

DNA G-quadruplex structures mould the DNA methylome

Shi-Qing Mao¹, Avazeh T. Ghanbarian¹, Jochen Spiegel¹, Sergio Martínez Cuesta¹, Dario Beraldi^{1,4}, Marco Di Antonio², Giovanni Marsico¹, Robert Hänsel-Hertsch¹, David Tannahill¹ and Shankar Balasubramanian^{*1,2,3}

¹ Cancer Research UK Cambridge Institute, Li Ka Shing Centre, Cambridge, UK

² Department of Chemistry, University of Cambridge, Cambridge, UK

³ School of Clinical Medicine, University of Cambridge, Cambridge, UK

⁴ Current address: Institute of Cancer Sciences, University of Glasgow, Glasgow, UK

* Correspondence should be addressed to S.B. (sb10031@cam.ac.uk)

Abstract

Control of DNA methylation level is critical for gene regulation, and the factors that govern hypomethylation at CpG islands (CGIs) are still being uncovered. Here, we provide evidence that G-quadruplex (G4) DNA secondary structures are genomic features that influence methylation at CGIs. We show that the presence of G4 structure is tightly associated with CGI hypomethylation in the human genome. Surprisingly, we find that these G4 sites are enriched for DNA methyltransferase 1 (DNMT1) occupancy, which is consistent with our biophysical observations that DNMT1 exhibits higher binding affinity for G4s as compared to duplex, hemi-methylated or single-stranded DNA. The biochemical assays also show that the G4 structure itself, rather than sequence, inhibits DNMT1 enzymatic activity. Based on these data, we propose that G4 formation sequesters DNMT1 thereby protecting certain CGIs from methylation and inhibiting local methylation.

Introduction

Methylation of cytosine at C-5 is a key DNA modification in development and disease^{1,2}. In mammals, cytosine methylation occurs predominantly at CpG dinucleotides and is installed and maintained by three DNA methyltransferase enzymes (DNMT1, 3A and 3B) that are essential for development³⁻⁵. CpGs occur less frequently than expected in the mammalian genome and show a bimodal distribution with respect to methylation^{6,7}. Sparsely distributed CpGs (~90%), found in genic and intergenic regions, tend to be highly methylated, while CpGs

35 found in dense GC-rich regions, so-called CpG Islands (CGIs), are largely depleted of
36 methylation and are prevalent at the promoters of house-keeping and developmental
37 genes^{8,9}. Outside of embryonic development, gross methylation patterns are generally stable
38 across different tissues^{10,11}. Nonetheless during key cellular events, methylation can be
39 dynamic at specific loci to modulate gene expression, such as *de novo* methylation of some
40 promoter CGIs with intermediate CpG density during lineage commitment¹².

41
42 General rules on maintenance of the default methylation state are being uncovered and
43 several studies suggest that CGI hypomethylation is sequence-dependent¹³⁻¹⁸. Furthermore,
44 DNMTs are reported to be actively and continuously excluded from CpG-poor distal
45 regulatory regions through competitive inhibition with DNA binding proteins, such as NRF1
46 and CTCF/REST, thus maintaining the hypomethylated state of regulatory regions^{19,20}. Lowly
47 methylated regions also co-localise with DNase I hypersensitivity sites marking accessible
48 chromatin regions²¹. The presence of enhancer chromatin marks, such as histone
49 modifications, also play an important role in forming the unique chromatin structure of
50 CGIs^{22,23}. In mouse embryonic stem cells, the CXXC finger protein 1, Cfp1, is believed to
51 promote CGI hypomethylation through binding unmethylated CpG and recruitment of H3K4
52 methyltransferases to promote H3K4me3^{24,25}. However, Cfp1 binding and/or H3K4me3 are
53 not required for the 'protection' of CGI from DNA methylation since Cfp1 knockout results in
54 a dramatic loss of H3K4me3 at CGIs without increasing DNA methylation²⁴. This suggests that
55 Cfp1 binding and/or H3K4me3 are not required to prevent CGIs from DNA methylation, thus
56 there may be other factors that are fundamental to impart the hypomethylated state.

57
58 Alternative DNA secondary structures, known as G-quadruplexes (G4s) are found within
59 certain G-rich sequences and arise through the self-association of guanine bases to form
60 stacked tetrads (**Fig. 1a**)²⁶. G4s are increasingly being recognised as important features in the
61 genome and over 700,000 G4s have been biophysically mapped in purified human genomic
62 DNA by high-throughput sequencing²⁷. G4 structures have been observed in human using
63 immunofluorescence with a G4-specific antibody (BG4)²⁸, and linked with transcriptional
64 regulation and are enriched in gene promoters including many oncogenes^{26,29}. Recently, G4-
65 chromatin immunoprecipitation sequencing (G4-ChIP-seq) has been developed to map G4

66 structures in human chromatin^{30,31}. Corroborating a link with transcription, the majority of
67 G4-ChIP-seq sites were found predominantly in regulatory, nucleosome-depleted chromatin,
68 particularly in gene promoters^{31,32}. As both G4s and hypomethylated CGIs are associated
69 with actively transcribed genes^{9,31}, this raises the question of whether there is an interplay
70 between G4 formation and DNA methylation.

71
72 Herein we present evidence that most G4 structures, as detected by G4-ChIP-seq, are formed
73 in regions comprising unmethylated CGIs in the human genome. We also uncover a striking
74 co-localisation of G4 structures and DNMT1 docked at CGIs, and we demonstrate that
75 DNMT1 methylation activity is inhibited by DNA G4 structures. Our data suggest a
76 mechanism for the 'protection' of CGIs from methylation by G-quadruplex structures that
77 locally sequester and inhibit DNMT1.

