause
57 abnormal gene expression and downstream effects in tissues. A previous study by Gentilini *et al*
58 [7] specifically defined epigenetic mutations as extreme outliers within a population, with
59 methylation levels exceeding three times interquartile ranges (IQR) of the first quartile ($Q1-3 \times$
60 IQR) or the third quartile ($Q3+3 \times$ IQR). They found that the total numbers of epigenetic
61 mutations increased exponentially with age. However, since this finding was based on a cross-
62 sectional study, it needs to be validated in a longitudinal setting, where the accumulation of
63 epigenetic mutations over time can be followed within the same individuals. Moreover, it is not
64 yet known what the clinical consequences of accumulated epigenetic mutations are, and if
65 individuals with a high burden of epigenetic mutations are prone to develop cancer as previously

66 suggested [6,8].

67 In this study, we used a Swedish twin cohort including 375 individuals sampled up to five times in
68 late life across 18 years (Table 1). We first validated the age-related increase of epigenetic
69 mutations from a longitudinal perspective. Next, we identified important factors associated with
70 the number of epigenetic mutations, including sex, cellular composition (CD19 B-cells), genetic
71 background and technical artifacts. In parallel, we analyzed the direction of change in
72 methylation level and characterized the epigenetic mutations as High- (HMO) and Low
73 Methylation Outliers (LMO). We also studied the association between epigenetic mutations and
74 cancer, as well as the genetic influence on epigenetic mutations using a twin approach. Last, we
75 validated a select set of epigenetic mutations using bisulfite pyrosequencing.

76 **Results**

77 **Longitudinal accumulation of epigenetic mutations is exponentially** 78 **associated with age**

79 To explore the longitudinal increase in number of epigenetic mutations, we measured DNA
80 methylation data (Illumina 450k array) repeatedly in whole blood samples (n=994) from
81 participants in the Swedish Adoption/Twin Study of Aging (SATSA; Table 1) [9]. To avoid
82 confounding by underlying genetic variation, we removed 20,660 CpGs that were associated with
83 at least one single nucleotide polymorphism (SNP) ($p < 1e-14$) within 1 Mbps (mega base pairs), i.e.
84 cis-methylation quantitative loci (cis-meQTLs). In the remaining 370,234 CpGs, the number of
85 epigenetic mutations ranged from 58 to 26,291 in each sample, using the definition in Gentilini *et*
86 *al* [7]. Across samples, the number of epigenetic mutations had a right-skewed distribution,

87 which was close to normal distribution after log₁₀-transformation (Figure S1).

88 After identifying epigenetic mutations in SATSA, we found that the log₁₀ total number of
89 epigenetic mutations increased with age ($p=1.22e-13$) longitudinally (Figure 1A). We also
90 identified additional factors and confounders associated with the number of epigenetic
91 mutations (Table 2). Women had a slightly higher average number of epigenetic mutations than
92 men ($p=6.33e-3$). Low sample quality, as defined by the log₁₀-transformed number of CpGs with
93 detection p-values over 0.01, was positively associated with the total number of epigenetic
94 mutations ($p=1.48e-117$). In general, unreliable samples tended to have more epigenetic
95 mutations, indicating that measurement errors could also be identified as epigenetic mutations.
96 However, after adjusting the mixed models for detection p-value, the effect of age on number of
97 epigenetic mutations remained unchanged. Using predicted cellular compositions, CD19+ B cell
98 composition was positively associated with the total number of epigenetic mutations ($p=5.06e-$
99 23). After removing cis-meQTLs, the first genetic principal component (PC) showed only a minor
100 effect on the total number of epigenetic mutation ($p=0.041$).

101 Out of all CpGs, 237,398 (64%) were defined as epigenetic mutations in at least one sample, but
102 only 1,185 (0.32%) CpGs were mutated in more than 50 samples; subsequently defined as
103 frequently mutated CpGs. Only two of the 1,185 frequently mutated CpGs were also identified to
104 be age-differentially methylated sites (aDMS) in our previous study [3]. The frequently mutated
105 CpGs were still significantly associated with age, sample quality, CD19+ B cell compositions and
106 genetic PC1, while sex was no longer significant (Table 2).

107 **High/Low Methylation Outliers**

108 Compared to normal methylation levels in the population, epigenetic mutations can be either

109 higher or lower in methylation level. Hence, we defined HMO and LMO as CpGs with abnormally
110 higher or lower methylation levels than the average (Figure S2). Of the defined epigenetic
111 mutation sites, almost half were identified as HMOs and the other half as LMOs (118,259 HMOs
112 and 119,175 LMOs). Thirty-six CpGs were defined as both HMOs and LMOs because those sites
113 had intermediate methylation levels and very small IQRs. However, among the frequently
114 mutated CpGs, there were significantly more HMOs than LMOs (969 and 216, $p < 1e-16$) (Figure 2).
115 Nevertheless, numbers of both sets of frequent mutations (log10-transformed) significantly
116 increased with age ($p = 2.09e-17$ for HMOs and $p = 1.14e-05$ for LMOs) (Figure 1B and C). Sex was
117 no longer a significant factor with either frequent HMOs or LMOs. The composition of CD19+ B
118 cell was still strongly associated with HMOs ($p = 2.25e-12$), but only marginally significant for
119 LMOs ($p = 0.046$). Sample quality, as measured by detection p-value, showed strong effects on
120 both frequent HMOs and LMOs, however LMOs were much more influenced ($p = 8.09e-30$) than
121 HMOs ($p = 3.58e-8$). Moreover, the first genetic principal component became a significant factor
122 ($p = 7.65e-5$) when analyzing frequent HMOs, while it had no effect on LMOs ($p = 0.92$) (Table 2).

123 **Functional annotation of epigenetic mutations**

124 To characterize HMO and LMOs, we examined their locations in relation to CpG island regions
125 and regulatory features. Compared to all CpGs analyzed, where 33.5% of CpGs locate in CpG
126 islands, HMOs were enriched within CpG islands (63% of CpGs, $p < 1e-16$) and frequent HMOs
127 even more so (88% of CpGs, $p < 1e-16$). On the other hand, LMOs were mostly located outside of
128 CpG islands (88% CpGs outside of CpG islands, $p < 1e-16$), but the opposite was true for frequent
129 LMOs, which were enriched in CpG islands (51% of CpGs, $p = 8.6e-8$) (Figure 3). We further
130 explored regulatory features of the frequent epigenetic mutations using the Ensembl database

131 [10], and found that frequent HMOs were enriched in promoter regions ($p=1.1e-10$), but less
132 likely to be found in CCCTC-Binding factor (CTCF) binding sites ($p=1.4e-09$) and regions of open
133 chromatin ($p=3.6e-07$) (Figure 4A). The frequent LMOs, on the other hand, were enriched in CTCF
134 ($p=7.7e-12$) and transcription factor binding sites ($p=3.9e-05$), open chromatin ($p=0.0012$), and
135 promoter flanking regions ($p=0.041$), while depleted in promoter regions ($p=6.9e-19$) (Figure 4B).

136 **Epigenetic mutation is associated with cancer diagnosis**

137 As aberrant DNA methylation levels in gene regulatory regions may cause abnormal gene
138 expression, which may be associated with cancer, we analyzed epigenetic mutations in relation to
139 cancer diagnosis in the SATSA participants. Cancer diagnosis date was retrieved using linkage to
140 The National Patient Registry (prior to May 2016) including ICD-codes for all cancer types (ICD7
141 codes 140-205, ICD8 codes 140-209, ICD9 codes 140-208, ICD10 codes C00-C97 and B21). The
142 SATSA participants included 29 prevalent cancer cases diagnosed already at baseline, and 79
143 incident cases that developed cancer during the follow-up period. Hence, information on
144 whether the participant was diagnosed with cancer by the end of the follow-up was tested in the
145 mixed model for associations with log₁₀-transformed numbers of epigenetic mutations. Samples
146 of individuals with cancer, including samples before and after cancer diagnosis, were observed to
147 have a higher number of frequent HMOs ($p=0.013$), but no associations were found for frequent
148 LMOs ($p=0.71$, Table 2). Furthermore, in the survival analysis, people with a higher number of
149 frequent HMOs had a higher risk of cancer incidence (Table S1).

150 **Epigenetic mutations are shared within twin pairs**

151 By applying a co-twin control design we could further study the genetic effect and the genetic-
152 age interaction in association with epigenetic mutations. We calculated the number of shared
153 epigenetic mutations within a twin pair sampled at the same time, and studied their association

154 with time and twin zygosity using a random effects model (Table 3). The numbers of shared
155 epigenetic mutations were normalized in order to compare the effect sizes from different sets of
156 CpGs. First, taking all CpGs into account ($n=390,894$), the number of shared epigenetic mutations
157 increased significantly with age ($\beta=0.019$, $p=0.026$) and MZ pairs shared more epigenetic
158 mutations than DZ pairs ($\beta=1.078$, $p=3.41e-18$). After excluding 20,660 cis-meQTL CpGs, the age
159 effect became stronger ($\beta=0.025$, $p=5.98e-3$) while the zygosity effect was smaller ($\beta=0.855$,
160 $p=1.05e-11$). Last, within the 20,660 cis-meQTL-CpGs, the number of shared epigenetic
161 mutations was not associated with age ($\beta=2.86e-4$, $p=0.969$), while the zygosity difference
162 ($\beta=1.461$, $p=8.34e-28$) was larger than in results from non-meQTL-CpGs. None of the three tests
163 showed significant twin zygosity-age interaction or sex effect.

