

1 **Title:** The emergence of the brain non-CpG methylation system in vertebrates

2

3 **Authors:** Alex de Mendoza^{1,2,3}, Daniel Poppe^{1,2}, Sam Buckberry^{1,2}, Jahnvi Pflueger^{1,2}, Caroline
4 B. Albertin^{4,5}, Tasman Daish⁶, Stephanie Bertrand⁷, Elisa de la Calle-Mustienes⁸, Jose Luis
5 Gomez-Skarmeta^{8,†}, Joseph R. Nery⁹, Joseph R. Ecker^{9,10}, Boris Baer¹¹, Clifton W. Ragsdale⁵,
6 Frank Grützner⁶, Hector Escriva⁷, Byrappa Venkatesh¹², Ozren Bogdanovic^{1,2,13,14}, Ryan Lister^{1,2,*}

7

8 **Affiliations**

9 ¹Australian Research Council Centre of Excellence in Plant Energy Biology, School of Molecular
10 Sciences, The University of Western Australia, Perth, Western Australia, Australia

11 ²Harry Perkins Institute of Medical Research, Perth, Western Australia, Australia

12 ³Queen Mary, University of London. School of Biological and Chemical Sciences, London, United
13 Kingdom

14 ⁴Eugene Bell Center for Regenerative Biology and Tissue Engineering, Marine Biological
15 Laboratory, Woods Hole, MA 02543, USA

16 ⁵Department of Neurobiology, University of Chicago, Chicago, Illinois 60637, USA

17 ⁶School of Biological Sciences, University of Adelaide, Adelaide, South Australia 5005, Australia

18 ⁷Sorbonne Université, CNRS, Biologie Intégrative des Organismes Marins (BIOM), Observatoire
19 Océanologique, Banyuls-sur-Mer, France.

20 ⁸Centro Andaluz de Biología del Desarrollo (CABD), CSIC-Universidad Pablo de Olavide-Junta
21 de Andalucía, Seville, Spain

22 ⁹Genomic Analysis Laboratory, Salk Institute for Biological Studies, La Jolla, California, USA

23 ¹⁰Howard Hughes Medical Institute, The Salk Institute for Biological Studies, La Jolla, California,
24 USA

25 ¹¹Center for Integrative Bee Research, Department of Entomology, The University of California
26 Riverside

27 ¹²Comparative Genomics Laboratory, Institute of Molecular and Cell Biology, A*STAR, Biopolis,
28 Singapore 138673

29 ¹³Genomics and Epigenetics Division, Garvan Institute of Medical Research, Sydney, New South
30 Wales, Australia

31 ¹⁴School of Biotechnology and Biomolecular Sciences, Faculty of Science, University of New
32 South Wales, Sydney, New South Wales, Australia

33 † Deceased September 16th, 2020.

34 * Corresponding author: Ryan Lister ryan.lister@uwa.edu.au

35

36

37 **Abstract**

38 Mammalian brains feature exceptionally high levels of non-CpG DNA methylation alongside the
39 canonical form of CpG methylation. Non-CpG methylation plays a critical regulatory role in
40 cognitive function, which is mediated by the binding of MeCP2, the transcriptional regulator that
41 when mutated causes Rett Syndrome. However, it is unclear if the non-CpG neural methylation
42 system is restricted to mammalian species with complex cognitive abilities or has deeper
43 evolutionary origins. To test this, we investigated brain DNA methylation across 12 distant animal
44 lineages, revealing that non-CpG methylation is restricted to vertebrates. We discovered that in
45 vertebrates, non-CpG methylation is enriched within a highly conserved set of developmental
46 genes transcriptionally repressed in adult brains, indicating that it demarcates a deeply conserved
47 regulatory program. Concomitantly, we found that the writer of non-CpG methylation, DNMT3A,
48 and the reader, MeCP2, originated at the onset of vertebrates as a result of the ancestral
49 vertebrate whole genome duplication. Together, we demonstrate how this novel layer of
50 epigenetic information assembled at the root of vertebrates and gained new regulatory roles
51 independent of the ancestral form of the canonical CpG methylation. This suggests the

52 emergence of non-CpG methylation may have fostered the evolution of sophisticated cognitive
53 abilities found in the vertebrate lineage.

54

55 **Main text**

56 **Introduction**

57 Cytosine DNA methylation (mC) is the most abundant base modification in animal genomes^{1,2}. In
58 vertebrates, most of the CpG dinucleotides (> 80%) in the genome are methylated³. In contrast,
59 most invertebrates show sparse methylation, where most CpG methylation accumulates on
60 transcribed gene bodies^{4,5}. However, cytosine methylation can also occur in the CpH (where H is
61 C, A, or T) dinucleotide context. In mammals, CpH methylation is mostly restricted to a few tissues
62 and cell types⁶, such as embryonic stem cells, neurons, and muscle. Embryonic stem cells display
63 CpH methylation enriched on transcribed gene bodies, while neural tissues accumulate high
64 levels of CpH methylation on transcriptionally silent genes⁷⁻¹². CpH methylation is deposited *de*
65 *novo* by the DNMT3A or DNMT3B methyltransferases, and unlike CpG methylation, is not
66 maintained after genome replication by the DNA methyltransferase DNMT1¹¹. Thus, post-mitotic
67 neurons can accumulate CpH methylation since they do not undergo genome replication. In
68 contrast to CpG methylation, CpH methylation is accumulated in the brain after birth, coinciding
69 with synaptogenesis and synaptic pruning^{7,13}. Furthermore, CpH methylation shows cell-type
70 specific patterns in distinct neurons and glia^{7,8,14}, and is the most abundant form of DNA
71 methylation in neurons. Most importantly, CpH methylation is bound by MeCP2, a highly
72 expressed transcriptional regulator that can cause Rett syndrome, a strong autistic phenotype,
73 when mutated^{15,16}. Similarly, mutations in DNMT3A and abnormal cytosine methylation are also
74 linked to neurological diseases¹⁷. Therefore, the role of DNA methylation and CpH methylation in
75 neural maturation and cognitive functions is well established in mammals. To date, CpH
76 methylation has been observed in the brain of human, mouse, and a songbird^{7,18,19}, thus the roles
77 of this unique epigenomic feature could potentially be linked to complex brain functions. However,

78 neither the evolutionary origin of CpH methylation nor the molecular basis that allowed the
79 emergence of this new methylation context to appear has so far been unraveled.

80 The morphology of the vertebrate brain is highly conserved, with a tripartite organization
81 that is found from lampreys to mammals²⁰. However, the homology between the vertebrate brain
82 and that of distantly related invertebrates remains uncertain^{21,22}. Notwithstanding this, all animal
83 brains are mainly composed of neurons and glia, ectodermal-derived neural cell types that have
84 deep evolutionary roots²³. Thus, to understand the evolution of neural CpG and CpH methylation
85 and its relationship to cognitive complexity, here we study the evolution of neural methylation
86 within and outside the vertebrate lineage.