78

79 **Results**

80 **G4 structures in active chromatin are found within hypomethylated CGIs**

81 To explore any potential relationship between G4 structures in chromatin and methylation
82 levels, we employed human K562 chronic myelogenous leukaemia cells in which methylation
83 has been comprehensively characterised at single base resolution using whole genome
84 bisulfite sequencing (WGBS) by the ENCODE project³³. We generated a genome-wide dataset
85 for G4 structures by G4 ChIP-Seq³¹ using the G4-specific antibody BG4²⁸ and compared the
86 BG4 peak overlap with CGIs³⁴. Strikingly, we found that the majority of BG4 peaks (79%,
87 7111/8952) overlapped with a CGI (covering 23% of all CGIs) (**Fig. 1b**). The majority of CpG
88 island regions span 200 to 1000 bp (median/mean, 569/775 bp), while BG4 peak regions
89 span 100 to 400 bp (median/mean, 205/226 bp) (**Fig. 1c**). 83% (5935/7111) of these CGIs
90 overlap with one BG4 peak (**Fig. 1d**). Furthermore, when the level of methylation at BG4
91 peaks was considered, we noted that there was a dramatic absence of methylation at BG4
92 peaks (mean 1%, median 0.5%), compared with average genome methylation (28.4%) (**Fig.**
93 **1e**). To rule out any effect of the cytosine methylation state on the ability of the BG4
94 antibody to recognise a G4 structure, an ELISA binding assay was used to show that BG4 can
95 bind to G4 structured DNA with equal affinity irrespective of the presence of cytosine
96 methylation (**SI Fig. 1**). DNase I hypersensitive sites (DHS), which mark open chromatin, are

97 also mainly hypomethylated (mean 11%, median 2.5%, **Fig. 1e**). Confirming our previous
98 observations³¹, the majority of BG4 peaks are found in open chromatin (DHSs, 97%,
99 8655/8952), and it is notable that these sites have the lowest methylation levels (**Fig. 1e**).
100 Overall CGI methylation (mean 27%, median 8%, **Fig. 1e**) shows a broader distribution than
101 BG4 regions, since some CGIs are associated with active hypomethylated promoters while
102 others with inactive genes or gene bodies and thus are more heavily methylated^{9,35}. This
103 prominent association between BG4 peaks, hypomethylation and particular CGIs is
104 suggestive of a functional link between G4 secondary structures and the establishment
105 and/or maintenance of low methylation status at these CGIs in active chromatin.

106
107 It has recently been concluded from work in mouse embryonic stem cells, that both high
108 CpG-density and high GC-richness are required to establish the hypomethylated state at
109 CGIs¹⁵. It is therefore notable that BG4 peaks have a similar level of GC richness to CGIs (**Fig.**
110 **1f**) with most (79%) being located in regions of CpG density comparable to that seen in CGIs
111 (**Fig. 1g**). It has been suggested that CpG density alone is only a minor determinant of the
112 unmethylated state, as dense CpG sequences embedded within an AT-rich context are
113 invariably highly methylated when inserted into the mouse genome^{15,16}. Indeed, when we
114 compare the average methylation of CGIs (**Fig. 1h**) to that of BG4 peaks at different CpG
115 dinucleotide densities, it is noteworthy that across a wide range of CpG densities, BG4 peak
116 regions are always largely devoid of methylation (**Fig. 1h**). We also confirmed these
117 observations using an alternative CGI definition set generated by CpGCluster algorithm³⁶ (**SI**
118 **Fig. 2a, 2b**). Furthermore, when methylation at CGIs is considered with respect to the
119 presence or absence of a BG4 peak, it is noteworthy that there is an almost a total lack of
120 methylation at CGIs with BG4 than CGIs without (**SI Fig. 2c**). This strongly suggests that CGIs
121 associated with the physical presence of a G4 structure generally have particularly low
122 methylation.

123
124 To explore low methylation in different CGI contexts, we calculated methylation levels
125 relative to BG4 presence in CGIs containing i) no promoter or DHS site, ii) a promoter alone,
126 iii) a DHS site alone and iv) both a promoter and DHS site. It is apparent that CGIs containing
127 a BG4 peak always have lower methylation in open (DHS +) or closed (DHS -) chromatin, or in
128 the presence or absence of a promoter. (**SI Fig. 2d**). CGIs with a DHS site and promoter but

129 without a BG4 peak (4500 sites) have higher methylation (mean 2%, median 2%) than those
130 CGIs (5567 sites) with a BG4 peak plus promoter and DHSs (mean 1%, median 0.5%) (**SI Fig.**
131 **2d**, right two panels). The lowest observed methylation states are found therefore at sites
132 carrying a G4 structure, suggesting that the physical presence of a G4 structure within CGI is
133 an important feature with respect to the hypomethylation state. This is illustrated in **Fig. 1i**
134 which shows the co-incidence of BG4 peaks with hypomethylated promoter CGIs for a
135 representative genome region.

136
137 In earlier work, we found that treatment of human epidermal keratinocytes (HaCaT) cells
138 with the HDAC inhibitor entinostat led to increased BG4 binding signal primarily located in
139 open chromatin promoter regions³⁷. We therefore generated WGBS datasets to examine
140 DNA methylation changes with respect to BG4 signal. Consistent with our observation in
141 K562 cells, BG4 peaks in HaCaT cells have lower methylation compared with open chromatin
142 and CGI regions (**SI Fig. 2e**). In open-chromatin promoter CGI regions, 307 had a significant
143 increase in BG4 signal (BG4 increase, > 1.5-fold change in signal and FDR < 0.05, see Online
144 Methods). No change in BG4 signal was seen in 3261 CGI promoter regions before and after
145 treatment (BG4 constant), or for 1504 G-rich CGI promoter regions that do not have a BG4
146 peak (BG4 negative) but have the potential to form a G4 *in vitro*²⁷. Despite open-chromatin
147 promoter CGI regions already being predisposed to low methylation, we see a statistically
148 significant additional drop in methylation levels at CGIs where BG4 peak size increases after
149 HDAC inhibition (**SI Fig. 2f**). Overall, these data support that formation of G-quadruplex
150 structures in CGIs is linked to lower methylation.

151
152 **DNMT1 is sequestered at G4 structures associated with low methylation**
153 Given that regions where G4 structures marked by BG4 peaks are generally observed to be
154 hypomethylated, we considered that the DNA methyltransferases might have some form of
155 physical interaction with G4 structures in the chromatin context. We focused on DNMT1
156 since DNMT1 knockout is lethal causing global DNA methylation loss in all dividing somatic
157 cells and human embryonic stem cells (ESCs)^{3,5,38,39}, whereas DNMT3A/B knockouts mainly
158 affect non-CpG methylation in human ESCs⁵. When we considered the distribution of DNMT1
159 binding sites in K562 cells, downloaded from ENCODE³³ (516,483 peaks in total across both
160 biological replicates), we found that 52% (4611/8952, Monte Carlo simulation's P-value

161 1.25e-04) of the G4 structures mapped by G4-ChIP (BG4 peaks) overlapped with at least one
162 DNMT1 binding site. Of the remaining 4341 BG4 peaks, 4003 were within 1 Kb of a DNMT1
163 binding site. The proximity of BG4 peaks to DNMT1 recruitment sites is illustrated graphically
164 in **Fig. 2a** for a representative genome region. Intriguingly, when the distribution of DNMT1
165 binding is plotted relative to high, intermediate or low methylated CGIs, we observe a
166 prominent enrichment of DNMT1 binding at lowly methylated CGIs which overlaps with
167 those regions with the highest BG4 peak density as well as DHS sites (**Fig. 2b**). A similar
168 profile is also seen using alternative CGI definition set generated by CpGCluster³⁶ (**SI Fig. 3**).
169 The observation that DNMT1 enrichment at G4 regions that lack methylation is, at first
170 glance, somewhat unexpected and counter-intuitive, given that DNMT1 installs methylation.
171 We therefore considered the possibility of a mechanism whereby DNMT1 protein is
172 sequestered at these sites in active chromatin but prevented from methylating CpGs in that
173 locality.

