

1 Non-canonical aberrant DNA hypermethylation in glioma

2

3 Agustin F. Fernandez^{1,†,*}, Gustavo F. Bayón^{1,†}, Marta I. Sierra¹, Rocio G. Urdinguio²,
4 Estela G. Toraño^{1,2}, Maria García^{1,2}, Antonella Carella^{1,2}, Virginia Lopez², Pablo
5 Santamarina^{1,2}, Thalia Belmonte^{1,2}, Juan Ramon Tejedor¹, Isabel Cobo^{1,3}, Pablo
6 Menendez^{3,4}, Cristina Mangas¹, Cecilia Ferrero¹, Luís Rodrigo⁵, Aurora Astudillo⁶,
7 Ignacio Ortea⁷, Sergio Cueto Díaz⁸, Pablo Rodríguez-Gonzalez⁹, J. Ignacio García
8 Alonso⁹, Manuela Mollejo¹⁰, Bárbara Meléndez¹⁰, Gemma Dominguez¹¹, Felix
9 Bonilla¹¹, and Mario F. Fraga^{2,*}

10

11 ¹Institute of Oncology of Asturias (IUOPA), ISPA-HUCA, Universidad de Oviedo, Oviedo,
12 Spain.

13 ²Nanomaterials and Nanotechnology Research Center (CINN-CSIC)-Universidad de Oviedo-
14 Principado de Asturias, Spain.

15 ³Josep Carreras Leukemia Research Institute and Department of Biomedicine, School of
16 Medicine, University of Barcelona, Barcelona, Spain

17 ⁴Institució Catalana de Recerca i Estudis Avançats (ICREA) and Centro de Investigacion
18 Biomedica en Red en Cancer CIBER-ONC, ISCIII, Barcelona, Spain.

19 ⁵Department of Gastroenterology, Hospital Universitario Central de Asturias (HUCA), Oviedo,
20 Spain.

21 ⁶Department of Pathology, Hospital Universitario Central de Asturias and Instituto Universitario
22 de Oncología del Principado de Asturias, Oviedo, Spain.

23 ⁷Proteomics Unit, IMIBIC, Maimonides Institute for Biomedical Research, Córdoba, Spain.

24 ⁸Mass Spectrometry Unit, University of Oviedo.

25 ⁹Department of Physical and Analytical Chemistry, University of Oviedo

26 ¹⁰Department of Pathology, Hospital Virgen de la Salud, Toledo, Spain. Avd. Barber 30, Toledo
27 45005

28 ¹¹Servicio de Oncología Médica, Hospital Universitario Puerta de Hierro. Majadahonda,
29 Facultad de Medicina, Universidad Autónoma de Madrid, Madrid, Spain

30

31 *Corresponding Authors:

32 Mario F. Fraga: mffraga@cinn.es

33 Agustin F Fernandez: affermandez@hca.es

34 [†]Same contribution.

35

36 **Abstract**

37 Aberrant DNA hypermethylation is a hallmark of cancer although the underlying
38 molecular mechanisms are still poorly understood. To study the possible role of 5-
39 hydroxymethylcytosine (5hmC) in this process we analyzed the global and locus-
40 specific genome-wide levels of 5hmC in primary samples from 54 gliomas and 72
41 colorectal cancer patients. Levels of 5hmC in colorectal cancer were very low and no
42 consistent changes were detected between control tissues and tumors. As expected,
43 levels of 5hmC in non-tumoral brain samples were high and significantly reduced at the
44 49,601 CpG sites in gliomas. Strikingly, hypo-hydroxymethylation at 4,627 (9.3%) of
45 these CpG sites was associated with aberrant DNA hypermethylation. The DNA regions
46 containing these CpG sites were enriched in H3K4me2, and presented a different
47 genuine chromatin signature to that characteristic of the genes classically aberrantly
48 hypermethylated in cancer. We conclude that this data identifies a novel 5hmC-
49 dependent non-canonical class of aberrant DNA hypermethylation in glioma.

50

51

52 **Introduction**

53 DNA methylation at the fifth position of cytosine (5mC) has been one of the most
54 studied epigenetic modifications in mammals to date. 5mC is involved in the regulation
55 of multiple physiological and pathological processes, including cancer, and when
56 located at gene promoters, it is usually linked to transcriptional repression.

57 As distinctive features of tumorigenesis, local DNA hypermethylation and global
58 hypomethylation have been attributed to changes in 5mC levels [10; 11]. However, the
59 discovery a few years ago, of 5-hydroxymethylcytosine (5hmC), a new epigenetic mark
60 resulting from 5mC oxidation, is reshaping our view of the cancer epigenome [29; 47].
61 This 5mC to 5hmC conversion in mammals is mediated by ten-eleven translocation
62 proteins (TET1, TET2, and TET3), a family of α -ketoglutarate (α KG) and Fe(II)-
63 dependent dioxygenases[21; 47]. Global levels of 5hmC in the genome fluctuate
64 considerably according to tissue type, and are consistently around 10-fold lower than
65 those of 5mC, though it is interesting that the highest levels of both marks are found in
66 brain [13; 16; 24; 32; 37; 44].

67 Several studies suggest that 5hmC is not only an intermediate of DNA demethylation,
68 but that it also plays a role in cancer biology [23; 33; 41; 46; 55]. In this vein, a broad
69 loss of 5hmC has been reported in different human cancers including melanoma,
70 glioma, breast, colon, gastric, kidney, liver, lung, pancreatic, and prostate cancers [17;
71 23; 27; 30; 32; 34; 55].

72 The fact that there are now methods available that distinguish 5mC and 5hmC positions
73 at single-base resolution within the genome prompted us to reassess the role of DNA
74 methylation status in tumorigenesis from a 5hmC perspective. The method used here
75 allowed us to describe global and genome-wide locus-specific 5mC and 5hmC patterns
76 in colon and brain samples, to identify a specific chromatin signature associated with
77 changes of these epigenetic marks in cancer and, most notably, to describe a novel non-
78 canonical type of aberrant DNA hypermethylation in cancer.

79 **Results**

80 *Global changes of 5mC and 5hmC in cancer*

81 To evaluate the role of 5hmC in the changes of DNA methylation observed in cancer,
82 we first analyzed the levels of 5hmC and 5mC at repetitive DNA in 84 normal and 123
83 tumor samples obtained from patients with colorectal cancer and glioma. Bisulfite
84 pyrosequencing was used to determine the level of both epigenetic modifications in 4
85 different types of repeated DNA: the retrotransposons LINE-1 and AluYb8, and the
86 pericentromeric tandem repeats Sat-alpha and NBL-2 [49]. These 4 DNA regions
87 contain most of the genomic methylation and, consequently, global DNA methylation
88 level is highly dependent on their 5mC content [51]. As expected, 5mC levels at most
89 repeated DNA in healthy tissue was high, and was reduced in tumor samples (**Figure**
90 **1A**). In contrast, the levels of 5hmC at repeated DNA in healthy tissue was very low,
91 and tumoral tissue showed even lower levels of 5hmC in these DNA regions (**Figure**
92 **1B**). However, the differences were very small and, consequently, they cannot explain
93 the global loss of this mark in cancer observed by mass spectrometry [23; 27; 28; 39].

94

95 *5mC and 5hmC profiling in colorectal and brain tissue*

96 As changes in 5hmC at repeated DNA were not able to explain the global changes
97 previously observed by mass spectrometry, we hypothesized that these changes
98 primarily occur at single copy sequences. To investigate this possibility in more detail,
99 we first used 450K Infinium methylation arrays to determine the level and genomic
100 distribution of 5mC and 5hmC at 479,423 CpG sites in 11 non-tumoral colorectal
101 samples and 5 healthy brain tissue samples, all from different donors. A preliminary
102 examination of the data revealed that the beta values of the oxidized samples (true 5mC)
103 were much lower than their non-oxidized counterparts (5mC+5hmC) in brain
104 (Wilcoxon rank sum test; $p < 0.001$; $W = 2.34e13$) than in colorectal (Wilcoxon rank sum
105 test; $p < 0.001$; $W = 5.52e13$) tissues (see Materials and Methods), which indicates that, as
106 expected, levels of 5hmC are higher in brain tissue than in the colon (**Figure 2A**). In
107 line with this, we identified 111,633 and 5,089 hydroxymethylated CpG sites (5hmC
108 sites) in brain and colorectal tissue respectively (**Figure 2A, and Supplementary**
109 **Table 1 and 2**) (see Materials and Methods).

110 The analysis of the genomic distribution of the 5hmC sites showed that, in both
111 colorectal and brain tissue, hydroxymethylation is enriched at the low CpG-density
112 regions interrogated by the array (Wilcoxon non-parametric test; $p < 0.001$, $D = -0.29$, and

113 $p < 0.001$, $D = -0.5$, respectively) (**Figure 2B**). Consequently, the 5hmC sites were
114 enriched in non-CpG islands (non-CGI) in both colon and brain (chi-square test;
115 $p < 0.001$; OR=1.93, and $p < 0.001$, OR=3.45, respectively) and infrequent in CGIs (chi-
116 square test; $p < 0.001$, OR=0.14, and $p < 0.001$, OR=0.13) (**Figure 2C**). With respect to
117 genes, 5hmC sites were enriched in introns in both brain and colorectal tissue (chi-
118 square test; $p < 0.001$, OR=1.82, and $p < 0.001$, OR=1.76, respectively), but were less
119 frequent than expected in intergenic regions in colorectal tissue and in gene promoters
120 in brain tissue (chi-square test; $p < 0.001$, OR=0.58, and $p < 0.001$, OR=0.6) (**Figure 2D**).
121 Moreover, hydroxymethylated CpG sites were farther away from centromeres in brain
122 (Wilcoxon non-parametric test, $p < 0.001$, $D = 0.01$) and telomeres in both colorectal
123 tissues and brain (Wilcoxon non-parametric test, $p < 0.001$, $D = 0.07$, and $p < 0.001$,
124 $D = 0.02$, respectively) than the median in terms of other background sites, although the
125 size of these shifts was rather small (**Figure 2-figure supplement 1A**).

126 To identify possible chromatin marks associated with 5hmC sites in colorectal and brain
127 tissue, we compared these CpG sites with previously published data on a range of
128 histone modifications and chromatin modifiers in 10 different cell types (see Materials
129 and Methods) (**Figure 2E**). This approach identified statistically significant associations
130 (Fisher's exact test; $p < 0.05$) between the 5hmC sites and the active histone marks
131 H3K4me1, H3K36me3, and H4K20me1, in both colorectal and brain tissue (**Figure**
132 **2E**). Interestingly, in colorectal tissue, 5hmC was also enriched in other activating
133 histone posttranslational modifications (PTMs) such as H3K79me2, and H3K4me2
134 (**Figure 2E**). Finally, a similar framework was used to test for the enrichment of our
135 selected probes over the computer-generated chromatin segmentation states from the
136 ENCODE ChromHMM project (see Materials and Methods). In total, fifteen states were
137 used to segment the genome, and these were then grouped and colored to highlight
138 predicted functional elements. This approach showed that the hmC sites were
139 significantly enriched in states associated with enhancers and transcription in both
140 colorectal and brain tissue (Fisher's exact test; $p < 0.05$) (**Figure 2-figure supplement**
141 **1B**).

142

143 *Locus-specific alterations of 5hmC in colorectal cancer (CRC) and glioma*

144 To identify differentially hydroxymethylated CpG sites (d5hmC) at single copy
145 sequences in cancer, we used 450K methylation arrays to analyze 11 additional
146 colorectal tumors and 9 primary tumors obtained from patients with glioma (see

147 Materials and Methods). A total of 49,601 CpG sites that were hypohydroxymethylated
148 were identified in gliomas, but almost no hyper-hydroxymethylated sites were found
149 (see Materials and Methods) (**Figure 3A and Supplementary Table 3**). In contrast, no
150 significant methylation changes were found in colorectal tumors (**Figure 3A**) and thus
151 subsequent stages of the study focused on glioma alone. Hierarchical clustering using
152 the differentially hydroxymethylated CpG sites showed the correct classification of
153 normal and tumor samples (**Figure 3B**). The analysis of the genomic distribution of the
154 hypo-hydroxymethylated CpG sites in gliomas showed an enrichment at low CpG
155 density regions (Wilcoxon rank sum test, $p < 0.001$, $D = -0.41$), and consequently at non-
156 CpG islands (chi-squared test, $p < 0.001$, $OR = 2.53$) (**Figure 3C**). With respect to gene
157 location, hypo-hydroxymethylation was more frequent in introns (chi-squared test,
158 $p < 0.001$, $OR = 1.77$) (**Figure 3C**).