87

88 **Results**

89 **Brain CpG methylation recapitulates differences between vertebrates and invertebrates**

90 To investigate the evolution of neural DNA methylation, we gathered forebrain samples from
91 representative species of major vertebrate lineages. We generated whole genome bisulfite
92 sequencing (WGBS) data from adult forebrain regions for six vertebrate species (Fig. 1a),
93 including opossum, platypus, chicken, zebrafish, elephant shark and arctic lamprey, and we
94 reanalysed previously published datasets from another four^{7,18,24}. For invertebrates, we generated
95 new data for two lineages with highly complex brains and behaviours. As representatives of
96 insects, we generated WGBS data for honeybee whole brains from a queen. As a cephalopod
97 representative, we obtained material from the California two-spot octopus, for which we sampled
98 and performed WGBS for both the supraesophageal and the subesophageal brains. As an out-
99 group to vertebrates, we generated new data for neural tube material from the European
100 amphioxus. The anterior neural tube is homologous to and shares many epigenomic similarities
101 with the vertebrate brain^{25,26}. Therefore, this dataset comprises the broadest assessment of adult
102 neural DNA methylation to date, encompassing major animal phyla with highly complex brains.

103 To understand major differences in methylation across species, we first analysed CpG
104 methylation, since it is the preferred context for animal DNA methyltransferases²⁷. As previously
105 reported, vertebrates show higher CpG methylation levels than invertebrates (Fig. 1a)^{1,4}. The high
106 global levels of CpG methylation in vertebrate genomes have been proposed to correlate with the
107 size of the genome or its high level of repetitive content^{4,28}. However, the octopus genome is
108 larger and has comparable repeat content to some vertebrate species²⁹. Still, the octopus genome
109 shows typical invertebrate global methylation levels (~10% mCpG/CpG) and most CpGs in the
110 genome are unmethylated (Fig. 1a,b), thus contradicting previous hypotheses regarding the
111 evolutionary origin of hypermethylation in vertebrates. Additionally, hypermethylation (global
112 mCpG/CpG > 70%) is not found in all vertebrate samples. The arctic lamprey and both bird
113 species show lower levels of global methylation than other vertebrate species (Fig. 1a). These
114 vertebrate lower global methylation levels are explained by an overwhelming majority of
115 intermediately methylated CpG positions (Fig. 1b). Intermediate methylation observed in the arctic
116 lamprey brain coincides with previous observations from sperm, muscle and heart methylation
117 levels in another species of lamprey³⁰. Interestingly, intermediate methylation levels correspond
118 to very heterogeneous methylation at the read level, suggesting noisy inheritance of methylation
119 after cell division (Extended Data 1). Given the phylogenetic position of lampreys, the intermediate
120 methylation levels in this lineage might represent a middle step in the transition between the
121 mosaic methylomes of invertebrates to the fully methylated genomes of jawed vertebrates³⁰.
122 However, avian intermediate methylomes represent a secondary reduction since all earlier
123 splitting lineages show hypermethylation. The evolutionary causes of such reduction in
124 methylation are unclear, since genome size does not explain methylation levels, even within
125 vertebrates, given that elephant shark has higher methylation and a smaller genome than birds
126 (Fig. 1a). Surprisingly, lampreys and other cyclostomes have genomes enriched in CpG
127 dinucleotides, unlike any other vertebrate (Fig. 1a, Extended Data 2). In sum, the CpG methylation

128 landscape in the brain reflects known differences between vertebrate and invertebrate genomes,
129 yet challenges prior assumptions about the evolution of hypermethylation in vertebrates.

130

131 **Lamprey genomes are not affected by methyl-CpG hypermutability**

132 Methylated cytosines are known to be prone to deaminate into thymines³¹. This tendency towards
133 deamination makes CpG sites hotspots of mutability and genetic variation³². Furthermore,
134 methylated CpG mutability is believed to be responsible for the global depletion of CpG sites in
135 vertebrate genomes^{1,4}. To explore these observations, we first gathered global CpG dinucleotide
136 content in the sampled species (Figure 1a). Whereas all jawed vertebrates show strong depletions
137 of CpG dinucleotides, lampreys and other cyclostomes do not show such depletions. In fact, the
138 ratio of CpG dinucleotides in lamprey genomes is similar to that of species that lack cytosine DNA
139 methylation (Figure 1a). To further investigate this anomaly, we used WGBS data to identify
140 Single Nucleotide Variants (SNV) in all sampled species (Extended Data 2). All jawed vertebrates
141 showed a higher frequency of variants at CpG dinucleotides with respect to other dinucleotides.
142 However, the arctic lamprey did not show such an enrichment. The intermediate methylation
143 levels found in lamprey genomes could explain why CpG dinucleotides are not disproportionately
144 affected by mutagenesis and depleted as seen in other vertebrate lineages. However, avian
145 genomes also have intermediate methylation levels and still show archetypal global CpG
146 depletion and disproportionate variants on CpG sites. Therefore, how lampreys avoid or
147 compensate for methylation-derived mutagenesis remains unclear, yet could be linked to somatic
148 DNA elimination in this lineage³³.

149

150 **Brain CpH methylation is restricted to vertebrates**

151 To avoid methylation mutability confounding our measurements of CpH methylation, we first
152 discarded all CpH positions that showed evidence of being CpG dinucleotide variants in the
153 sequenced WGBS reads. We then measured global genomic methylation levels at CpA, CpT and

154 CpC dinucleotides for each species and compared these to the bisulfite non-conversion rates in
155 the unmethylated lambda DNA spike in control for each WGBS experiment (Fig. 2a). All
156 vertebrates showed CpA and CpT global methylation above non-conversion levels, whereas
157 invertebrates did not. As previously reported in mammals, CpA is the preferred context for non-
158 CpG methylation in all vertebrates, while CpC is rarely methylated^{7,34}. We next interrogated
159 whether there is a wider sequence context in which CpH methylation gets preferentially deposited,
160 as it occurs in mammals^{7,35}. We gathered the neighbouring positions from the 10,000 most highly
161 methylated CpH sites in each species, finding that the trinucleotide CAC and additional bases
162 conform to an overrepresented motif conserved across vertebrates (Fig. 2b). The flanking bases
163 surrounding the CpH sites coincide with the flanking sequence preference reported for
164 DNMT3A³⁶. This CpH flanking motif was not detectable in non-neural samples for elephant shark,
165 zebrafish or *Xenopus*, confirming that mCpH is not a bisulfite sequencing bias and mCpH neural
166 specificity extends beyond mammals (Extended Data 3). Similarly, the CpH flanking motif was not
167 detectable in invertebrates (Fig. 2b). Furthermore, methylation levels on the highest methylated
168 CpH sites were lower in amphioxus, honeybee, and octopus (mC/C < 20%) compared to any
169 vertebrate brain (Fig. 2c). Thus, invertebrate CpH methylation is likely to be a rare off-target
170 consequence of DNMT activity. In contrast, the robust mammalian neural CpH methylation levels
171 are conserved across the vertebrate lineage.

