



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Transition to naïve human pluripotency mirrors pan-cancer DNA hypermethylation

Citation for published version:

Patani, H, Rushton, MD, Higham, J, Teijeiro, SA, Oxley, D, Cutillas, P, Sproul, D & Ficz, G 2020, 'Transition to naïve human pluripotency mirrors pan-cancer DNA hypermethylation', *Nature Communications*.
<https://doi.org/10.1038/s41467-020-17269-3>

Digital Object Identifier (DOI):

<https://doi.org/10.1038/s41467-020-17269-3>

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Nature Communications

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 Transition to naïve human pluripotency mirrors pan-cancer DNA hypermethylation

2

3 Patani H¹, Rushton MD¹, Higham J², Teijeiro SA¹, Oxley D³, Cutillas P¹, Sproul D², Ficz G^{1*}

4 ¹ Centre for Haemato-Oncology, Barts Cancer Institute, Queen Mary University of London, London, UK, EC1M 6BQ

5 ² MRC Human Genetics Unit and Edinburgh Cancer Research Centre, MRC Institute of Genetics & Molecular
6 Medicine, University of Edinburgh, Western General Hospital, Crewe Road, Edinburgh, UK, EH4 2XU

7 ³ Mass Spectrometry Facility, Babraham Institute, Cambridge, UK, CB22 3AT

8 * Correspondence: g.ficz@qmul.ac.uk (G.F.)

9

10

11 Abstract

12

13 Epigenetic reprogramming is a cancer hallmark but how it unfolds during early neoplastic
14 events and its role in carcinogenesis and cancer progression is not fully understood. Here
15 we show that resetting from primed to naïve human pluripotency results in acquisition of a
16 DNA methylation landscape mirroring the cancer DNA methylome, with gradual
17 hypermethylation of bivalent developmental genes. We identify a dichotomy between
18 bivalent genes that do and do not become hypermethylated, which is also mirrored in
19 cancer. We find that loss of H3K4me3 at bivalent regions is associated with gain of
20 methylation. Additionally, we observe that promoter CpG island hypermethylation is not
21 restricted solely to emerging naïve cells, suggesting that it is a feature of a heterogeneous
22 intermediate population during resetting. These results indicate that transition to naïve
23 pluripotency and oncogenic transformation share common epigenetic trajectories, which
24 implicates reprogramming and the pluripotency network as a central hub in cancer
25 formation.

26

27 Introduction

28

29 Disruption of DNA methylation patterns is a hallmark of human cancers, typically
30 characterised by loss of global genomic DNA methylation accompanied by site-specific
31 hypermethylation¹⁻⁴. DNA hypomethylation is typically associated with genomic
32 instability^{5,6}, while site-specific DNA hypermethylation occurs at promoter CpG islands
33 (CGIs) and can be associated with repression of tumour suppressor genes in cancer cells^{7,8}.
34 The majority of these observations have been made in cancer cell lines or primary cancer
35 cells, but they are not fully representative of the processes occurring during the transition of
36 normal cells into malignant cells. The underlying mechanisms that give rise to these
37 opposing patterns of genomic DNA methylation in early stages of human cancer
38 development remain elusive, as does the timing and biological function of such events in
39 cancer initiation and progression. To this end, a recent study demonstrated that ageing and
40 cancer associated DNA hypermethylation accelerates cellular transformation in a *Braf*^{V600E}
41 mouse colon organoid system, through suppression of Wnt signalling regulators in a
42 progressive manner⁹. This study functionally links promoter CGI hypermethylation with
43 oncogenic transformation, demonstrating a causal relationship. Nevertheless, how *de novo*
44 DNA methyltransferase activity is preferentially targeted to specific regions of the genome
45 in the context of aberrant cancer methylation remains largely a mystery.

46 It has been hypothesised that cancer cells follow an evolutionary trajectory towards a stem
47 cell state, which allows both self-renewal and differentiation¹⁰, and more recently, cancer-
48 related mutations have been identified in naïve human embryonic stem cells (hESCs)¹¹.

49

50 Here we identify cancer-like DNA methylation changes during primed to naïve hESC
51 resetting using the recently developed NANOG/KLF2 overexpression + 2iL+Gö method
52 (comprising two small-molecule inhibitors of MEK and GSK3 β , human recombinant
53 leukaemia inhibitory factor, and a pan-PKC inhibitor)¹². Our system provides a unique
54 opportunity to investigate the mechanism of DNA hypermethylation in human cells in a
55 temporal manner, and sheds light on the role of the transcription factor and pluripotency
56 networks in driving cancer-like DNA hypermethylation.

57

58 **Results**

59

60 **Naïve resetting induces CGI promoter hypermethylation**

61

62 To investigate the kinetics of the changing DNA methylation landscape between primed and
63 naïve hESCs, we transitioned primed hESCs to the naïve state as previously described¹², by
64 inducing *NANOG/KLF2* transgenes with doxycycline. We also captured the two intermediary
65 states, termed 'early transition' and 'late transition' when the cells are in 2iL+dox or 2iL+Gö,
66 respectively (Fig. 1a). We see global DNA demethylation of the genome in naïve cells as
67 reported previously¹², measured by the Infinium MethylationEPIC array and mass
68 spectrometry (Supplementary Figs 1a-1c). The loss of 5-methylcytosine (5mC) is gradual and
69 is accompanied by the loss of its oxidation product, 5-hydroxymethylcytosine (5hmC)
70 (Supplementary Fig. 1a). Interestingly, while the majority of the genome is demethylated,
71 we observe hypermethylation of a subset of CpGs (an increase of >10% methylation
72 compared to primed hESCs), exemplified by the *HOXA* cluster (Figs 1b, 1c, Supplementary
73 Fig. 1d). This gain in methylation is evident as cells go through the early transition of
74 resetting, with a peak of hypermethylated CpGs as the cells go through the late transition of
75 resetting (Figs 1b, 1c). Although the peak of hypermethylation coincides with the cells being
76 transitioned into 2iL+ Gö conditions, the abundance of hypermethylation is independent of
77 the addition of Gö (Supplementary Fig. 1e), indicating a time-dependent accrual of DNA
78 methylation instead. As the cells stabilise in the naïve state, we observe maintenance of a
79 proportion of hypermethylated sites, while some CpGs show only a transient gain in
80 methylation (Fig. 1b). The reproducibility of the hypermethylation during the resetting
81 process is apparent from the strong overlap between hypermethylated sites across
82 biologically independent MethylationEPIC arrays (with 2 or 3 cell populations assayed within
83 each array) and when compared to published whole genome bisulfite sequencing (WGBS)
84 data, suggesting that the site-specific gain in methylation is not random, and likely has a
85 biological function (Supplementary Figs 1f, 1g). Moreover, as primed hESCs and hESCs
86 during the early transition of resetting proliferate and cycle at comparable rates as
87 measured by loss of bromodeoxyuridine (BrdU), the site-specific gain in methylation upon
88 resetting is the result of an active process rather than the selection of an existing
89 subpopulation of cells (Supplementary Fig 1h).

90

91 We next sought to investigate the genomic context with which hypermethylation occurs.
92 We utilised the Encyclopaedia of DNA elements (ENCODE) ChIP-seq datasets for the H1
93 primed hESC cell line and overlapped them with resetting-associated hypermethylated
94 probes. We observed that hypermethylated probes are enriched within regions marked by
95 H3K4me1/2/3 and H3K27me3 in primed hESCs (Fig. 1d). The majority of these fall within
96 regions marked by bivalent histone modifications, defined by co-occurrence of H3K4me3
97 and H3K27me3 (Fig. 1e). Bivalency typically marks regulatory regions (promoters and
98 enhancers with overlapping transcription factor binding sites)¹³, and as expected, we saw a
99 striking overlap with CGIs and regulatory regions, which is not the case for hypomethylated
100 probes (Figs 1f, 1g). In addition, the majority of the hypermethylated probes reside within
101 ChromHMM predicted poised promoters (Fig. 1h). To validate our results, we also re-
102 analysed published whole-genome bisulfite sequencing data for primed and naïve hESCs¹²,
103 and identified 26,625 regions (300bp each) that were hypermethylated (>10% increase in
104 naïve vs primed) in naïve hESCs. In agreement with our Illumina MethylationEPIC array data,
105 we observed enrichment of these regions within loci marked by bivalency in primed hESCs,
106 with a bias towards regulatory regions and CGIs (Supplementary Figs 2a-2e). This strongly
107 reinforces the highly reproducible nature of the DNA hypermethylation that occurs upon
108 resetting of primed hESCs to naïve pluripotency.

109

110 To eliminate the possibility that the hypermethylation is simply an artefact of the
111 NANOG/KLF2 + 2iLGö *in vitro* resetting system, we compared 2iLGö naïve hypermethylated
112 regions to hypermethylated regions identified in naïve cells generated using the alternative
113 methods of naïve hESC generation and those in the human inner cell mass (ICM)¹⁴⁻¹⁷. We
114 saw a significant overlap between the hypermethylated regions in each of these data sets
115 compared to the 2iLGö naïve cells, in both cases also enriched within H3K4me1/2/3 and
116 H3K27me3 regions in primed hESCs (Supplementary Figs 3a-3g). Interestingly, we do not
117 detect any hypermethylation in mouse embryonic stem cells cultured in 2i compared to
118 those cultured in serum¹⁸, and the hypermethylation present in the mouse ICM¹⁹ is far less
119 extensive and does not enrich at bivalent regions (Supplementary Figs 3h-3j). Overall from
120 these results, we can conclude that the hypermethylation pattern is a feature of *in vitro*
121 naïve human pluripotency and recapitulates the *in vivo* relationship between the ICM and
122 the post-implantation embryo.