159 To identify possible chromatin signatures associated with DNA hypo-
160 hydroxymethylation in gliomas, we compared our list of hypo-hydroxymethylated CpG
161 sites with previously published data on a range of histone modifications and chromatin
162 modifiers in 11 different cell types (see Materials and Methods) (**Figure 3D**).
163 Interestingly, this approach showed an enrichment of hypo-hydroxymethylation at
164 chromatin regions marked with the activating histone PTMs H3K4me1, H3K36me3,
165 H4K20me1 and H3K79me2 (Fisher's exact test, $p < 0.05$) (**Figure 3D**), but not with the
166 repressive histone modification H3K27me3, which has been previously shown to be
167 associated with aberrant DNA hypermethylation in cancer [38; 52] (**Figure 3D**). A
168 similar framework was used to test for the enrichment of our selected probes over the
169 computer-generated chromatin segmentation states from the ENCODE ChromHMM
170 project. Using this approach, we found that hypohydroxymethylated CpG sites were
171 significantly associated with transcription regulation and enhancers (Fisher's exact test;
172 $p < 0.05$) (**Figure 3E**).

173

174 *DNA hypo-hydroxymethylation identifies a novel type of non-canonical aberrant* 175 *DNA hyper-methylation in glioma*

176 To study the relationship between changes in 5mC and 5hmC in glioma, we first
177 identified aberrantly methylated CpG (d5mC) sites. The comparison of the methylation
178 data between tumoral and control samples (see Materials and Methods) identified 2,727
179 hypo- and 12,050 hyper-methylated CpG sites in gliomas (**Supplementary Tables 4**
180 **and 5**). Next, we compared these d5mC sites with the previously identified hypo-

181 hydroxymethylated CpG sites (**Figure 3A, Supplementary Table 3**). Interestingly, this
182 approach showed that 4,627 (38.4%) of the CpG sites aberrantly hypermethylated in
183 gliomas also lose 5hmC (**Figure 4A, Supplementary Table 6**).

184 To investigate, at a functional genomic level, the characteristics of these two classes of
185 aberrantly hypermethylated CpG sites in gliomas we first analyzed their genomic
186 distribution in relation to density of CpG sites and we found that the hypermethylated
187 CpG sites that lose 5hmC (hyper5mC-hypo5hmC) were enriched in low density CpG
188 regions (Wilcoxon rank sum test, $p < 0.001$, $D = -0.11$) as compared with the
189 hypermethylated CpG sites that showed no changes in 5hmC (hyper5mC) (Wilcoxon
190 rank sum test, $p < 0.001$, $D = -0.23$) (**Figure 4B, Supplementary Tables 6 and 7**). In line
191 with this, hyper5mC-hypo5hmC sites were strongly depleted from CGIs (chi-squared
192 test, $p < 0.001$, $OR = 0.42$) (**Figure 4B**). Hierarchical clustering using the differentially
193 methylated CpG sites showed that the hyper5mC-hypo5hmC sites were slightly more
194 methylated in control brain samples than the hyper5mC sites, and that they were more
195 uniformly hypermethylated in glioma (**Figure 4C**). To further corroborate our results,
196 we took advantage of recently published data on the whole-genome bisulfite sequencing
197 (WGBS) in glioma [42]. We found that, in addition to a large percentage of CpGs (n:
198 4,051; 88%) showing the same patterns of change as in our methylation arrays, the
199 WGBS analysis identified more than 10^6 new hyper5mC-hypo5hmC sites, thus
200 confirming that this is a frequent event in glioma (**Figure 4-figure supplement 1**).

201 Next, to identify possible chromatin signatures associated with the two classes of
202 aberrantly hypermethylated CpG sites in gliomas, we compared our data with
203 previously published data on a range of histone modifications and chromatin modifiers
204 in 11 different cell types (see Materials and Methods) (**Figure 5A**). This approach
205 confirmed the association between hyper5mC and the repressive histone marks
206 H3K9me3 and H3K27me3 (Fisher's exact test, $p < 0.05$) [36; 38; 52]. The hyper5mC-
207 hypo5hmC sites showed a completely different chromatin signature, with enrichment in
208 the activating histone PTMs H3K4me1, H3K36me3, H3K79me2 and H4K20me1
209 (Fisher's exact test, $p < 0.05$) (**Figure 5A**). Notably, as compared with the chromatin
210 signature of the whole set of hypo-hydroxymethylated CpGs in glioma, these CpG sites
211 were particularly enriched at the H3K4me2 histone mark (Fisher's exact test, $p < 0.001$,
212 OR in [1.19, 1.78] for all cell lines in the Broad Histone project) (**Figure 5B**).

213 These results indicate that the hyper5mC sites behave like the aberrantly
214 hypermethylated canonical CpG sites in cancer, whilst the hyper5mC-hypo5hmC sites

215 represent a novel and functionally different non-canonical type of aberrantly methylated
216 DNA sequence in glioma (**Figure 5A, 5B, Supplementary Tables 6 and 7**). In support
217 of this notion, experiments focused on the computational prediction of functional
218 elements confirmed the enrichment of canonical aberrant hypermethylation in
219 promoters and repressed sequences and revealed a completely different pattern for non-
220 canonical hypermethylation, one which is more closely associated with enhancers and
221 transcriptional regulation (Fisher's exact test; $p < 0.05$) (**Figure 5-figure supplement**
222 **1**).

223

224 *Distinct functional role of canonical and non-canonical aberrant hypermethylation in* 225 *glioma*

226 To identify possible differences between the functional role of canonical and non-
227 canonical aberrant DNA hypermethylation in glioma we first ascribed CpG sites to
228 specific genes and then used HOMER to carry out gene ontology analyses of each
229 group of genes (see methods). Using this approach, we identified 1,921 genes
230 displaying canonical hypermethylation, 2,042 displaying non-canonical
231 hypermethylation and 938 displaying both types of aberrant hypermethylation (**Figure**
232 **6A, Supplementary Tables 8, 9 and 10**). As expected, GO analyses showed an
233 enrichment of development and differentiation processes in canonical genes [6] (**Figure**
234 **6A, Supplementary Table 11**). In contrast, non-canonical genes were strongly enriched
235 in cell signaling and protein processing pathways (**Figure 6A, Supplementary Table**
236 **12**).

237 To further investigate the functional role of canonical and non-canonical
238 hypermethylation in cancer, we compared our methylation data with previously
239 published gene expression data in the same type of tumor (see Materials and Methods).
240 Results showed that 681 (31%) of the canonical and 585 (24%) of the non-canonical
241 aberrantly hypermethylated genes were repressed in gliomas (**Figure 6B**).

242 Genomic distribution analysis of both types of aberrant hypermethylation confirmed the
243 enrichment of canonical hypermethylation in exons (chi-squared test, $p < 0.001$,
244 $OR = 1.79$ for general exons, $OR = 2.01$ for first exons), while non-canonical
245 hypermethylation was more frequent in introns (chi-squared test, $p < 0.001$, $OR = 1.7$)
246 (**Figure 6C**). The genes frequently downregulated in glioma, *SLC14A* and the *SMAD7*,
247 represent two bona fide examples of this pattern of non-canonical aberrant
248 hypermethylation (**Figure 6D, Figure 6-figure supplement 1**).

249 Taken as a whole, these results indicate that both types of aberrant hypermethylation
250 have a similar effect on gene expression, but that they affect different types of genes and
251 gene regions.

252 Discussion

253 During recent decades, it has largely been accepted that aberrant genomic DNA
254 methylation is a hallmark of cancer [10; 11] and the best-known DNA methylation
255 alterations in tumors were the aberrant hypermethylation of CpG island promoters, and
256 global DNA hypomethylation. In both cases, the alterations were mostly attributed to
257 changes in the overall content and genomic distribution of 5mC [10; 11].

258 The vast majority of studies on DNA methylation and cancer have been based on the
259 sodium bisulfite modification of the genomic DNA, a chemical reaction that allows C
260 and 5mC to be distinguished by polymerase chain reaction [19]. However, this approach
261 cannot distinguish between 5mC and 5-hydroxymethylcytosine (5hmC), the latter being
262 a chemical modification of the cytosine first identified in bacteriophages in 1952 [54],
263 and which has recently been found to be quite abundant in specific mammalian tissue
264 [29]. 5hmC is synthesized from 5mC by the Ten-eleven Translocation (Tet) Enzymes, a
265 family of proteins that can also catalyze the successive conversion of 5hmC to 5-
266 formylcytosine and then to 5-carboxylcytosine, both of which can be transformed to
267 unmodified C [40]. Although 5hmC was originally described as simply a demethylation
268 intermediate of C [16; 25], recent data suggest that this may be an epigenetic mark in its
269 own right [3; 22]. Thus, as most previous studies did not distinguish between 5mC and
270 5hmC, and it appears that DNA hydroxymethylation might play a specific role in
271 cancer, in this work we aimed to re-evaluate changes in DNA methylation in cancer,
272 paying special attention to the specific contribution of 5hmC.

273 To identify the DNA regions affected by hydroxymethylation changes in cancer, we
274 first focused on four types of repeated DNA (LINE1, Sat α , NBL2 and ALUYB8).
275 Among them, the LINE1 repeat is of particular interest because it contains almost 20%
276 of the genomic 5mC, and it has been proposed to be a surrogate of global DNA
277 methylation [51]. Our results confirmed that tumors lose 5mC at repeated DNA [7].
278 However, the level of 5hmC at repeated DNA in healthy samples was very low and no
279 significant differences were observed compared to tumors, which indicates that the
280 global DNA hypo-hydroxymethylation previously observed in cancer [23; 27; 28; 39]
281 does not principally occur at repeated DNA. As changes in 5hmC at repeated DNA
282 could not explain the global differences previously observed by mass spectrometry, we
283 decided to study the possible contribution of single copy sequences. Genome-wide
284 profiling of 5mC and 5hmC of healthy tissue identified a 10-fold increase in abundance
285 of CpG sites frequently hydroxymethylated in brain compared to colorectal tissue,

286 providing evidence that the level of this epigenetic mark is highly tissue type-
287 dependent, and also that it is very abundant in the brain [16; 24; 29; 32; 37; 44].
288 Moreover, 5hmC was enriched in regions with low CpG density and in introns in both
289 colorectal and brain tissue. As the 5hmC is enriched in different genomic regions, these
290 results support the notion that 5hmC is not simply a demethylation intermediate [1].
291 Interestingly, 5hmC co-localized in regions marked with the activating histone PTM
292 H3K4me1. This histone mark has been previously associated with gene enhancers [20;
293 48], which suggests that DNA hydroxymethylation might play a role in gene regulation
294 in trans. Moreover, we have recently found an association between H3K4me1 and DNA
295 hypomethylation during aging in stem and differentiated cells [12], which may
296 represent an interesting link between aging and cancer at these genomic regions.
297 Colorectal tumors showed more changes with respect to 5mC than to 5hmC. However,
298 in contrast, glioma presented more changes in 5hmC than in 5mC, suggesting that the
299 dynamics of DNA methylation and hydroxymethylation in cancer is highly tumor-type
300 dependent. Moreover, the great number of hypo-hydroxymethylated single CpG sites in
301 glioma could explain the global differences previously observed by mass spectrometry
302 [23; 27; 28; 39] and suggests that, in contrast to 5mC, most DNA hypo-
303 hydroxymethylation in brain tumors occurs at single copy sequences.

304 The behavior of 5hmC led us to next identify two types of CpG sites aberrantly
305 hypermethylated in glioma: aberrantly hypermethylated CpG sites that showed no
306 changes in 5hmC; and hypermethylated CpG sites that lose 5hmC. The first of these
307 sites display similar chromatin signatures to previously described genes aberrantly
308 hypermethylated in cancer (i.e. enrichment in the repressive histone marks H3K9me3
309 and H3K27me3) [36; 38; 52]. In contrast, the second type of aberrantly
310 hypermethylated CpG sites were enriched in the activating histone PTMs H3K4me1,
311 H3K36me3, H3K79me2, H4K20me1 and H3K4me2. As these CpG sites present a
312 genuine chromatin signature which is different to the repressive chromatin signature of
313 the classical genes aberrantly hypermethylated in cancer [36; 38; 52], we conclude that
314 they represent a novel 5hmC-dependent non-canonical class of aberrant DNA
315 hypermethylation in glioma. As this gain in 5mC is inversely correlated with loss of
316 5hmC, it was not possible to identify this significant alteration in previous studies using
317 the classical sodium bisulfite-based technologies, since they are not able to distinguish
318 between the two chemical modifications.