172

173 **CpH methylation is functionally decoupled from CpG methylation across vertebrates**

174 In mammalian brains, CpH methylation deposition does not fully recapitulate CpG methylation^{6,7}.
175 While CpG methylation is found on transcribed and silent gene bodies alike, CpH methylation is
176 depleted on transcriptionally active gene bodies in neurons. To test whether brain CpH
177 methylation anti-correlates with transcription, we classified genes in deciles of expression for each
178 species and assessed the corresponding gene body CpG and CpA methylation levels (Extended
179 Data 4). A clear anti-correlation pattern between transcription and CpA methylation was observed

180 for mammals, birds and the frog (Spearman's r). However, this anti-correlation was not evident in
181 opossum, zebrafish, elephant shark and lamprey. This lack of anti-correlation in these species
182 might respond to different cell-type compositions biasing the measurements. The proportion of
183 neurons versus glia depends on the exact brain region and varies in a species-specific manner³⁷,
184 and species with smaller brains might display higher cell-type heterogeneity in similarly sized
185 samples, as for instance birds have higher neuron densities than mammals³⁸. In fact, all four
186 species not showing CpA methylation anti-correlation with transcription show lower levels of CpH
187 methylation on the highest methylated CpH sites (Fig. 2C), which suggests a lower ratio of
188 neurons to glia. Another possible explanation is that the anti-correlation with transcription evolved
189 in tetrapods, and was secondarily lost in opossum. In contrast, CpG methylation also shows some
190 degree of anti-correlation with expression levels in most vertebrate brain samples, whereas
191 invertebrates show the typical positive correlation between CpG methylation and transcription.

192 Despite the existence of differences in cell type composition across the brains of different
193 species, we reasoned that common methylation patterns should be observable across species if
194 similar pathways are regulated in a similar manner across most neural cells. In fact, distinct brain
195 regions show similar CpH methylation patterns in mammals³⁹. Consistently, transcriptional and
196 enhancer landscapes at the organ and tissue level are conserved across vertebrates^{40,41}. To test
197 if methylation patterns are conserved, we classified all genes in each species into 10 deciles
198 based on the weighted average of CpG and CpA methylation along the gene body (Extended
199 Data 5). For each hypermethylated and hypomethylated gene subset (top and bottom decile), we
200 obtained Gene Ontology (GO) enrichments (Fig. 3, Supplementary Table 1). Hypomethylated
201 genes in the CpG context largely represent developmental genes, predominantly transcription
202 factors. Such genes are found in methylation canyons or valleys, where lack of methylation in the
203 gene body and surrounding regions is mediated by histone modifications such as H3K27me3 and
204 H3K4me3^{42,43}. These same GOs appear enriched in non-brain samples, suggesting that CpG
205 methylation valleys are shared across tissues (Extended Data 3). In contrast, highly methylated

206 genes in the CpG context did not show deeply conserved GO patterns, and the few GOs that
207 appear in more than one species have housekeeping functions. On the contrary, hypermethylated
208 genes in the CpA context belong to developmental functions across all vertebrates (Fig. 3), and
209 many are related to signaling pathways, cell adhesion, or cell differentiation. On the other hand,
210 genes with the lowest levels of CpA methylation have housekeeping functions. Unlike with CpG
211 methylation, non-brain samples do not recapitulate any of these CpA enrichments (Extended Data
212 3). However, CpA and CpG methylation patterns are not completely unlinked, since there is a
213 high degree of overlap between genes found in both the lowly methylated categories (Extended
214 Data 5), which suggests that methylation protection on hypomethylated genes occurs through
215 restricting access of DNA methyltransferases^{44,45}. However, the developmental genes that are
216 CpG hypomethylated and CpA hypermethylated show very little overlap, which is indicative of
217 differential removal or deposition of methylated cytosines occurring in these regions. Invertebrates
218 do not exhibit conservation of these patterns. Surprisingly, birds show higher conservation of GOs
219 for genes methylated in the CpA context than for the CpG context (Fig. 3), suggesting that CpG
220 methylation state is not maintained yet CpA methylation is deposited in a conserved set of genes.

221 To corroborate the functional patterns gathered by GO analysis, we measured the CpG
222 and CpA methylation levels of genes classified by gene family or function. Methylation levels on
223 transcription factors, signaling molecules, synaptic genes and ribosomal proteins (Supplementary
224 Fig. 1), showed overall consistent patterns with the GO analysis approach. Among the
225 orthologues found in the highly methylated CpA category across species (≥ 7 species,
226 Supplementary Table 2) there are signaling molecules (WNT16, BMP7) and transcription factors
227 (*FOXP2*, *EOMES/TBR2*, *GLI3*, *PROX1*, *SOX6*, *SALL1*) that have been previously shown to be
228 involved in neural progenitor cell maintenance and differentiation. Furthermore, these sets of
229 conserved CpA methylated genes show declining gene expression in adult stages in the brains
230 of mammals and birds compared to earlier developmental stages (Extended Data 6). Therefore,
231 CpA methylation accumulates on a conserved subset of developmental genes across the

232 vertebrate lineage, likely marking and contributing to silencing genes no longer required in the
233 fully developed adult brain.

234

235 **DNMT3A is the ancestral writer of neural CpH methylation in vertebrates**

236 Given that the establishment of CpH methylation coincided with the origin of vertebrates, a new
237 “writer” able to deposit CpH methylation should have also evolved concomitantly. In mammals
238 DNMT3A is responsible for neural CpH methylation^{10,13}, whereas CpH methylation in stem cells
239 is mediated by DNMT3B⁴⁶. To gain an evolutionary perspective on the distribution and origin of
240 these genes, we performed a phylogenetic analysis of DNMT3 enzymes in animals (Fig. 4a,
241 Extended Data 7). While invertebrate genomes typically contain a single DNMT3 gene, DNMT3A
242 and DNMT3B evolved at the root of vertebrates. DNMT3A and DNMT3B are located in syntenic
243 regions (Supplementary Fig. 2), confirming that they represent ohnologues: the paralogues
244 product of the ancestral two rounds of whole genome duplication (WGD) in vertebrates, as
245 previously reported^{47,48}. More unexpectedly, we found that DNMT3L, a degenerate paralogue with
246 non-catalytic methyltransferase domain⁴⁹, is present in two lamprey genomes and non-avian
247 reptiles, suggesting it might be the third ohnologue derived from the WGD (Fig. 4a, Extended
248 Data 7). However, not all DNMT3 ohnologues are widely retained across vertebrates; lampreys
249 and amphibians do not encode a DNMT3B copy (Fig. 4c). Given that both species have neural
250 CpH methylation, only DNMT3A orthologues can have a role as writers of CpH methylation in
251 these species. This, in turn, would support an ancestral role of DNMT3A in neural CpH
252 methylation. Consistently, zebrafish DNMT3A orthologues have been shown to be expressed in
253 brain tissues⁵⁰, and we detect DNMT3A transcripts in all vertebrate brain samples (Extended Data
254 7). Furthermore, the differential deposition patterns of CpH methylation in neural and stem cells
255 seems to have been mediated by changes in the PWWP domain in DNMT3A and DNMT3B
256 ohnologues after gene duplication (Supplementary Fig. 3). In summary, phylogeny and

257 distribution of DNMT3 paralogues suggests that DNMT3A was the ancestral “writer” of neural
258 CpH methylation in vertebrates.