123

124 **Developmental genes are hypermethylated and repressed**

125

126 To identify the genes targeted by hypermethylation, we performed gene ontology (GO)
127 analysis of the genes and classified a gene as hypermethylated if it possessed a
128 hypermethylated probe/ region within 1500bp upstream of its transcription start site (TSS).
129 GO analysis revealed an extensive enrichment of hypermethylated genes in developmental
130 pathways, particularly neuronal development, while hypomethylated genes show much
131 weaker enrichment in pathways involved in cell cycle and metabolism (Fig. 2a,
132 Supplementary Fig. 4a). To investigate whether hypermethylation is associated with gene
133 silencing, we performed temporal transcriptome analysis of cells during the transition from
134 primed to naïve pluripotency and observed thousands of genes differentially expressed at
135 each stage of the transition compared to primed hESCs, with an enrichment of
136 downregulated genes in developmental pathways (Supplementary Figs 4b, 4c). We

137 observed that the average expression of genes which undergo hypermethylation is
138 attenuated in naïve compared to primed hESCs, but that hypermethylated genes are
139 characterised by low average expression in primed hESCs (Figs 2b, 2c). Hypermethylation
140 may play a functional role in these cells by contributing to downregulation of developmental
141 pathways, perhaps enhancing pluripotency and blocking differentiation by providing a more
142 stable gene repression mechanism¹³. Intriguingly, we observe a subset of genes, notably the
143 HOX gene family, that become hypermethylated and are upregulated upon hESC resetting
144 (Fig. 2d), despite no evidence of 5mC oxidation to 5hmC at their promoters (Fig. 2e).
145 Interestingly, the methylation dynamics of these HOX genes do not differ from genes that
146 are hypermethylated and downregulated (Fig. 2c, 2d, Supplementary Fig. 4d). This suggests
147 that there may be population heterogeneity during resetting, whereby a subpopulation of
148 cells which do not undergo HOX promoter hypermethylation exhibit upregulation of the
149 genes.
150

151 **DNMT3A controls early *de novo* methylation**

152
153 We next sought to identify the epigenetic regulators that are responsible for the deposition
154 of *de novo* DNA methylation. Of the *de novo* DNMT3 family of DNA methyltransferases,
155 *DNMT3B* is highly expressed but is transiently downregulated upon resetting. The mRNA
156 level of *DNMT3A* is transiently upregulated (Supplementary Fig. 5a), though this is not
157 reflected in the protein level (Supplementary Fig. 5g). The catalytically inactive *DNMT3L* is
158 upregulated (Supplementary Fig. 5a) and considered a marker of naïve pluripotency²⁰. We
159 generated constitutive knockdown primed hESC cell lines using two short hairpin RNAs
160 (shRNAs) targeting each of the three genes, as well as one specifically targeting the long
161 isoform of *DNMT3A*, known as *DNMT3A1* (Supplementary Fig. 5b). We subjected each of
162 the cell lines to resetting until the early transition, at which stage hypermethylation is
163 already detectable, and thereafter to the late transition. In the early transition, knockdown
164 of *DNMT3B* and *DNMT3L* had little impact on the level of methylation (Fig. 3a,
165 Supplementary Fig. 5c). Knockdown of *DNMT3A*, however, was able to abolish
166 hypermethylation (Fig. 3a). A recent study demonstrated isoform-specific recruitment of
167 DNMT3A1 to bivalent CGIs in mouse embryonic stem cells²¹. However, specific knockdown
168 of *DNMT3A1* had no impact on the level of hypermethylation (Supplementary Fig. 5d),
169 suggesting that the more dominantly expressed *DNMT3A2* carries out *de novo* methylation
170 early during resetting, independently of DNMT3L. It is worth noting, however, that during
171 the late transition of resetting, both DNMT3A and DNMT3B knockdown cells show a partial
172 reduction in methylation compared to the control cells (Fig. 3b). This indicates that
173 DNMT3B, which by this stage is transcriptionally expressed at a higher level than during the
174 early transition (Supplementary Fig. 5a), contributes to hypermethylation along with
175 DNMT3A. Additionally, when we reset either DNMT3A or DNMT3B knock down cells to the
176 naïve state, we see a reduction in the pluripotency of the cells as measured by their alkaline
177 phosphatase activity as well as reduced expression levels of naïve pluripotency genes in
178 DNMT3A but not DNMT3B knock down cells (Figs 3c, 3d, Supplementary Fig. 5e). As both
179 DNMT3A and DNMT3B seem to contribute to hypermethylation during the late transition of
180 resetting, we additionally generated a double knock down cell line of the two genes in
181 primed hESCs and subjected the cells to resetting to the naïve state (Supplementary Fig 5f).
182 The DNMT3A/DNMT3B double knock down cells show significantly reduced levels of

183 hypermethylation during the early transition as well as in naïve hESCs, and they exhibit
184 reduced pluripotency as measured by their alkaline phosphatase activity and expression
185 levels of naïve pluripotency genes, though TFCP2L1 and KLF4 do not change significantly
186 (Figs 3e-g). Collectively, this indicates a putative role of the *de novo* methyltransferases in
187 stabilization of the naïve pluripotent state.

188

189 Aside from the *de novo* methyltransferases, we also hypothesised that loss of ten-eleven
190 translocation (TET) enzymes may be the cause of bivalent CGI hypermethylation^{22,23}. TET1 is
191 expressed in primed hESCs and subsequently downregulated at the protein level as hESCs
192 progress through the early transition of resetting (Supplementary Figs 6a, 6b). We
193 generated TET1-overexpressing primed hESCs (Supplementary Figs 6c, 6d) and subjected
194 them to resetting until the early transition. We measured DNA methylation at selected
195 target loci but observed no change in the levels of DNA methylation, though TET1 is
196 modestly overexpressed (Supplementary Fig. 6e). These data indicate that
197 hypermethylation of bivalent loci upon resetting is independent of TET1 loss.

198

199 To determine whether *de novo* methylation is strictly correlated with the acquisition of
200 naïve pluripotency during resetting, we used a published cell-surface marker, SUSD2, that
201 has been shown to separate naïve from primed hESCs during resetting²⁴, and is able to
202 identify increasing numbers of naïve hESCs during the transition to the naïve state
203 (Supplementary Fig. 7a). We sorted SUSD2+ and SUSD2- cells during the early and late
204 transition of resetting and used a targeted approach to measure DNA methylation at
205 selected resetting -associated hypermethylated loci (Fig. 4a). We observed no significant
206 difference in the level of hypermethylation in SUSD2+ and SUSD2- fractions, though SUSD2+
207 cells displayed higher expression of naïve pluripotency genes (Figs. 4b, 4c, Supplementary
208 Fig. 7b). In conjunction with the observation that reduced expression of *de novo*
209 methyltransferases during resetting impacts naïve pluripotency, this indicates that while *de*
210 *nov*o methylation may be important for the stability of the naïve pluripotent state, it is
211 insufficient for the acquisition of naïve pluripotency. Collectively, this suggests that the
212 hypermethylation is more a feature of the heterogeneous intermediate population of
213 partially reset cells.

214

215 **Bivalent CGIs that lose H3K4me3 gain DNA methylation**

216

217 To investigate the mechanism of bivalent CGI hypermethylation, we first classified bivalent
218 regions as those possessing H3K4me3 and H3K27me3 peaks in primed hESCs utilising the
219 ENCODE ChIP-seq data-set. We observed that only 41% of these loci gain DNA methylation
220 upon resetting (Fig. 5a), indicating that the presence of bivalent chromatin is not sufficient
221 for acquiring *de novo* DNA methylation upon resetting. We therefore divided all bivalent
222 regions in primed hESCs into those that do and do not gain methylation during resetting
223 (hypermethylated regions defined as >10% increase in naïve vs primed hESC). Taking the
224 nearest gene to each region (within 1500bp of the TSS), GO analysis of the bivalent
225 hypermethylated group showed a strong enrichment for developmental pathways, while
226 the bivalent non-hypermethylated group showed much lower enrichment of other
227 biological processes (Fig. 5a). This points towards common regulation of the developmental
228 genes that exhibit hypermethylation during resetting.

229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248

We hypothesized that there are intrinsic differences between bivalent regions that do and do not gain methylation, which both begin with similar chromatin states. As DNA methylation and H3K4me3 are known to be mutually exclusive²⁵, we performed CHIP-qPCR of H3K4me3 at bivalent DNA regions, across the time course of resetting. We observed a loss of H3K4me3 at bivalent regions that become hypermethylated, whilst bivalent non-hypermethylated regions retain their levels of H3K4me3 (Fig. 5b). In contrast, the levels of H3K27me3 exhibit little change, despite the presence of DNA methylation (Fig. 5c). While H3K27me3 and DNA methylation are considered to be mutually exclusive at CpG rich regions during development²⁶, co-existence of the two modifications has previously been reported²⁷. It is known that loss of H3K4me3 is permissive to the gain of DNA methylation, but this on its own cannot explain the specific gain of methylation at these regions, as both DNMT3A and DNMT3B possess an ADD domain capable of mediating the interaction of the enzymes with unmethylated H3K4^{25,28-30}, and several loci that have been shown to lose H3K4me3 in naïve hESCs do not undergo hypermethylation¹⁴. Additionally, despite comparable absolute protein levels of DNMT3A and DNMT3B, as measured by mass spectrometry (Supplementary Fig. 5g), only DNMT3A deposits DNA methylation during the early transition of hESC resetting. Moreover, the strong bias towards developmental genes suggests that the gradual hypermethylation is not a stochastic process.

249 **Transcription factors influence hypermethylation**

250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274

Our data led us to hypothesize that an additional player, likely a DNA-binding factor, facilitates DNMT3A-mediated hypermethylation in the early stages of hESC resetting. To understand the unique properties of the bivalent hypermethylated group, we performed differential motif analysis of these regions, with the bivalent non-hypermethylated regions as a control set. We identified a number of motifs corresponding to DNA-binding transcription factors enriched at regions that undergo hypermethylation (Fig. 6a). To test whether these proteins are expressed, we performed total proteomics of primed and early transition hESCs and identified proteins that were upregulated during the early transition compared to primed hESCs (Fig. 6b). Through cross-comparison of the two analyses, we short-listed two candidate transcription factors, SOX15 and NFKB1, which are upregulated during early resetting and show an enrichment of binding sites at hypermethylated regions (Figs 6a, 6b). We identified an additional two candidate transcription factors, FOXC1 and ZFH3 (Fig. 6a), which were transcriptionally upregulated during the early transition based on RNA-seq data but not detected in any samples by proteomic analysis, likely due to technical limitations of the method in detecting nuclear transcription factors³¹. We generated constitutive knockdown cell lines using two shRNAs targeting each of the four candidate genes (Supplementary Fig. 8a) and subjected each of the cell lines to resetting until the early transition. We measured the expression of naïve pluripotency genes to test whether the knock down cells undergo resetting similarly to control cells, and found that their expression is not significantly altered in the knock downs compared to control cells (Supplementary Fig. 8b). Strikingly, upon resetting, knock down of each of the transcription factors reduced the level of hypermethylation at target loci analysed, suggesting that the transcription factor network active during the early transition of resetting is involved in bivalent promoter CGI hypermethylation (Fig. 6c). The impact of each of the knockdowns on

275 DNA methylation is higher in regions where at least one of the highly expressed
276 transcription factors are predicted to bind compared to sites that are not bound by any of
277 the four transcription factors (Fig. 6c). Interestingly, the reduction in methylation is
278 observed in each of the transcription factor knock downs at SOX15 and NFKB1 predicted
279 binding sites, indicating that the impact is not limited to the specific binding sites for each
280 transcription factor (Fig. 6c). This points to a network synergy in preferentially mediating *de*
281 *novo* methylation at these sites. As the reduction in DNA hypermethylation upon
282 transcription factor knockdown is only partial, however, this suggests that additional
283 mechanisms are also at play.

284

285 To test whether signalling changes associated with factors required for induction of the
286 naïve state induction could explain such a mechanism, we conducted hESC resetting until
287 the early transition, each time removing one of these factors. Resetting in the absence of
288 the MEK inhibitor or GSK3 β inhibitor or concomitant removal of both inhibitors still resulted
289 in hypermethylation at target loci analysed (Fig. 6d), suggesting that hypermethylation may
290 be coordinated by the overexpressed *NANOG* and *KLF2* or the associated pluripotency
291 network. Collectively, these data indicate that upon reprogramming hESCs to the naïve
292 state, hypermethylation is driven by the transcription factor network that becomes active
293 upon resetting, and that this is synchronised by the core pluripotency network.