174

175 **DNMT1 selectively binds to and is inhibited by G4 structures**

176 To address whether DNMT1 binds G4 structures directly, we carried out biophysical
177 measurements using an enzyme-linked immunosorbent assay (ELISA) to measure the binding
178 of recombinant human FLAG-tagged full-length DNMT1 protein to immobilized target DNA
179 structures (see Online Methods). Biotinylated single-stranded oligonucleotides of sequence
180 based on the promoters of BCL2, KIT2 and MYC were chosen as these fold into well-
181 characterised G4 structures^{40–42}. Mutated versions (BCL2-mut, KIT2-mut and MYC-mut) that
182 are unable to form a G4 structure were also used as controls. The presence or absence of G4
183 folded structure with G4 oligonucleotides and mutated controls was confirmed by circular
184 dichroism (CD) spectroscopy and ultraviolet (UV) thermal melting spectroscopic analysis (**SI**
185 **Fig. 4a-f**). We found that DNMT1 binds to all three G4 structures with low nanomolar affinity
186 ($K_d[\text{BCL2}] = 9.6 \pm 0.3 \text{ nM}$; $K_d[\text{KIT2}] = 15.2 \pm 0.4 \text{ nM}$, $K_d[\text{MYC}] = 25.3 \pm 0.4 \text{ nM}$, $n=3$; **Fig. 3a-c**).
187 DNMT1 showed a lower binding affinity for unmethylated duplex DNA (BCL2, $107 \pm 5 \text{ nM}$; **Fig.**
188 **3d**) and there was no specific binding observed for the single stranded mutated
189 oligonucleotide controls. Notably, DNMT1 generally showed a greater affinity for G4
190 structures than known DNMT1 substrates such as a hemi-methylated duplex DNA (BCL2, $85 \pm$
191 7 nM , **Fig. 3e**), or a synthetic poly(dI-dC)₅₀ substrate ($75 \pm 2 \text{ nM}$, **Fig. 3f**). DNMT1 binding to
192 G4 structures does not appear to depend on CpG dinucleotides, since the absence of CpG in

193 the MYC G4 did not preclude DNMT1 binding (**Fig. 3c**). To begin to dissect the binding mode
194 of DNMT1 to G4s, we used a competition ELISA assay. 50 nM immobilized BCL2 G4 (Kd for
195 DNMT1 = 9.6 nM, **Fig. 3a**) was incubated with DNMT1 protein in the presence of increasing
196 concentrations of competitors. Even with 100-fold excess (5 μ M) of DNA duplex (Kd for
197 DNMT1 = 107 nM, **Fig. 3d**) or poly-dIdC (Kd for DNMT1 = 75 nM, **Fig. 3f**), there was no
198 inhibition (**Fig. 3g**). This suggests that G4s and duplex DNA occupy different binding sites and
199 that the catalytic domain is not involved in G4 binding. A similar, non-overlapping G4 and
200 duplex binding has also been observed in other proteins such as TRF2⁴³ and Rap1⁴⁴.

201

202 The relatively high binding affinity and selectivity of DNMT1 for DNA G4 structures is
203 consistent with the observation that DNMT1 shows some localisation to G4 structures in
204 K562 cells (**Fig. 2a, b**). To validate the association with low methylation in the locality G4 sites
205 in the genome, we evaluated whether G4 DNA could actually inhibit DNA methylation on a
206 standard assay using poly(dI-dC)_n as substrate⁴⁵ of DNMT1 using a fluorometric biochemical
207 assay (Abcam, see Online Methods). Specifically, we evaluated different concentrations of
208 folded G4-structured oligonucleotides or mutated non-G4 controls, where the presence or
209 absence of G4 structure had been confirmed by CD spectroscopy (**SI Fig. 4g-i**). We indeed
210 found that each of three G4 structures resulted in significant inhibition of DNMT1
211 methyltransferase activity whereas the mutated control oligonucleotides did not (**Fig 3h-j**).
212 Gratifyingly, the potency of inhibition by each G4 was related to the binding affinity for
213 DNMT1, as determined by ELISA with BCL2 being the most potent inhibitor (50% inhibition at
214 ~25 nM), MYC being least potent (50% inhibition at ~1 μ M) and KIT2 being intermediate (50%
215 inhibition at 90 nM). No inhibition of activity was seen with mutated controls ranging from
216 400 nM-8 μ M concentration. C-rich oligonucleotides complementary to the G4 sequences
217 (BCL2-CCC, KIT2-CCC and MYC-CCC) or corresponding duplex DNA also had no effect on
218 DNMT1 activity (**SI Fig. 4j-l**). We also tested G4 oligonucleotides that were able to fold into a
219 G4 structure but carried a reduced number CpGs (BCL2, KIT) or had a number of artificially
220 introduced CpGs (MYC). In all cases, changes in the number of CpG sites only had minor
221 effects on DNMT1 inhibition (**SI Fig. 4j-l**). Taken together, these results indicate a novel and
222 unexpected feature of G4 structures as potential genomic regions that promote the
223 unmethylated state through recruitment and inhibition of DNMT1 activity.