319 Aberrant DNA hypermethylation in cancer was discovered more than 30 years ago, but
320 the underlying molecular mechanisms are still poorly understood. For example, it has
321 been proposed that genes enriched in bivalent histone modifications (H3K4me3 and
322 H3K27me3) and polycomb group proteins during embryo development are prone to
323 become aberrantly hypermethylated in cancer [36; 38; 52] but the molecular basis of
324 this is unknown. Our data suggest that tumor cells might in fact acquire aberrant DNA
325 methylation through various different pathways. Moreover, in the case of the non-
326 canonical hypermethylation, the previous loss of 5hmC suggests that aberrant
327 hypermethylation at these DNA regions could be due to an attempt by the cell to reverse
328 or repair the loss of 5hmC at functionally sensible loci. This possibility is supported by
329 the fact that the non-canonical aberrant hypermethylation described here seems to play
330 an important role in gene regulation. Intriguingly, 5hmC at gene promoters has also
331 been proposed to protect from aberrant hypermethylation in colorectal cancer [50].
332 Thus, although it seems that 5hmC plays an important role in the regulation of the DNA
333 methylation changes in cancer, more research is needed to fully understand its role.
334 The non-canonical aberrant hypermethylation described here seems to have a similar
335 overall effect on gene expression as classical canonical hypermethylation, although the
336 type of genes and the genomic regions affected are very different. Previous research has
337 shown that the repression of developmental genes affected by canonical aberrant
338 hypermethylation promotes tumorigenesis [6]. However, the possible functional role of
339 disruption of cell signaling and protein processing pathways affected by the non-
340 canonical hypermethylation described in this study remains to be elucidated. Future
341 research is thus needed to address this issue, and to determine whether the two types of
342 aberrant DNA hypermethylation have distinct functional roles in cancer.
343

344 **Materials and methods**

345 *Normal samples and primary tumors*

346 The colon and brain samples analyzed in this study were collected at the Hospital
347 Universitario Central de Asturias (HUCA), the Hospital Virgen de la Salud, Toledo, and
348 the Hospital Universitario Puerta de Hierro, Madrid. The samples studied comprised 72
349 normal colons, 13 normal brains, 72 colorectal primary tumors and 54 glioblastomas.
350 The study was approved by the Clinical Research Ethics Committee and all the
351 individuals involved provided written informed consent.

352

353 *Pyrosequencing assays*

354 5mC and 5hmC patterns at repetitive sequences (LINE1, ALUBY8, Sat α and NBL2)
355 were analyzed by pyrosequencing using previously described primers [49]. To calculate
356 5hmC levels, each sample was analyzed using two methods performed in parallel; an
357 oxidative bisulfite conversion (oxBS) and a bisulfite-only conversion (BS), in
358 accordance with the TrueMethyl® Array Kit User Guide (CEGX, Version 2) with some
359 modifications. Briefly, DNA samples were cleaned using Agencourt AMPure XP
360 (Beckman Coulter) then oxidated with 1 μ L of a KRuO₄ (Alpha Aesar) solution (375
361 mM in 0.3 M NaOH), after which bisulfite conversion was performed using EpiTect
362 bisulfite kit (Qiagen®).

363 After PCR amplification of the region of interest in oxBS and BS samples,
364 pyrosequencing was performed using PyroMark Q24 reagents, and vacuum prep
365 workstation, equipment and software (Qiagen®). 5hmC levels were obtained when
366 methylation values of oxBS samples (represents true 5mC) were subtracted from their
367 corresponding BS treated pairs (the latter representing 5mC+5hmC).

368

369 *Genome-wide DNA methylation analysis with high-density arrays*

370 Microarray-based DNA methylation profiling was performed with the
371 HumanMethylation 450 BeadChip [2]. Oxidative bisulfite (oxBS) and bisulfite-only
372 (BS) conversion was performed using the TrueMethyl® protocol for 450K analysis
373 (Version 1.1, CEGX) following the manufacturer's recommended procedures.
374 Processed DNA samples were then hybridized to the BeadChip (Illumina), following
375 the Illumina Infinium HD Methylation Protocol. Genotyping services were provided by
376 the Spanish Centro Nacional de Genotipado (CEGEN-ISCI) (www.cegen.org). DNA

377 methylation data were downloaded from ArrayExpress accession numbers E-MTAB-
378 6003 (brain) and E-MTAB-xxx (colon).

379

380 *HumanMethylation450 BeadChip data preprocessing*

381 Raw IDAT files were processed using the R/Bioconductor package minfi [15] (version
382 1.14.0), implementing the SWAN algorithm [35] to correct for differences in the
383 microarray probe designs. No background correction or control probe normalization
384 was applied. Probes where at least two samples had detection p-values > 0.01 , and
385 samples where at least 5500 probes had detection p-values > 0.01 were filtered out. M-
386 values and beta values were computed as the final step in the preprocessing procedure.
387 In line with a previously published methodology [5], M-values were used for the
388 statistical analyses and beta values for effect size thresholding, visualization and report
389 generation.

390

391 *Batch effect correction*

392 In order to detect whether there was any batch effect associated with technical factors,
393 the visualization technique of multidimensional scaling (MDS) was employed to
394 highlight any strange interaction affecting the different samples. Where necessary,
395 posterior adjustment of the samples was performed by means of the SVA method [31]
396 implemented in the R/Bioconductor sva package (version 3.14.0).

397

398 *Computation of hydroxymethylation levels*

399 Beta values from oxBS samples were subtracted from their corresponding BS treated
400 pairs, generating an artificial dataset representing the level of 5hmC for each probe and
401 sample as per a previously published methodology [45]. One further dataset was created
402 to represent the 5mC levels using beta values from oxBS samples.

403

404 *Detection of differentially methylated probes*

405 Differential methylation and hydroxymethylation of an individual probe was determined
406 by a moderated t-test implemented in the R/Bioconductor package limma [43]. A linear
407 model, with methylation or hydroxymethylation levels as response and the sample
408 group (normal/tumoral) as the principal covariate of interest, was then fitted to the
409 methylation or hydroxymethylation data. Surrogate Variables generated using SVA
410 were also included in the model definition, but excluding those found to be correlated to

411 the phenotype of interest. P values were corrected for multiple testing using the
412 Benjamini-Hochberg method for controlling false discovery rate (FDR). An FDR
413 threshold of 0.001 was employed to determine differentially methylated and
414 hydroxymethylated probes. Additionally, these probes were filtered according to their
415 effect size, keeping only those probes with methylation or hydroxymethylation changes
416 between-groups which exceeded the median of all differences for the same comparison.
417 The probes without no significant 5hmC signal on control samples were filtered out
418 from the set of hypo-hydroxymethylated probes in glioma.

419

420 *Identification of hydroxymethylated probes*

421 In order to identify those probes representing the regions where the 5hmC mark is
422 located, a differential hydroxymethylation analysis was performed as described
423 previously [14] using a dataset containing both oxBS and BS versions of the control
424 samples. Probes with significant differences in beta values between the BS and oxBS
425 samples were considered to be enriched for the 5hmC mark. An FDR threshold of 0.001
426 was employed. No filtering on effect size was applied in this case.

427

428 *Histone enrichment analysis*

429 In order to analyze the enrichment of histone marks on a subset of probes, we used the
430 information contained in the UCSC Genome Browser Broad Histone track from the
431 ENCODE Project. Histone mark peaks were downloaded for every combination of cell
432 line and antibody. For each track, a 2x2 contingency table was built to represent the
433 partition of the whole set of possible probes in the microarray with respect to the
434 membership of the subset of interest and the overlap between the probes and the histone
435 peaks. A Fisher's exact test was used to determine whether there was significant
436 enrichment of the selected histone mark for the subset of interest. P-values were
437 adjusted for multiple comparisons using the Benjamini-Hochberg method for
438 controlling FDR. A significance level of 0.05 was used to determine whether the given
439 combination of histone mark and cell line presented a significant change in proportion.
440 Additionally, the base-2 logarithm of the Odds Ratio (OR) was used as a measure of
441 effect size.

442

443 *Chromatin segment enrichment analysis*

444 Data from the BROAD ChromHMM Project were downloaded from the UCSC Genome
445 Browser site. Each of the tracks comprising this dataset represents a different
446 segmentation generated by a Hidden Markov Model (HMM) using Chip-Seq signals
447 from the Broad Histone Project as inputs. The segmentations were later curated and
448 labelled according to their functional status [8; 9]. In order to detect any significant
449 enrichment in the proportion of probes in a given subset of interest belonging to one
450 functional category, an analysis strategy similar to the one employed for the detection of
451 histone enrichment was performed. In this case, a 2x2 contingency table was built using
452 segments of a given functional status rather than antibodies. A Fisher's exact test was
453 employed, and significant combinations were detected using a FDR threshold of 0.05
454 (Benjamini-Hochberg procedure). Again, the base-2 logarithm of the OR was used as a
455 measure of effect size.

456

457 *Genomic region analysis*

458 The probes in the microarray were assigned to a genomic region according to their
459 position relative to the transcript information extracted from the R/Bioconductor
460 package TxDb.Hsapiens.UCSC.hg19.knownGene (package version 3.1.2). A probe was
461 said to be in a promoter region if it was located in a region up to 2kb upstream of the
462 transcription start site (TSS) of any given transcript. Similarly, a set of mutually
463 exclusive regions were defined inside the transcripts, namely 5UTR, 3UTR, First Exon,
464 Exon and Intron. A probe could only belong to one category, hence if the location of a
465 probe overlapped with two or more regions in different transcripts, it was assigned to
466 the region with a higher level of precedence (i.e. in the order stated above, earlier
467 mention indicates higher precedence). If a probe was not assigned to any of these
468 special regions, it was labelled by default as Intergenic. A contingency table was built
469 for each of the subsets, partitioning the whole set of probes according to membership to
470 a given category and the subset of interest. A Pearson's χ^2 test was used to determine
471 whether there was any significant change in proportion between the number of probes
472 marked as belonging to a given region inside and outside the subset of interest. A
473 significance level of 0.05 was employed, and effect size measured by OR.

474

475 *CGI status analysis*

476 Similar to the genomic region analysis, probes were labelled according to their relative
477 position to CpG-islands (CGIs), the locations of which were obtained from the

478 R/Bioconductor package FDb.InfiniumMethylation.hg19 (package version 2.2.0). The
479 generation procedure of these CGIs is described by [53], i.e. ‘CpG shores’ were defined
480 as the 2kbp regions flanking a CGI. ‘CpG shelves’ were defined as the 2kbp regions
481 either upstream of or downstream from each CpG shore. Probes not belonging to any of
482 the regions thus far mentioned were assigned to the special category ‘non-CGI’ with
483 each probe being assigned to only one of the categories. A 4x2 contingency table was
484 constructed for each subset of probes in order to study the association between the given
485 subset and the different CGI categories. A χ^2 test was used to determine whether any of
486 the categories had a significant association with the given subset. For each of the CGI
487 status levels, a 2x2 contingency table was defined and another χ^2 test used to
488 independently evaluate the association of the given subset with each status level, a
489 significance level of 0.05 being employed for all tests. Effect size was reported as the
490 OR for each of the individual tests.

491

492 *Analysis of CpG density*

493 For each of the probes in the HumanMethylation450 microarray, CpG density was
494 measured as the number of CG 2-mers present divided by the number which would be
495 theoretically possible in a 2kbp window with the CpG under study at its centre. A
496 Wilcoxon non-parametric test was used to determine if any significant difference
497 existed between the CpG density of each subset of interest and that of the array probes
498 in the background. A significance level of 0.05 was employed for all tests. Effect size
499 was measured using Cliff's Delta (D).

500

501 *Gap distance analysis*

502 Distance to both the centromere and telomere was measured for each of the probes in
503 the HumanMethylation450 microarray. In order to find significant differences between
504 the probes within the subset of interest and those in the background, a Wilcoxon non-
505 parametric test was used. Once again, a significance level of 0.05 was employed for all
506 tests, and Cliff's Delta (D) was used as a measure of effect size.

507

508 *Microarray background correction*

509 Although it is sometimes referred to as a genome-wide solution, the
510 HumanMethylation450 BeadChip only covers a fraction of the entire genome. In its
511 27K predecessor, the probes were mainly located at gene promoter regions, while the

512 newer HumanMethylation450 BeadChip additionally includes probes located inside
513 genes and in intergenic regions [4].

514 The irregular distribution of probes can however lead to unwanted biases when studying
515 whether a selected subset of probes is enriched with respect to any functional or clinical
516 mark. For this reason, here a reference to the background distribution of features was
517 included in all statistical tests performed in order to prevent our conclusions from being
518 driven by the irregular distribution of probes. In qualitative tests (CGI status, genomic
519 region, or histone mark enrichment), the contingency matrix was built to represent the
520 background distribution of the microarray. In quantitative tests (CpG density, distance
521 to centromeres and telomeres) the corresponding metric was compared between the
522 subset of interest and the remaining probes in the microarray. Thus, any significant
523 result would indicate a departure from the fixed background distribution and ignore any
524 bias inherent in the test.

525

526 *Gene ontology analysis and annotation*

527 Probe sets were converted to gene sets by using the annotation information from the
528 R/Bioconductor package TxDb.Hsapiens.UCSC.hg19.knownGene (version 3.1.2). A
529 probe was assigned to a gene if the probe was contained within the overlap of all the
530 genomic regions represented by the different transcripts belonging to that gene, or in a
531 2kbp region upstream of the corresponding TSS. Probes converted this way can be
532 assigned to one or more genes, or to zero (i.e. intergenic probes).