294

295 **Resetting-associated hypermethylation is mirrored in cancer**

296

297 *De novo* DNA methylation of bivalent chromatin in the context of a hypomethylated
298 genome has also been reported in cancer cell lines and primary tumours^{32,33}. To investigate
299 the link between the hypermethylation patterns associated with hESC resetting and the re-
300 emergence of such patterns in cancer, we compared hypermethylated CpGs at each stage of
301 resetting with regions previously identified as hypermethylated in B-cell chronic lymphocytic
302 leukaemia (B-CLL)³⁴. We observed that the most substantial overlap is found between
303 hypermethylated CpGs associated with the late transition of resetting and B-CLL
304 hypermethylation, corresponding to the peak of hypermethylation we observe during
305 reprogramming (Supplementary Fig. 9a). Interestingly, we also see a more significant
306 overlap between hypermethylated regions in B-CLL or colon cancer^{34,35} and resetting -
307 associated hypermethylated CpGs with a low basal methylation level (<5%) in primed hESCs
308 (Supplementary Figs 9b, 9c). These data demonstrate a substantial overlap between naïve
309 stem cell and cancer hypermethylation.

310 We further hypothesised that cancer cells may recapitulate the dichotomy between bivalent
311 genes that do and do not become hypermethylated. To test this, we compared resetting -
312 associated bivalent hypermethylated and non-hypermethylated CpGs with data from the
313 cancer genome atlas (TCGA) pan-cancer atlas^{36,37}. We found a significant gain in methylation
314 between normal and cancer tissue for bivalent CpGs identified as hypermethylated during
315 the resetting process, compared to unmethylated bivalent CpGs (Fig. 7a, Supplementary Fig.
316 10a). This was consistent across all cancer types analysed. In addition, when we defined
317 CpGs that were hypermethylated compared to normal tissues across cancer types, we found
318 that they were significantly enriched in hESC hypermethylated bivalent CpGs (Fisher's test *p*-
319 value < 2x10⁻¹⁶, percentage fold-change 11.13, Fig. 7b). Similar results were also observed
320 when defining hypermethylated CpGs for each of the 9 individual tumour types (Fig. 7b).

321 Together these results indicate that resetting -associated hypermethylation parallels pan-
322 cancer hypermethylation, though the individual CpGs methylated in each cancer type vary³⁸.
323 We see no further enrichment of resetting -associated hypermethylated CpGs with more
324 advanced stages of cancer (Supplementary Fig. 9d), or with datasets derived from
325 metastatic tissues (Supplementary Fig. 9e).

326 To further investigate the parallels between hypermethylation in naïve hESCs and cancer,
327 we asked whether the regions hypermethylated in naïve hESCs were also marked by
328 H3K27me3 in normal human tissues that give rise to cancers. By examining H3K27me3 ChIP-
329 seq from 8 normal tissues, we found that the regions surrounding hypermethylated CpGs
330 had significantly higher levels of H3K27me3 than non-hypermethylated CpGs from naïve
331 hESCs (Fig. 7c, Supplementary Fig. 10b). This included the HOX loci which we had seen to be
332 strongly hypermethylated in naïve hESCs as well as loci which are previously reported as
333 gaining DNA methylation in cancer such as SFRP1 (Fig. 7d)⁸.

334 Together, these results demonstrate that resetting of primed hESCs to naïve pluripotency
335 and cancer follow similar epigenetic trajectories, both involving the acquisition of DNA
336 methylation at developmental gene promoters marked by H3K27me3. There is growing
337 evidence in the literature regarding the acquisition of stem-like properties and expression of
338 pluripotency genes in cancers^{10,39}. This makes it intriguing to speculate that the pluripotency
339 network plays a key role in establishing cancer hypermethylation. Considering the causal
340 relationship between hypermethylation and acceleration to one-step transformation⁹, this
341 warrants further investigation of means to disrupt *de novo* accumulation of DNA
342 methylation at CGI promoters.

343

344

345 Discussion

346

347 We propose the concept that reprogramming events during primed to naïve resetting are a
348 fundamental feature of human cancers and possibly a very early step in cancer evolution.
349 Although naïve stem cells and cancer cells seem developmentally distant, our data supports
350 the hypothesis that cancers follow an evolutionary trajectory towards an embryonic stem
351 cell state, which allows both self-renewal and differentiation¹⁰.

352

353 We demonstrate the dynamic acquisition of DNA methylation, primarily at CGI promoters,
354 upon the transition from primed to naïve pluripotency in hESCs. We observed the highest
355 levels of hypermethylation in the heterogeneous population of cells present during the late
356 transition of resetting. It is worth noting that during the late transition of resetting, there is
357 population heterogeneity as demonstrated by the study of HOX gene promoters, as well as
358 heterogeneous expression of naïve stem cell markers as exemplified by the SUSD2
359 expression, but the hypermethylation and expression of naïve stem cell markers do not
360 strictly co-occur. It seems that hypermethylation is a feature of this heterogeneous
361 intermediate population of partially reset cells, and it is partially maintained as the cells
362 stabilise in the naïve state.

363

364 We demonstrate that the hypermethylation that we observed upon the transition from
365 primed to naïve pluripotency mirrors the frequently observed aberrant hypermethylation in
366 human cancers. Such parallels have been drawn previously in other mammalian species and
367 developmental contexts⁴⁰. However, the data we present here demonstrates

368 hypermethylation conserved across *in vitro* and *in vivo* human pluripotency¹⁷, strengthened
369 by its reproducibility across multiple *in vitro* resetting methods^{14,16}. Moreover, it is notable
370 that we do not observe comparable hypermethylation in the mouse ICM or in *in vitro* mouse
371 ESCs cultured with 2i inhibitors. This observation has potential implications for making
372 inferences with regards to epigenetic processes between species, both in development and
373 in the study of cancer, as has been noted previously⁴¹.

374
375 We have observed that bivalent loci are almost exclusively susceptible to DNA
376 hypermethylation during the transition to naïve pluripotency. However, not all bivalent loci
377 become hypermethylated, suggesting that the presence of bivalent chromatin is not the
378 only factor required for acquiring *de novo* DNA methylation upon resetting. Several studies
379 have demonstrated that the aberrant gains of DNA methylation observed in cancer also
380 occur at H3K27me3-marked loci⁴²⁻⁴⁴. Loci that gain DNA methylation in cancer are also
381 enriched in sets of loci that are bivalent in embryonic stem cells^{33,42,45}. The mechanistic basis
382 of this relationship is unclear. Although DNMTs can interact with EZH2 which is responsible
383 for deposition of H3K27me3⁴⁶, recruitment of DNMT3A by PRC2 is not sufficient to trigger
384 *de novo* DNA methylation⁴⁷. It is particularly noteworthy that bivalent loci that undergo
385 hypermethylation both upon resetting and in the context of cancer belong to
386 developmental pathways³³, distinguishing them from other bivalent loci that do not gain
387 methylation despite having a comparable starting chromatin configuration.

388
389 Our data indicates that DNMT3A is responsible for DNA hypermethylation during the early
390 transition of primed to naïve resetting, and that both DNMT3A and DNMT3B contribute to
391 hypermethylation during the late transition of resetting. At this stage, the reduction in
392 methylation is only partial upon knockdown of either enzyme, therefore alternative factors
393 such as the putative *de novo* activity of DNMT1 cannot be excluded⁴⁸. A knockdown of
394 either DNMT3A or DNMT3B impacted the stability of the naïve state as measured by
395 alkaline phosphatase activity, with a knockdown of DNMT3A additionally disrupting the
396 naïve transcriptional network. Interestingly, a double knockdown of DNMT3A and DNMT3B
397 had an even greater impact on the naïve state as measured by alkaline phosphatase activity,
398 suggesting a synergy between the two *de novo* methyltransferases in influencing either the
399 transition to the naïve state or naïve cell stability. It remains to be deduced whether
400 DNMT3A/DNMT3B impact the naïve state through the hypermethylation they contribute to
401 or through a non-catalytic role. Additionally, in order to better understand the functional
402 impact of resetting-associated hypermethylation, its potential impact on the differentiation
403 potential of the naïve cells remains to be investigated independently of the known roles of
404 *de novo* methyltransferases in stem cell differentiation.

405
406 Our data points towards the transcription factor network established upon resetting playing
407 a role in the targeting or recruitment of DNMT3A to loci that gain methylation. Whilst we
408 cannot currently differentiate between a direct interaction of DNMT3A with transcription
409 factors or an indirect network-driven effect on targeting of the enzyme, loci-specific
410 recruitment of DNMT3A via transcription factors has been previously demonstrated⁴⁹ and *in*
411 *vitro* data supports the ability of DNMT3A to interact directly with numerous transcription
412 factors⁵⁰. Our data are also indicative of the overexpressed NANOG and KLF2 coordinating
413 *de novo* methylation, however studies have shown that KLF2 is not expressed *in vivo* in the
414 human inner cell mass^{51,52}, where we also observe hypermethylation. Additionally, we

415 observe comparable hypermethylation in naïve hESCs generated using two transgene-
416 independent methods of resetting. This collectively suggests that the core pluripotency
417 network, to which NANOG belongs, is likely responsible for coordinating the transcriptional
418 changes that drive DNA hypermethylation. There is growing evidence in the literature
419 regarding the acquisition of stem-like properties and expression of pluripotency genes in
420 cancers^{10,39}. This makes it intriguing to speculate that a transcriptional programme
421 associated with the naïve pluripotency network could drive a shared mechanism of
422 hypermethylation during resetting of primed hESCs to naïve pluripotency and in cancer
423 development, either preceding or in conjunction with genetic mutations.

424

425 Hypermethylation during primed to naïve resetting affects developmental genes whose
426 expression is generally low and is further attenuated upon hypermethylation, as is often
427 observed in cancer^{33,38}. The function of hypermethylation in cancer remains a topic of
428 debate. While several studies have shown clear repressive roles of hypermethylation for
429 individual tumour suppressor genes^{7,53}, it remains less well understood what the purpose of
430 hypermethylation of a large number of loci might be. It has been proposed that aberrant
431 hypermethylation in cancer may act to block cellular differentiation, thus enabling cancer
432 cells to continue to propagate in their more primitive states^{33,44,54}, and this has been
433 experimentally demonstrated in a recent study⁹. During reprogramming of somatic cells to
434 induced pluripotent stem cells (iPS), global DNA demethylation occurs late⁵⁵ and is a
435 bottleneck for efficient reprogramming⁵⁶. Our data indicates that in addition to global DNA
436 demethylation which is efficiently erased in further naïve resetting, gain of DNA methylation
437 in bivalent developmental gene promoters will lock cells in a primitive state. The
438 commonality in methylation patterns across cancer types, each harbouring different driver
439 mutations, suggests that these methylation changes may be regulated by a common
440 overarching mechanism and occur early in tumourigenesis, as has been demonstrated
441 previously in one of the few models of early cancer development⁵⁷. In line with this, the
442 notion that cancer cells follow an evolutionary trajectory towards a stem cell state^{10,39}
443 makes the transition from primed to naïve pluripotency an interesting model to study
444 biological processes such as DNA methylation that likely occur early during cancer initiation,
445 and may be analogous to dedifferentiation. Additional molecular features of the primed to
446 naïve state transition appear analogous to cancer hallmarks⁵⁸, such as altered metabolism¹²,
447 loss of imprints¹⁵, loss of DNA hydroxymethylation⁵⁹ and genomic instability^{11,15,60}. Whether
448 they are related to the changing epigenetic landscape remains unexplored, but further use
449 of this model system may shed light on the emergence of these characteristics during
450 cellular transformation. We propose that naïve resetting may provide a good model system
451 to understand whether other molecular processes associated with cellular reprogramming
452 play a role in tumourigenesis.