224

225 **Recruitment of DNMT to G4 structures shapes the methylome**

226 The above data suggest that there is a striking lack of methylation (**Fig. 1e, h, SI Fig. 2a-d**) in
227 chromatin regions where G4 structure formation is observed. To rigorously question whether
228 this observation was related to the detectable formation of a G4 structure in chromatin (i. e.
229 a BG4 peak), or merely the G-rich sequences *per se* with potential to form a G4 structure, the
230 methylation profile for BG4 peak regions was compared to those of G-rich sequences that
231 can physically form a G4 structure in an in vitro sequencing assay²⁷ (here called Sequences
232 with potential to form G4s, **Fig. 4a**). As the majority of BG4 peaks (8,210) are found in open
233 chromatin, only sequences with potential to form G4s located in open chromatin (36,015)
234 were considered. The mean and median length is 226/205 bp for BG4 peaks, and 383/285 bp
235 for the latter. G-rich sequences with the potential to form a G4 are largely hypomethylated
236 (12%), with methylation levels rising in the flanking regions (45%), whereas BG4 peaks have
237 substantially lower methylation (down to 1%) and flanking regions being more methylated
238 (60%). The contrast between lowest methylation at BG4 sites with highest methylation at
239 distal flanking regions is also exemplified in the genome browser view in **Fig. 1i**. While G-
240 richness as defined by G4 sequence without structure is a feature that correlates with the
241 lack of methylation, there is a further dramatic loss of methylation due to the physical
242 presence of a G4 structure with these regions also being marked by a greater methylation
243 flanking the G4 structure. Regions with a G4 structure also correspond to CGIs that mark
244 particular active genes, and is in keeping with our previous data showing that G4s are
245 associated with particular chromatin states to promote elevated transcription³¹. R-loops
246 (three-stranded DNA-RNA hybrids) have also been linked to reduced methylation in
247 transcribed CpG island promoters^{46,47}. As R-loops form in a similar genomic context to G4s,
248 we tested the correlation of R-loops, BG4 peaks and methylation. Using the K562 R-loop
249 dataset⁴⁶, we found that 5685 BG4 peaks overlap with a R-loop, while 3267 BG4 peaks do not
250 and that BG4 peaks are depleted of methylation independent of R-loop presence (**SI Fig. 5**).
251 This suggests that G4 structure is strongly linked to hypomethylation, irrespective of the
252 presence of R-loop.

253

254 **Discussion**

255 Here we have provided evidence for a link between a DNA secondary structure formation
256 and epigenetic status. We have uncovered a unique chromatin context whereby certain CGIs

257 in active chromatin are depleted in methylation but carry a G4 structure and also the
258 surrounding flanking regions display higher than average methylation. This suggests that G4s
259 may impart a previously unknown and important function in the establishment of epigenome.
260 We propose a model (**Fig. 4b**) in which G4 formation, together with transcription factor
261 binding^{19,20}, contributes to loss of methylation at key genomic loci by sequestering DNMT1,
262 via G4 recognition, and locally inhibiting DNMT1 function at CpG islands. It is noteworthy
263 that this mechanism resembles a recently proposed model for recruitment and inhibition of
264 PRC2 complex by a RNA G-quadruplex present in the HOTAIR lncRNA^{48,49}. This suggests there
265 may be other mechanisms for epigenetic regulation that operate by the sequestration and
266 inhibition of epigenetic modifiers mediated by high affinity interactions with nucleic acid
267 secondary structures.

268

269 **Accession Codes**

270 K562 datasets for DHS (ENCSR000EPC), DNMT1 ChIP-seq (ENCSR987PBI) and whole genome
271 bisulfate sequencing (ENCSR765JPC) were downloaded from ENCODE. G4-ChIP-seq data sets
272 for K562 and WGBS datasets for entinostat-treated and untreated HaCaT cells are available
273 at the NCBI GEO repository under accession number GSE107690. G4-ChIP-seq data in
274 entinostat-treated and untreated HaCaT cells were taken from GSE76688.

275

276 **Acknowledgements**

277 This work is supported by a core CRUK award (C14303/A17197). S.B. is a Senior Investigator
278 of the Wellcome Trust (grant no. 099232/z/12/z). JS is a Marie Curie Fellow of the European
279 Union (747297-QAPs-H2020-MSCA-IF-2016).

280

281 **Author Contributions**

282

283 The project was conceived by SM and SB. SM designed and carried out all the experiments
284 with discussions from DB, JS, RHH, DT & SB. SM designed the analysis strategies with
285 discussions from AG, DT & SB. JS performed G4-ChIP-seq experiments. AG & SMC carried out
286 all computational analysis with discussions from SM, DT, DB, RHH & GM. MD carried out the
287 CD spectroscopy and UV melting experiments. All authors interpreted the results. SM, DT &
288 SB wrote the paper with input from all authors.

289

290 **Competing interests**

291 SB is an advisor and shareholder of Cambridge Epigenetix limited.

292 **References**

- 293 1. Smith, Z. D. & Meissner, A. DNA methylation: roles in mammalian development. *Nat.*
294 *Rev. Genet.* **14**, 204–20 (2013).
- 295 2. Feinberg, A. P. & Tycko, B. The history of cancer epigenetics. *Nat. Rev. Cancer* **4**, (2004).
- 296 3. Li, E., Bestor, T. H. & Jaenisch, R. Targeted mutation of the DNA methyltransferase
297 gene results in embryonic lethality. *Cell* **69**, 915–926 (1992).
- 298 4. Okano, M., Bell, D. W., Haber, D. A. & Li, E. DNA methyltransferases Dnmt3a and
299 Dnmt3b are essential for de novo methylation and mammalian development. *Cell* **99**,
300 247–257 (1999).
- 301 5. Liao, J. *et al.* Targeted disruption of DNMT1, DNMT3A and DNMT3B in human
302 embryonic stem cells. *Nat. Genet.* **47**, 469–478 (2015).
- 303 6. Bird, A. DNA methylation patterns and epigenetic memory. *Genes Dev.* **16**, 6–21
304 (2002).
- 305 7. Lister, R. *et al.* Human DNA methylomes at base resolution show widespread
306 epigenomic differences. *Nature* **462**, 315–322 (2009).
- 307 8. Illingworth, R. S. & Bird, A. P. CpG islands - 'A rough guide'. *FEBS Lett.* **583**, 1713–1720
308 (2009).
- 309 9. Deaton, A. & Bird, A. CpG islands and the regulation of transcription. *Genes Dev.* **25**,
310 1010–1022 (2011).
- 311 10. Reik, W., Dean, W. & Walter, J. Epigenetic reprogramming in mammalian development.
312 *Science* **293**, 1089–93 (2001).
- 313 11. Li, E. Chromatin modification and epigenetic reprogramming in mammalian
314 development. *Nat. Rev. Genet.* **3**, 662–673 (2002).
- 315 12. Reik, W. Stability and flexibility of epigenetic gene regulation in mammalian
316 development. *Nature* **447**, 425–432 (2007).
- 317 13. Long, H. K., King, H. W., Patient, R. K., Odom, D. T. & Klose, R. J. Protection of CpG
318 islands from DNA methylation is DNA-encoded and evolutionarily conserved. *Nucleic*
319 *Acids Res.* **44**, 6693–6706 (2016).
- 320 14. Lienert, F. *et al.* Identification of genetic elements that autonomously determine DNA
321 methylation states. *Nat. Genet.* **43**, 1091–1097 (2011).
- 322 15. Wachter, E. *et al.* Synthetic CpG islands reveal DNA sequence determinants of
323 chromatin structure. *Elife* **3**, 1–16 (2014).
- 324 16. Krebs, A. R., Dessus-Babus, S., Burger, L. & Schübeler, D. High-throughput engineering
325 of a mammalian genome reveals building principles of methylation states at CG rich
326 regions. *Elife* **3**, e04094 (2014).
- 327 17. Takahashi, Y. *et al.* Integration of CpG-free DNA induces de novo methylation of CpG
328 islands in pluripotent stem cells. *Science* **356**, 503–508 (2017).
- 329 18. Quante, T. & Bird, A. Do short, frequent DNA sequence motifs mould the epigenome?
330 *Nat. Rev. Mol. Cell Biol.* **17**, 257–62 (2016).
- 331 19. Domcke, S. *et al.* Competition between DNA methylation and transcription factors
332 determines binding of NRF1. *Nature* **528**, 575–579 (2015).
- 333 20. Stadler, M. B. *et al.* DNA-binding factors shape the mouse methylome at distal
334 regulatory regions. *Nature* **480**, 490–5 (2011).
- 335 21. Thurman, R. E. *et al.* The accessible chromatin landscape of the human genome.
336 *Nature* **489**, 75–82 (2012).
- 337 22. Ooi, S. K. T. *et al.* DNMT3L connects unmethylated lysine 4 of histone H3 to de novo
338 methylation of DNA. *Nature* **448**, 714–717 (2007).