533 After gene conversion, each subset of interest was analyzed using the HOMER software
534 tool [18]. The software was configured to use the whole set of genes represented in the
535 HumanMethylation450 architecture as a background. HOMER tested the genes in each
536 subset of interest against 21 different databases, including the Gene Ontology (GO)
537 Biological Process, Molecular Function and Cellular Component ontologies, as well as
538 KEGG and Reactome pathway databases, among many others.

539

540 *Circular visualization and track smoothing*

541 In order to plot the CpG and histone peak information on the circular genome-wide and
542 example graphs, smoothing was applied to the data. CpG enrichment information for
543 canonical and non-canonical hypermethylation was generated by partitioning the
544 genome into intervals of 10kbp and assigning to each a score corresponding to the
545 average coverage of the selected CpGs in the interval.

546

547 *Whole-genome bisulfite sequencing (WGBS) datasets*

548 Supplementary data referenced in [42] was used as a validation dataset in glioblastoma.
549 Previously processed data in the form of quantified methylation for each CpG measured
550 in both strands of the genome was downloaded and filtered. Only methylation measures
551 from CpGs having a total read count higher than 10 were retained.

552 The resulting dataset comprised only two samples (normal and tumoral), so a
553 descriptive strategy was used to distinguish the different types of probes according to
554 their methylation status. Hydroxymethylated probes were identified as those having a
555 5hmC measure higher than 0.1. Differentially methylated probes were defined as those
556 having an absolute difference in their methylation values, between the control and
557 tumor samples, higher than a given threshold (0.2 for 5mC and 0.1 for 5hmC).

558 The validation datasets may contain either one or two methylation measures for each
559 CpG in the genome as they measure methylation in both strands. Strand-agnostic CpG
560 regions representing the CpG dinucleotides with at least one measure were defined in
561 order to compute the degree of intersection between the WGBS and methylation arrays
562 results.

563

564 *TCGA expression datasets*

565 In order to analyze changes in gene expression, samples of glioblastoma multiforme
566 (GBM) were selected from among the data generated by the TCGA Research Network
567 (<http://cancergenome.nih.gov>). Expression Level-3 pre-processed data was obtained for
568 572 GBM samples (10 controls and 562 tumors). The moderated t-test approach in the
569 R/Bioconductor package *limma* was used to assess the differential expression status of
570 each gene in the TCGA datasets. The normalized expression ratio in the TCGA datasets
571 was used as the response variable, and the sample group (normal/tumoral) as the
572 covariate of interest. No adjustment for possible confounders was performed in this
573 case. An FDR threshold of 0.001 was used to correct for multiple hypotheses. No
574 filtering on effect size was applied in this case.

575

576 *Data analysis workflow*

577 All the necessary steps for upstream and downstream analyses were defined and
578 implemented using the Snakemake tool [26], which helps data scientists to generate a

579 reproducible and inherently parallel processing pipeline. Individual workflow tasks
580 were implemented in R (version 3.2.2) and Python (version 3.4.3).
581

582 **Acknowledgments**

583 We thank Ronnie Lendrum for editorial assistance. We also thank the Tumor Bank of
584 the Hospital Virgen de la Salud (BioB-HVS, Toledo, Spain) for providing tumor
585 samples. This work has been financially supported by: the Plan Nacional de I+D+I
586 2013-2016/FEDER (PI15/00892 to M.F.F. and A.F.F.); the ISCIII-Subdirección
587 General de Evaluación y Fomento de la Investigación, and the Plan Nacional de I+D+I
588 2008-2011/FEDER (CP11/00131 to A.F.F.); IUOPA (to G.F.B. and M.S); the
589 Fundación Científica de la AECC (to R.G.U.); the Fundación Ramón Areces (to M.F.F);
590 FICYT (to E.G.T., M.G.G. and A.C.); and the Asturias Regional Government
591 (GRUPIN14-052 to M.F.F.). Work in P.M. lab is supported by the European Research
592 Council (CoG-2014-646903), the Spanish Ministry of Economy-Competitiveness (SAF-
593 SAF2013-43065), the Obra Social La Caixa-Fundació Josep Carreras, and the
594 Generalitat de Catalunya. P.M. is an investigator in the Spanish Cell Therapy
595 cooperative network (TERCEL). The IUOPA is supported by the Obra Social Cajastur-
596 Liberbank, Spain.

597

598

599 **Competing interests:** The authors declare that no competing interests exist.

600

601 **References**

602

- 603 [1] M. Bachman, S. Uribe-Lewis, X. Yang, M. Williams, A. Murrell, and S.
604 Balasubramanian, 5-Hydroxymethylcytosine is a predominantly stable DNA
605 modification. *Nat Chem* 6 (2014) 1049-55.
- 606 [2] M. Bibikova, B. Barnes, C. Tsan, V. Ho, B. Klotzle, J.M. Le, D. Delano, L. Zhang,
607 G.P. Schroth, K.L. Gunderson, J.B. Fan, and R. Shen, High density DNA
608 methylation array with single CpG site resolution. *Genomics* 98 (2011) 288-95.
- 609 [3] C.G. Chapman, C.J. Mariani, F. Wu, K. Meckel, F. Butun, A. Chuang, J. Madzo,
610 M.B. Bissonette, J.H. Kwon, and L.A. Godley, TET-catalyzed 5-
611 hydroxymethylcytosine regulates gene expression in differentiating colonocytes
612 and colon cancer. *Sci Rep* 5 (2015) 17568.
- 613 [4] S. Dedeurwaerder, M. Defrance, E. Calonne, H. Denis, C. Sotiriou, and F. Fuks,
614 Evaluation of the Infinium Methylation 450K technology. *Epigenomics* 3 (2011)
615 771-84.
- 616 [5] P. Du, X. Zhang, C.C. Huang, N. Jafari, W.A. Kibbe, L. Hou, and S.M. Lin,
617 Comparison of Beta-value and M-value methods for quantifying methylation
618 levels by microarray analysis. *BMC Bioinformatics* 11 (2010) 587.
- 619 [6] H. Easwaran, S.E. Johnstone, L. Van Neste, J. Ohm, T. Mosbrugger, Q. Wang, M.J.
620 Aryee, P. Joyce, N. Ahuja, D. Weisenberger, E. Collisson, J. Zhu, S.
621 Yegnasubramanian, W. Matsui, and S.B. Baylin, A DNA hypermethylation
622 module for the stem/progenitor cell signature of cancer. *Genome Res* 22 (2012)
623 837-49.
- 624 [7] M. Ehrlich, DNA hypomethylation in cancer cells. *Epigenomics* 1 (2009) 239-59.
- 625 [8] J. Ernst, and M. Kellis, Discovery and characterization of chromatin states for
626 systematic annotation of the human genome. *Nat Biotechnol* 28 (2011) 817-25.
- 627 [9] J. Ernst, P. Kheradpour, T.S. Mikkelsen, N. Shores, L.D. Ward, C.B. Epstein, X.
628 Zhang, L. Wang, R. Issner, M. Coyne, M. Ku, T. Durham, M. Kellis, and B.E.
629 Bernstein, Mapping and analysis of chromatin state dynamics in nine human cell
630 types. *Nature* 473 (2011) 43-9.
- 631 [10] M. Esteller, Aberrant DNA methylation as a cancer-inducing mechanism. *Annu*
632 *Rev Pharmacol Toxicol* 45 (2005) 629-56.
- 633 [11] A.P. Feinberg, and B. Tycko, The history of cancer epigenetics. *Nat Rev Cancer* 4
634 (2004) 143-53.
- 635 [12] A.F. Fernandez, G.F. Bayon, R.G. Urdinguio, E.G. Torano, M.G. Garcia, A.
636 Carella, S. Petrus-Reurer, C. Ferrero, P. Martinez-Cambor, I. Cubillo, J. Garcia-
637 Castro, J. Delgado-Calle, F.M. Perez-Campo, J.A. Riancho, C. Bueno, P.
638 Menendez, A. Mentink, K. Mareschi, F. Claire, C. Fagnani, E. Medda, V.
639 Toccaceli, S. Brescianini, S. Moran, M. Esteller, A. Stolzing, J. de Boer, L.
640 Nistico, M.A. Stazi, and M.F. Fraga, H3K4me1 marks DNA regions
641 hypomethylated during aging in human stem and differentiated cells. *Genome*
642 *Res* 25 (2014) 27-40.
- 643 [13] G. Ficz, M.R. Branco, S. Seisenberger, F. Santos, F. Krueger, T.A. Hore, C.J.
644 Marques, S. Andrews, and W. Reik, Dynamic regulation of 5-
645 hydroxymethylcytosine in mouse ES cells and during differentiation. *Nature* 473
646 (2011) 398-402.
- 647 [14] S.F. Field, D. Beraldi, M. Bachman, S.K. Stewart, S. Beck, and S.
648 Balasubramanian, Accurate measurement of 5-methylcytosine and 5-

- 649 hydroxymethylcytosine in human cerebellum DNA by oxidative bisulfite on an
650 array (OxBS-array). *PLoS One* 10 (2015) e0118202.
- 651 [15] J.P. Fortin, A. Labbe, M. Lemire, B.W. Zanke, T.J. Hudson, E.J. Fertig, C.M.
652 Greenwood, and K.D. Hansen, Functional normalization of 450k methylation
653 array data improves replication in large cancer studies. *Genome Biol* 15 (2014)
654 503.
- 655 [16] D. Globisch, M. Munzel, M. Muller, S. Michalakis, M. Wagner, S. Koch, T.
656 Bruckl, M. Biel, and T. Carell, Tissue distribution of 5-hydroxymethylcytosine
657 and search for active demethylation intermediates. *PLoS One* 5 (2010) e15367.
- 658 [17] M.C. Haffner, A. Chaux, A.K. Meeker, D.M. Esopi, J. Gerber, L.G. Pellakuru, A.
659 Toubaji, P. Argani, C. Iacobuzio-Donahue, W.G. Nelson, G.J. Netto, A.M. De
660 Marzo, and S. Yegnasubramanian, Global 5-hydroxymethylcytosine content is
661 significantly reduced in tissue stem/progenitor cell compartments and in human
662 cancers. *Oncotarget* 2 (2011) 627-37.
- 663 [18] S. Heinz, C. Benner, N. Spann, E. Bertolino, Y.C. Lin, P. Laslo, J.X. Cheng, C.
664 Murre, H. Singh, and C.K. Glass, Simple combinations of lineage-determining
665 transcription factors prime cis-regulatory elements required for macrophage and
666 B cell identities. *Mol Cell* 38 (2010) 576-89.
- 667 [19] J.G. Herman, J. Jen, A. Merlo, and S.B. Baylin, Hypermethylation-associated
668 inactivation indicates a tumor suppressor role for p15INK4B. *Cancer Res* 56
669 (1996) 722-7.
- 670 [20] G.C. Hon, C.X. Song, T. Du, F. Jin, S. Selvaraj, A.Y. Lee, C.A. Yen, Z. Ye, S.Q.
671 Mao, B.A. Wang, S. Kuan, L.E. Edsall, B.S. Zhao, G.L. Xu, C. He, and B. Ren,
672 5mC oxidation by Tet2 modulates enhancer activity and timing of transcriptome
673 reprogramming during differentiation. *Mol Cell* 56 (2014) 286-97.
- 674 [21] S. Ito, A.C. D'Alessio, O.V. Taranova, K. Hong, L.C. Sowers, and Y. Zhang, Role
675 of Tet proteins in 5mC to 5hmC conversion, ES-cell self-renewal and inner cell
676 mass specification. *Nature* 466 (2010) 1129-33.
- 677 [22] D. Jiang, Y. Zhang, R.P. Hart, J. Chen, K. Herrup, and J. Li, Alteration in 5-
678 hydroxymethylcytosine-mediated epigenetic regulation leads to Purkinje cell
679 vulnerability in ATM deficiency. *Brain* 138 (2015) 3520-36.
- 680 [23] S.G. Jin, Y. Jiang, R. Qiu, T.A. Rauch, Y. Wang, G. Schackert, D. Krex, Q. Lu,
681 and G.P. Pfeifer, 5-Hydroxymethylcytosine is strongly depleted in human
682 cancers but its levels do not correlate with IDH1 mutations. *Cancer Res* 71
683 (2011) 7360-5.
- 684 [24] T. Khare, S. Pai, K. Koncivicius, M. Pal, E. Kriukiene, Z. Liutkeviciute, M. Irimia,
685 P. Jia, C. Ptak, M. Xia, R. Tice, M. Tochigi, S. Morera, A. Nazarians, D.
686 Belsham, A.H. Wong, B.J. Blencowe, S.C. Wang, P. Kapranov, R. Kustra, V.
687 Labrie, S. Klimasauskas, and A. Petronis, 5-hmC in the brain is abundant in
688 synaptic genes and shows differences at the exon-intron boundary. *Nat Struct
689 Mol Biol* 19 (2012) 1037-43.
- 690 [25] M. Klug, S. Schmidhofer, C. Gebhard, R. Andreesen, and M. Rehli, 5-
691 Hydroxymethylcytosine is an essential intermediate of active DNA
692 demethylation processes in primary human monocytes. *Genome Biol* 14 (2013)
693 R46.
- 694 [26] J. Koster, and S. Rahmann, Snakemake--a scalable bioinformatics workflow
695 engine. *Bioinformatics* 28 (2012) 2520-2.
- 696 [27] T.F. Kraus, D. Globisch, M. Wagner, S. Eigenbrod, D. Widmann, M. Munzel, M.
697 Muller, T. Pfaffeneder, B. Hackner, W. Feiden, U. Schuller, T. Carell, and H.A.
698 Kretzschmar, Low values of 5-hydroxymethylcytosine (5hmC), the "sixth base,"