453

454 **Methods**

455

456 **Cell lines**

457 WA09/H9 NK2 primed hESCs were kindly provided by Austin Smith¹² with permission from
458 WiCell. All hESCs were cultured on irradiated mouse embryonic fibroblasts (iMEF). iMEFs
459 were seeded at a density of 1×10^6 cells per 6-well plate, in 5% O₂, 7% CO₂ at 37°C in a
460 humidified incubator.

461

462 **Cell culture**

463 Primed H9-NK2 cells containing doxycycline-inducible *KLF2* and *NANOG* transgenes coupled
464 to Venus were maintained in conventional medium (KSR/FGF) comprised of DMEM/F-12
465 (Sigma Aldrich) with 20% KSR (ThermoFisher Scientific) and 10ng per ml basic fibroblast
466 growth factor (bFGF; Peprotech), supplemented with 2mM L-glutamine (ThermoFisher
467 Scientific), 100uM 2-mercaptoethanol (2ME) (ThermoFisher Scientific), 1% MEM non-
468 essential amino acids (ThermoFisher Scientific), and 50mg per ml Penicillin-Streptomycin.
469 Cultures were passaged every 5-6 days as small clumps by dissociation with a buffer
470 containing 1mg per ml Collagenase IV (ThermoFisher Scientific), 0.025% Trypsin
471 (ThermoFisher Scientific), 1mM CaCl₂ and KSR at a final concentration of 20% in PBS.
472 Medium was changed daily.

473 Resetting to the naïve state was carried out as previously described¹². Conventional human
474 embryonic stem cells (hESCs) were dissociated to single cells with trypsin and re-plated in
475 the presence of 10µM Rho-associated kinase inhibitor (ROCKi [Y-27632]; Sigma Aldrich).
476 After 24 hours, media was changed to primed media with 1µM doxycycline (Sigma Alrich).
477 The following day, media was changed to 2iL+dox media composed of 50% DMEM/F12 and
478 50% Neurobasal (ThermoFisher Scientific) supplemented with 2mM L-glutamine, 100µM
479 2ME, N2 (ThermoFisher Scientific), B27 (ThermoFisher Scientific), 1µM PD0325901
480 (StemCell Technologies), 1µM CHIR99021 (StemCell Technologies), human recombinant LIF
481 (Peprotech), 1x Penicillin-Streptomycin and 1µM doxycycline. Media was changed daily.
482 Cells were split every 4-5 days after dissociation to single cells using Accutase (Sigma
483 Aldrich). After 2 weeks, doxycycline was withdrawn and PKC inhibitor Gö6983 (Sigma
484 Aldrich) was added at a concentration of 5µM. Cells in 2iL+Gö were split every 4-5 days after
485 dissociation to single cells using Accutase.

486

487 **Stable knock down or overexpression cell line generation**

488 Short hairpin RNA (shRNA) constructs were obtained from Dharmacon in the TRC pLKO.1
489 lentiviral vector. Sequences are listed in Supplementary Table 1. Knock down of DNMT3A2
490 was not carried out due to an inability to design an effective shRNA against its single unique
491 exon. Entry clone for overexpression of TET1 was obtained from Harvard PlasmID
492 Repository (TET1 in pENTR223; HsCD00399189) and a recombination reaction was
493 performed with the pLenti CMV puro DEST (w118-1) destination vector with Gateway LR
494 clonase II (ThermoFisher Scientific) to generate expression vectors. To generate lentiviral
495 particles, HEK293T cells were transfected with the shRNA plasmid or expression vector for a
496 target gene, the packaging construct pCMV Δ8.91, and a vesicular stomatitis virus
497 glycoprotein (VSV-G) containing envelope expressing plasmid pMD2.G, using jetPrime
498 (Polyplus) at a ratio of 1:2. Primed hESCs were treated with 6µg per ml polybrene (Sigma
499 Aldrich), transduced with filtered lentiviral particles, and stable hESC knock down or
500 overexpression cell lines were generated by puromycin selection (1µg per ml) of
501 successful integrants.

502

503 **Western Blotting**

504 Whole cell lysates were extracted in RIPA buffer (Sigma Aldrich), with protease inhibitor
505 cocktail (Sigma Aldrich). Proteins (concentration determined by BCA assay) were separated
506 by electrophoresis on a 4-12% Bis-Tris gel in MOPS running buffer (ThermoFisher Scientific)
507 and then transferred to polyvinylidene difluoride (PVDF) membranes (Merck Millipore).

508 Membranes were blocked with 5% skimmed milk for 45 min at room temperature and
509 incubated overnight at 4°C with primary antibody (TET1: Source Bioscience; GTX124207 at
510 1:1000 and GAPDH: Cell Signalling Technologies; 2118S at 1:2500) in blocking buffer.
511 Membranes were incubated for 1 hour at room temperature with horseradish peroxidase-
512 conjugated secondary antibodies sheep-anti-mouse IgG or sheep-anti-rabbit IgG (1:5,000;
513 GE Healthcare; NA931, NA934). Membranes were washed in 0.1% Tween-20 in PBS (PBST),
514 and detection was performed with enhanced chemiluminescence (ThermoFisher Scientific),
515 with visualisation on the Amersham Imager 600 (GE Healthcare). Uncropped blots are
516 provided in Supplementary Fig. 11.

517

518 **qPCR**

519 Cells were dissociated to single cells using Accutase and serially plated for 2 hours to
520 eliminate excess iMEFs. Total RNA was isolated from pelleted hESCs using the Direct-zol RNA
521 mini-prep kit (Zymo) and treated with the DNA-free™ DNA removal kit (ThermoFisher
522 Scientific). Complementary DNA (cDNA) was made using a high-capacity RNA to cDNA kit
523 (ThermoFisher Scientific). Real-time PCR was carried out using one-step Sybr green reaction
524 mix (Bio-Rad) on the CFX384 Touch™ Real Time PCR detection system (Bio-Rad). An
525 endogenous control (GAPDH) was used to normalise expression. Primer sequences are listed
526 in Supplementary Table 2.

527

528 **Chromatin Immunoprecipitation**

529 Cells were cross-linked with 1% formaldehyde for 10 minutes at room temperature with
530 gentle rocking, after which the formaldehyde was quenched with 1.25M glycine. Chromatin
531 was then extracted from the cross-linked cells using the chromatin extraction kit (Abcam),
532 as per the manufacturer's instructions. The extracted chromatin was then fractionated by
533 sonication at 4° (12 cycles of 15s on, 60s off; Diagenode Bioruptor® Plus). The size of the
534 sonicated chromatin was then checked by agarose gel electrophoresis. Chromatin
535 immunoprecipitation was carried out using the ChIP – One Step kit (Abcam) with a starting
536 total of 5µg of chromatin. Immunoprecipitation was carried out as per the manufacturer's
537 instructions and the following quantities of antibody were used for immunoprecipitation;
538 H3K4me3 0.5µg (Abcam; ab8580), H3K27me3 2µg (Abcam; ab195477). As a loading control
539 for assessing immunoprecipitation we isolated input DNA for each sample, which represent
540 the starting quantity of chromatin prior to immunoprecipitation. Input and
541 immunoprecipitated DNA were quantified by real-time PCR, and data is shown as the %
542 enrichment relative to the input for each sample. Primer sequences are listed in Table 3.

543

544 **Bisulfite Sequencing Analysis**

545 Bismark coverage files downloaded from GEO were uploaded into SeqMonk (v1.41.0),
546 where the genomes were binned into 300bp probe windows. Methylation quantitation
547 was carried out using the 'Bisulphite methylation over features' pipeline in SeqMonk
548 (v1.41.0), with a 300bp probe carried forward if it contained at least 5 CpGs each with at
549 least 3 counts. Motif enrichment analysis was performed using the analysis of motif
550 enrichment (AME) tool on the MEME suite (v5.0.4)⁶¹, searching against the human
551 HOCOMOCO (v11 FULL) database. Sequences were scored using the average odds score
552 and motif enrichment calculated using Fisher's exact test.

553 **Overlap analysis**

554 Overlap analysis was performed in R using the package regioneR (Version 3.8:
555 <https://bioconductor.org/packages/release/bioc/html/regioneR.html>). Overlap was
556 performed using the 'overlapPermTest' function with 1000 permutations. Random
557 regions were generated for the hg19 genome using the 'circularRandomizeRegions'
558 function. Random loci generation was restricted to loci present in the Illumina EPIC array
559 (for overlaps performed with Illumina EPIC array probes) or to regions with a (G+C)
560 fraction >0.55 and a CpG observed-to-expected ratio >0.6 (for overlaps performed with
561 bisulfite sequencing data). ENCODE and ChromHMM data for the H1 hESC cell line were
562 downloaded from the UCSC genome browser. For ENCODE data, StdPk files were
563 downloaded for each histone modification and genomic coordinates extracted (as BED
564 files) for use in the overlap analysis.

565 **Mass Spectrometry of Nucleosides**

566 Genomic DNA was digested using DNA Degradase Plus (Zymo Research) according to the
567 manufacturer's instructions and nucleosides were analysed by LC-MS/MS on a Q-Exactive
568 mass spectrometer (Thermo Scientific) fitted with a nanoelectrospray ion-source (Proxeon).
569 All samples and standards had a heavy isotope-labelled nucleoside mix added prior to mass
570 spectral analysis (2'-deoxycytidine-¹³C₁, ¹⁵N₂ (Santa Cruz), 5-(methyl-²H₃)-2'-deoxycytidine
571 (Santa Cruz), 5-(hydroxymethyl)-2'-deoxycytidine-²H₃ (Toronto Research Chemicals). MS2
572 data for 5hmC, 5mC and C were acquired with both the endogenous and corresponding
573 heavy-labelled nucleoside parent ions simultaneously selected for fragmentation using a 5
574 Th isolation window with a 1.5 Th offset. Parent ions were fragmented by Higher-energy
575 Collisional Dissociation (HCD) with a relative collision energy of 10%, and a resolution setting
576 of 70,000 for MS2 spectra. Peak areas from extracted ion chromatograms of the relevant
577 fragment ions, relative to their corresponding heavy isotope labelled internal standards,
578 were quantified against a six-point serial 2-fold dilution calibration curve, with triplicate
579 runs for all samples and standards.