259

260 **MeCP2 evolved as CpH reader from an ancestral DNA repair protein**

261 In mammals, the silencing capacity of CpH methylation has been attributed to the methylation
262 “reader” MeCP2³⁴. MeCP2 is a Methyl-CpG Binding Domain (MBD) containing protein, capable
263 of binding both methylated CpG and CpA dinucleotides^{34,51}. Furthermore, MeCP2 has been
264 shown to bind methylated CAC *in vitro* and *in vivo*, the most common context of CA methylation
265 in the brain^{51,52}. To better understand if CpH methylation co-evolved with MeCP2, we performed
266 a phylogenetic analysis of MBD proteins in animals (Fig. 4b, Extended Data 8). We found that
267 MeCP2 is deeply conserved in all vertebrates, including lampreys and chondrichthyans. MeCP2
268 branches as a sister group to the MBD4 family, as reported previously⁵³. MBD4 is conserved
269 across vertebrates, however, it is associated with DNA repair and not gene regulation⁵⁴ implying
270 that MeCP2 evolved as a duplication of an ancestral MBD4-like gene.

271 Besides the conserved MBD domain, MeCP2 has vastly diverged from the ancestral
272 invertebrate MBD4-like family. Whereas MBD4 contains a C-terminal glycosylase domain,
273 involved in mismatch repair of CpG dinucleotides, MeCP2 harbors a transcriptional repression
274 domain (TRD) and a C-terminal domain. The TRD domain is known to interact with multiple
275 histone modifying complexes associated with transcriptional silencing, such as Sin3, CoREST
276 and N-CoR^{55–57}. Most surprisingly, we found that many parts of the TRD are conserved beyond
277 vertebrates, being found in amphioxus MBD4/MECP2 orthologue (Fig. 4d, Extended Data 9),
278 which represents an intermediate step between MBD4 and MeCP2. Moreover, we found that
279 amphioxus transcribes a longer MBD4/MECP2 isoform that includes the glycosylase domain
280 involved in DNA repair and a shorter isoform lacking this domain. When assessing the isoform
281 usage across developmental stages and tissues in amphioxus, we found that the longer
282 MBD4/MECP2 isoform is preferentially expressed in developmental samples, whereas the short

283 version is predominant in adult tissues (Fig. 4e, Extended Data 10). This suggests that
284 MBD4/MECP2 in amphioxus has DNA repair functions predominantly during development, and
285 gene regulatory activities in adult tissues, and thus a dual function achieved using alternative
286 isoforms. In vertebrates, MeCP2 could have evolved and specialised as a consequence of gene
287 duplication linked to WGD, in which one of the MBD4-like duplicated loci lost the glycosylase
288 domain and gained a new C-terminal domain restricting it to gene regulation, whereas the other
289 copy lost the TRD domain and maintained the glycosylase domain, reverting to the pre-chordate
290 MBD4 domain architecture specialised in DNA repair.

291 These changes in protein structure and function must have imposed new functional
292 constraints on MeCP2. Since MeCP2 protein is expressed at histone levels and proposed to
293 partially substitute H1 in neurons⁵⁸, high levels of conservation in MeCP2 would be expected.
294 Consistently, we found that the MBD had 70% identity between lamprey and human MeCP2
295 orthologues, but only ~40% identity between MBD4 orthologues (Fig. 4b). In contrast, the MBD
296 domain in amphioxus MBD4/MECP2 is quite divergent from both MBD4 and MeCP2 (Extended
297 Data 9), suggesting that it does not have the capacity to bind CpH methylation like MeCP2, which
298 is consistent with the lack of CpH methylation in amphioxus neural tube (Fig. 2). Also influencing
299 DNA binding specificity, MeCP2 harbours two AT-hook motifs^{51,59}, which are conserved across
300 vertebrates and amphioxus MBD4/MECP2 (Extended Data 9). Thus, the binding specificities of
301 MeCP2 evolved in a stepwise manner, first gaining the AT-Hooks in the MBD4-like chordate
302 ancestor, and then acquiring the vertebrate MBD CpH methylation binding capacity that became
303 fixed after the subfunctionalization of MeCP2.

304

305 **Discussion**

306 Here we show how a functionally conserved new layer of epigenomic regulation was assembled
307 at the origin of the vertebrate lineage (Fig. 5). Neural CpH methylation evolved from gene
308 machinery ancestrally involved in CpG methylation. Despite CpH methylation having non-

309 overlapping distribution patterns with CpG methylation, CpH methylation is not fully independent
310 of CpG methylation, as it is deposited by DNMT3 enzymes able to methylate both sequence
311 contexts. Furthermore, CpH methylation is read by MeCP2, which also binds CpG methylation.
312 This scenario contrasts with that of plants, in which the different contexts of cytosine methylation
313 are fully uncoupled. Specialised DNMTs are responsible for CpG and CpH methylation deposition
314 and maintenance, and CpH methylation is largely restricted to transposable elements⁶⁰.
315 Nevertheless, there is extensive cross-talk between CpH and CpG methylation in plants, since
316 CpG gene body methylation is lost in species that have lost CMT3, a DNMT that methylates the
317 CHG context⁶¹. Instead, such a dual readout of CpG and CpH methylation seems to be absent
318 from invertebrate genomes, as CpH methylation is very scarce. Here we show how brain DNA
319 methylation in amphioxus, honeybee, and octopus are depleted of CpH methylation, as the low
320 levels of CpH methylation cannot be distinguished from non-conversion rates. Furthermore,
321 invertebrates lack a functionally consistent pattern of deposition of CpH on gene bodies as
322 observed in vertebrates. Therefore, it is likely that previous reports of CpH methylation in
323 invertebrate genomes are due to off-target activity of DNMT3^{62,63}, suggestive of CpH methylation
324 in invertebrates not being fully constituted into an autonomous epigenomic layer.

325 We hypothesize that the evolution of MeCP2 was instrumental in the fixation of CpH
326 methylation as a regulatory mark in the brain. CpH methylation could have originally accumulated
327 in neurons simply as a by-product of the lack of DNA replication. However, the capacity of MeCP2
328 to specifically read CpH methylation could have enabled and reinforced the silencing roles of CpH
329 methylation as a hub for chromatin silencing in a pathway partially independent of CpG
330 methylation. In fact, mice that preserve neural CpG methylation patterns but lack CpH methylation
331 recapitulate the transcriptional deregulation caused by MeCP2 loss³⁹, suggesting that CpH
332 methylation is what drives the specific roles of MeCP2 in the brain. Furthermore, mice encoding
333 a modified MeCP2 version lacking the ability to bind to methylated CpA (while still preserving the
334 capacity to bind to methylated CpGs) show Rett syndrome-like phenotypes⁵². Our finding that

335 CpH methylation and MeCP2 evolved concomitantly argues in favour of a key role of this
336 epigenomic layer in neural functions across the whole vertebrate lineage. Despite the fact that we
337 do not know at which developmental time point CpH methylation is deposited in most vertebrate
338 lineages, or the MeCP2 binding patterns in most species, we speculate that the CpH roles in
339 neural maturation and memory formation described in mammals could extend to all vertebrates.