164 **Epigenetic mutations were validated using pyrosequencing**

165 To verify epigenetic mutations identified from 450k array, we selected four frequently mutated
166 CpGs (One HMO: cg05270750, and three LMOs: cg17338133, cg25351353, cg05124918) in 93
167 samples from 26 individuals for validation with pyrosequencing. In general, the pyrosequencing
168 results were well correlated with methylation data measured by the 450k array (cg05270750:
169 $r=0.84$; cg17338133: $r=0.59$; cg25351353: $r=0.80$; cg05124918: $r=0.77$). In addition, we compared
170 methylation levels of mutated samples to the normal group using results from the 450k array and
171 pyrosequencing respectively. In pyrosequencing data, significant differences were observed
172 between mutated samples and normal ones, using the same definition of a mutated sample as
173 that for the 450k array data (Table 4). Hence, pyrosequencing technically validated epigenetic
174 mutations identified from the 450k array. Although the agreement between the two methods
175 was generally good, we still observed large differences between pyrosequencing and 450k data in

176 some samples, where four samples in cg17338133 and six samples in cg 05124918 showed over
177 15% methylation level differences between 450k array and pyrosequencing data after centering
178 their mean methylation levels. This indicates that we might wrongly-detect or fail to detect
179 epigenetic mutations from 450k chip data. In general, pyrosequencing data were more stable and
180 changes in methylation levels were smoother than that from 450k array (Figure 5). For example,
181 in cg05270750 measured by the 450k array (Figure 5E), one participant was identified to have
182 epigenetic mutations in the first three measures, but the methylation level turned back to normal
183 status in the last two measures. However, pyrosequencing data showed the methylation levels of
184 the five measures from this individual were consistently defined as epigenetic mutations.

185 **Functional validation of epigenetic mutations in cancer tissues**

186 To further verify the overabundance of epigenetic mutations in cancer tissues, we picked a gene
187 PR/SET domain 7 (*PRDM7*) which was the only gene related to CpGs tested in pyrosequencing
188 (cg05270750), and analyzed DNA methylation and gene expression data of the gene in tumor
189 tissues and normal adjacent tissues using The Cancer Genome Atlas (TCGA) [11] data
190 downloaded from Wanderer [12]. We selected the four most common cancer types in both sexes
191 combined: lung cancer, breast cancer, colorectal cancer and prostate cancer [13]. The total
192 numbers of tumor and normal adjacent samples were 2,209 and 261 respectively, all cancer
193 types combined. On average, the expression levels of *PRDM7* were higher in tumor tissues than
194 normal adjacent tissues in all cancer types, but the difference was only statistically significant for
195 lung cancer ($p=1.83e-09$, Table S2). For DNA methylation data, the tumor tissues had significantly
196 lower methylation levels than normal adjacent tissues in the gene body (Figure 6A). However, for
197 CpGs in the *PRDM7* promoter (from cg06295223 to cg26935333), there was no significant

198 difference between the mean methylation levels of cancer and normal adjacent tissues (Figure
199 6A). To quantify and compare epigenetic mutations in both tissues, we used the distribution of
200 normal adjacent samples to determine epigenetic mutation cutoffs. By calculating the number of
201 epigenetic mutations in tissue samples, tumor tissues had higher proportions of epigenetic
202 mutations in the gene body, while epigenetic mutations were not observed in normal adjacent
203 tissues. In the gene promoter, tumor and normal adjacent tissues had similar and relatively low
204 proportions of epigenetic mutations (Figure 6B).

205 **Discussion**

206 In this study, we analyzed age-related accumulation of epigenetic mutations from a longitudinal
207 perspective in old Swedish twins. Apart from being exponentially associated with age, epigenetic
208 mutations were also associated with sex, CD19+ B cell count, genetic background, cancer
209 incidence and technical factors. We further analyzed frequent HMOs and LMOs separately and
210 found that biological factors, including B cell compositions and genetic factors, were more
211 strongly associated with frequent HMOs than LMOs, while LMOs were more influenced by
212 technical factors. Moreover, cancer diagnosis was significantly associated with the increase of
213 epigenetic mutations, especially among frequent HMOs, while the same was not true for LMOs.
214 Emerging evidence indicate that epigenetic mutations could be related to cancer [6], as
215 epigenetic mutations may cause abnormal gene expression, which could contribute to the
216 development of cancer. On the other hand, mutated DNA sequences and abnormal epigenetic
217 regulation in tumor cells may in turn cause more epigenetic mutations. In this study, we found
218 that the number of epigenetic mutations was significantly higher in samples of individuals who

219 were diagnosed with cancer by the end of follow-up. Therefore, we conclude that the number of
220 epigenetic mutations may accumulate long before the diagnosis of cancer. The survival analysis
221 further showed that a higher number of frequent HMOs could be a risk factor for cancer
222 incidence. These results support a previous finding where the number of epigenetic mutations
223 were higher in tumor tissues than in normal tissues [8]. Follow-up studies with more participants
224 are needed to better establish the possible relationship between epigenetic mutations and
225 cancer.

226 In this study, DNA methylation data were corrected for cellular compositions predicted by the
227 Houseman method [14], yet imputed CD19+ B cell count was significantly associated with
228 epigenetic mutations, but not other cell types. A possible explanation could be that B cells have a
229 unique methylation pattern compared to other lymphocytes [15]. Also, B cell composition was
230 still a strong factor for frequent HMOs but the effect became very weak for frequent LMOs,
231 probably because cell specific CpGs are enriched in promoter regions [15] where HMOs are
232 mostly found.

233 When studying functional annotations associated with the epigenetic mutations, we found that
234 the location and regulatory features were different for frequent HMOs and LMOs. The observed
235 enrichment of HMOs in CpG islands and promoter regions indicated that HMOs were more
236 related to biological function than LMOs, which is in line with the fact that technical bias was
237 significant in LMOs.

238 The concept of epigenetic mutations should be discussed in relation to methylation variability, as
239 they both describe methylation divergence between individuals. However, epigenetic mutations
240 refer to more extreme methylation levels carried by a small number of individuals, while

241 methylation variability is considered to be a population pattern. In contrast to traditional
242 association studies on methylation levels, where CpGs of higher variances are more likely to have
243 statistical power, CpGs of high variances could have too large inter quartile ranges to be identified
244 as epigenetic mutations by definition. Therefore, the identified frequent epigenetic mutations
245 were different from the age-associated CpGs or age-varied CpGs reported prior to this study
246 using the same data [3,5], and thus may contribute to the aging processes by other ways than
247 through the epigenetic drift.

248 Even after excluding cis-meQTL CpGs, a small genetic effect captured by the first genetic PC was
249 associated with epigenetic mutations, especially in frequent HMOs. To further explore how
250 genetic background and age affected the accumulation of epigenetic mutations, we studied the
251 number of shared epigenetic mutations between twins over time. Here we did not simply
252 exclude cis-meQTL CpGs, but considered them as epigenetic mutations caused by genetic variants
253 inherited at birth. For all CpGs and non-meQTL CpGs, we observed both age and genetic effect
254 associated with the number of shared epigenetic mutations within the twin pair. To isolate the
255 genetic effect, we specifically analyzed cis-meQTL CpGs and found that in this selection, the
256 number of shared epigenetic mutations did not change with age. This result was consistent with
257 a previous paper showing that meQTL-CpG associations are stable over time [16]. Additionally,
258 we failed to detect an interaction between genetic factors and age, indicating that the increase of
259 epigenetic mutations with age was not dependent on the genetic background. Therefore, the
260 remaining genetic effect observed after removing cis-meQTL CpGs was probably due to trans-
261 meQTLs or unidentified cis-meQTLs. In conclusion, the age effect on the accumulation of
262 epigenetic mutations is independent of genetic background. However, we might not have enough

263 statistical power to detect a significant age-genetic interaction on shared epigenetic mutations,
264 since the age effect estimated for MZ twins was larger than for DZ twins. Moreover, due to the
265 limit of the age range in this study (48 to 98 years), we could not exclude the possibility of
266 genetic-associated development of epigenetic mutations in early ages, which remains to be
267 examined by future studies.

268 Technical artifacts and poor sample quality could lead to erroneous measures that interfere with
269 identifying true biological methylation outliers. Although sample quality control based on
270 detection p-value was applied in the pre-processing pipe-line of the methylation data, it was still
271 found to strongly influence the identification of the epigenetic mutations. Although the technical
272 effect was strong and hard to avoid, the effect of age on epigenetic mutations was not biased as
273 we randomized samples on microarrays. Another important technical artifact is the batch effect
274 from different arrays, but we adjust for batches both in data pre-processing and as a random
275 effect in the mixed effect model. Hence, despite the confounding issues from different technical
276 biases when analyzing methylation outliers, the underlying biological phenomenon of increasing
277 number of epigenetic mutations with age still holds.