- 339 23. Du, J., Johnson, L. M., Jacobsen, S. E. & Patel, D. J. DNA methylation pathways and
340 their crosstalk with histone methylation. *Nat. Rev. Mol. Cell Biol.* **16**, 519–532 (2015).
- 341 24. Clouaire, T. *et al.* Cfp1 integrates both CpG content and gene activity for accurate
342 H3K4me3 deposition in embryonic stem cells. *Genes Dev.* **26**, 1714–1728 (2012).
- 343 25. Thomson, J. P. *et al.* CpG islands influence chromatin structure via the CpG-binding
344 protein Cfp1. *Nature* **464**, 1082–1086 (2010).
- 345 26. Hänsel-Hertsch, R., Di Antonio, M. & Balasubramanian, S. DNA G-quadruplexes in the
346 human genome: Detection, functions and therapeutic potential. *Nat. Rev. Mol. Cell*
347 *Biol.* **18**, 279–284 (2017).
- 348 27. Chambers, V. S. *et al.* High-throughput sequencing of DNA G-quadruplex structures in
349 the human genome. *Nat. Biotechnol.* **33**, 1–7 (2015).
- 350 28. Biffi, G., Tannahill, D., McCafferty, J. & Balasubramanian, S. Quantitative visualization
351 of DNA G-quadruplex structures in human cells. *Nat. Chem.* **5**, 182–6 (2013).
- 352 29. Rhodes, D. & Lipps, H. J. G-quadruplexes and their regulatory roles in biology. *Nucleic*
353 *Acids Res.* **43**, 8627–8637 (2015).
- 354 30. Hänsel-Hertsch, R., Spiegel, J., Marsico, G., Tannahill, D. & Balasubramanian, S.
355 Genome-wide mapping of endogenous G-quadruplex DNA structures by chromatin
356 immunoprecipitation and high-throughput sequencing. *Nat. Protoc.* **13**, 551–564
357 (2018).
- 358 31. Hänsel-Hertsch, R. *et al.* G-quadruplex structures mark human regulatory chromatin.
359 *Nat. Genet.* **48**, 1267–1272 (2016).
- 360 32. De, S. & Michor, F. DNA secondary structures and epigenetic determinants of cancer
361 genome evolution. *Nat. Struct. Mol. Biol.* **18**, 950–5 (2011).
- 362 33. Dunham, I. *et al.* An integrated encyclopedia of DNA elements in the human genome.
363 *Nature* **489**, 57–74 (2012).
- 364 34. Gardiner-Garden, M. & Frommer, M. CpG Islands in vertebrate genomes. *J. Mol. Biol.*
365 **196**, 261–282 (1987).
- 366 35. Jones, P. A. Functions of DNA methylation: islands, start sites, gene bodies and beyond.
367 *Nat. Rev. Genet.* **13**, 484–92 (2012).
- 368 36. Hackenberg, M. *et al.* CpGcluster: A distance-based algorithm for CpG-island detection.
369 *BMC Bioinformatics* **7**, 1–13 (2006).
- 370 37. Chen, L. *et al.* R-ChIP Using Inactive RNase H Reveals Dynamic Coupling of R-loops with
371 Transcriptional Pausing at Gene Promoters. *Mol. Cell* **68**, 745–757.e5 (2017).
- 372 38. Fan, G. *et al.* DNA hypomethylation perturbs the function and survival of CNS neurons
373 in postnatal animals. *J. Neurosci.* **21**, 788–797 (2001).
- 374 39. Jackson-Grusby, L. *et al.* Loss of genomic methylation causes p53-dependent apoptosis
375 and epigenetic deregulation. *Nat. Genet.* **27**, 31–39 (2001).
- 376 40. Dai, J. *et al.* An intramolecular G-quadruplex structure with mixed parallel/antiparallel
377 G-strands formed in the human BCL-2 promoter region in solution. *J. Am. Chem. Soc.*
378 **128**, 1096–1098 (2006).
- 379 41. Kuryavyi, V., Phan, A. T. & Patel, D. J. Solution structures of all parallel-stranded
380 monomeric and dimeric G-quadruplex scaffolds of the human c-kit2 promoter. *Nucleic*
381 *Acids Res.* **38**, 6757–6773 (2010).
- 382 42. Ambrus, A., Chen, D., Dai, J., Jones, R. A. & Yang, D. Solution structure of the
383 biologically relevant G-quadruplex element in the human c-MYC promoter.
384 Implications for G-quadruplex stabilization. *Biochemistry* **44**, 2048–2058 (2005).
- 385 43. Biffi, G., Tannahill, D. & Balasubramanian, S. An intramolecular G-quadruplex structure
386 is required for binding of telomeric repeat-containing RNA to the telomeric protein