340 Recent evidence suggests that the ancestral whole genome duplication may not have had
341 an impact on the evolution of CpG hypermethylation in vertebrates⁶⁴, however, it allowed the
342 emergence of neural CpH methylation. DNMT3 paralogues that are specialised in different
343 functions emerged after duplication, as exemplified by DNMT3A methylating CAC trinucleotides
344 in neural tissues whereas DNMT3B methylates CAG trinucleotides in stem cells⁴⁶. In the case of
345 MeCP2 and MBD4, the duplication allowed the specialisation of both copies to perform unique
346 functions, which was only partially attained in amphioxus through differential usage of isoforms,
347 as previously observed for a vertebrate neural-specific splicing factor⁶⁵. Therefore, our work
348 unveils the stepwise assembly of a critical regulatory novelty in vertebrate brains. This novelty
349 likely had an impact on the complexity of behaviours and cognitive processes found across the
350 vertebrate lineage.

351

352 **Methods**

353

354 **Brain DNA collection**

355 Arctic lamprey (*Lethenteron camtschaticum*) and elephant shark (*Callorhynchus milii*) forebrains
356 were collected from frozen samples, belonging to adult animals collected in Hokkaido, Japan and
357 Queenscliff, Victoria, Australia respectively. Chicken (*Gallus gallus*) and zebrafish (*Danio rerio*)
358 forebrains were collected from adult individuals reared in the CABD, Spain, approved by the
359 Ethical Committees from the University Pablo de Olavide, CSIC and the Andalucían government.
360 The platypus (*Ornithorhynchus anatinus*) frontal lobe cortex and the gray short-tailed opossum

361 (*Monodelphis domestica*) brain samples were obtained from adult male frozen samples according
362 to the University of Adelaide biosafety and ethics committee regulations (Institutional Biosafety
363 Committee, Dealing ID 12713, permits ID1111998.2, NPWS A193 and ID1814535.1).
364 Mediterranean amphioxus (*Branchiostoma lanceolatum*) neural tubes were dissected from 6
365 adults collected in Argeles-sur-Mer, France with special permission provided by the Prefect of
366 Region Provence Alpes Côte d'Azur. For the honeybee (*Apis mellifera*), a whole brain from an
367 adult egg laying queen was collected at the University of Western Australia. California two-spot
368 octopus (*Octopus bimaculoides*) samples were obtained from a single adult female octopus in
369 compliance with the EU Directive 2010/63/EU guidelines on cephalopod use and the University
370 of Chicago Animal Care and Use Committee. Both the supraesophageal and subesophageal
371 brains from the octopus were dissected as previously described²⁹. To purify genomic DNA,
372 DNeasy Blood and tissue Kit (Qiagen) and phenol-chloroform DNA extraction methods were
373 used.

374

375 **Whole Genome Bisulfite Sequencing**

376 We followed the MethylC-seq protocol for library preparation⁶⁶. In brief, for each species, 500 ng
377 to 1 µg of brain genomic DNA was mixed with 0.1% to 0.5% (w/w) of unmethylated lambda phage
378 genomic DNA. The mixed DNA was sheared into 200 bp fragments using a Covaris Sonicator
379 S220. Then methylated Illumina adaptors (Nextflex Bisulfite-seq adaptors, BIOO scientific) were
380 ligated to sheared DNA, and bisulfite conversion was performed using EZ DNA Methylation-Gold
381 kit (Zymo Research) following the manufacturer's instructions. After bisulfite treatment, DNA was
382 purified and amplified using universal Illumina primers and KAPA HiFi HotStart Uracil+ DNA
383 polymerase (Kapa Biosystems). The honeybee library was obtained using the same protocol with
384 minor modifications, MethylCode Bisulfite Conversion Kit (Thermo Fisher) was used for bisulfite
385 conversion and the PfuTurbo Cx Hotstart DNA Polymerase (Agilent) was used for library
386 amplification. All libraries but the honeybee and amphioxus samples were sequenced in a Illumina

387 HiSeq 1500 instrument in single-end mode, with reads spanning 100 bp. The honeybee samples
388 were sequenced with an Illumina Genome Analyzer Ix in single-end mode, with reads spanning
389 84 bp, and amphioxus were sequenced in a NovaSeq 6000 in a paired-end 28-87 bp format.

390

391 **Methylation analysis**

392 The newly generated WGBS datasets were complemented with available data from previous
393 studies^{7,18,24,67}, corresponding to the NCBI Sequence Read Archive (SRA) accessions
394 SRX314948 for 6 week old mouse frontal cortex, SRX306585 for 25 year old human frontal cortex,
395 SRX1002603 for zebrafish (*Danio rerio*) adult brain, SRX1162705 for *Xenopus tropicalis* adult
396 brain, SRX2645741 for elephant shark liver and SRX1064224 for great tit (*Parus major*) adult
397 whole brain. All WGBS reads were trimmed using fastp⁶⁸ with default parameters and mapped to
398 the reference genomes using BS-Seeker2⁶⁹ specifying Bowtie 2⁷⁰ as the aligner in end-to-end
399 mode. Duplicated reads were discarded using Sambamba⁷¹, unconverted reads were filtered out
400 using the XS:i:1 sam flag from BS-Seeker2, and methylation calls were obtained using
401 CGmapTools⁷². Previously processed WGBS datasets for *Ciona intestinalis* and the sea anemone
402 *Nematostella vectensis* were obtained from Gene Expression Omnibus (GEO) GSE19824⁷³ and
403 GSE124016⁶⁴.

404 Since methylated CpG sites are prone to deamination, after the deamination of a symmetric CpG
405 site it becomes a non-symmetric CpA site. Therefore, some CpA positions in the reference
406 genomes are likely to represent genetic variants in which individuals might have CpG
407 dinucleotides. Distinguishing those sites is crucial to accurately measure CpH methylation, to
408 avoid confounding variant hypermethylated CpG sites for CpA positions. Therefore, the
409 ATCGmap file resulting from CGmapTools was parsed with AWK to identify CpH sites with $\geq 20\%$
410 of reads supporting a guanine in the downstream position of a methylated cytosine. Those
411 positions were discarded from the final CGmap file.

412 Single Nucleotide Variants were obtained using CGmapTools 'snv' function (-m bayes --bayes-
413 dynamicP parameters) from the WGBS ATCGmap file. For each SNV position, the upstream and
414 downstream dinucleotides based on the reference genome were obtained using BEDTools⁷⁴.