278 Validation of the epigenetic mutations identified in 450k data was done by pyrosequencing,
279 which also detected aberrant methylation levels proving that they were true biological outliers
280 and not simply technical errors. However, some samples showed very different results between
281 the two methods suggesting measurement errors existed. When comparing results from the two
282 methods, pyrosequencing data were more stable and better indicated that epigenetic mutations
283 were persistent over time, which supported the accumulation of epigenetic mutations as a factor
284 of aging.

285 The HMO site cg05270750 validated by pyrosequencing is located in the promoter region of the
286 gene *PRDM7*, which encodes a Histone-Lysine Trimethyltransferase involved in histone
287 modification. To further explore the potential consequence of epigenetic mutations, we analyzed
288 DNA methylation and gene expression of gene *PRDM7* in data on tumor and normal adjacent
289 tissues from TCGA. The expression of *PRDM7* in normal adjacent tissues were very low, as
290 previously seen [17]. Nevertheless, we observed higher expression of *PRDM7* in tumor tissues,
291 especially in lung cancers, suggesting the abnormal expression of *PRDM7* could be related to the
292 dysregulation of histone modification in tumor. On the other hand, we observed similar
293 proportions of epigenetic mutations between tumor and normal adjacent tissues in the gene
294 promoter, but more epigenetic mutations in the gene body for tumor tissues. Since normal
295 adjacent tissue can be regarded as an intermediate state between healthy and tumor tissues, it is
296 suggested that, in the process of cancer development, epigenetic mutations were likely to first
297 accumulate in gene promoters and then spread to the whole epigenome.

298 **Conclusions**

299 In summary, using longitudinal DNA methylation data, we showed that the accumulation of
300 epigenetic mutations is exponentially associated with age in old adults, and once mutations are
301 established, they are stable over time. Furthermore, epigenetic mutations are enriched in
302 important regulatory sites, e.g. promoter regions of genes involved in histone modification
303 processes, which could potentially be an explanation to why people who develop cancer have
304 more epigenetic mutations than others do. In addition, we showed that the burden of
305 accumulation associated with the human aging process is unlikely to be driven by underlying

306 genetic background. Hence, accumulation of epigenetic mutations is an underexplored area in
307 the field of aging, and warrants further studies to enhance our understanding of this
308 phenomenon.

309 **Methods**

310 **Study population**

311 Twins as participants in this study were enrolled in the SATSA longitudinal cohort study [18]. After
312 quality control, a total of 994 blood samples obtained from 375 individuals in five longitudinal
313 waves (1992-2012) were used in the analyses. The 375 participants had a mean age of 68.9 years
314 (SD=9.7) at their first measurement, and 223 (59.5%) were women. Of the 375 participants, 197
315 contributed samples in three or more waves. Phenotype data were collected through
316 comprehensive questionnaires and physical testing at each sampling wave. Phenotypes used in
317 this study include chronological age, sex, zygosity, smoking status and cancer diagnosis.

318 **DNA methylation data**

319 DNA methylation data were obtained from DNA extracted from whole blood samples measured
320 by Infinium HumanMethylation450 BeadChips. In total 485,512 CpG sites were measured for
321 each sample. The quality control and preprocessing methods of the methylation data were
322 described in a previous study [3]. Samples from individuals lacking genetic data were removed,
323 retaining a total of 994 samples for analyses. Blood cellular compositions were estimated by the
324 Houseman method [14] using a reference panel [15]. The methylation data were adjusted by
325 cellular compositions using a linear regression before the analyses. Additionally, batch effects,
326 which were detected as slides on the 450k chip, were adjusted using the Combat method from

327 the sva package [19].

328 **Genotype data and imputation**

329 Genetic data were measured by Infinium PsychArray (Illumina Inc., San Diego, CA, USA) with
330 588,454 SNPs detected for every individual. After quality control, data were imputed to the 1000
331 Genomes Project phase 1 version 3 reference [20] using IMPUTE2 version 2.3.2 [21,22] with
332 default parameters. The first 10 PCs were calculated based on a linkage disequilibrium pruned set
333 of directly genotyped autosomal SNPs.

334

335 **Identifying epigenetic mutations**

336 The definition of an epigenetic mutation was consistent with Gentilini *et al* [7]. For each CpG, the
337 quartiles of methylation levels were calculated for every CpG using the first observation available
338 from each individual, and were calculated separately for men and women to avoid the sex effect
339 on methylation levels. Samples having methylation levels three times the inter quartile range
340 higher than the third quartile or lower than the first quartile were identified as mutated outliers.
341 Methylation levels were presented in beta-values, which indicate the methylation proportions.
342 CpGs associated with cis-meQTLs (<1 Mbps) were removed from further epigenetic mutation
343 analyses. For the rest of the CpGs, outlier samples were identified as epigenetic mutations, and
344 the total number of epigenetic mutations was counted for every sample. Identified epigenetic
345 mutations were classified into HMOs and LMOs according to whether they exceed the upper or
346 lower boundary of normal methylation levels (defined as 3 times IQR higher than the third
347 quantile or lower than the first quantile).

348 **Statistical analysis**

349 A mixed effect model was fitted to measure the association of the number of epigenetic
350 mutations on age and other factors (Equation 1). A log-10 transformation was applied to the
351 number of epigenetic mutations to form a distribution closer to a normal distribution. For each
352 sample, the log10-transformed number of CpGs with detection p-values over 0.01 was used to
353 indicate the sample quality. In the formula, i , j and k denote individual, slide batch and waves; β_0 ,
354 β_1 , β_2 , β_3 , β_4 , β_5 , β_6 denote fixed intercepts, fixed coefficient of age, sex, CD19 B cell
355 composition, first genetic principal component, detection p-value and whether the individual
356 developed cancer; u_0 , u_1 and ε denotes random intercept of individual, slide batch and random
357 error.

358

$$\begin{aligned} Mut_{i,j,k} = & \beta_0 + \beta_1 Age_{i,j,k} + \beta_2 Sex_i + \beta_3 Bcell_{i,j,k} + \beta_4 PC1_i + \beta_5 Dpval_{i,j,k} + \beta_6 Cancer_i + u_{0i} \\ & + u_{1j} + \varepsilon_{i,j,k} \quad (Eq. 1) \end{aligned}$$

359 The survival analysis of cancer diagnosis and epigenetic mutations was performed using a Cox
360 model. The model included sex, current smoking as baseline exposure, number of epigenetic
361 mutations as a time-varying covariate, and attained age as the time scale. The model was further
362 adjusted for twin pair and batch effect using robust standard error.

363 In twin analysis, a mixed effect model was used to study the number of exact same epigenetic
364 mutations between paired twins measured at the same time in association with age, sex and twin
365 zygosity (Equation 2),

$$\log_{10} N_{i,j} = \beta_0 + \beta_1 Age_{i,j} + \beta_2 Sex_i + \beta_3 Zyg_i + \beta_4 Zyg_i \times Age_{i,j} + u_{0i} + \varepsilon_{i,j} \quad (Eq. 2)$$

366 where i and j denote individual and longitudinal measure; β_0 , β_1 , β_2 , β_3 , β_4 denote fixed
367 intercept, fixed coefficient of age, sex, zygosity and zygosity-age interaction; u_{0i} , and ε denote

368 random intercept of individual and random error.

369 All statistical analyses were performed in R version 3.4.3.

370 **Pyrosequencing**

371 In total, 93 samples from 26 individuals were measured by pyrosequencing to validate epigenetic
372 mutations in 4 CpGs (cg05270750, cg17338133, cg25351353, cg05124918). The samples were
373 selected to present 4 to 5 longitudinal measures for every individual. The selection of CpGs was
374 based on their primer quality, and having large numbers of mutated samples. The primers of the
375 four CpGs were designed using the software PyroMark Assay Design by QIAGEN. DNA samples
376 were converted by bisulfite reaction performed on EZ-96 DNA Methylation-Gold™ MagPrep kit
377 provided by ZYMO RESEARCH CORP. Converted samples were randomized in a 96-well plate and
378 sequenced for each CpG on PyroMark Q96 ID provided by QIAGEN. The raw data were processed
379 in PyroMark Q24 Software v2.5.8 by QIAGEN.

380 **Declarations**

381 **Fundings**

382 This study was supported by NIH grants R01 [AG04563, AG10175, AG028555], the MacArthur
383 Foundation Research Network on Successful Aging, the European Union's Horizon 2020 research
384 and innovation programme [No. 634821], the Swedish Council for Working Life and Social
385 Research (FAS/FORTE) [97:0147:1B, 2009-0795, 2013-2292], the Swedish Research Council [825-
386 2007-7460, 825-2009-6141, 521-2013-8689, 2015-03255, 2015-06796], the Karolinska Institutet
387 delfinansiering (KID) grant for doctoral students (YW), the KI Foundation, the Strategic Research

388 Area in Epidemiology at Karolinska Institutet and by Erik Rönnerbergs donation for scientific studies

389 in aging and age-related diseases.