580

581 **Targeted bisulfite sequencing**

582 Bisulfite PCR primers were designed against an *in silico* bisulfite converted reference
583 sequence using the Bisulfite Primer Seeker software (Zymo) or Methprimer (Urogene), and
584 universal Illumina adapter sequences were added to the 5' end of each primer. Cells were
585 dissociated to single cells using Accutase and serially plated for 2 hours to eliminate excess
586 iMEFs. DNA was isolated from pelleted cells using the PureLink Genomic DNA mini kit
587 (ThermoFisher Scientific). Bisulfite conversion of DNA was carried out using the Imprint®
588 DNA Modification kit (Sigma Aldrich), following the manufacturer's instructions. The
589 modified DNA was amplified using the loci specific bisulfite PCR primers (listed in
590 Supplementary Table 4) and HotStar Taq DNA Polymerase (Qiagen). The PCR conditions
591 were as follows: 95 °C for 15 min; 94 °C for 30 seconds; 56 °C for 30 seconds; 72 °C for 1
592 min; Repeat steps 2-4 29X; 72 °C for 10 min; Hold 12°C. PCR products were purified using
593 SPRI beads (Agencourt AMPure XP, Beckman Coulter). Amplicons were PCR amplified
594 with 8 cycles using a universal Illumina forward primer and an indexed reverse primer
595 and quantified with the Kapa Library quantification kit for Illumina (Roche). For larger
596 experiments, multiplex targeted bisulfite sequencing was performed using the 48x48
597 layout on the Fluidigm C1 system (Fluidigm), coupled with Illumina MiSeq sequencing.
598 Fluidigm primers are listed in Supplementary Table 5. Amplicons from a single sample
599 were pooled and sequencing was performed on an Illumina MiSeq with 150bp paired-end

600 reads, using v3 chemistry, at Barts and the London Genome Centre (London, UK). Reads
601 were quality trimmed and mapped to a personalised human genome composed of
602 amplicon sequences, using Bismark (v.0.19.0), followed by extraction of methylation calls.
603

604 **Infinium MethylationEPIC BeadChip assay**

605 Genomic DNA was extracted using the PureLink Genomic DNA mini kit (ThermoFisher
606 Scientific). Bisulfite conversion of DNA was carried out using the Imprint[®] DNA Modification
607 kit (Sigma Aldrich), following the manufacturer's instructions. Infinium MethylationEPIC
608 BeadChip assay (Illumina) was performed according to manufacturer instructions by Barts
609 and the London Genome Centre (London, UK). The Bioconductor package ChAMP (version
610 2.11.3: <https://bioconductor.org/packages/release/bioc/html/ChAMP.html>) was used to
611 process raw Infinium idat files using the GRCh37 human genome manifest file.

612 **TCGA Analysis**

613 Illumina 450K DNA Methylation data spanning 396965 CpGs and 9664 samples was
614 downloaded from the Pan Cancer Atlas ([https://gdc.cancer.gov/about-](https://gdc.cancer.gov/about-data/publications/pancanatlas)
615 [data/publications/pancanatlas](https://gdc.cancer.gov/about-data/publications/pancanatlas)). All samples from individuals without both a tumour and
616 normal tissue sample were removed. Samples from tumour types with less than 30
617 individuals were removed. In order to assess only CpGs deemed "bivalent", CpGs outside of
618 regions that showed a peak of H3K27me3 and H3K4me3 in ENCODE H1 hESCs were
619 removed. For this analysis, raw infinium IDAT files from the hESC resetting experiment were
620 processed using minfi and normalised via the singlesample Noob method⁶². CpGs used for
621 analysis were filtered for those that are unmethylated in primed hESCs (mean beta < 0.3).
622 Unmethylated probes were restricted to those CpGs with mean Beta < 0.3 during the
623 primed to naïve transition. Hypermethylated hESC probes were defined using ChAMP and
624 restricted to those CpGs with $\Delta\text{Beta} > 0.1$ in either the early transition, late transition or
625 naïve state. Probes hypermethylated in cancer were also defined using ChAMP, and similarly
626 restricted to those CpGs with $\Delta\text{Beta} > 0.1$ between normal tissue and tumour samples. The
627 overlap enrichment of cancer hypermethylated and hESC hypermethylated regions was
628 determined via Fisher's exact test. For the creation of heatmaps, data was first ordered by
629 sample based on mean methylation of all CpGs, and then by CpG based on mean
630 methylation across all samples of every cancer type. Statistical significance was calculated
631 using a paired Wilcoxon test.
632

633 **Human Tissue H3K27me3 Analysis**

634 H3K27me3 and control read alignments were downloaded as BAM files for each normal
635 human tissue examined from ENCODE (<https://www.encodeproject.org/>)^{63,64} that
636 corresponded to the cancers profiled by TCGA, except HSNC for which no obvious
637 corresponding normal sample was available. Bigwig files for genome browser visualisation
638 and peaks of H3K27me3 for comparison were obtained in the same way. ENCODE IDs for
639 each experiment and data file can be found in Supplementary Data 1. To examine
640 H3K27me3 levels around naïve hypermethylated CpGs, windows of 500bp was defined
641 centred around each CpG. ChIP-seq read counts/window were calculated using BEDtools'
642 coverage function. Read counts were scaled to counts per 10 million based on total number
643 of mapped reads/sample and divided by the input read count to provide a normalised read
644 count. To prevent windows with zero reads in the input sample generating a normalised
645 count of infinity, an offset of 0.5 was added to all windows prior to scaling and input

646 normalisation. Regions where coverage was 0 in all samples were removed from the
647 analysis. A similar procedure was used to generate heatmaps of H3K27me3 levels, using
648 multiple 250bp windows to span 5kb on either side of each CpG. Colour scales for ChIP-seq
649 heatmaps range from the minimum to the 90% quantile of the normalised read count.

650

651 **Analysis of 5hmC by glucMS-qPCR**

652 Genomic DNA was treated with T4 Phage β -glucosyltransferase (T4-BGT; NEB) according to
653 the manufacturer's instructions. Glucosylated genomic DNA was digested with 10 U of
654 either HpaII, MspI or no enzyme (mock digestion) at 37 °C overnight, followed by
655 inactivation for 20 min at 80 °C. The HpaII- and MspI-resistant fraction was quantified by
656 qPCR using primers designed around at least one HpaII/MspI site, and normalizing to the
657 mock digestion control. Resistance to MspI directly translates into percentage of 5hmC,
658 whereas 5mC levels were obtained by subtracting the 5hmC contribution from the total
659 HpaII resistance. Primers used are listed in Supplementary Table 6.

660

661 **Mass Spectrometry-based proteomics**

662 Cells from 3 independent biological replicates per condition were washed twice with ice
663 cold PBS supplemented with 1 mM Na_3VO_4 and 1 mM NaF, lysed in urea buffer (8M urea in
664 20 mM in HEPES pH 8.0, 1 mM Na_3VO_4 , 1 mM NaF, 1mM $\text{Na}_4\text{P}_2\text{O}_7$ and 1 mM sodium β -
665 glycerophosphate) for 30 min and homogenized by sonication (15 cycles of 30s on 30s off;
666 Diagenode Bioruptor[®] Plus). Insoluble material was removed by centrifugation at 20,000 xg
667 and protein levels in the cell extracts were quantified by bicinchoninic acid (BCA) analysis.
668 For trypsin digestion, 100 μg of protein was reduced and alkylated by sequential incubation
669 with 10 mM DTT and 16.6 mM iodoacetamide for 1h and 30min, respectively. Urea
670 concentration was diluted to 2 M with 20 mM HEPES (pH 8.0), 80 μL of preconditioned
671 trypsin beads [(50% slurry of TLCK-trypsin (Thermo-Fisher Scientific; Cat. #20230)] were
672 added and samples were incubated for 16h at 37C with agitation. Peptide solutions were
673 desalted using 10 mg OASIS-HLB cartridges (Waters, Manchester, UK). Briefly, OASIS
674 cartridges were accommodated in a vacuum manifold (-5 mmHg), activated with 1 mL ACN
675 and equilibrated with 1.5 mL washing solution (1% ACN, 0.1% TFA). Peptides were loaded
676 into the cartridges, washed with 1 mL of washing solution, eluted with 500 μL of ACN
677 solution (30% ACN, 0.1% TFA), dried in a speed vac (RVC 2-25, Martin Christ
678 Gefriertrocknungsanlagen) and stored at -80C. Dried peptides were dissolved in 0.1% TFA
679 and analysed by nanoflow ultimate 3,000 RSL nano instrument coupled to a Q Exactive plus
680 mass spectrometer (Thermo Fisher Scientific). Gradient elution was from 3% to 35% solvent
681 B in 120 min at a flow rate 300 nL/min with solvent A being used to balance the mobile
682 phase (buffer A was 0.1% formic acid in water and B was 0.1% formic acid in acetonitrile) .
683 The spray voltage was 1.95 kV and the capillary temperature was set to 255C. The Q-
684 Exactive plus was operated in data dependent mode with one survey MS scan followed by
685 15 MS/MS scans. The full scans were acquired in the mass analyser at 375-1500 m/z with
686 the resolution of 70,000 and the MS/MS scans were obtained with a resolution of 17,500.
687 Overall duty cycle generated chromatographic peaks of approximately 30s at the base,
688 which allowed the construction of extracted ion chromatograms (XICs) with at least 10 data
689 points. Mascot Daemon 2.5.0 was used to automate peptide identification from MS data.
690 Peak list files (MGFs) from RAW data were generated with Mascot Distiller v2.5.1 and
691 loaded into the Mascot search engine (v2.5) in order to match MS/MS data to peptides.
692 Searches were performed against the SwissProt Database (release December 2015) with a

693 FDR of ~1% and restricted to the human entries. Mass tolerance of ± 10 ppm for the MS
694 scans and ± 25 mmu for the MS/MS scans, 2 trypsin missed cleavages, carbamidomethyl Cys
695 as a fixed modification and PyroGlu on N-terminal Gln and oxidation of Met as variable
696 modifications were allowed. The in-house developed Pescal software was used for label-
697 free peptide quantification as described before⁶⁵, XICs for all the peptides identified across
698 all samples were constructed with ± 2 min and ± 7 ppm retention time and mass windows,
699 respectively. Peak areas from all XICs were calculated. The maximum intensity value for the
700 2 technical replicates was selected and used for further analysis. Intensity values for each
701 peptide were normalized to total sample intensity. Statistical significance was calculated
702 using two tail unpaired Student's t-test. Multiplicity correction was performed by applying
703 the Benjamini-Hochberg method on the p-values, to control the false discovery rate (FDR).
704 Differences were considered significant when $FDR < 0.05$. Proteins with a Mascot score > 40
705 were used for analysis. Data are available in Supplementary Data 2.

706 **RNA-sequencing**

707 Total RNA was extracted using Direct-zol RNA mini kit (Zymo) and DNase treated
708 (ThermoFisher Scientific), before mRNA was isolated from 500ng of total RNA using
709 Dynabeads mRNA DIRECT purification kit (ThermoFisher Scientific) and fragmented with
710 RNA fragmentation reagent (ThermoFisher Scientific). First strand cDNA synthesis was
711 performed with SuperScript III First-Strand Synthesis System and $3 \mu\text{g } \mu\text{l}^{-1}$ random
712 hexamers (ThermoFisher Scientific) followed by second strand synthesis with DNA
713 polymerase I and RNase H. After purification using SPRI beads, the double stranded cDNA
714 was ligated to in house designed adapters (based on TruSeq Indexed adapters (Illumina))
715 using NEBNext Ultra II (NEB) followed by 15 cycles of amplification and library
716 purification. Library size distribution and molarity was assessed by the DNA 1000 assay on
717 the 2100 Bioanalyzer (Agilent), and libraries were quantified with the Kapa Library
718 quantification kit for Illumina (Roche). Sequencing was performed on an Illumina NextSeq
719 with 75bp paired-end sequencing at Barts and the London Genome Centre (London, UK).
720 Read quality was determined using FASTQC. Genomic mapping of short reads was
721 performed using hisat2 (v. 2.1.0) to the human genome (GRCh38). Reads were counted
722 for each sample using FeatureCounts (Subread, v. 1.6.3)⁶⁶. RNA-sequencing analysis was
723 performed using the R package EdgeR (v3.18.1)⁶⁷. Upregulated and downregulated genes
724 were called as those with Benjamini-Hochberg corrected $FDR < 0.05$ and a \log_2 fold
725 change > 1 . Pathway enrichment analysis was performed using DAVID Bioinformatics
726 Resources^{68,69}.