415 To estimate methylation heterogeneity in each sample, we followed the Proportion of Discordant
416 Reads (PDR) measure previously proposed for heterogeneous tumour samples⁷⁵. We first
417 selected CpG positions for which coverage was ≥ 10 , and filtered for those that had at least 3 CpG
418 ± 40 bp around them. Then we selected 100,000 of these CpGs randomly in every genome
419 (sample function in R) and obtained the per read methylation levels on the reads that overlapped
420 these positions. We only retained CpGs that had at least 5 reads covering ≥ 4 CpGs. Fully
421 methylated and unmethylated reads were counted as concordant, whereas intermediate
422 methylation was counted as discordant.

423 CGmap files were imported into R using the bsseq package⁷⁶, and all methylation calculations
424 were performed using in-built functions getCoverage and getMeth. CpH methylation was initially
425 calculated for each dinucleotide context to obtain the global levels (mC/C), however, gene body
426 level calculations were restricted to CpA dinucleotides since it is the predominant context.

427 For each species, CpH positions were sorted by methylation level (mC/C), and the top 10,000
428 were selected to have a comparable number across species. The neighbouring regions were
429 obtained using BEDTools in a strand-specific manner, and collapsed into sequence motifs with
430 ggseqlogo in R⁷⁷.

431 Protein-coding genes were classified into 10 deciles according to CpA and CpG methylation levels
432 along the gene body. Gene body methylation level measurements were obtained from the
433 weighted average of all cytosine calls in a given region divided by the total amount of coverage
434 in the C positions. Genes without enough covered CpG positions (≥ 30) and mean coverage ($\geq 4\times$)
435 were discarded.

436

437 **Gene Ontology enrichments**

438 Gene Ontology (GO) enrichments were obtained using g:Profiler⁷⁸ gProfileR R package, using
439 ensembl gene ids. For the arctic lamprey and the elephant shark, which were not present on the
440 g:Profiler database, OrthoFinder⁷⁹ was used to obtain orthology relationships with human genes.
441 Then, gene ids from both species and each decile were converted to human gene ids, which were
442 used to obtain GO enrichments using g:Profiler with 'hsapiens', limiting the background to all the
443 human genes detected in each orthology search. Significance was corrected with the g:Profiler
444 inbuilt g:SCS algorithm. The final set of GOs shown in Fig. 3 represent GOs that are enriched in
445 the maximum number of species and are not non-redundant. The full list of GOs and KEGG
446 pathways for each species and comparison are found in Supplementary Table 1.

447

448 **RNA-seq analysis**

449 Brain RNA-seq reads from previous publications^{7,18,25,29,40,65,80,81} were downloaded from SRA.
450 SRX314972 was used for human adult frontal cortex, SRX314992 was used for mouse adult
451 prefrontal cortex, SRX081894 for opossum brain, SRX081882 for the platypus brain, SRX081869
452 for the chicken brain, SRX904626 for the great tit brain, SRX191164 for *Xenopus tropicalis* brain,
453 SRX4184230 for zebrafish adult forebrain, SRX154851 for elephant shark brain, SRX2267405
454 for the Arctic lamprey brain, SRX1045432 for the octopus supraesophageal brain, and
455 PRJNA416866 for all amphioxus tissues. For the honeybee brain, we extracted matched DNA
456 and RNA samples from workers and queens, using a Trizol extraction protocol and prepared
457 Illumina stranded TruSeq RNA-seq libraries, which were sequenced on an Illumina Genome
458 Analyzer IIX.

459 Kallisto⁸² was then used to quantify gene expression, based on the canonical isoform for each
460 gene as per ENSEMBL annotations. For genomes without ENSEMBL annotation, we used the
461 isoform that encoded the longest open reading frame.

462 Developmental time-series from human, mouse, opossum and chicken were downloaded from
463 <https://apps.kaessmannlab.org/evodevoapp/>⁸³, gene expression was standardized for each gene
464 dividing the RPKM value against the maximum level of expression of that given gene.

465 To determine isoform usage in amphioxus MBD4 locus, we gathered the non-overlapping regions
466 between the short and the long MBD4 isoforms, added 100 padding N bases (to allow paired-end
467 sequencing mapping) and made a transcriptome index using Kallisto⁸², which was also used to
468 quantify isoform abundance without using reads from the common sequence between isoforms.

469

470 **Gene search and phylogeny**

471 MBD family genes were searched using HMMER3⁸⁴ with the PFAM PF01429 model against the
472 proteomes of a representative subset of animal genomes (Supplementary Table 3). Hits were
473 extracted and aligned with MAFFT⁸⁵ in LINS-I mode, and an initial pruning of the alignment was
474 performed to avoid members of the SETDB1/2 and BAZ2A/B families, since the MBD domain in
475 these family is derived and accumulates an excess of mutations. The resulting alignment was
476 then trimmed manually, to maximize the number of positions on the MBD domain and avoiding
477 spurious aligned regions. The resulting alignment was then used in IQ-TREE⁸⁶ to obtain maximum
478 likelihood phylogenetic reconstruction, letting the software to choose the best fitting substitution
479 model (-m TEST) and obtaining 100 non-parametric bootstrap replications to compute nodal
480 supports. Protein domain architectures for each sequence were obtained using HMMER3 with
481 the PFAM A database using the “hmmscan” program. MECP2 domains not defined in PFAM were
482 obtained from previous publications describing the TRD and CTD domains^{15,55}. TRD and CTD
483 alignments spanning all vertebrate major lineages were used to generate HMM models with
484 HMMER3 hmmbuild program, and were searched using hmmsearch against the selected animal
485 proteomes.

486 For obtaining DNMT3 sequences, we used BLASTP search using human DNMT3A as query
487 against the proteomes of all species, selecting the best hits for each species. For species where

488 we could not find a specific ohnologue, we searched in NCBI against the whole clade using
489 BLASTP (e.g. DNMT3B in amphibians) to certify that absence is not due to genome assembly
490 incompleteness. Similarly, DNMT3L was searched using BLASTP in NCBI against all lineages
491 except mammals, to detect ohnologues in all reptilian lineages (turtles, crocodylians and
492 squamates) except birds. The resulting sequences were aligned with MAFFT in EINS-I mode and
493 trimmed using TrimAL (-automated1). The phylogenetic tree was computed as for MBDs.
494 PWWP alignments were obtained from a subset of full length DNMT3 sequences, using one
495 representative species for each lineage. The sequences were aligned using MAFFT LINS-I mode
496 and the sequence logos were obtained using ggseqlogo in R. The alignments were visualised
497 using Geneious software.