390 **Availability of data and material**

391 The datasets generated and analyzed during the current study are available in Array Express

392 database of EMBL-EBL (www.ebi.ac.uk/arrayexpress) under accession number E-MTAB-7309.

393 **Authors' contributions**

394 SH, NP and YW conceived and designed this study. YW performed data processing, statistical

395 analysis and drafted the manuscript. YW and MA conducted pyrosequencing for validation. SH,

396 ÅH, RK, JJ, IK and MA contributed to the manuscript writing. All authors read and approved the

397 final manuscript.

398 **Ethics approval and consent to participate**

399 All participants in SATSA have provided written informed consents. This study was approved by

400 the ethics committee at Karolinska Institutet with Dnr 2015/1729-31/5.

401 **Competing interests**

402 The authors declare that they have no competing interests.

403

404

References

405

406

1. López-Otín C, Blasco MA, Partridge L, Serrano M, Kroemer G. The Hallmarks of Aging. *Cell*. 2013 Jun 6;153(6):1194–217.

407

408

2. Bell JT, Tsai P-C, Yang T-P, Pidsley R, Nisbet J, Glass D, et al. Epigenome-Wide Scans Identify Differentially Methylated Regions for Age and Age-Related Phenotypes in a Healthy Ageing Population. *PLoS Genet*. 2012 Apr 19;8(4):e1002629.

409

410

411

3. Wang Y, Karlsson R, Lampa E, Zhang Q, Hedman ÅK, Almgren M, et al. Epigenetic influences on aging: a longitudinal genome-wide methylation study in old Swedish twins. *Epigenetics*. 2018 Sep 2;13(9):975–87.

412

413

414

4. Fraga MF, Ballestar E, Paz MF, Ropero S, Setien F, Ballestar ML, et al. Epigenetic differences arise during the lifetime of monozygotic twins. *Proc Natl Acad Sci U S A*. 2005 Jul 26;102(30):10604–9.

415

416

417

5. Wang Y, Pedersen NL, Hägg S. Implementing a method for studying longitudinal DNA methylation variability in association with age. *Epigenetics*. 2018 Aug 3;13(8):866–74.

418

419

6. Feinberg AP, Koldobskiy MA, Göndör A. Epigenetic modulators, modifiers and mediators in cancer aetiology and progression. *Nat Rev Genet*. 2016 May;17(5):284–99.

420

421

7. Gentilini D, Garagnani P, Pisoni S, Bacalini MG, Calzari L, Mari D, et al. Stochastic epigenetic mutations (DNA methylation) increase exponentially in human aging and correlate with X chromosome inactivation skewing in females. *Aging*. 2015 Aug;7(8):568–78.

422

423

424

8. Gentilini D, Scala S, Gaudenzi G, Garagnani P, Capri M, Cescon M, et al. Epigenome-wide association study in hepatocellular carcinoma: Identification of stochastic epigenetic mutations through an innovative statistical approach. *Oncotarget*. 2017 Jun 27;8(26):41890–902.

425

426

427

428

9. Magnusson PKE, Almqvist C, Rahman I, Ganna A, Viktorin A, Walum H, et al. The Swedish Twin Registry: establishment of a biobank and other recent developments. *Twin Res Hum Genet Off J Int Soc Twin Stud*. 2013 Feb;16(1):317–29.

429

430

431

10. Zerbino DR, Wilder SP, Johnson N, Juettemann T, Flicek PR. The Ensembl Regulatory Build. *Genome Biol*. 2015;16:56.

432

433

11. The Cancer Genome Atlas Home Page [Internet]. The Cancer Genome Atlas - National Cancer Institute. 2011 [cited 2019 Feb 6]. Available from: <https://cancergenome.nih.gov/>

434

435

12. Díez-Villanueva A, Mallona I, Peinado MA. Wanderer, an interactive viewer to explore DNA methylation and gene expression data in human cancer. *Epigenetics Chromatin*. 2015 Jun 23;8(1):22.

436

437

- 438 13. BW S, CP W. World Cancer Report 2014 [Internet]. [cited 2019 Mar 1]. Available from:
439 [http://publications.iarc.fr/Non-Series-Publications/World-Cancer-Reports/World-Cancer-](http://publications.iarc.fr/Non-Series-Publications/World-Cancer-Reports/World-Cancer-Report-2014)
440 [Report-2014](http://publications.iarc.fr/Non-Series-Publications/World-Cancer-Reports/World-Cancer-Report-2014)
- 441 14. Houseman EA, Accomando WP, Koestler DC, Christensen BC, Marsit CJ, Nelson HH, et al. DNA
442 methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinformatics*.
443 2012 May 8;13(1):86.
- 444 15. Reinius LE, Acevedo N, Joerink M, Pershagen G, Dahlén S-E, Greco D, et al. Differential DNA
445 Methylation in Purified Human Blood Cells: Implications for Cell Lineage and Studies on
446 Disease Susceptibility. *PLoS ONE*. 2012 Jul 25;7(7):e41361.
- 447 16. Gaunt TR, Shihab HA, Hemani G, Min JL, Woodward G, Lyttleton O, et al. Systematic
448 identification of genetic influences on methylation across the human life course. *Genome*
449 *Biol*. 2016 Mar 31;17(1):61.
- 450 17. GTEx Portal [Internet]. [cited 2019 Mar 1]. Available from:
451 <https://gtexportal.org/home/gene/PRDM7>
- 452 18. Finkel D, Pedersen NL. Processing Speed and Longitudinal Trajectories of Change for
453 Cognitive Abilities: The Swedish Adoption/Twin Study of Aging. *Aging Neuropsychol Cogn*.
454 2004 Jun 1;11(2–3):325–45.
- 455 19. Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD. The sva package for removing batch
456 effects and other unwanted variation in high-throughput experiments. *Bioinformatics*. 2012
457 Mar 15;28(6):882–3.
- 458 20. The 1000 Genomes Project Consortium. A global reference for human genetic variation.
459 *Nature*. 2015 Oct 1;526(7571):68–74.
- 460 21. Howie BN, Donnelly P, Marchini J. A Flexible and Accurate Genotype Imputation Method for
461 the Next Generation of Genome-Wide Association Studies. *PLOS Genet*. 2009 Jun
462 19;5(6):e1000529.
- 463 22. Howie B, Marchini J, Stephens M. Genotype Imputation with Thousands of Genomes. *G3*
464 *Genes Genomes Genet*. 2011 Nov 1;1(6):457–70.

465

466

467

Figure legends

468

Figure 1. The number of epigenetic mutations (log₁₀-transformed) increased longitudinally with age in a longitudinal perspective using genome-wide DNA methylation data from repeated whole blood samples collected in the Swedish Adoption/Twin Study of Aging (SATSA; n=375 participants). The numbers of epigenetic mutations of samples were counted from: A) total epigenetic mutations (n=370,234, p=1.22e-13 for association with age), B) frequent high methylation outliers (HMO) (n=969, p=2.09e-17 for association with age), and C) frequent low methylation outliers (LMO) (n=216, p=1.14e-05 for association with age).

475

476

Figure 2. The distribution of mutated samples for high methylation outliers (HMOs) and low methylation outliers (LMOs). For most CpGs, epigenetic mutations only occurred in a small number of samples, but HMOs were more likely to appear in a large number of samples (n>50) than LMOs (969 HMOs and 216 LMOs, p<1e-16).

477

478

479

480

481

Figure 3. Proportions of high methylation outliers (HMOs) and low methylation outliers (LMOs) in different CpG island regions. HMOs are enriched in CpG islands (p<1e-16) while LMOs are more distributed outside of CpG islands (p<1e-16), especially in open sea regions. However, both frequent HMOs and LMOs are enriched in CpG islands (p<1e-16 and p=8.6e-8).

482

483

484

485

486

Figure 4. The distribution of regulatory features of frequent high methylation outliers (HMOs) and low methylation outliers (LMOs). Compared to the background distribution of the 450k array design, frequent HMOs were enriched in promoter regions (A), while the opposite was true for LMOs (B).

487

488

489

490

Figure 5. The longitudinal change of four CpGs in 93 samples from 26 individuals measured by 450k array (left panel) and pyrosequencing (Pyroseq, right panel) techniques. Methylation levels of A) cg05270750 from 450k-chip, B) cg05270750 from Pyroseq, C) cg17338133 from 450k-chip, D) cg17338133 from Pyroseq, E) cg25351353 from 450k-chip, F) cg25351353 from Pyroseq, G) cg05124918 from 450k-chip, H) cg05124918 from Pyroseq. Samples are shown as points colored by their mutation status defined by the 450k data and lines links longitudinal samples collected in the same individual.