- 387 TRF2. *J. Am. Chem. Soc.* **134**, 11974–11976 (2012).
- 388 44. Giraldo, R. & Rhodes, D. The yeast telomere-binding protein RAP1 binds to and
389 promotes the formation of DNA quadruplexes in telomeric DNA. *EMBO J.* **13**, 2411–
390 2420 (1994).
- 391 45. Bacolla, A., Pradhan, S., Roberts, R. J. & Wells, R. D. Recombinant Human DNA
392 (Cytosine-5) Methyltransferase. *J. Biol. Chem.* **274**, 33002–33010 (1999).
- 393 46. Sanz, L. A. *et al.* Prevalent, Dynamic, and Conserved R-Loop Structures Associate with
394 Specific Epigenomic Signatures in Mammals. *Mol. Cell* **63**, 167–178 (2016).
- 395 47. Ginno, P. A., Lott, P. L., Christensen, H. C., Korf, I. & Chédin, F. R-Loop Formation Is a
396 Distinctive Characteristic of Unmethylated Human CpG Island Promoters. *Mol. Cell* **45**,
397 814–825 (2012).
- 398 48. Wang, X. *et al.* Targeting of Polycomb Repressive Complex 2 to RNA by Short Repeats
399 of Consecutive Guanines. *Mol. Cell* **65**, 1056–1067.e5 (2017).
- 400 49. Wang, X. *et al.* Molecular analysis of PRC2 recruitment to DNA in chromatin and its
401 inhibition by RNA. *Nat. Struct. Mol. Biol.* (2017). doi:10.1038/nsmb.3487
402

403 **Figure 1. G4 formation is associated with hypomethylation at CGIs**
404 **a)** A G-tetrad stabilized by Hoogsteen hydrogen bonding and a central monovalent cation
405 (left). Schematic representations of a three-tetrad G4 structure (Right). **b)** Venn diagram
406 illustrating the overlap of G4 structure formation (BG4 peaks) and CGIs. **c)** Violin plot
407 showing size distribution of BG4 peaks and CGIs. **d)** Count of BG4 peaks overlapping a CGI. **e)**
408 Box and whisker plot showing the average methylation for BG4 peaks (n = 8,210), DHSs (n =
409 142,115) and CGIs (n = 27,073). Centre line represents the median value separating upper
410 and lower quartiles in the box, whiskers indicate 1.5× interquartile range (IQR), points are
411 actual values of outliers. Note that methylation level at CpG sites with less than 5x coverage
412 is considered unreliable and discarded. **f)** Histogram showing the distribution of BG4 peaks
413 and CGIs relative to percentage of GC. **g)** Histogram showing the distribution of BG4 peaks
414 and CGIs relative to percentage of CpGs per 100 bp. **h)** Box and whisker plot showing the
415 methylation levels for BG4 peaks and CGIs at different CpG densities. Note that by definition
416 there are no CGIs at a CpG density < 5 CpGs/100bp and that at > 20 CpGs/100bp there are
417 few CGIs (1) and BG4 (36) peaks to consider. The number of CGI regions and BG4 peaks in
418 each category are presented on top of the plot. **i)** An IGV screen shot illustrating the co-
419 incidence of BG4 peaks (blue) with hypomethylated promoter CGIs (green) and DHSs (orange)
420 for a representative genome region from Chr 7. Shown are normalised signal. Whole genome
421 bisulfite sequence tracks are in black (top). RefGene tracks are in grey (bottom).

422

423 **Figure 2. DNMT1 is recruited to BG4 peaks associated with low methylation**

424 **a)** An IGV screen shot showing the co-incidence (blue-masked) of BG4 peaks (blue) with
425 DNMT1 ChIP-seq peaks (red) and CGIs (green) at hypomethylated region from Chr 6. Orange-
426 masked regions are hypermethylated and enriched with DNMT1 presence, but not BG4
427 signal. Whole genome bisulfite sequence tracks in black (top). **b)** Binding profile of DNMT1
428 shown across CGIs with low (less than 20%, n = 16,523), intermediate (between 20% and 80%,
429 n = 6,042) and high (more than 80%, n = 4,266) methylation. Y-axis shows the number of
430 reads in the ChIP normalised to 1 of sequencing depth (also known as Reads Per Genomic
431 Content (RPGC), more details in computational methods). Replicate 1 and 2 are indicated in
432 red and blue respectively. Above each plot is a heat map showing the enrichment of BG4
433 peaks and DHSs across the respective regions. The heat maps show RPGC of active chromatin
434 marks (DHSs) and BG4 peaks on these three classes of CGIs.

435

436 **Figure 3. DNMT1 selectively binds and is inhibited by G4 structures**

437 **a-f)** ELISA assays testing the binding of recombinant DNMT1 to G4 structure and control
438 oligonucleotides. Binding curves for: **a)** BCL2 G4 and non-G4-forming control (BCL2-mut); **b)**
439 KIT2 G4 and non-G4-forming control (KIT2-mut); **c)** MYC G4 and non-G4-forming control
440 (MYC-mut); **d)** BCL2 duplex DNA; **e)** BCL2 hemi-methylated duplex DNA; **f)** poly(dI-dC), 100 nt.
441 Absorbance was measured at 450 nm. a.u., arbitrary unit. Sequences of oligonucleotides are
442 given below the graphs. **g)** Binding curve of BCL2 G4 in presence of different concentration of
443 BCL2 duplex or poly(dIdC)_n. **h-j)** Relative methylation activity of recombinant DNMT1 in
444 presence of G4 structure and control oligonucleotides: **h)** BCL2 G4 and BCL2-mut; **i)** KIT2 G4
445 and KIT2-mut; **j)** MYC G4 and MYC-mut. Shown are mean ± s.d., n = 3 independent
446 experiments in all plots but **g** (n = 2).

447

448 **Figure 4. Recruitment of DNMT1 by G4 structures shapes the methylome in G-rich regions**

449 **a)** Plot showing the average methylation profile centred around G4 forming regions (red and
450 blue are replicates 1 and 2 respectively, n = 7,491) or G4 sequences without structure
451 (orange and green are replicates 1 and 2 respectively, n = 36,015). The plot extends ± 5 Kb
452 from the centre. The dotted line denotes the lowest methylation level of G4 sequence
453 without structure. **b)** Proposed model for potential involvement of G4 structures and
454 methylation control at CGIs: i) G4 structures sequester DNMT1 due to high affinity binding; ii)
455 G4 structures inhibit the methylation activity of DNMT1. Together with the binding of
456 transcription factors, G4 structures contribute to protection of CGIs from methylation.