498

499

500 **References**

- 501 1. Schübeler, D. Function and information content of DNA methylation. *Nature* **517**, 321–326
502 (2015).
- 503 2. Luo, C., Hajkova, P. & Ecker, J. R. Dynamic DNA methylation: In the right place at the right
504 time. *Science* **361**, 1336–1340 (2018).
- 505 3. Bird, A. DNA methylation patterns and epigenetic memory. *Genes Dev.* **16**, 6–21 (2002).
- 506 4. Suzuki, M. M. & Bird, A. DNA methylation landscapes: provocative insights from
507 epigenomics. *Nat. Rev. Genet.* **9**, 465–476 (2008).
- 508 5. de Mendoza, A., Lister, R. & Bogdanovic, O. Evolution of DNA Methylome Diversity in
509 Eukaryotes. *J. Mol. Biol.* (2019) doi:10.1016/j.jmb.2019.11.003.
- 510 6. He, Y. & Ecker, J. R. Non-CG Methylation in the Human Genome. *Annu. Rev. Genomics*
511 *Hum. Genet.* **16**, 55–77 (2015).
- 512 7. Lister, R. *et al.* Global epigenomic reconfiguration during mammalian brain development.

- 513 *Science* **341**, 1237905–1237905 (2013).
- 514 8. Mo, A. *et al.* Epigenomic Signatures of Neuronal Diversity in the Mammalian Brain. *Neuron*
515 **86**, 1369–1384 (2015).
- 516 9. Lister, R. *et al.* Hotspots of aberrant epigenomic reprogramming in human induced
517 pluripotent stem cells. *Nature* **471**, 68–73 (2011).
- 518 10. Guo, J. U. *et al.* Distribution, recognition and regulation of non-CpG methylation in the adult
519 mammalian brain. *Nat. Neurosci.* **17**, 215–222 (2014).
- 520 11. Ziller, M. J. *et al.* Genomic distribution and inter-sample variation of non-CpG methylation
521 across human cell types. *PLoS Genet.* **7**, e1002389 (2011).
- 522 12. Gabel, H. W. *et al.* Disruption of DNA-methylation-dependent long gene repression in Rett
523 syndrome. *Nature* (2015) doi:10.1038/nature14319.
- 524 13. Stroud, H. *et al.* Early-Life Gene Expression in Neurons Modulates Lasting Epigenetic
525 States. *Cell* **171**, 1151–1164.e16 (2017).
- 526 14. Luo, C. *et al.* Single-cell methylomes identify neuronal subtypes and regulatory elements in
527 mammalian cortex. *Science* **357**, 600–604 (2017).
- 528 15. Amir, R. E. *et al.* Rett syndrome is caused by mutations in X-linked MECP2, encoding
529 methyl-CpG-binding protein 2. *Nat. Genet.* **23**, 185–188 (1999).
- 530 16. Lyst, M. J. & Bird, A. Rett syndrome: a complex disorder with simple roots. *Nat. Rev.*
531 *Genet.* **16**, 261–275 (2015).
- 532 17. Tatton-Brown, K. *et al.* Mutations in the DNA methyltransferase gene DNMT3A cause an
533 overgrowth syndrome with intellectual disability. *Nat. Genet.* **46**, 385–388 (2014).
- 534 18. Laine, V. N. *et al.* Evolutionary signals of selection on cognition from the great tit genome
535 and methylome. *Nat. Commun.* **7**, 10474 (2016).
- 536 19. Derks, M. F. L. *et al.* Gene and transposable element methylation in great tit (*Parus major*)
537 brain and blood. *BMC Genomics* **17**, 332 (2016).
- 538 20. Sugahara, F. *et al.* Evidence from cyclostomes for complex regionalization of the ancestral

- 539 vertebrate brain. *Nature* (2016) doi:10.1038/nature16518.
- 540 21. Roth, G. Convergent evolution of complex brains and high intelligence. *Philos. Trans. R.*
541 *Soc. Lond. B Biol. Sci.* **371**, (2015).
- 542 22. Holland, L. Z. *et al.* Evolution of bilaterian central nervous systems: a single origin?
543 *Evodevo* **4**, 27 (2013).
- 544 23. Arendt, D., Tosches, M. A. & Marlow, H. From nerve net to nerve ring, nerve cord and
545 brain--evolution of the nervous system. *Nat. Rev. Neurosci.* **17**, 61–72 (2016).
- 546 24. Bogdanović, O. *et al.* Active DNA demethylation at enhancers during the vertebrate
547 phylotypic period. *Nat. Genet.* **48**, 417–426 (2016).
- 548 25. Marlétaz, F. *et al.* Amphioxus functional genomics and the origins of vertebrate gene
549 regulation. *Nature* **564**, 64–70 (2018).
- 550 26. Albuixech-Crespo, B. *et al.* Molecular regionalization of the developing amphioxus neural
551 tube challenges major partitions of the vertebrate brain. *PLoS Biol.* **15**, e2001573 (2017).
- 552 27. Lyko, F. The DNA methyltransferase family: a versatile toolkit for epigenetic regulation. *Nat.*
553 *Rev. Genet.* **19**, 81–92 (2018).
- 554 28. Mugal, C. F., Arndt, P. F., Holm, L. & Ellegren, H. Evolutionary consequences of DNA
555 methylation on the GC content in vertebrate genomes. *G3* **5**, 441–447 (2015).
- 556 29. Albertin, C. B. *et al.* The octopus genome and the evolution of cephalopod neural and
557 morphological novelties. *Nature* (2015) doi:10.1038/nature14668.
- 558 30. Zhang, Z. *et al.* Genome-wide and single-base resolution DNA methylomes of the Sea
559 Lamprey (*Petromyzon marinus*) Reveal Gradual Transition of the Genomic Methylation
560 Pattern in Early Vertebrates. *bioRxiv* 033233 (2015) doi:10.1101/033233.
- 561 31. Shen, J. C., Rideout, W. M., 3rd & Jones, P. A. The rate of hydrolytic deamination of 5-
562 methylcytosine in double-stranded DNA. *Nucleic Acids Res.* **22**, 972–976 (1994).
- 563 32. Pfeifer, G. P. Mutagenesis at methylated CpG sequences. *Curr. Top. Microbiol. Immunol.*
564 **301**, 259–281 (2006).

- 565 33. Smith, J. J. *et al.* The sea lamprey germline genome provides insights into programmed
566 genome rearrangement and vertebrate evolution. *Nat. Genet.* (2018) doi:10.1038/s41588-
567 017-0036-1.
- 568 34. Kinde, B., Gabel, H. W., Gilbert, C. S., Griffith, E. C. & Greenberg, M. E. Reading the
569 unique DNA methylation landscape of the brain: Non-CpG methylation, hydroxymethylation,
570 and MeCP2. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 6800–6806 (2015).
- 571 35. Xie, W. *et al.* Base-resolution analyses of sequence and parent-of-origin dependent DNA
572 methylation in the mouse genome. *Cell* **148**, 816–831 (2012).
- 573 36. Wienholz, B. L. *et al.* DNMT3L modulates significant and distinct flanking sequence
574 preference for DNA methylation by DNMT3A and DNMT3B in vivo. *PLoS Genet.* **6**,
575 e1001106 (2010).
- 576 37. Herculano-Houzel, S. The glia/neuron ratio: how it varies uniformly across brain structures
577 and species and what that means for brain physiology and evolution. *Glia* **62**, 1377–1391
578 (2014).
- 579 38. Olkowitz, S. *et al.* Birds have primate-like numbers of neurons in the forebrain. *Proc. Natl.*
580 *Acad. Sci. U. S. A.* **113**, 7255–7260 (2016).
- 581 39. Clemens, A. W. *et al.* MeCP2 Represses Enhancers through Chromosome Topology-
582 Associated DNA Methylation. *Mol. Cell* **77**, 279–293.e8 (2020).
- 583 40. Brawand, D. *et al.* The evolution of gene expression levels in mammalian organs. *Nature*
584 **478**, 343–348 (2011).
- 585 41. Villar, D. *et al.* Enhancer Evolution across 20 Mammalian Species. *Cell* **160**, 554–566
586 (2015).
- 587 42. Jeong, M. *et al.* Large conserved domains of low DNA methylation maintained by Dnmt3a.
588 *Nat. Genet.* **46**, 17–23 (2014).
- 589 43. Xie, W. *et al.* Epigenomic analysis of multilineage differentiation of human embryonic stem
590 cells. *Cell* **153**, 1134–1148 (2013).