491

492

493

494

495

496

497

498

Figure 6. Comparing the DNA methylation and epigenetic mutation patterns of gene *PRDM7* between tumor and normal adjacent tissues. Data were downloaded from TCGA through Wanderer. The cancer types included lung cancer, breast cancer, colorectal cancer and prostate cancer. A) The location of CpGs related to gene *PRDM7* in UCSC genome browser. B) The methylation levels of CpGs in gene *PRDM7*. Tumor and normal adjacent tissues had similar methylation levels in the gene promoter, while the methylation levels of tumor tissues in the gene body were significantly lower than normal adjacent tissues. C) The proportion of epigenetic mutations in tumor and normal adjacent tissues. Tumor tissues had higher proportions of epigenetic mutations in the gene body, while both tumor and normal adjacent tissues had similar but low proportion of epigenetic mutations in the gene promoter.

499

500

501

502

503

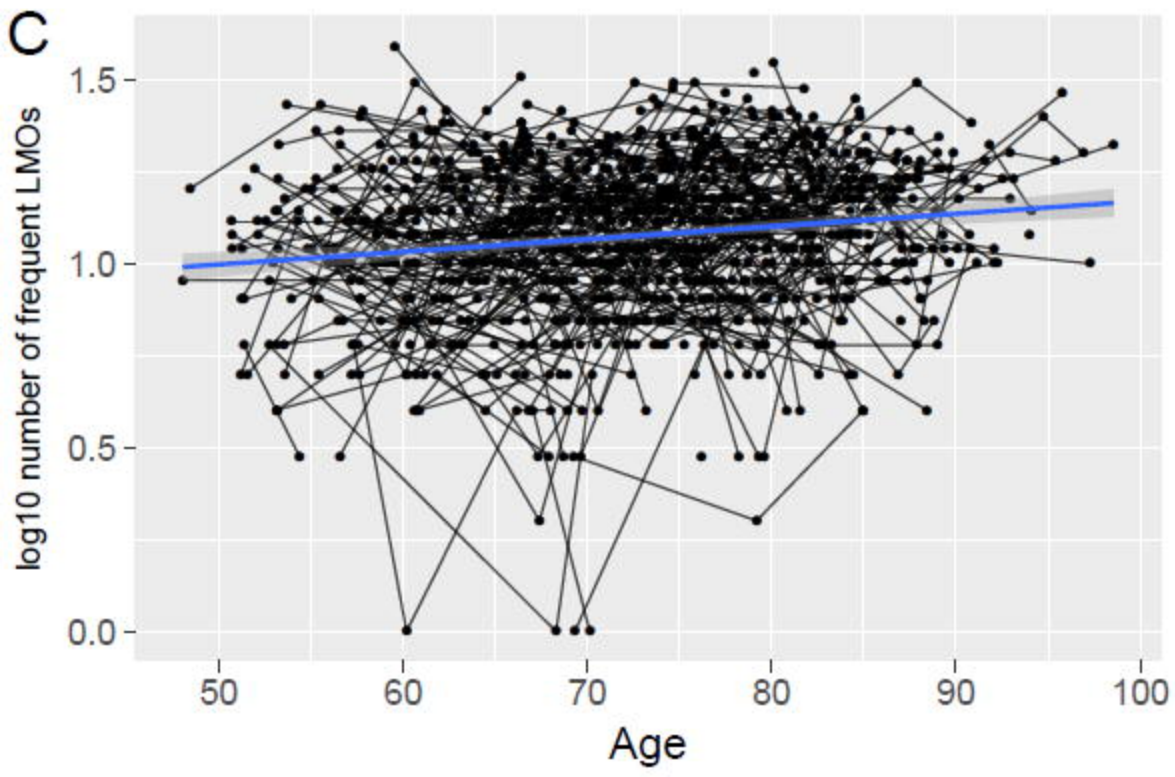
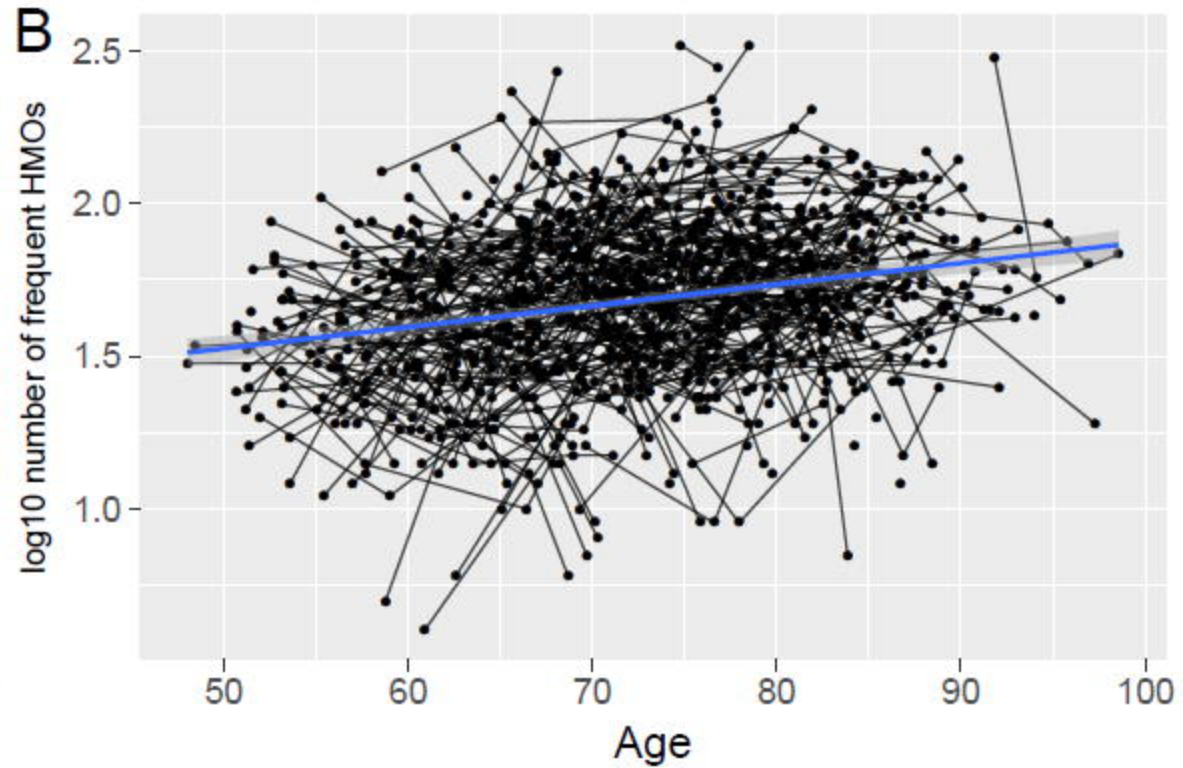
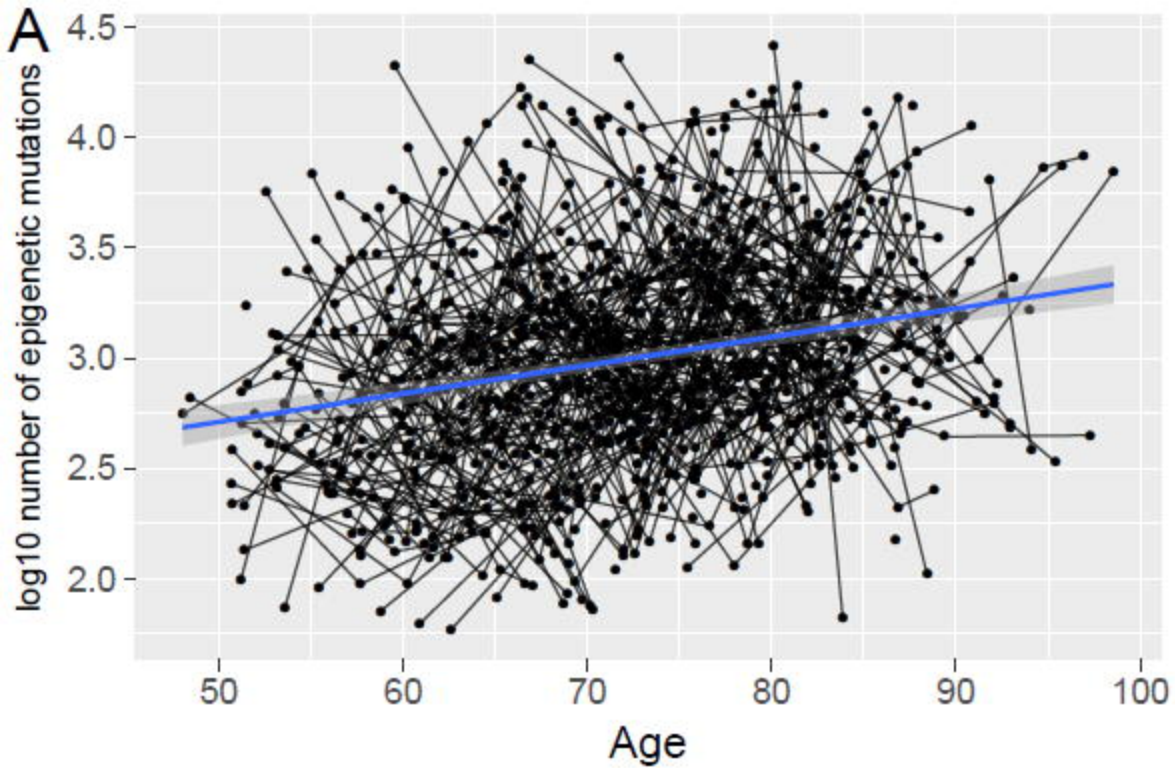
504