457

458 **Online Methods**

459

460 **Cell culture**

461 Mycoplasma-free human chronic myelogenous leukaemia K-562 cells (CCL-243) were
462 purchased from ATCC and grown in RPMI1640 (Glutamine plus, Life Technologies)
463 supplemented with 10% of fetal bovine serum and 100 U/ml penicillin-streptomycin (Life
464 Technologies). All cell stocks were regularly tested for mycoplasma contamination.

465

466 **G-quadruplex CHIP-seq**

467 ChIP-seq for G-quadruplex structures (G4-ChIP-seq) was performed using the G4-specific
468 antibody BG4 essentially as described previously³¹.

469

470 **Oligonucleotide annealing**

471 All oligonucleotides were PAGE purification quality (Sigma). For G4 formation, 10 μ M DNA
472 oligonucleotide was annealed in 10 mM Tris HCl, pH 7.4, 100 mM KCl by heating at 95 °C for
473 5 min followed by gradually cooling to 21 °C. For double-stranded DNA, 10 μ M forward and
474 reverse strand oligonucleotides were mixed and annealed in 10 mM Tris HCl, pH 7.4, 100 mM
475 NaCl in the same manner. 20 μ M poly(dI-dC)₅₀ was annealed as for double-stranded DNA.

476

477 **Enzyme-linked immunosorbent assay (ELISA)**

478 ELISAs for binding affinity and specificity were performed as described previously²⁸ with
479 minor modifications. Briefly, biotinylated oligonucleotides were bound to Pierce™
480 Streptavidin Coated High Capacity Plates (ThermoFisher) followed by blocking with 1.5% BSA
481 and incubation with recombinant full-length human FLAG-tagged DNMT1 protein (Active
482 Motif, Cat. No: 31404) in ELISA buffer (100 mM KCl, 50 mM KH₂PO₄, pH7.4). After three
483 washes with ELISA buffer, detection was achieved with an anti-FLAG horseradish peroxidase
484 (HRP)-conjugated antibody (ab1238, Abcam) and TMB (3,3',5,5'-tetramethylbenzidine) ELISA
485 Substrate (Fast Kinetic Rate, ab171524, Abcam). Signal intensity was measured at 450 nm on
486 a PHERAstar microplate reader (BMG Labtech). Dissociation constants (Kd) were calculated
487 from saturation binding curves assuming one-site binding using Prism (GraphPad Software
488 Inc.). Standard error of mean (s.e.m.) values were calculated from three replicates.

489

490 **In vitro DNA methylation assay**

491 DNA methylation assays were performed using a DNMT Activity Assay Kit (Fluorometric,
492 ab113468, Abcam) as per manufacturer's instructions. Briefly, 100ng recombinant DNMT1
493 was incubated with substrate assay wells in presence of different concentrations of G4 or
494 non-G4 oligonucleotides at 37 °C for 90 min. Methylation levels were quantified from the
495 binding of an anti-5-methylcytosine antibody detected by fluorescent secondary antibody.
496 Fluorescence signal was measured using a PHERAstar microplate reader (530 nm excitation,
497 590 nm emission). DNMT enzyme activity is proportional to the fluorescence intensity (RFU,

498 relative fluorescence unit) measured. Relative methylation activity is then calculated against
499 mock control.

500

501 **Circular dichroism spectroscopy**

502 CD spectra were recorded on an Applied Photo-physics Chirascan circular dichroism
503 spectropolarimeter using a 1 mm path length quartz cuvette. CD measurements were
504 performed at 298 K over a range of 220-300 nm using a response time of 0.5 s, 1 nm pitch
505 and 0.5 nm bandwidth. The recorded spectra represent a smoothed average of three scans,
506 zero-corrected at 300 nm (Molar ellipticity θ is quoted in 105 deg cm² dmol⁻¹). The
507 absorbance of the buffer was subtracted from the recorded spectra. Oligonucleotides were
508 dissolved in lithium cacodylate buffer (100 mM, pH 7.2) containing 100 mM of KCl and 1 mM
509 EDTA to the concentration of 10 μ M. 200 μ L of the oligonucleotides were annealed prior
510 measurement by warming up to 90 °C and slowly cooling down at room temperature.

511

512 **UV Melting**

513 For UV melting experiments, measurements were collected using a Varian Cary 100-Bio
514 UV-visible spectrophotometer by following absorbance at 295 nm. Samples (200 μ l) with
515 final concentration of 2 μ M were measured in black, small window, 1 cm path-length quartz
516 cuvettes, covered with a layer of mineral oil (50 μ l). Samples were equilibrated at 5 °C for 10
517 min, heated to 95 °C and cooled back to 5 °C at a rate of 0.5 °C/min. The samples were held
518 for a further 10 min and then the 5 °C to 95 °C ramp was repeated. Data were recorded every
519 1 °C during both the melting and cooling steps. 200 μ L of oligonucleotides were annealed
520 prior measurement by warming up to 90 °C and slowly cooling down at room temperature.

521

522 **Bioinformatics Software and Scripts**

523 Bioinformatic data analyses and processing were performed using Perl, Bash, Python and R
524 programming languages. The following tools were also used: cutadapt (1.15)⁵⁰, BWA
525 (0.7.15)⁵¹, Picard (2.8.3), (<http://broadinstitute.github.io/picard>), MACS (2.1.1)⁵², Bedtools
526 (2.26.0), (<http://bedtools.readthedocs.io/en/latest/content/overview.html>), Deeptools
527 (2.5.1)⁵³ and Bismark (v0.19.0)⁵⁴.