- 699 are associated with anaplasia in human brain tumors. *Int J Cancer* 131 (2015)
700 1577-90.
- 701 [28] T.F. Kraus, G. Kolck, A. Greiner, K. Schierl, V. Guibourt, and H.A. Kretzschmar,
702 Loss of 5-hydroxymethylcytosine and intratumoral heterogeneity as an
703 epigenomic hallmark of glioblastoma. *Tumour Biol* 36 (2012) 8439-46.
- 704 [29] S. Kriaucionis, and N. Heintz, The nuclear DNA base 5-hydroxymethylcytosine is
705 present in Purkinje neurons and the brain. *Science* 324 (2009) 929-30.
- 706 [30] Y. Kudo, K. Tateishi, K. Yamamoto, S. Yamamoto, Y. Asaoka, H. Ijichi, G.
707 Nagae, H. Yoshida, H. Aburatani, and K. Koike, Loss of 5-
708 hydroxymethylcytosine is accompanied with malignant cellular transformation.
709 *Cancer Sci* 103 (2012) 670-6.
- 710 [31] J.T. Leek, and J.D. Storey, Capturing heterogeneity in gene expression studies by
711 surrogate variable analysis. *PLoS Genet* 3 (2007) 1724-35.
- 712 [32] W. Li, and M. Liu, Distribution of 5-hydroxymethylcytosine in different human
713 tissues. *J Nucleic Acids* 2011 (2011) 870726.
- 714 [33] C.G. Lian, Y. Xu, C. Ceol, F. Wu, A. Larson, K. Dresser, W. Xu, L. Tan, Y. Hu, Q.
715 Zhan, C.W. Lee, D. Hu, B.Q. Lian, S. Kleffel, Y. Yang, J. Neiswender, A.J.
716 Khorasani, R. Fang, C. Lezcano, L.M. Duncan, R.A. Scolyer, J.F. Thompson, H.
717 Kakavand, Y. Houvras, L.I. Zon, M.C. Mihm, Jr., U.B. Kaiser, T. Schatton,
718 B.A. Woda, G.F. Murphy, and Y.G. Shi, Loss of 5-hydroxymethylcytosine is an
719 epigenetic hallmark of melanoma. *Cell* 150 (2012) 1135-46.
- 720 [34] C. Liu, L. Liu, X. Chen, J. Shen, J. Shan, Y. Xu, Z. Yang, L. Wu, F. Xia, P. Bie, Y.
721 Cui, X.W. Bian, and C. Qian, Decrease of 5-hydroxymethylcytosine is
722 associated with progression of hepatocellular carcinoma through downregulation
723 of TET1. *PLoS One* 8 (2013) e62828.
- 724 [35] J. Maksimovic, L. Gordon, and A. Oshlack, SWAN: Subset-quantile within array
725 normalization for illumina infinium HumanMethylation450 BeadChips. *Genome*
726 *Biol* 13 (2012) R44.
- 727 [36] K.M. McGarvey, J.A. Fahrner, E. Greene, J. Martens, T. Jenuwein, and S.B.
728 Baylin, Silenced tumor suppressor genes reactivated by DNA demethylation do
729 not return to a fully euchromatic chromatin state. *Cancer Res* 66 (2006) 3541-9.
- 730 [37] C.E. Nestor, R. Ottaviano, J. Reddington, D. Sproul, D. Reinhardt, D. Dunican, E.
731 Katz, J.M. Dixon, D.J. Harrison, and R.R. Meehan, Tissue type is a major
732 modifier of the 5-hydroxymethylcytosine content of human genes. *Genome Res*
733 22 (2012) 467-77.
- 734 [38] J.E. Ohm, K.M. McGarvey, X. Yu, L. Cheng, K.E. Schuebel, L. Cope, H.P.
735 Mohammad, W. Chen, V.C. Daniel, W. Yu, D.M. Berman, T. Jenuwein, K.
736 Pruitt, S.J. Sharkis, D.N. Watkins, J.G. Herman, and S.B. Baylin, A stem cell-
737 like chromatin pattern may predispose tumor suppressor genes to DNA
738 hypermethylation and heritable silencing. *Nat Genet* 39 (2007) 237-42.
- 739 [39] B.A. Orr, M.C. Haffner, W.G. Nelson, S. Yegnasubramanian, and C.G. Eberhart,
740 Decreased 5-hydroxymethylcytosine is associated with neural progenitor
741 phenotype in normal brain and shorter survival in malignant glioma. *PLoS One*
742 7 (2012) e41036.
- 743 [40] N. Plongthongkum, D.H. Diep, and K. Zhang, Advances in the profiling of DNA
744 modifications: cytosine methylation and beyond. *Nat Rev Genet* 15 (2014) 647-
745 61.
- 746 [41] E.L. Putiri, R.L. Tiedemann, J.J. Thompson, C. Liu, T. Ho, J.H. Choi, and K.D.
747 Robertson, Distinct and overlapping control of 5-methylcytosine and 5-

- 748 hydroxymethylcytosine by the TET proteins in human cancer cells. *Genome*
749 *Biol* 15 (2014) R81.
- 750 [42] E.-A. Raiber, D. Beraldi, S. Martinez Cuesta, G.R. McInroy, Z. Kingsbury, J.
751 Becq, T. James, M. Lopes, K. Allinson, S. Field, S. Humphray, T. Santarius, C.
752 Watts, D. Bentley, and S. Balasubramanian, Base resolution maps reveal the
753 importance of 5-hydroxymethylcytosine in a human glioblastoma. *npj Genomic*
754 *Medicine* 2 (2017) 6.
- 755 [43] M.E. Ritchie, B. Phipson, D. Wu, Y. Hu, C.W. Law, W. Shi, and G.K. Smyth,
756 limma powers differential expression analyses for RNA-sequencing and
757 microarray studies. *Nucleic Acids Res* 43 (2015) e47.
- 758 [44] C.X. Song, K.E. Szulwach, Y. Fu, Q. Dai, C. Yi, X. Li, Y. Li, C.H. Chen, W.
759 Zhang, X. Jian, J. Wang, L. Zhang, T.J. Looney, B. Zhang, L.A. Godley, L.M.
760 Hicks, B.T. Lahn, P. Jin, and C. He, Selective chemical labeling reveals the
761 genome-wide distribution of 5-hydroxymethylcytosine. *Nat Biotechnol* 29
762 (2011) 68-72.
- 763 [45] S.K. Stewart, T.J. Morris, P. Guilhamon, H. Bulstrode, M. Bachman, S.
764 Balasubramanian, and S. Beck, oxBS-450K: a method for analysing
765 hydroxymethylation using 450K BeadChips. *Methods* 72 (2015) 9-15.
- 766 [46] M. Sun, C.X. Song, H. Huang, C.A. Frankenberger, D. Sankarasharma, S. Gomes,
767 P. Chen, J. Chen, K.K. Chada, C. He, and M.R. Rosner, HMGA2/TET1/HOXA9
768 signaling pathway regulates breast cancer growth and metastasis. *Proc Natl*
769 *Acad Sci U S A* 110 (2013) 9920-5.
- 770 [47] M. Tahiliani, K.P. Koh, Y. Shen, W.A. Pastor, H. Bandukwala, Y. Brudno, S.
771 Agarwal, L.M. Iyer, D.R. Liu, L. Aravind, and A. Rao, Conversion of 5-
772 methylcytosine to 5-hydroxymethylcytosine in mammalian DNA by MLL
773 partner TET1. *Science* 324 (2009) 930-5.
- 774 [48] E.G. Torano, G.F. Bayon, A. Del Real, M.I. Sierra, M.G. Garcia, A. Carella, T.
775 Belmonte, R.G. Urdinguio, I. Cubillo, J. Garcia-Castro, J. Delgado-Calle, F.M.
776 Perez-Campo, J.A. Riancho, M.F. Fraga, and A.F. Fernandez, Age-associated
777 hydroxymethylation in human bone-marrow mesenchymal stem cells. *J Transl*
778 *Med* 14 (2016) 207.
- 779 [49] R.G. Urdinguio, G.F. Bayon, M. Dmitrijeva, E.G. Torano, C. Bravo, M.F. Fraga,
780 L. Bassas, S. Larriba, and A.F. Fernandez, Aberrant DNA methylation patterns
781 of spermatozoa in men with unexplained infertility. *Hum Reprod* 30 (2015)
782 1014-28.
- 783 [50] S. Uribe-Lewis, R. Stark, T. Carroll, M.J. Dunning, M. Bachman, Y. Ito, L. Stojic,
784 S. Halim, S.L. Vowler, A.G. Lynch, B. Delatte, E.J. de Bony, L. Colin, M.
785 Defrance, F. Krueger, A.L. Silva, R. Ten Hoopen, A.E. Ibrahim, F. Fuks, and A.
786 Murrell, 5-hydroxymethylcytosine marks promoters in colon that resist DNA
787 hypermethylation in cancer. *Genome Biol* 16 (2015) 69.
- 788 [51] D.J. Weisenberger, M. Campan, T.I. Long, M. Kim, C. Woods, E. Fiala, M.
789 Ehrlich, and P.W. Laird, Analysis of repetitive element DNA methylation by
790 MethyLight. *Nucleic Acids Res* 33 (2005) 6823-36.
- 791 [52] M. Widschwendter, H. Fiegl, D. Egle, E. Mueller-Holzner, G. Spizzo, C. Marth,
792 D.J. Weisenberger, M. Campan, J. Young, I. Jacobs, and P.W. Laird, Epigenetic
793 stem cell signature in cancer. *Nat Genet* 39 (2007) 157-8.
- 794 [53] H. Wu, B. Caffo, H.A. Jaffee, R.A. Irizarry, and A.P. Feinberg, Redefining CpG
795 islands using hidden Markov models. *Biostatistics* 11 (2010) 499-514.
- 796 [54] G.R. Wyatt, and S.S. Cohen, A new pyrimidine base from bacteriophage nucleic
797 acids. *Nature* 170 (1952) 1072-3.

798 [55] H. Yang, Y. Liu, F. Bai, J.Y. Zhang, S.H. Ma, J. Liu, Z.D. Xu, H.G. Zhu, Z.Q.
799 Ling, D. Ye, K.L. Guan, and Y. Xiong, Tumor development is associated with
800 decrease of TET gene expression and 5-methylcytosine hydroxylation.
801 *Oncogene* 32 (2013) 663-9.

802

803

804

805 **Figure legends**

806

807 **Figure 1. 5mC and 5hmC levels at repetitive DNA sequences in glioma and CRC.**

808 5mC (A) and 5hmC (B) values of several repetitive regions (AluYb8, LINE-1, NBL-2,
809 and Sat-alpha) measured by bisulfite pyrosequencing in controls and glioma (left
810 panels) and CRC (right panels). Individual CpG site values for each repeat are
811 displayed. P-values are shown.

812

813 **Figure 2. Characterization of DNA 5hmC in normal brain and colon samples. (A)**

814 Box plots showing differences between average Beta values of 5mC+5hmC (BS) and
815 true 5mC (OxBS) in both normal brain and colon. On the right are Hilbert curves
816 showing the amount and genomic distribution of 5hmC in brain and colon. (B)
817 Associations between 5hmC and CpG density. (C) Distribution of 5hmC CpG sites
818 relative to CpG island status and compared to the array background (450K). (D)
819 Distribution of 5hmC CpG sites relative to different genomic regions. (E) Heatmaps
820 showing significant enrichment of the 5hmC CpG sites, identified in brain and colon,
821 with different histone marks contained in the UCSC Browser Broad Histone track from
822 the ENCODE project. Color code indicates the significant enrichment based on log2
823 odds ratio (OR).