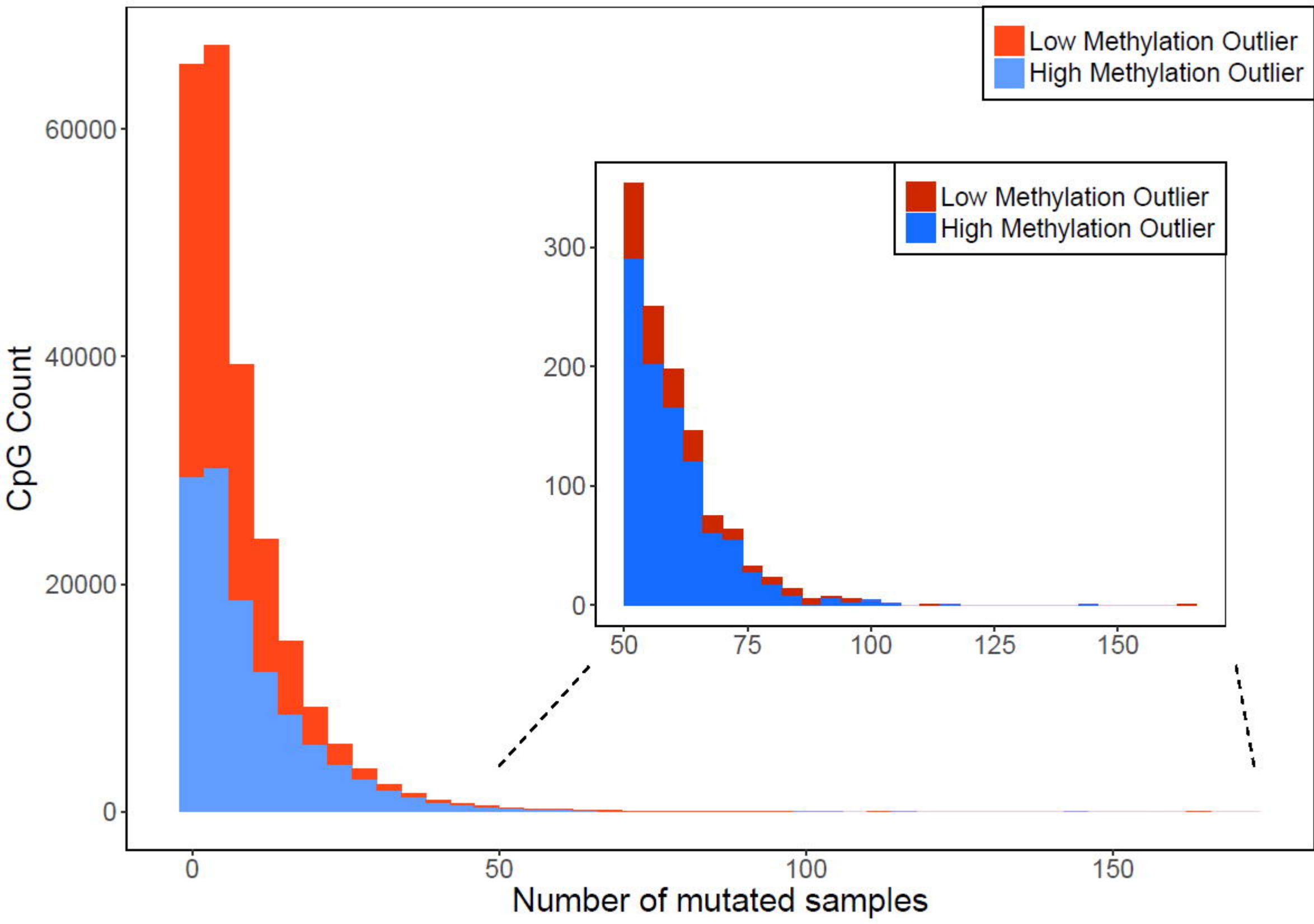
505

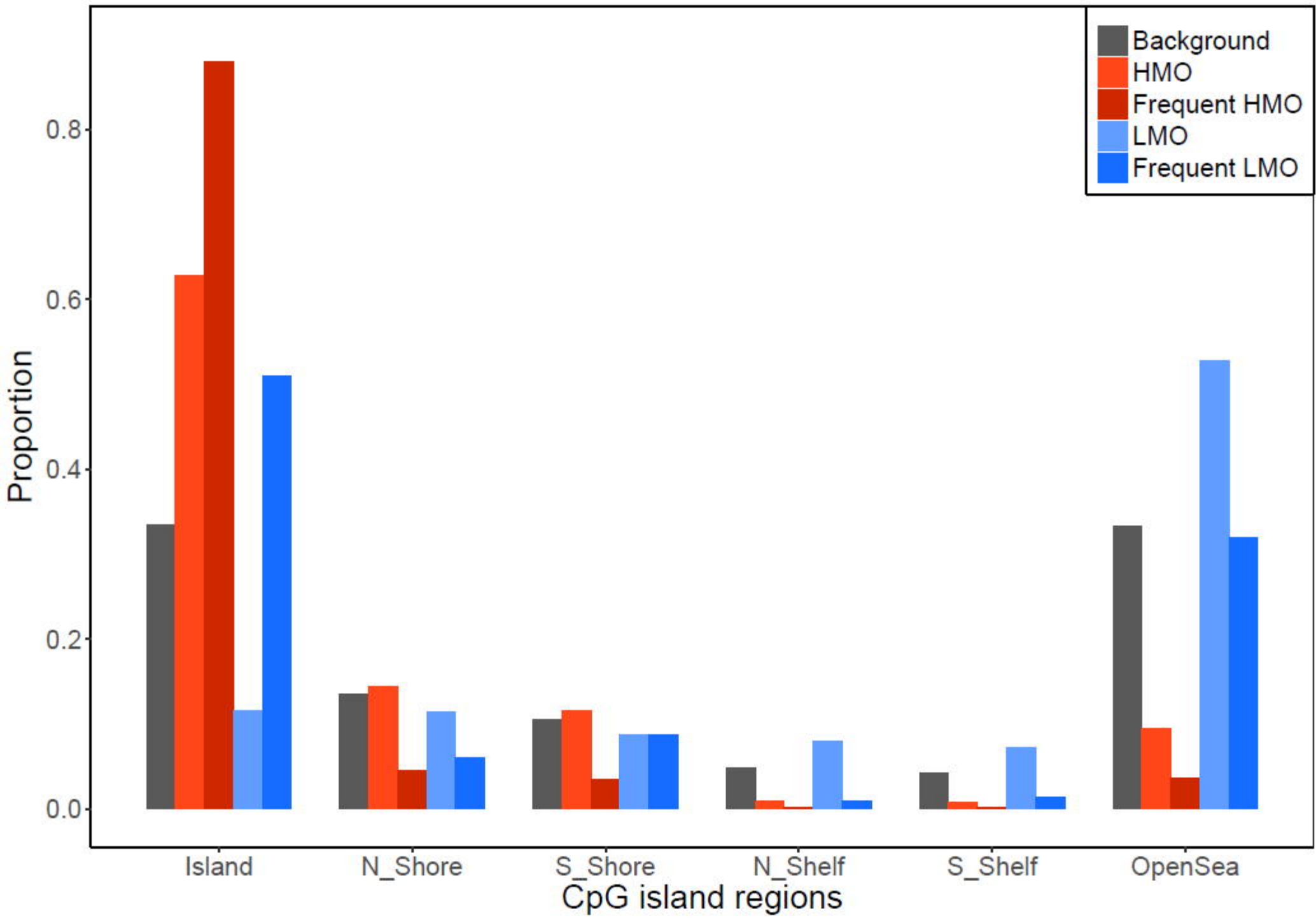