727

728 **Alkaline Phosphatase Assay**

729 Cells were seeded a 96-well plate. After 24 hours, the Amplite™ Colorimetric Alkaline
730 Phosphatase Assay kit (Strattech) was used to measure alkaline phosphatase activity
731 according to manufacturer's instructions.

732

733 **Flow cytometry, fluorescence-activated cell sorting (FACS)**

734 Cells were dissociated to single cells with Accutase and washed with 3%FCS/PBS, before
735 being blocked in 10% FCS/PBS. Cells were resuspended in 2.5ul SUSD2-PE antibody
736 (Biolegend; 327406) and 5ul anti-feeder-APC antibody (Miltenyi Biotec; 130-120-802) for 15
737 minutes at 4°C in the dark. Alternatively, following 1 hour pulse labelling with 10um BrdU,
738 cells were fixed, permeabilised, blocked and stained using the APC BrdU Flow Kit (BD

739 Pharminogen) following manufacturer's instructions, with the addition of 5ul of anti-feeder-
740 PE antibody (Miltenyi Biotec; 130-120-166). Cells were wash twice in 3% FCS/PBS and then
741 stained with DAPI for 15 minutes at 4°C in the dark. Samples were either analysed on an LSR
742 Fortessa cell analyser (BD Biosciences) or FACS sorted on the BD FACS Aria Fusion cell sorter.
743 Flow cytometry data analysis was carried out using FlowJo Version 10 software.

744

745 **Statistical Analysis**

746 Significance testing was performed using Prism (v.7.04, v.8.4.2) and Student's T-test, one-
747 way ANOVA or two-way ANOVA with Bonferroni post-hoc tests as specified in the Figure
748 legends. Where applicable, data are plotted as mean ± SEM. Representative data are
749 shown where experiments were repeated at least twice with similar results.

750 **Data availability**

751 All datasets have been deposited in the Gene Expression Omnibus and are accessible
752 under [<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE128130>] GSE128130.
753 Additional data used include ENCODE, ChromHMM and TCGA pan-cancer data,
754 HOCOMOCO (v11 FULL), SwissProt (Dec 2015 release). Data for human naive resetting
755 methods was downloaded from GSE60945, GSE76970, GSE90168, data for mouse 2i and
756 serum ESCs from GSE42923, and data for human and mouse in vivo development from
757 GSE34864 and GSE49828. The source data for Figs. 2e, 3c-g, 4b-c, 5b-c, 6c-d, and
758 Supplementary Figs 1a, 1h, 5a-b, 5d-g, 6a-e, 7b, 8a-b are provided as a source data file.

759

760 Correspondence and requests for materials should be addressed to g.ficz@qmul.ac.uk.

761

762

763

764

765

766

767

768 **References**

769

- 770 1 Jones, P. A. & Baylin, S. B. The epigenomics of cancer. *Cell* **128**, 683-692,
771 doi:10.1016/j.cell.2007.01.029 (2007).
- 772 2 Baylin, S. B. & Jones, P. A. A decade of exploring the cancer epigenome - biological
773 and translational implications. *Nat Rev Cancer* **11**, 726-734, doi:10.1038/nrc3130
774 (2011).
- 775 3 Esteller, M., Corn, P. G., Baylin, S. B. & Herman, J. G. A Gene Hypermethylation
776 Profile of Human Cancer. *Cancer Research* **61**, 3225-3229 (2001).
- 777 4 Feinberg, A. P., Ohlsson, R. & Henikoff, S. The epigenetic progenitor origin of human
778 cancer. *Nat Rev Genet* **7**, 21-33, doi:10.1038/nrg1748 (2006).

779 5 Eden, A., Gaudet, F., Waghmare, A. & Jaenisch, R. Chromosomal Instability and
780 Tumors Promoted by DNA Hypomethylation. *Science* **300** (2003).

781 6 Gaudet, F. *et al.* Induction of Tumors in Mice by Genomic Hypomethylation. *J Biol*
782 *Chem* **280**, 17986-17991 (2003).

783 7 Jones, P. A. & Baylin, S. B. The fundamental role of epigenetic events in cancer. *Nat*
784 *Rev Genet* **3**, 415-428, doi:10.1038/nrg816 (2002).

785 8 Sproul, D. & Meehan, R. R. Genomic insights into cancer-associated aberrant CpG
786 island hypermethylation. *Brief Funct Genomics* **12**, 174-190, doi:10.1093/bfpg/els063
787 (2013).

788 9 Tao, Y. *et al.* Aging-like Spontaneous Epigenetic Silencing Facilitates Wnt Activation,
789 Stemness, and Braf(V600E)-Induced Tumorigenesis. *Cancer Cell* **35**, 315-328 e316,
790 doi:10.1016/j.ccell.2019.01.005 (2019).

791 10 Chen, H. & He, X. The Convergent Cancer Evolution toward a Single Cellular
792 Destination. *Mol Biol Evol* **33**, 4-12, doi:10.1093/molbev/msv212 (2016).

793 11 Avior, Y., Eggan, K. & Benvenisty, N. Cancer-Related Mutations Identified in Primed
794 and Naive Human Pluripotent Stem Cells. *Cell Stem Cell* **25**, 456-461,
795 doi:10.1016/j.stem.2019.09.001 (2019).

796 12 Takashima, Y. *et al.* Resetting transcription factor control circuitry toward ground-
797 state pluripotency in human. *Cell* **158**, 1254-1269, doi:10.1016/j.cell.2014.08.029
798 (2014).

799 13 Bernstein, B. E. *et al.* A bivalent chromatin structure marks key developmental genes
800 in embryonic stem cells. *Cell* **125**, 315-326, doi:10.1016/j.cell.2006.02.041 (2006).

801 14 Theunissen, T. W. *et al.* Systematic identification of culture conditions for induction
802 and maintenance of naive human pluripotency. *Cell Stem Cell* **15**, 471-487,
803 doi:10.1016/j.stem.2014.07.002 (2014).

804 15 Pastor, W. A. *et al.* Naive Human Pluripotent Cells Feature a Methylation Landscape
805 Devoid of Blastocyst or Germline Memory. *Cell Stem Cell* **18**, 323-329,
806 doi:10.1016/j.stem.2016.01.019 (2016).

807 16 Guo, G. *et al.* Epigenetic resetting of human pluripotency. *Development* **144**, 2748-
808 2763, doi:10.1242/dev.146811 (2017).

809 17 Guo, H. *et al.* The DNA methylation landscape of human early embryos. *Nature* **511**,
810 606-610, doi:10.1038/nature13544 (2014).

811 18 Ficiz, G. *et al.* FGF signaling inhibition in ESCs drives rapid genome-wide
812 demethylation to the epigenetic ground state of pluripotency. *Cell Stem Cell* **13**, 351-
813 359, doi:10.1016/j.stem.2013.06.004 (2013).

814 19 Smith, Z. D. *et al.* A unique regulatory phase of DNA methylation in the early
815 mammalian embryo. *Nature* **484**, 339-344, doi:10.1038/nature10960 (2012).

816 20 Ware, C. B. Concise Review: Lessons from Naive Human Pluripotent Cells. *Stem Cells*,
817 doi:10.1002/stem.2507 (2016).

818 21 Manzo, M. *et al.* Isoform-specific localization of DNMT3A regulates DNA methylation
819 fidelity at bivalent CpG islands. *EMBO J* **36**, 3421-3434,
820 doi:10.15252/embj.201797038 (2017).

821 22 Xu, Y. *et al.* Genome-wide regulation of 5hmC, 5mC, and gene expression by Tet1
822 hydroxylase in mouse embryonic stem cells. *Mol Cell* **42**, 451-464,
823 doi:10.1016/j.molcel.2011.04.005 (2011).

824 23 Verma, N. *et al.* TET proteins safeguard bivalent promoters from de novo
825 methylation in human embryonic stem cells. *Nat Genet* **50**, 83-95,
826 doi:10.1038/s41588-017-0002-y (2018).

827 24 Bredenkamp, N., Stirparo, G. G., Nichols, J., Smith, A. & Guo, G. The Cell-Surface
828 Marker Sushi Containing Domain 2 Facilitates Establishment of Human Naive
829 Pluripotent Stem Cells. *Stem Cell Reports* **12**, 1212-1222,
830 doi:10.1016/j.stemcr.2019.03.014 (2019).

831 25 Ooi, S. K. *et al.* DNMT3L connects unmethylated lysine 4 of histone H3 to de novo
832 methylation of DNA. *Nature* **448**, 714-717, doi:10.1038/nature05987 (2007).

833 26 Brinkman, A. B. *et al.* Sequential CHIP-bisulfite sequencing enables direct genome-
834 scale investigation of chromatin and DNA methylation cross-talk. *Genome Res* **22**,
835 1128-1138, doi:10.1101/gr.133728.111 (2012).

836 27 Gao, F. *et al.* Direct CHIP-bisulfite sequencing reveals a role of H3K27me3 mediating
837 aberrant hypermethylation of promoter CpG islands in cancer cells. *Genomics* **103**,
838 204-210, doi:10.1016/j.ygeno.2013.12.006 (2014).

839 28 Otani, J. *et al.* Structural basis for recognition of H3K4 methylation status by the DNA
840 methyltransferase 3A ATRX-DNMT3-DNMT3L domain. *EMBO Rep* **10**, 1235-1241,
841 doi:10.1038/embor.2009.218 (2009).

842 29 Guo, X. *et al.* Structural insight into autoinhibition and histone H3-induced activation
843 of DNMT3A. *Nature* **517**, 640-644, doi:10.1038/nature13899 (2015).

844 30 Zhang, Y. *et al.* Chromatin methylation activity of Dnmt3a and Dnmt3a/3L is guided
845 by interaction of the ADD domain with the histone H3 tail. *Nucleic Acids Res* **38**,
846 4246-4253, doi:10.1093/nar/gkq147 (2010).

847 31 Simicevic, J. & Deplancke, B. Transcription factor proteomics-Tools, applications, and
848 challenges. *Proteomics* **17**, doi:10.1002/pmic.201600317 (2017).

849 32 Bernhart, S. H. *et al.* Changes of bivalent chromatin coincide with increased
850 expression of developmental genes in cancer. *Sci Rep* **6**, 37393,
851 doi:10.1038/srep37393 (2016).

852 33 Easwaran, H. *et al.* A DNA hypermethylation module for the stem/progenitor cell
853 signature of cancer. *Genome Res* **22**, 837-849, doi:10.1101/gr.131169.111 (2012).

854 34 Kushwaha, G. *et al.* Hypomethylation coordinates antagonistically with
855 hypermethylation in cancer development: a case study of leukemia. *Hum Genomics*
856 **10 Suppl 2**, 18, doi:10.1186/s40246-016-0071-5 (2016).