824

825 **Figure 3. Alterations of 5hmC in CRC and glioma. (A)** Bar plot showing the number

826 of dh5mC sites in CRC and glioma. (B) Unsupervised hierarchical clustering and
827 heatmap including CpG sites with 5hmC loss in glioma. (C) Associations between
828 5hmC loss in glioma and density of CpGs (upper panel), CpG island status (middle
829 panel), and different genomic regions (lower panel). (D) Heatmaps showing significant
830 enrichment of hypo 5hmC CpGs identified in glioma with different histone marks
831 contained in the UCSC Browser Broad Histone track from the ENCODE project. (E)
832 Heatmaps showing significant enrichment of hypo 5hmC CpGs in gliomas with fifteen
833 “chromatin states” generated by a Hidden Markov Model (HMM) (right panel). Color
834 codes indicate the significant enrichment based on log2 odds ratio (OR).

835

836 **Figure 4. Relationships between changes in 5mc and 5hmc in glioma. (A)** Euler

837 diagram illustrating overlap of CpGs that lose 5hmC (hypo 5hmC) and gain 5mC (hyper
838 5mC) in glioma. (B) Associations between hypermethylated CpG sites that lose (or not)
839 5hmC and CpG density and CpG island status, compared to the array background

840 (450K). (C) Unsupervised hierarchical clustering and heatmap including CpG sites with
841 5mC changes (hyper- and hypomethylation) in glioma. Hypo- (purple) and non-hypo
842 (orange) 5hmC overlapped CpGs are indicated by colored lines on the annexed track.
843 Average beta methylation values are displayed from 0 (blue) to 1 (yellow).

844

845 **Figure 5. Canonical and non-canonical hypermethylation in glioma.** (A) Heatmaps
846 showing significant enrichment of CpG sites in glioma which exclusively gain 5mC
847 (canonical hypermethylation) (upper panel), and both lose 5hmC and gain 5mC (non-
848 canonical hypermethylation) (lower panel), with different histone marks contained in
849 the UCSC Browser Broad Histone track from the ENCODE project. Histone PTMs
850 related to activation and repression are distinguished by colors as indicated in the key.
851 (B) Circular representation of two representative chromosomes (12 and 17), indicating
852 genomic location of canonical (orange) and non-canonical (purple) hypermethylation in
853 glioma. Inner tracks display chromatin marks (H3K9me3, H3K27me3, and H3K4me2),
854 generated for NH-A cells. Two examples of genes showing canonical and non-canonical
855 hypermethylation associated with specific chromatin signatures are displayed below.

856

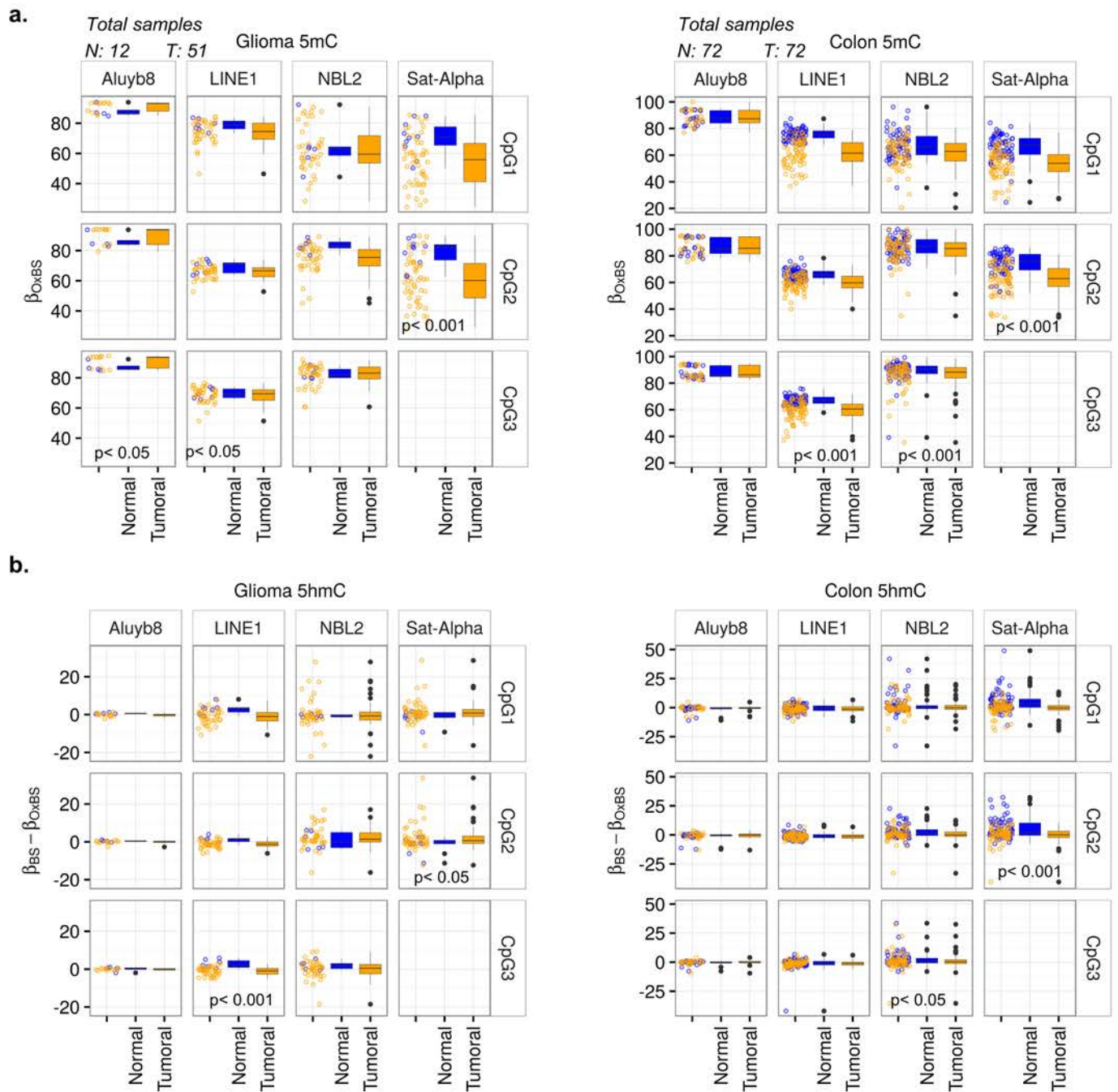
857 **Figure 6. Functional role of canonical and non-canonical hypermethylation in**
858 **glioma.** (A) Euler diagrams showing number of genes associated with canonical
859 hypermethylation, non-canonical hypermethylation, or both. On the right are
860 representative gene ontology terms (Biological process) of genes associated with
861 canonical (orange) and non-canonical (purple) hypermethylation, ranked by Q-value,
862 and enrichment score (relative risk). (B) Euler diagram showing overlap of canonical
863 and non-canonical hypermethylated genes with down-regulation. (C) Associations of
864 canonical and non-canonical hypermethylation in glioma with different genomic
865 regions. (D) Representative example of one gene (*SLC1A4*) showing non-canonical
866 hypermethylation in glioma (orange frame). Organization of the gene, locations of
867 CpGs included in the methylation array (black dots), and transcription start site (TSS)
868 are shown below. 5mC hypermethylation (blue to yellow) and 5hmC loss (gray to blue)
869 in glioma are shown above. Whole genome bisulfite sequencing (WGBS) data [42]
870 including all the CpG sites in the same region are shown on the right. The associated
871 change in gene expression is displayed below.