528 All scripts and software developed are released in the following GitHub page:

529 <https://github.com/sblab-bioinformatics/dna-g4-methylation-dnmt1>

530

531 **G4-ChIP-seq analysis**

532 Raw fastq reads from G4-ChIP-seq in K562 cells were trimmed with cutadapt⁵⁰ to remove
533 adapter sequences and low-quality reads (mapping quality < 10). Reads were aligned to the
534 human genome (version hg19) with BWA⁵¹ and duplicates were removed using Picard. Peaks
535 were called by MACS2⁵² ($p < 10^{-5}$) following our previous work³¹: [https://github.com/sblab-](https://github.com/sblab-bioinformatics/dna-secondary-struct-chrom-lands/blob/master/Methods.md)
536 [bioinformatics/dna-secondary-struct-chrom-lands/blob/master/Methods.md](https://github.com/sblab-bioinformatics/dna-secondary-struct-chrom-lands/blob/master/Methods.md)

537 Peaks were merged from different replicates with bedtools multiIntersect. Only peaks
538 overlapping in 3 out 5 replicates were considered high-confidence. K562 datasets for DHSs

539 (ENCSR000EPC), DNMT1 ChIP-seq (ENCSR987PBI), whole genome bisulfite sequencing
540 (ENCSR765JPC) were downloaded from ENCODE. Promoters were defined as 1 kb (+/-) from
541 the transcription start sites of 31,239 hg19 transcripts. Methylation levels at CpG sites with
542 less than 5x coverage were discarded. If not otherwise specified, CGI³⁴ were downloaded
543 using the UCSC's table browser and then ported to human genome release hg38 using the
544 batch coordinate conversion (liftover) tool of the UCSC. The alternative CGI sets were
545 generated using CpGCluster³⁶.

546
547 **Enrichment analysis**
548 ENCODE DHS and ChIP-seq data sets were normalised to sequencing depth of 1 (i.e. RPGC,
549 Reads Per Genomic Content). Sequencing depth is defined as: (total number of mapped
550 reads * fragment length) / effective genome size. The effective genome size was set to be
551 3,209,286,105 and enrichment values for DHSs and BG4 peaks over CGIs and their flanking
552 sequences were visualised in R using ggplot2 library. Enrichment values for DNMT1 over CGIs
553 and their flanks were visualised with DeepTools⁵³.

554
555 **Monte Carlo Simulation**
556 Monte Carlo simulation was used to calculate the significance of overlap between BG4 peaks
557 and high confidence DNMT1 peaks, defined by Irreproducible Discovery Rate (IDR) in
558 ENCODE's ChIP pipeline. We first counted how many BG4 peaks overlapped with OQs in
559 open chromatin (defined as all OQs seen potassium and/or PDS conditions (749,339
560 sequences)²⁷, which overlap at least one DHS region (43,506 sequences)). We then randomly
561 selected the same number of OQs from all OQs in open chromatin and counted how many
562 overlapped with at least one high confidence DNMT1 peak. The Monte Carlo P-value was
563 then calculated as $(N+1)/(M+1)$, where M is the number of iterations and N is the number of
564 times the same or more overlaps were observed between randomised OQs and high
565 confidence DNMT1 peaks (compared to the number of overlaps observed between BG4
566 peaks and high confidence DNMT1 peaks). Randomisation was repeated for 8000 times and
567 on average the number of overlaps between the shuffled OQs and DNMT1 were two-fold
568 less than those observed between BG4 and DNMT1 peaks.

569
570 **Differential methylation and BG4 binding analysis of entinostat treated HaCaT cells**
571 HaCaT cells were treated with 10 μ M entinostat for 48 hours as we previously described³¹.
572 Genomic DNA from untreated and treated cells were extracted with phenol/chloroform. 50
573 ng DNA were used to generate whole genome bisulfite sequencing libraries using Pico
574 Methyl-Seq Library Prep Kit from Zymo research. Libraries were sequenced using the pair-
575 end 150 bp high-output kit on Illumina Next-seq platform. Data from 4 runs were pooled
576 together. After quality assessment using FastQC
577 (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>), reads were processed to
578 remove adaptors and low-quality bases using cutadapt³. Options $-u 6 -u -1 -U 6 -U -1$ were
579 used to trim the initial six and last nucleotide bases. High-quality reads were aligned using

580 Bismark in paired-end mode with options --non_directional, --unmapped and -N 0 to hg19
581 reference genome. Reads were then de-duplicated and methylation was extracted. To
582 increase mapping efficiency and following previous work⁵⁵, unmapped reads resulting from
583 the paired-end alignments were then re-aligned in single-end mode with options --
584 non_directional and -N 0, and then deduplicated. Methylation was extracted for paired-end
585 and single-end alignments separately and then aggregated. Technical replicates for each
586 condition (before and after entinostat treatment) were merged and methylation counts were
587 aggregated by CpG site. A threshold was then applied to keep CpG sites with more than 5X
588 bisulfite sequencing depth both before and after treatment. This resulted in 21,106,307 CpG
589 sites, 75% of all CpG sites in hg19.

590 Differential BG4 binding analysis was done as previously reported³¹. Analysis focused on
591 open chromatin promoter regions (5351) which have at least one G-rich sequence²⁷ and
592 ATAC-seq peak unaltered in size (log₂ fold change = -0.6 to 0.6, FDR < 0.05) between
593 untreated and entinostat-treated HaCaT cells. Depending on differential BG4 signal, these
594 regions were categorized into BG4 gain (> 1.5-fold change in signal and FDR < 0.05) and BG4
595 constant and BG4 negative. 95% (5072/5351) of these regions are overlapping with CGIs.
596 Difference of the percentage methylation of the overlapping CGIs in each of the categories
597 were calculated. Statistic test was done with Mann–Whitney U test. Plotting of methylation
598 data were performed in the R programming language.

599
600 **Data availability.** K562 datasets for DHS (ENCSR000EPC), DNMT1 ChIP-seq (ENCSR987PBI)
601 and whole genome bisulfate sequencing (ENCSR765JPC) were downloaded from ENCODE.
602 G4-ChIP-seq data sets for K562 and WGBS datasets for entinostat-treated and untreated
603 HaCaT cells are available at the NCBI GEO repository under accession number GSE107690.
604 G4-ChIP-seq data in entinostat-treated and untreated HaCaT cells were taken from
605 GSE76688. Source data for figure 1d, e, h and Figure 3 are available with the paper online.

606
607 **References**

608 50. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing
609 reads. *EMBnet.journal* **17**, 10 (2011).

610 51. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
611 *Prepr. <https://arxiv.org/abs/1303.3997>* **00**, 1–3 (2013).

612 52. Zhang, Y. *et al.* Model-based Analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).

613 53. Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data
614 analysis. *Nucleic Acids Res.* **44**, W160–W165 (2016).

615 54. Krueger, F. & Andrews, S. R. Bismark: A flexible aligner and methylation caller for
616 Bisulfite-Seq applications. *Bioinformatics* **27**, 1571–1572 (2011).

617 55. Peat, J. R. *et al.* Genome-wide Bisulfite Sequencing in Zygotes Identifies Demethylation
618 Targets and Maps the Contribution of TET3 Oxidation. *Cell Rep.* **9**, 1990–2000 (2014).