506

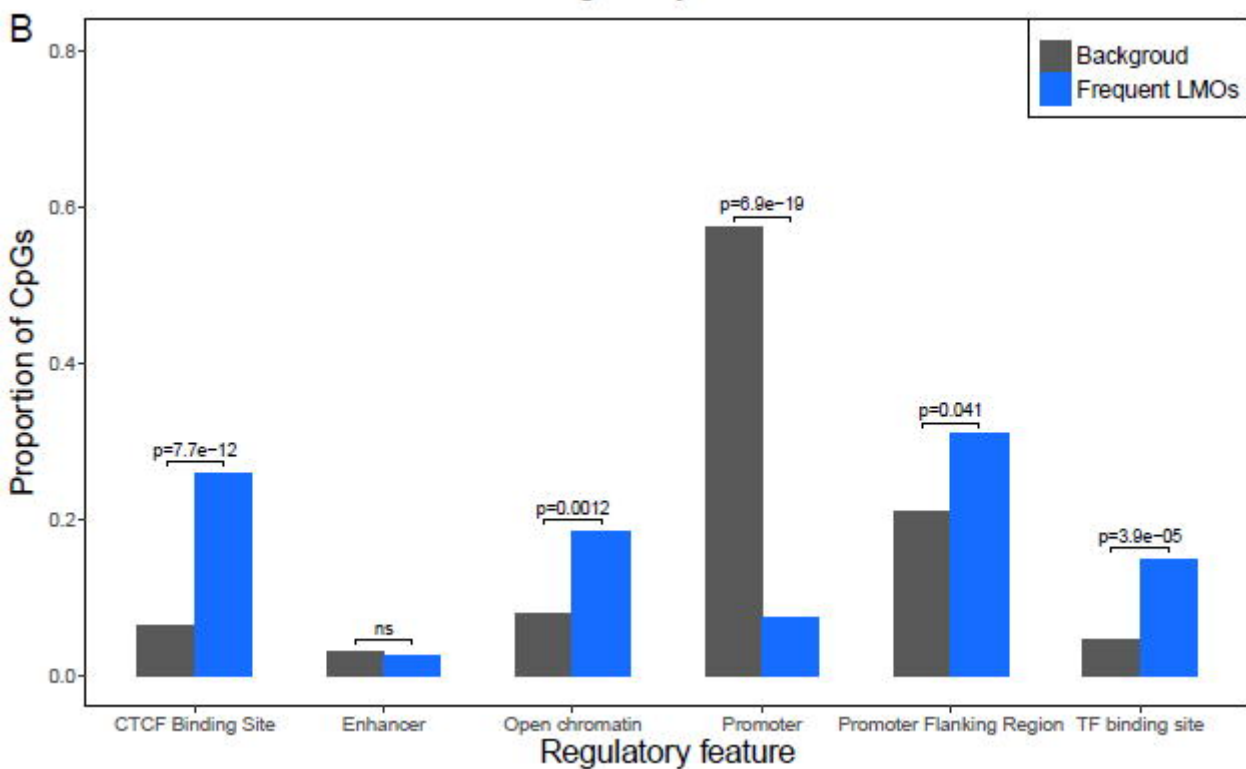
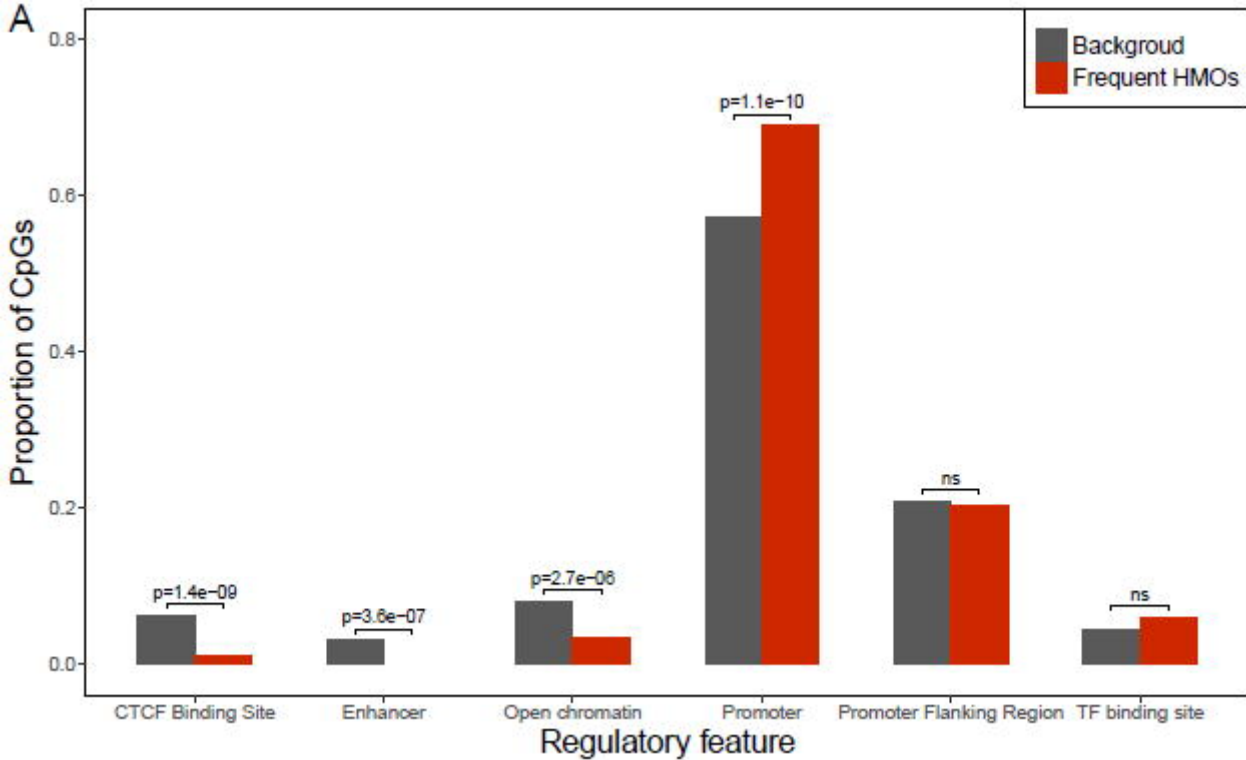
507