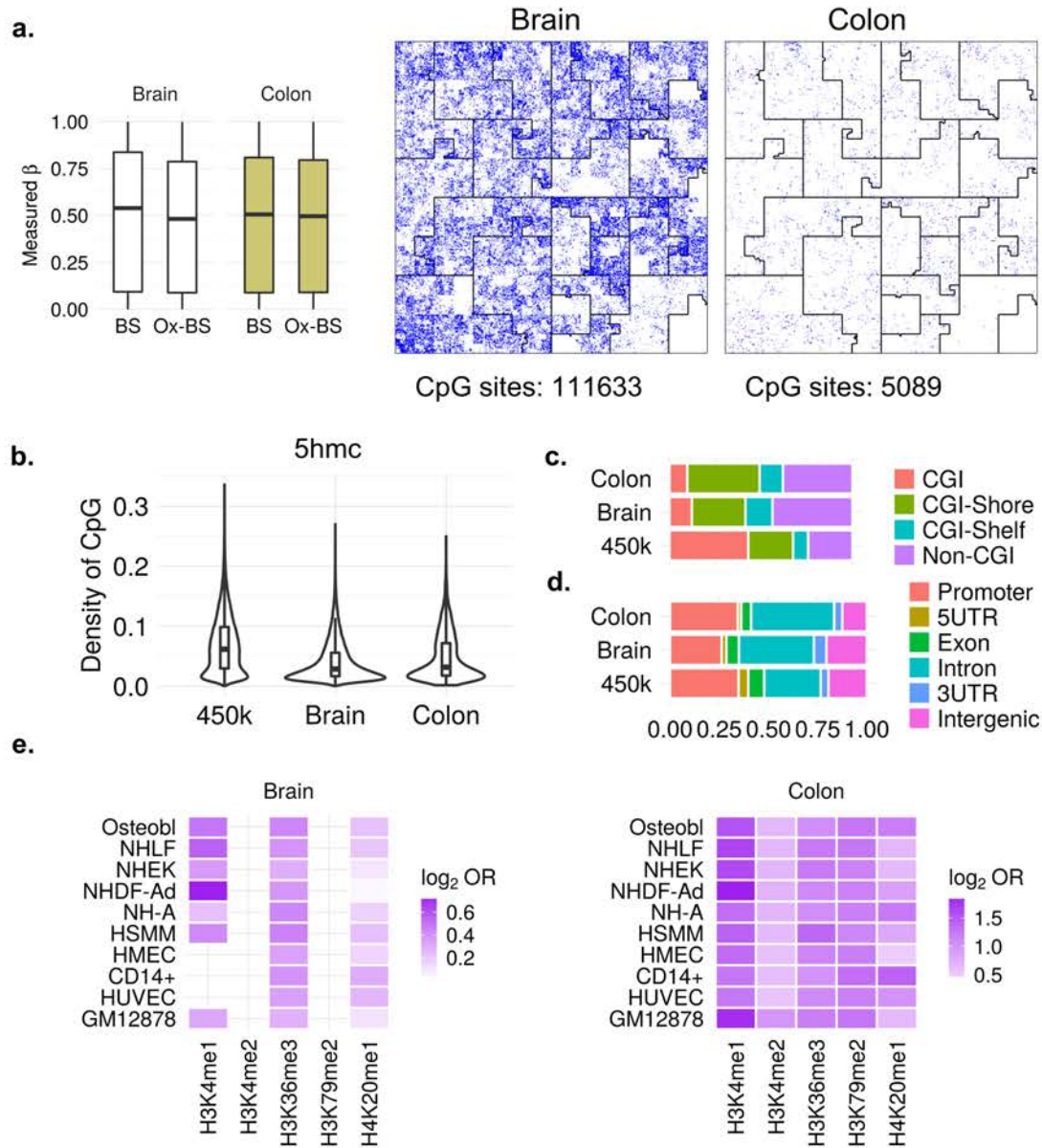
872

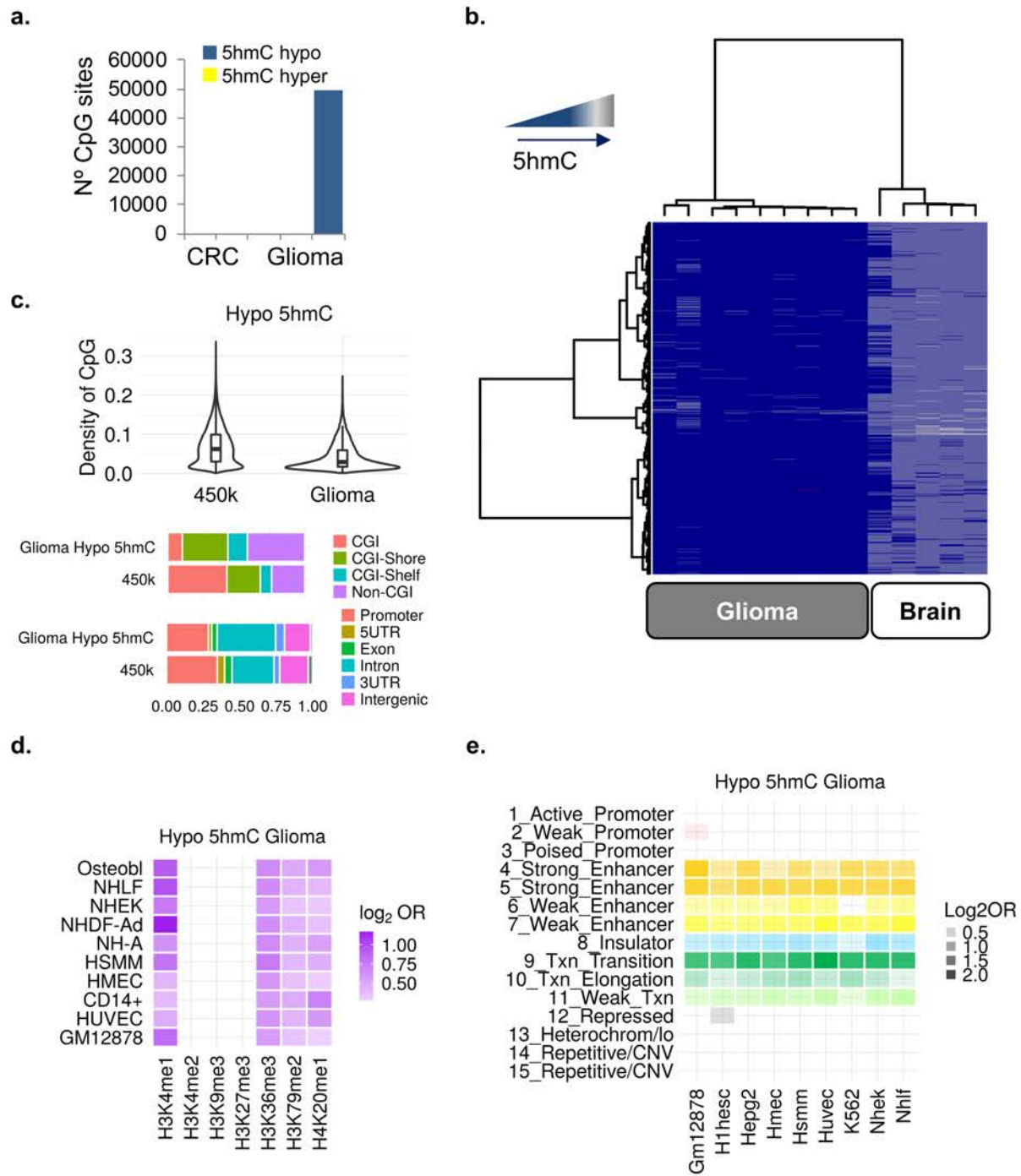
873 **Supplementary files**

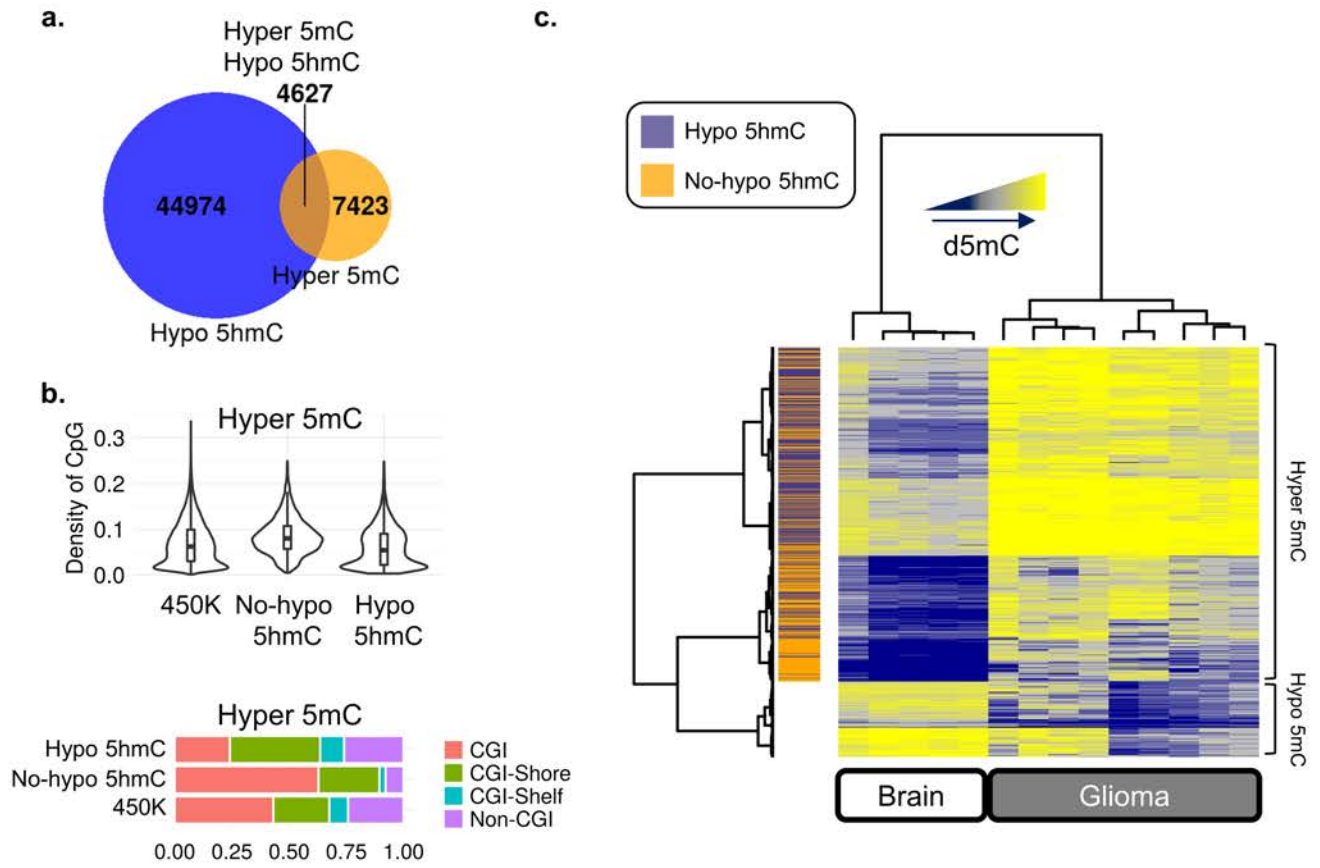
874 Supplementary Tables 1-12

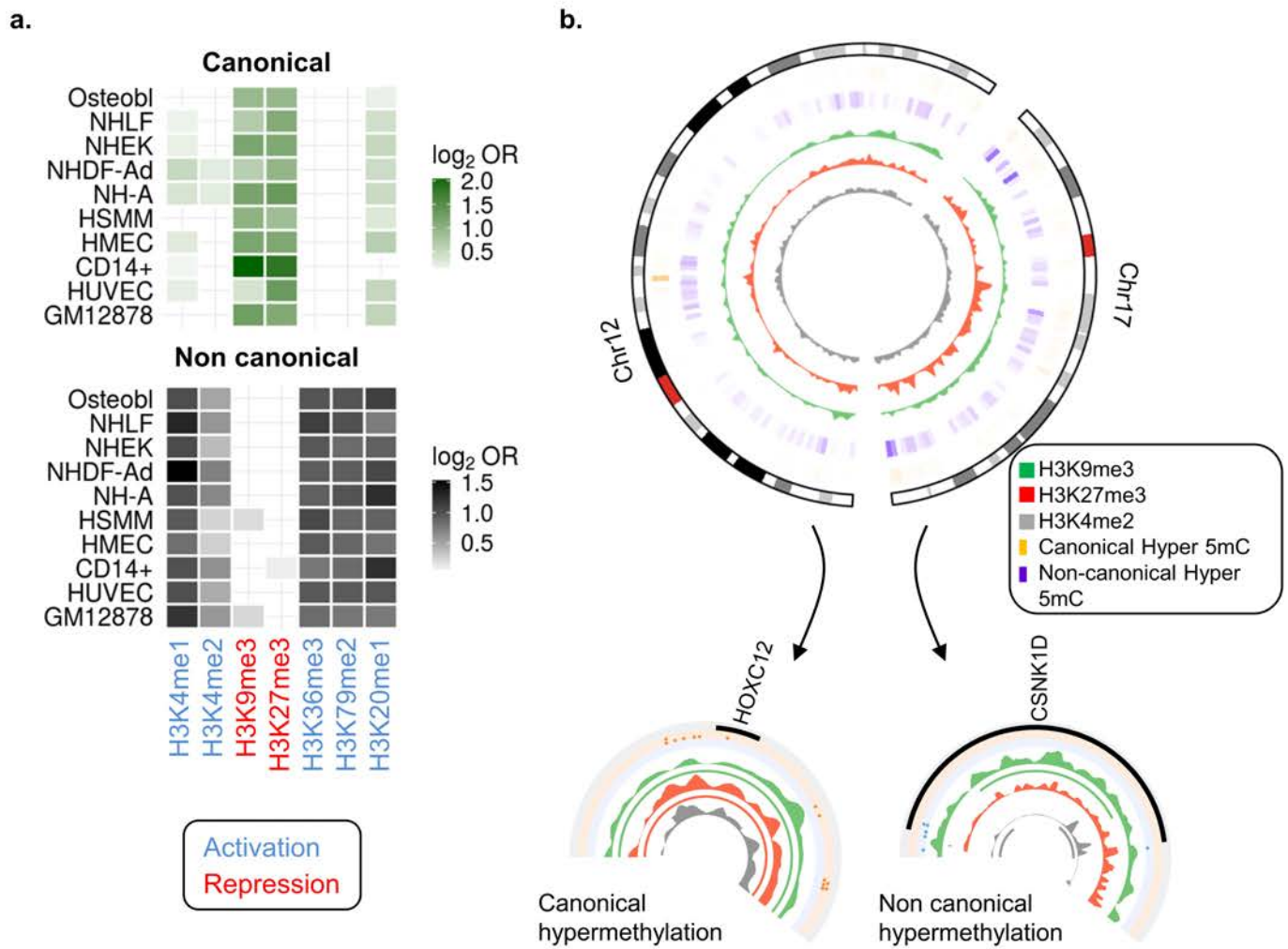
875 Supplementary figures 1-4











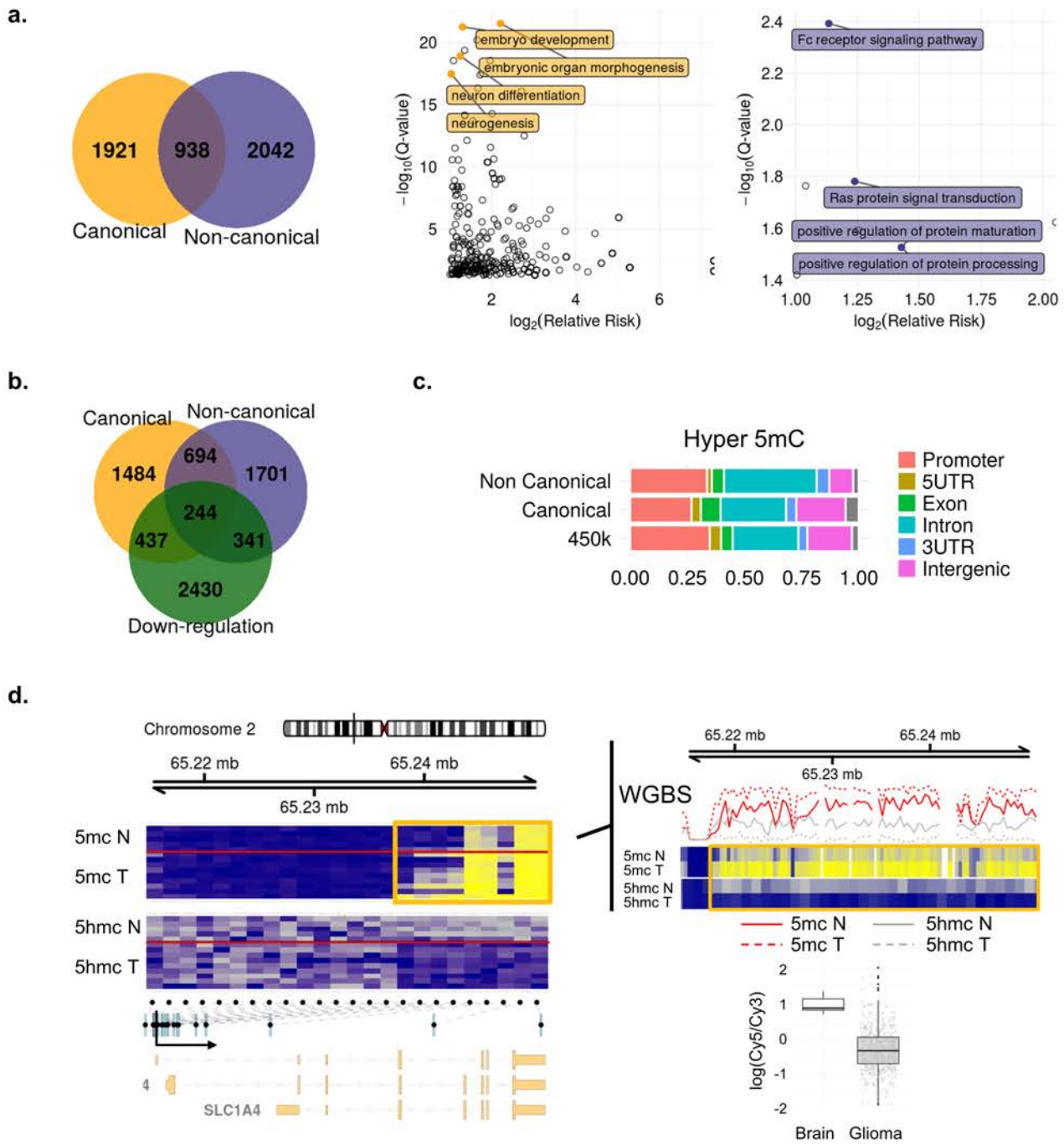


Figure 4-figure supplement 1

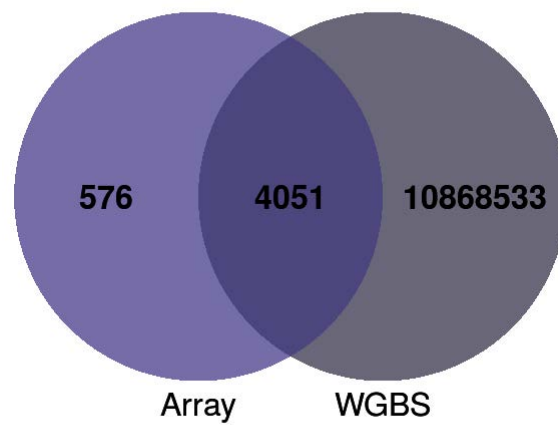


Figure 4-figure supplement 1. Venn diagram showing the overlap of hyper5mC-hypo5hmC sites in glioma obtained by methylation arrays and whole genome bisulfite sequencing (WGBS).