857 35 Hansen, K. D. *et al.* Increased methylation variation in epigenetic domains across
858 cancer types. *Nat Genet* **43**, 768-775, doi:10.1038/ng.865 (2011).

859 36 Cancer Genome Atlas Research, N. *et al.* The Cancer Genome Atlas Pan-Cancer
860 analysis project. *Nat Genet* **45**, 1113-1120, doi:10.1038/ng.2764 (2013).

861 37 Hoadley, K. A. *et al.* Cell-of-Origin Patterns Dominate the Molecular Classification of
862 10,000 Tumors from 33 Types of Cancer. *Cell* **173**, 291-304 e296,
863 doi:10.1016/j.cell.2018.03.022 (2018).

864 38 Sproul, D. *et al.* Tissue of origin determines cancer-associated CpG island promoter
865 hypermethylation patterns. *Genome Biol* **13** (2012).

866 39 Ben-Porath, I. *et al.* An embryonic stem cell-like gene expression signature in poorly
867 differentiated aggressive human tumors. *Nat Genet* **40**, 499-507, doi:10.1038/ng.127
868 (2008).

869 40 Smith, Z. D. *et al.* Epigenetic restriction of extraembryonic lineages mirrors the
870 somatic transition to cancer. *Nature* **549**, 543-547, doi:10.1038/nature23891 (2017).

871 41 Diede, S. J. *et al.* Fundamental differences in promoter CpG island DNA
872 hypermethylation between human cancer and genetically engineered mouse models
873 of cancer. *Epigenetics* **8**, 1254-1260, doi:10.4161/epi.26486 (2013).

874 42 Ohm, J. E. *et al.* A stem cell-like chromatin pattern may predispose tumor suppressor
875 genes to DNA hypermethylation and heritable silencing. *Nat Genet* **39**, 237-242,
876 doi:10.1038/ng1972 (2007).

877 43 Schlesinger, Y. *et al.* Polycomb-mediated methylation on Lys27 of histone H3 pre-
878 marks genes for de novo methylation in cancer. *Nat Genet* **39**, 232-236,
879 doi:10.1038/ng1950 (2007).

880 44 Widschwendter, M. *et al.* Epigenetic stem cell signature in cancer. *Nat Genet* **39**,
881 157-158, doi:10.1038/ng1941 (2007).

882 45 Ohm, J. E. *et al.* Cancer-related epigenome changes associated with reprogramming
883 to induced pluripotent stem cells. *Cancer Res* **70**, 7662-7673, doi:10.1158/0008-
884 5472.CAN-10-1361 (2010).

885 46 Vire, E. *et al.* The Polycomb group protein EZH2 directly controls DNA methylation.
886 *Nature* **439**, 871-874, doi:10.1038/nature04431 (2006).

887 47 Rush, M. *et al.* Targeting of EZH2 to a defined genomic site is sufficient for
888 recruitment of Dnmt3a but not de novo DNA methylation. *Epigenetics* **4**, 404-414,
889 doi:10.4161/epi.4.6.9392 (2014).

890 48 Jair, K. W. *et al.* De novo CpG island methylation in human cancer cells. *Cancer Res*
891 **66**, 682-692, doi:10.1158/0008-5472.CAN-05-1980 (2006).

892 49 Brenner, C. *et al.* Myc represses transcription through recruitment of DNA
893 methyltransferase corepressor. *The EMBO Journal* **24**, 336-346, doi:10.1038/ (2005).

894 50 Hervouet, E., Vallette, F. M. & Cartron, P.-F. Dnmt3/transcription factor interactions
895 as crucial players in targeted DNA methylation. *Epigenetics* **4**, 487-499,
896 doi:10.4161/epi.4.7.9883 (2014).

897 51 Yan, L. *et al.* Single-cell RNA-Seq profiling of human preimplantation embryos and
898 embryonic stem cells. *Nat Struct Mol Biol* **20**, 1131-1139, doi:10.1038/nsmb.2660
899 (2013).

900 52 Blakeley, P. *et al.* Defining the three cell lineages of the human blastocyst by single-
901 cell RNA-seq. *Development* **142**, 3613, doi:10.1242/dev.131235 (2015).

902 53 Saunderson, E. A. *et al.* Hit-and-run epigenetic editing prevents senescence entry in
903 primary breast cells from healthy donors. *Nat Commun* **8**, 1450, doi:10.1038/s41467-
904 017-01078-2 (2017).

905 54 Pfeifer, G. P. Defining Driver DNA Methylation Changes in Human Cancer. *Int J Mol*
906 *Sci* **19**, doi:10.3390/ijms19041166 (2018).

907 55 Polo, J. M. *et al.* A molecular roadmap of reprogramming somatic cells into iPS cells.
908 *Cell* **151**, 1617-1632, doi:10.1016/j.cell.2012.11.039 (2012).

909 56 Mikkelsen, T. S. *et al.* Dissecting direct reprogramming through integrative genomic
910 analysis. *Nature* **454**, 49-55, doi:10.1038/nature07056 (2008).

911 57 Hanley, M. P. *et al.* Genome-wide DNA methylation profiling reveals cancer-
912 associated changes within early colonic neoplasia. *Oncogene* **36**, 5035-5044,
913 doi:10.1038/onc.2017.130 (2017).

914 58 Hanahan, D. & Weinberg, R. A. Hallmarks of cancer: the next generation. *Cell* **144**,
915 646-674, doi:10.1016/j.cell.2011.02.013 (2011).

916 59 Ficzb, G. & Gribben, J. G. Loss of 5-hydroxymethylcytosine in cancer: cause or
917 consequence? *Genomics* **104**, 352-357, doi:10.1016/j.ygeno.2014.08.017 (2014).

- 918 60 Liu, X. *et al.* Comprehensive characterization of distinct states of human naive
919 pluripotency generated by reprogramming. *Nat Methods*, doi:10.1038/nmeth.4436
920 (2017).
- 921 61 Bailey, T. L. *et al.* MEME SUITE: tools for motif discovery and searching. *Nucleic Acids*
922 *Res* **37**, W202-208, doi:10.1093/nar/gkp335 (2009).
- 923 62 Fortin, J. P., Triche, T. J., Jr. & Hansen, K. D. Preprocessing, normalization and
924 integration of the Illumina HumanMethylationEPIC array with minfi. *Bioinformatics*
925 **33**, 558-560, doi:10.1093/bioinformatics/btw691 (2017).
- 926 63 Davis, C. A. *et al.* The Encyclopedia of DNA elements (ENCODE): data portal update.
927 *Nucleic Acids Res* **46**, D794-D801, doi:10.1093/nar/gkx1081 (2018).
- 928 64 Consortium, E. P. An integrated encyclopedia of DNA elements in the human
929 genome. *Nature* **489**, 57-74, doi:10.1038/nature11247 (2012).
- 930 65 Alcolea, M. P., Casado, P., Rodríguez-Prados, J. C., Vanhaesebroeck, B. & Cutillas, P.
931 R. Phosphoproteomic Analysis of Leukemia Cells under Basal and Drug-treated
932 Conditions Identifies Markers of Kinase Pathway Activation and Mechanisms of
933 Resistance. *Molecular and Cellular Proteomics* **11**, 453-466, doi:10.1074/ (2012).
- 934 66 Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose program
935 for assigning sequence reads to genomic features. *Bioinformatics* **30**, 923-930,
936 doi:10.1093/bioinformatics/btt656 (2014).
- 937 67 Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for
938 differential expression analysis of digital gene expression data. *Bioinformatics* **26**,
939 139-140, doi:10.1093/bioinformatics/btp616 (2010).
- 940 68 Huang, D. W., Sherman, B. T. & Lempicki, R. A. Bioinformatics enrichment tools:
941 paths toward the comprehensive functional analysis of large gene lists. *Nucleic Acids*
942 *Res* **37**, 1-13, doi:10.1093/nar/gkn923 (2009).
- 943 69 Huang, D. W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative analysis of
944 large gene lists using DAVID bioinformatics resources. *Nature Protocols* **4**, 44,
945 doi:10.1038/nprot.2008.211
946

947 **Acknowledgements**

948 We would like to thank Prof. Austin Smith at the WT–MRC Cambridge Stem Cell Institute
949 for providing H9 NK2 cells. We thank Dr. Charles Mein and team from Barts and The
950 London Genome Centre for the Infinium Methylation array and Sequencing services, and
951 we thank Dr. Vinothini Rajeeve, Dr Pedro Casado and Dr. Arran Dokal for the Mass
952 Spectrometry Proteomics service. We thank Dr. Miguel Branco and Dr. Emily Saunderson
953 for their advice on preparation of the manuscript. This work was supported by MRC grant
954 MR/M01892X/1; H.P. is supported by an MRC studentship [Ref: 1650326] and The Greg
955 Wolf fund. D.S. is a Cancer Research UK Career Development fellow (Ref: C47648/A20837),
956 and work in his laboratory is also supported by a MRC university grant to the MRC Human
957 Genetics Unit.

958 **Author Contributions**

959 Project was conceived by G.F. and developed by H.P. and M.D.R. Design of experimental
960 work, bioinformatics analysis and data interpretation was performed by H.P. and M.D.R.
961 Analysis of TCGA pan-cancer data was performed by J.H. Proteomics and analysis of
962 proteomics data was performed by S.A.T. Mass spectrometry of nucleosides was
963 performed by D.O. Research was supervised by G.F., D.S., and P.C. Manuscript was

964 written by H.P., M.D.R and G.F. All authors contributed to the editing of the manuscript
965 and approved its final version.
966

967 The authors declare no competing interests.

968

969 **Figure Legends**

970

971 **Figure 1. Primed to naïve resetting induces bivalent CGI promoter hypermethylation.** a)
972 Schematic detailing the model system and time points used in the study. 2iL+dox: 2 small
973 molecule inhibitors of MEK1/2 and GSK3 β (2i), human recombinant leukaemia inhibitory
974 factor (hLIF; collectively 2iL) and doxycycline. 2iL+Gö: 2iL and a pan-protein kinase C
975 inhibitor (PKCi), Gö. hESCs, human embryonic stem cells; b) Heatmap showing methylation
976 levels of the top 10,000 CpG probes that are differentially methylated ($\Delta\beta > 0.1$, adjPval <
977 0.05) in the early transition, late transition and naïve hESCs compared to primed hESCs.
978 Methylation β -value is indicated by the colour key. adjPval based on Benjamini-Hochberg
979 adjustment. c) Genome browser tracks for Infinium MethylationEPIC data capturing a
980 representative hypermethylated locus. The heatmap shows the raw methylation β -values
981 per CpG for each sample, while the subsequent rows show the per-probe difference in
982 methylation for each time point of resetting compared to primed hESCs. CGIs are
983 highlighted in green. d) Overlap of hypermethylated probes (n = 91119) with regions of
984 histone modification enrichment (obtained from the ENCODE ChIP-seq data for hESC cell
985 line H1: H3K4me1 n = 139971; H3K4me2 n = 73086; H3K4me3 n = 33270; H3K9me3 n =
986 86122; H3K27me3 n = 25909; H3K36me3 n = 35877; H3K79me2 n = 33205). e) The
987 proportion of hypermethylated probes (n = 46844 early transition, n = 91119 late transition,
988 n = 20297 naïve) that are marked by H3K4me1/2/3, H3K4me1/2/3 alone, or bivalency
989 (H3K4me3 and H3K27me3). f) Overlap of late transition hypermethylated probes (n = 91119)
990 and hypomethylated probes (n = 392,875) with CpG islands (n = 30344). g) Overlap of
991 hypermethylated and hypomethylated probes (as in 1f) with ENCODE regulatory regions
992 (promoters or enhancers). h) Proportion of late transition hypermethylated probes (n =
993 91119) that overlap with ENCODE predicted promoters and enhancers (as defined by
994 ChromHMM in the hESC cell line H1).
995 DMP, differentially methylated probes. For overlap analysis, data is presented as the log₂
996 corrected fold increase in the observed overlap compared to the mean overlap of 1000
997 randomly generated loci. ***P < 0.001.