508









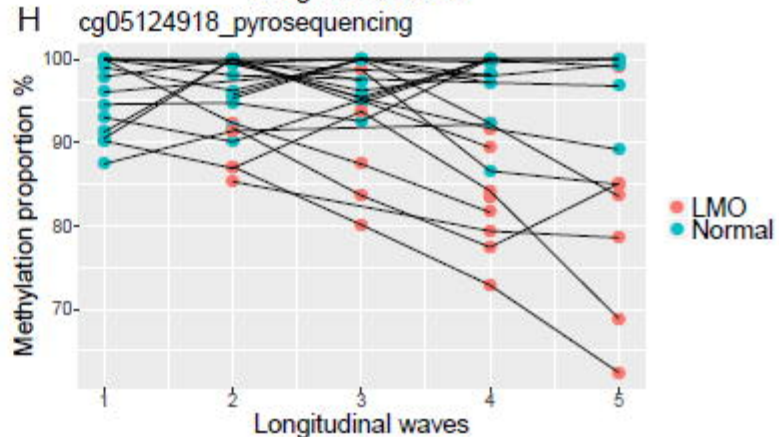
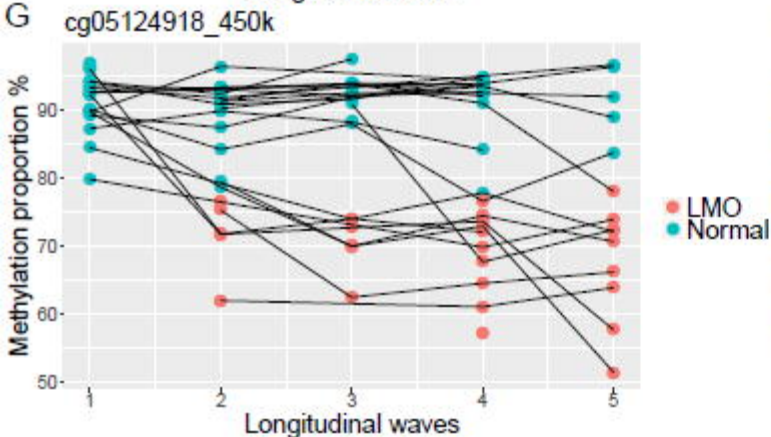
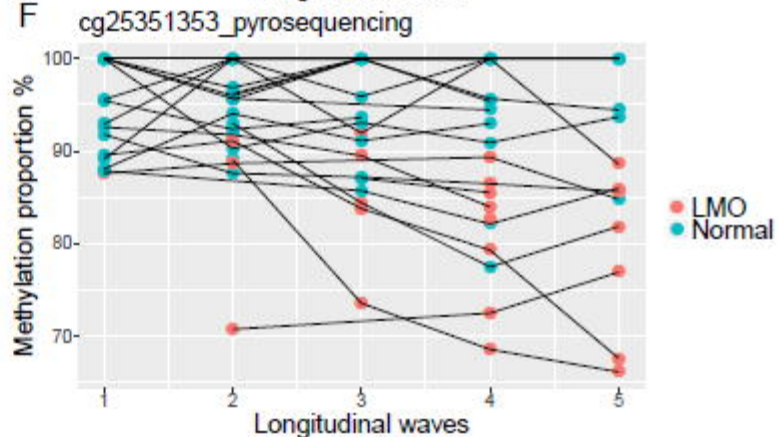
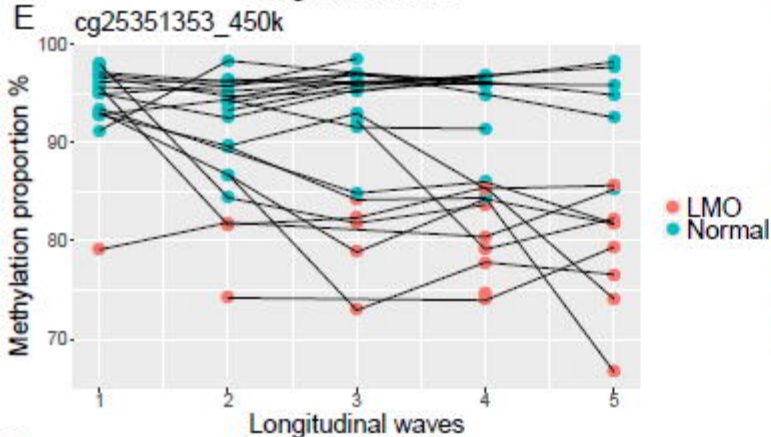
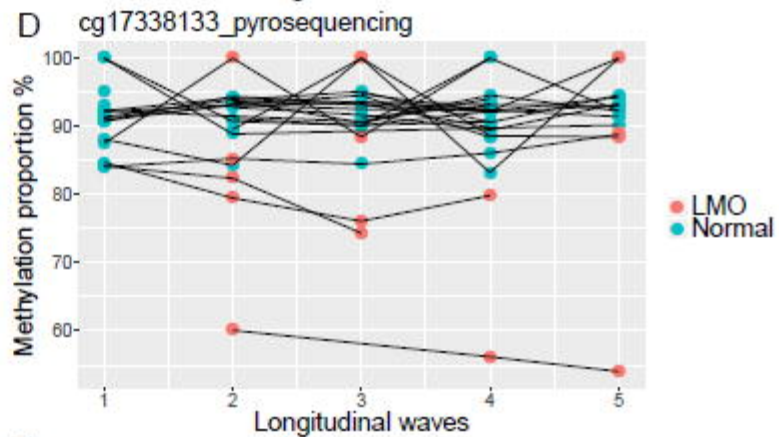
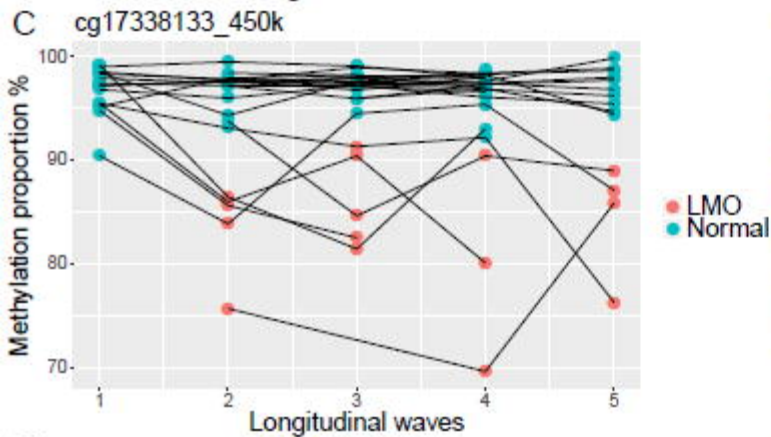
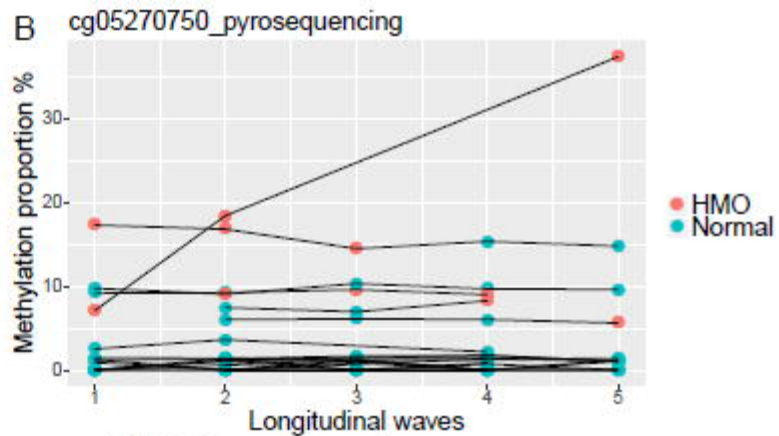
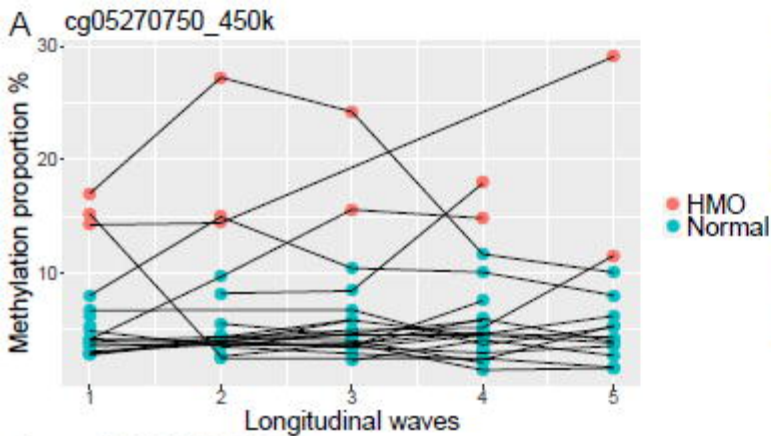


Table 1. Characteristics of study participants in SATSA (n=375 unique individuals).

Longitudinal wave	Year of sample collection	Number of Participants (new recruits)	Female Proportion	Age mean (SD)
1	1992-1994	232	58%	68.5 (9.1)
2	1999-2001	239 (101)	63%	71.1 (10.1)
3	2002-2004	186 (25)	54%	72.1 (9.1)
4	2008-2010	183 (14)	61%	76.2 (8.5)
5	2010-2012	154 (3)	66%	77.0 (8.4)

SATSA: The Swedish Adoption/Twin Study of Aging

Table 2. The association between number of epigenetic mutations (log10-transformed) and age from mixed models with confounders.

Number of epigenetic mutations	Effect sizes; (p-values)					
	Age (year)	Sex (Female to male)	CD19+ B cells (proportion)	1st genetic principal component	Sample quality*	Cancer diagnosis
Total epigenetic mutations	8.29e-03 (1.22e-13)	0.0722 (6.33e-03)	4.21 (5.06e-23)	0.445 (0.0413)	0.369 (1.48e-117)	0.0697 (0.0139)
Frequent epigenetic mutations	6.03e-03 (2.17e-19)	-0.0180 (0.33)	1.76 (1.37e-12)	0.595 (1.28e-04)	0.0573 (5.84e-13)	0.0478 (0.0164)
Frequent high methylation outliers	6.81e-03 (2.09e-17)	-0.0314 (0.16)	2.09 (2.25e-12)	0.750 (7.65e-05)	0.0512 (3.58e-08)	0.0602 (0.0130)
Frequent low methylation outliers	2.82e-03 (1.14e-05)	0.0340 (0.057)	0.474 (0.046)	0.0186 (0.92)	0.0888 (8.09e-30)	-6.99e-03 (0.71)

* Sample quality was indicated by the log10-transformed number of CpGs with a detection p-value over 0.01.

Table 3. The results of the scaled number of shared epigenetic mutations calculated from different sets of CpGs in association with age, sex, twin zygosity and zygosity-age interaction.

	Covariates	Estimate	Standard Error	P-value
All CpGs (390,894)	Age	0.019	8.59e-3	0.026
	Sex	0.208	0.107	0.055
	Zygoty (DZ)	-1.078	0.105	3.41e-18
	Zygoty (DZ)×Age	-0.012	0.011	0.284
Non-cis-meQTL CpGs (370,234)	Age	0.025	9.17e-3	5.98e-03
	Sex	0.183	0.116	0.117
	Zygoty (DZ)	-0.855	0.114	1.05e-11
	Zygoty (DZ)×Age	-0.013	0.012	0.263
Cis-meQTL CpGs (20,660)	Age	2.86e-4	7.61e-3	0.969
	Sex	0.194	0.107	0.071
	Zygoty (DZ)	-1.461	1.105	8.34e-28
	Zygoty (DZ)×Age	-3.77e-3	9.64e-3	0.696

meQTL: methylation quantitative trait loci

Table 4. Results from t-tests comparing methylation levels in samples with epigenetic mutations to normal samples using data from the 450k array and pyrosequencing.

Data	Number of samples		Mean difference (Methylation level, %)	p-value
	Normal	Mutation		
cg05270750, 450k-chip	81	12	13.39	4.34e-6
cg05270750, Pyroseq			10.79	2.01e-3
cg17338133, 450k-chip	76	17	13.11	6.39e-8
cg17338133, Pyroseq			9.35	0.02
cg25351353, 450k-chip	67	26	14.58	7.93e-17
cg25351353, Pyroseq			12.70	9.20e-8
cg05124918, 450k-chip	63	30	21.87	3.22e-20
cg05124918, Pyroseq			11.08	3.76e-07