Figure 5-figure supplement 1

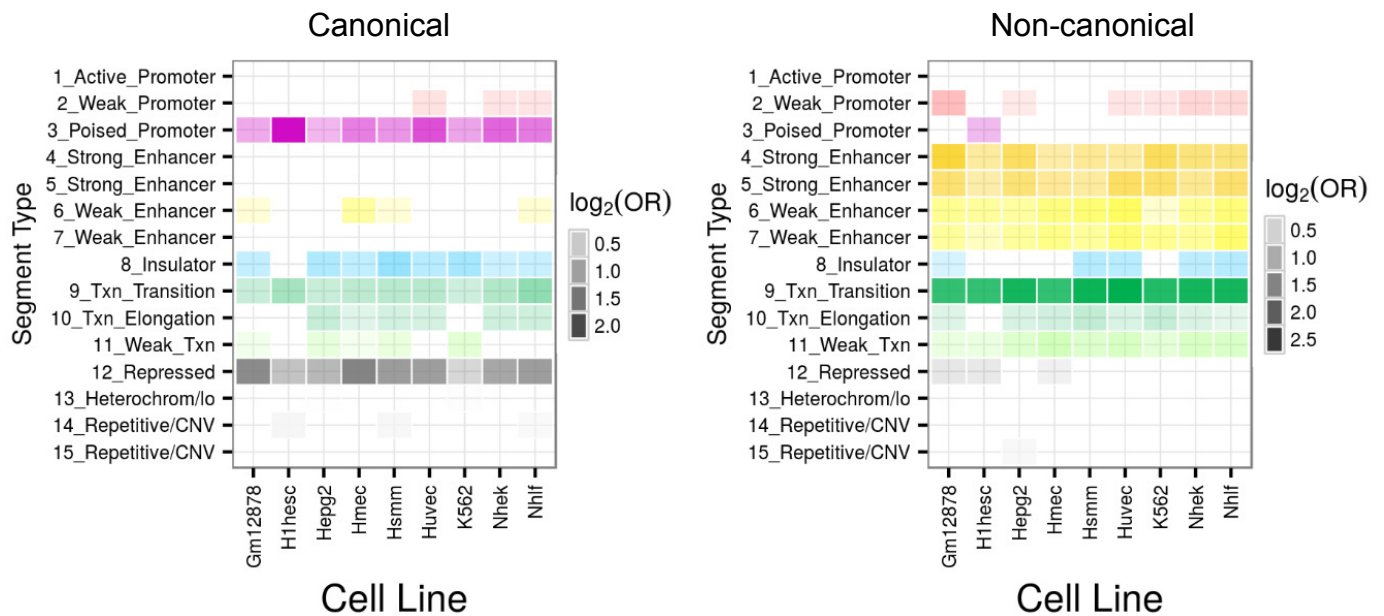


Figure 5-figure supplement 1. Canonical and non-canonical hypermethylation in glioma. Heatmaps showing significant enrichment of canonical (left panel) and non-canonical (right panel) hypermethylated CpG sites with fifteen “chromatin states” generated by a Hidden Markov Model (HMM). Colour codes indicate the significant enrichment based on \log_2 odds ratio (OR).

Figure 6-figure supplement 1

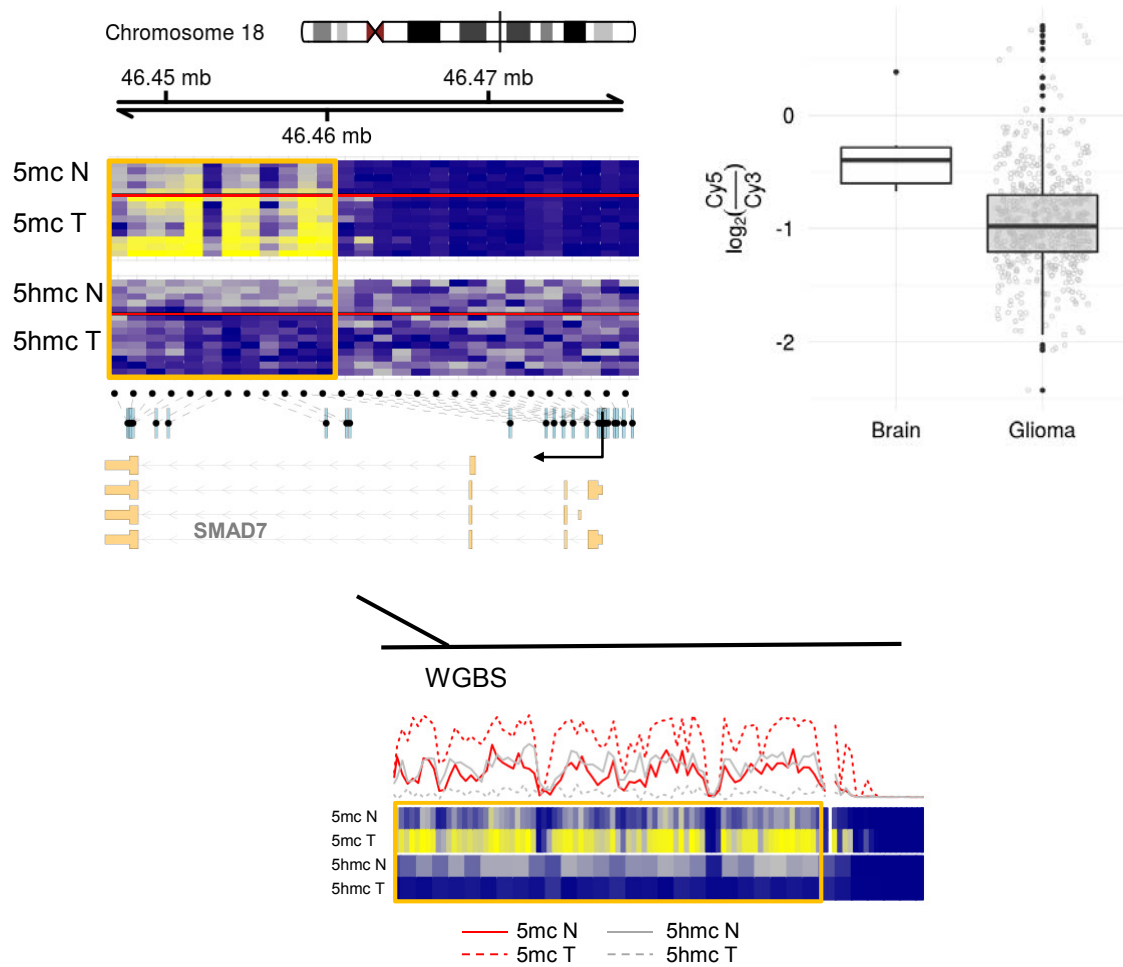


Figure 6-figure supplement 1. Functional role of non-canonical hypermethylation in glioma. Representative example of one gene (*SMAD7*) showing non-canonical hypermethylation in glioma (orange frame). Organization of the gene, location of CpGs included in the methylation array (black dots), and transcription start site (TSS) are shown below. 5mC hypermethylation (blue to yellow) and 5hmC loss (gray to blue) in gliomas are shown above. Lower panel shows the full genome bisulfite sequencing (WGBS) data including all the CpG sites in the same region. The associated change in gene expression is shown on the right.

Figure 2-figure supplement 1

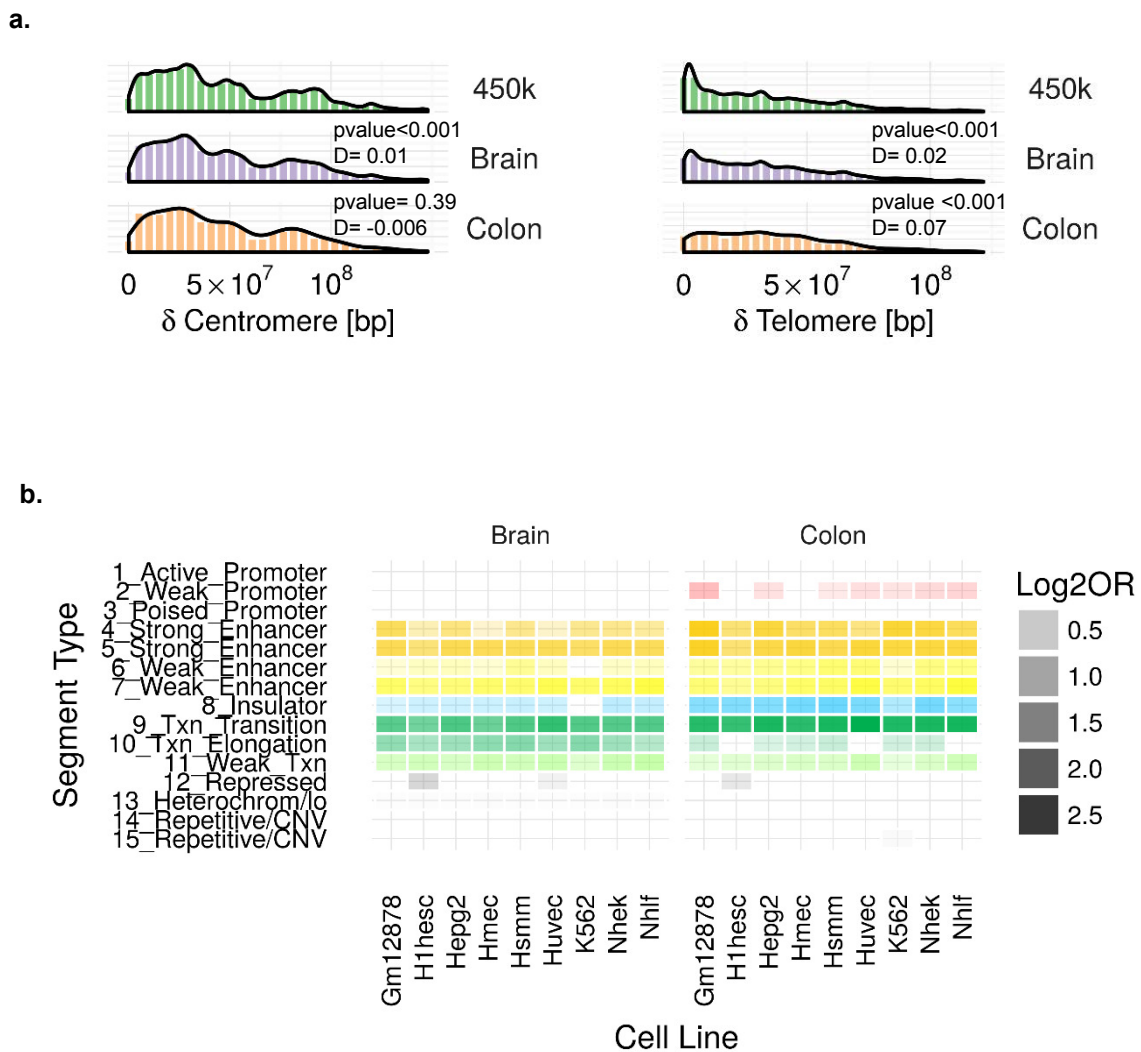


Figure 2-figure supplement 1. Characterization of DNA 5hmC in normal brain and colon samples. (a) Histograms and density plots showing the associations between 5hmC CpGs and distance to centromere (left) and telomeres (right). **(b)** Heatmaps showing significant enrichment of 5hmC CpG sites with fifteen “chromatin states” generated by a Hidden Markov Model (HMM). Colour codes indicate the significant enrichment based on \log_2 odds ratio (OR).