998

999 **Figure 2. Resetting results in hypermethylation and repression of developmental genes.** a)
1000 GO term analysis of hypermethylated and hypomethylated genes at the late transition of
1001 resetting compared to primed hESCs. A gene was classified as hypermethylated if a
1002 hypermethylated probe was present within 1500bp upstream of the TSS. b) Average gene
1003 expression for genes that are hypermethylated (average promoter methylation $\Delta\beta > 0.1$) or
1004 hypomethylated (average promoter methylation $\Delta\beta < 0.1$) in naïve compared to primed
1005 hESCs. Boxes represent the median and interquartile range and error bars represent the
1006 maximum and minimum values. Statistical significance determined via two-tailed paired

1007 Wilcoxon test. *** $P < 0.001$, **** $P < 0.0001$. CPM, Counts per Million. c) Scatter plot of genes
1008 that are hypermethylated and downregulated showing the average promoter methylation
1009 (average β -value of CpG probes within 1500bp of TSS) versus the log₂ CPM (counts per
1010 million) for each gene from RNA-seq data. Data for each individual time point is indicated by
1011 the colour key. d) Scatter plot showing the average promoter methylation of 21 HOX genes
1012 (average β -value of CpG probes within 1500bp of TSS) versus the log₂ CPM for each gene
1013 from RNA-seq data. Data for each individual time point is indicated by the colour key. e)
1014 GlucMS-qPCR in primed and late transition hESCs showing the percentage of 5mC (blue) and
1015 5hmC (orange) per CpG. Bars represent the mean of three biological replicates and error
1016 bars represent the SEM. Source data are provided as a Source Data file.

1017
1018 **Figure 3. Early *de novo* methylation is dependent on DNMT3A.** a&b) Heatmaps showing
1019 methylation levels in control and DNMT3A/B knock down samples during resetting.
1020 Heatmap shows the top 17,000 CpG differentially methylated probes (DMP; $\Delta\beta > 0.1$, $p <$
1021 0.05) in the early transition (Fig 3a) and late transition (Fig 3b) compared to primed hESCs
1022 (in wild type early transition (Fig 3a) or late transition (Fig 3b) compared to primed hESCs
1023 identified in analysis shown in Fig 1b). Methylation β -value is indicated by the colour key. c)
1024 qRT-PCR for naïve pluripotency genes in control and knock-down cells, in naïve hESCs. Bars
1025 represent mean of three biological replicates and error bars represent the SEM. Statistical
1026 difference between samples was calculated by a one-way ANOVA with a Bonferroni post-
1027 hoc test compared to the control. Human GAPDH was used to normalise expression. d)
1028 Alkaline phosphatase activity in knock down and control naïve hESCs. Data shown are the
1029 mean of 2 biological replicates (independent shRNA KD), each with 5 technical replicates.
1030 Error bars represent SEM. Statistical difference between samples was analysed by a one-
1031 way ANOVA with a Bonferroni post-hoc test compared to the control. e) Alkaline
1032 phosphatase activity in DNMT3A/DNMT3B double knock down and control naïve hESCs.
1033 Data shown are the mean of 5 replicates. Error bars represent SEM. Statistical difference
1034 between samples was analysed by a two-tailed Student's unpaired t-test. f) qRT-PCR for
1035 naïve pluripotency genes in control and DNMT3A/DNMT3B double knock-down naïve hESCs.
1036 Bars represent the mean of three biological replicates and error bars represent the SEM.
1037 Statistical difference between samples was calculated by a two-tailed Student's t-test.
1038 Human GAPDH was used to normalise expression. g) Plot showing % methylation in
1039 DNMT3A/DNMT3B double knock down and control cells during resetting. Each dot
1040 represents % methylation of single CpGs ($n = 57$) from 4 genomic regions analysed by
1041 targeted bisulfite sequencing. Red bars represent mean methylation for each sample.
1042 Statistical difference between samples was analysed by a Kruskal-Wallis test, with Dunn's
1043 multiple comparisons post-hoc test.
1044 For all panels, * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$, ****($P < 0.0001$), N.S. denotes not significant
1045 ($p > 0.05$). Source data are provided as a Source Data file.

1046
1047 **Figure 4. Hypermethylation is a feature of resetting and is not restricted to emerging naïve**
1048 **cells.** a) Flow cytometry dot plots showing SUSD2-PE staining on the x-axis against anti-
1049 feeder-APC staining on the y-axis for hESCs during the early and late transition. Boxes
1050 indicate SUSD2+ and SUSD2- cell populations that were sorted. b) Plot showing the %
1051 methylation in primed, early transition and late transition hESCs, for SUSD2+ and SUSD2-
1052 cell populations. Each dot represents the methylation % of single CpGs ($n = 57$) from 4
1053 genomic regions analysed by targeted bisulfite sequencing, and the red bars represent the

1054 mean methylation level for each sample. Statistical difference between samples was
1055 analysed by a Kruskal-Wallis test, with Dunn's multiple comparisons post-hoc test. N.S.
1056 denotes not significant ($p > 0.05$). c) qRT-PCR for naïve pluripotency genes in SUSD2+ and
1057 SUSD2- cell populations early and late transition hESCs. Bars represent the mean of three
1058 technical replicates and error bars represent SEM. Statistical significance between samples
1059 was analysed with a one-way ANOVA with a Bonferroni post-hoc test comparing all samples
1060 to each other. Human GAPDH was used to normalise expression. * $P < 0.05$, ** $P < 0.01$,
1061 *** $P < 0.001$, **** $P < 0.0001$, N.S. denotes not significant. Source data are provided as a
1062 Source Data file.

1063

1064 **Figure 5. Bivalent CGIs that lose H3K4me3 gain DNA methylation.** a) Pie chart showing the
1065 proportion of bivalent sites that do and do not gain methylation in naïve hESCs, and gene
1066 ontology analysis of these bivalent hypermethylated and non-hypermethylated genes
1067 respectively. b) ChIP-qPCR enrichment of H3K4me3 and c) H3K27me3 are shown for 6
1068 candidate bivalent regions (possessing both H3K4me3 and H3K27me3 histone modifications
1069 in primed cells) that become hypermethylated, and for 3 candidate bivalent regions that fail
1070 to become hypermethylated during resetting (bottom row). Data is show as the signal
1071 enrichment relative to the input sample with bars representing the mean of 2 independent
1072 experiments. Statistical difference between samples was analysed by a one-way ANOVA
1073 test, with Bonferroni post-hoc test of each time point compared to primed hESC. * indicates
1074 $p < 0.05$. Source data are provided as a Source Data file.

1075

1076 **Figure 6. Transcription and pluripotency factors influence hypermethylation.** a) A selection
1077 of the transcription factors with motifs enriched in bivalent hypermethylated regions, with
1078 bivalent non-hypermethylated regions used as a background control. Motif analysis was
1079 performed using the analysis of motif enrichment (AME) tool on the MEME suite. b)
1080 Volcano plots showing the difference in protein expression in early transition (72h and 1W)
1081 hESCs compared to primed hESCs (Supplementary Table 6). Each dot represents the \log_2
1082 fold change based on three biological replicates. Statistical difference between samples was
1083 analysed by a two-tailed student's t-test, corrected for multiple testing. Red dots indicate
1084 statistically significant changes ($\text{adj}P < 0.05$). Proteins of interest are highlighted with
1085 coloured and labelled symbols. c) Plot showing the % methylation for 4 different
1086 transcription factor knock downs and a non-silencing control. Data for each sample are an
1087 average of 2 independent shRNA knock downs. Each dot represents the methylation % of a
1088 single CpG analysed by targeted bisulfite sequencing, and the red lines represent the mean
1089 % methylation level for each sample from 5 genomic regions (SOX15 targets; $n = 64$ CpGs),
1090 or 4 genomic regions (NFKB1 regions; $n = 67$ CpGs, or regions without TF binding sites for
1091 any of the 4 TFs; $n = 66$ CpGs). Statistical difference between samples was analysed by a
1092 two-way ANOVA test, with Bonferroni post-hoc test of each TF knock down compared to the
1093 control. * $P < 0.05$, ** $P < 0.01$, and *** $P < 0.001$, **** $P < 0.0001$. N.S. denotes not significant. d)
1094 Targeted bisulfite-sequencing of 3 genomic regions. Each square represents the methylation
1095 % indicated by the colour key of a single CpG. The first column represents data from primed
1096 hESCs, and the subsequent columns represent data from early transition hESCs cultured in a
1097 variety of culture conditions indicated by the +/- symbols above. Source data are provided
1098 as a Source Data file.

1099

1100 **Figure 7. Resetting-associated hypermethylation is mirrored in cancer.** a) Differences in
1101 mean methylation level between matched normal tissue and tumour samples of bivalent
1102 CpGs identified as hypermethylated (N=23123) or not hypermethylated (N=25977) during
1103 the transition to the naïve state in hESCs. Data is presented for 592 individuals, separated by
1104 tumour location. P-values determined via paired Wilcoxon test (two-sided). CpGs used for
1105 analysis were filtered for those that are unmethylated in primed hESCs ($\beta < 0.3$).
1106 Lines=median; Box=25th–75th percentile; whiskers=1.5× interquartile range from box. b)
1107 Overlap of CpGs hypermethylated in cancers with bivalent CpGs hypermethylated during
1108 hESC resetting. Values for all cancers were generated by testing all cancer samples against
1109 all normal samples. P-values are determined using Fisher’s exact tests (two-sided). c)
1110 Heatmap of H3K27me3 distribution in 250bp windows around each hESC Hypermethylated
1111 CpG in normal colon (N=26180), ordered by total H3K27me3 reads within 5kb. Peaks are
1112 taken from ENCODE (Supplementary Table 8). Scale is in normalised reads per million. d)
1113 Selected genomic regions containing hESC Hypermethylated CpGs marked by H3K27me3 in
1114 the normal colon (highlighted light orange). Mean methylation in normal colon and
1115 colorectal tumours from TCGA, as well as the difference between them are shown in black.
1116 hESC Hypermethylated CpGs are shown in blue, and H3K27me3 ChIP reads from ENCODE
1117 are shown in orange (Supplementary Table 8). Scale for methylation data is 0 to 100%. Scale
1118 for methylation difference is -50% to 50%. Scale for H3K27me3 is 0 to 30 normalised reads
1119 per million reads.
1120
1121
1122