- 591 44. Sendžikaitė, G., Hanna, C. W., Stewart-Morgan, K. R., Ivanova, E. & Kelsey, G. A DNMT3A
592 PWWP mutation leads to methylation of bivalent chromatin and growth retardation in mice.
593 *Nat. Commun.* **10**, 1884 (2019).
- 594 45. Heyn, P. *et al.* Gain-of-function DNMT3A mutations cause microcephalic dwarfism and
595 hypermethylation of Polycomb-regulated regions. *Nat. Genet.* **51**, 96–105 (2019).
- 596 46. Lee, J.-H., Park, S.-J. & Nakai, K. Differential landscape of non-CpG methylation in
597 embryonic stem cells and neurons caused by DNMT3s. *Sci. Rep.* **7**, 11295 (2017).
- 598 47. Albalat, R., Martí-Solans, J. & Cañestro, C. DNA methylation in amphioxus: from ancestral
599 functions to new roles in vertebrates. *Brief. Funct. Genomics* **11**, 142–155 (2012).
- 600 48. Liu, J., Hu, H., Panserat, S. & Marandel, L. Evolutionary history of DNA methylation related
601 genes in chordates: new insights from multiple whole genome duplications. *Sci. Rep.* **10**,
602 970 (2020).
- 603 49. Greenberg, M. V. C. & Bourc'his, D. The diverse roles of DNA methylation in mammalian
604 development and disease. *Nat. Rev. Mol. Cell Biol.* **20**, 590–607 (2019).
- 605 50. Smith, T. H. L., Collins, T. M. & McGowan, R. A. Expression of the dnmt3 genes in
606 zebrafish development: similarity to Dnmt3a and Dnmt3b. *Dev. Genes Evol.* **220**, 347–353
607 (2011).
- 608 51. Lager, S. *et al.* MeCP2 recognizes cytosine methylated tri-nucleotide and di-nucleotide
609 sequences to tune transcription in the mammalian brain. *PLoS Genet.* **13**, e1006793
610 (2017).
- 611 52. Tillotson, R. *et al.* Neuronal non-CG methylation is an essential target for MeCP2 function.
612 *bioRxiv* 2020.07.02.184614 (2020) doi:10.1101/2020.07.02.184614.
- 613 53. Albalat, R. Evolution of DNA-methylation machinery: DNA methyltransferases and methyl-
614 DNA binding proteins in the amphioxus *Branchiostoma floridae*. *Dev. Genes Evol.* **218**,
615 691–701 (2008).
- 616 54. Millar, C. B. *et al.* Enhanced CpG mutability and tumorigenesis in MBD4-deficient mice.

- 617 *Science* **297**, 403–405 (2002).
- 618 55. Lyst, M. J. *et al.* Rett syndrome mutations abolish the interaction of MeCP2 with the
619 NCoR/SMRT co-repressor. *Nat. Neurosci.* **16**, 898–902 (2013).
- 620 56. Jones, P. L. *et al.* Methylated DNA and MeCP2 recruit histone deacetylase to repress
621 transcription. *Nat. Genet.* **19**, 187–191 (1998).
- 622 57. Nan, X. *et al.* Transcriptional repression by the methyl-CpG-binding protein MeCP2
623 involves a histone deacetylase complex. *Nature* **393**, 386–389 (1998).
- 624 58. Skene, P. J. *et al.* Neuronal MeCP2 is expressed at near histone-octamer levels and
625 globally alters the chromatin state. *Mol. Cell* **37**, 457–468 (2010).
- 626 59. Klose, R. J. *et al.* DNA binding selectivity of MeCP2 due to a requirement for A/T
627 sequences adjacent to methyl-CpG. *Mol. Cell* **19**, 667–678 (2005).
- 628 60. Law, J. A. & Jacobsen, S. E. Establishing, maintaining and modifying DNA methylation
629 patterns in plants and animals. *Nat. Rev. Genet.* **11**, 204–220 (2010).
- 630 61. Bewick, A. J. *et al.* On the origin and evolutionary consequences of gene body DNA
631 methylation. *Proc. Natl. Acad. Sci. U. S. A.* **113**, 9111–9116 (2016).
- 632 62. Bonasio, R. *et al.* Genome-wide and caste-specific DNA methylomes of the ants
633 *Camponotus floridanus* and *Harpegnathos saltator*. *Curr. Biol.* **22**, 1755–1764 (2012).
- 634 63. Harris, K. D., Lloyd, J. P. B., Domb, K., Zilberman, D. & Zemach, A. DNA methylation is
635 maintained with high fidelity in the honey bee germline and exhibits global non-functional
636 fluctuations during somatic development. *Epigenetics Chromatin* **12**, 62 (2019).
- 637 64. de Mendoza, A. *et al.* Convergent evolution of a vertebrate-like methylome in a marine
638 sponge. *Nat Ecol Evol* **3**, 1464–1473 (2019).
- 639 65. Torres-Méndez, A. *et al.* A novel protein domain in an ancestral splicing factor drove the
640 evolution of neural microexons. *Nat Ecol Evol* **3**, 691–701 (2019).
- 641 66. Urich, M. a., Nery, J. R., Lister, R., Schmitz, R. J. & Ecker, J. R. MethylC-seq library
642 preparation for base-resolution whole-genome bisulfite sequencing. *Nat. Protoc.* **10**, 475–

- 643 483 (2015).
- 644 67. Peat, J. R., Ortega-Recalde, O., Kardailsky, O. & Hore, T. A. The elephant shark
645 methylome reveals conservation of epigenetic regulation across jawed vertebrates.
646 *F1000Res.* **6**, 526 (2017).
- 647 68. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor.
648 *Bioinformatics* **34**, i884–i890 (2018).
- 649 69. Guo, W. *et al.* BS-Seeker2: a versatile aligning pipeline for bisulfite sequencing data. *BMC*
650 *Genomics* **14**, 774–774 (2013).
- 651 70. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*
652 **9**, 357–359 (2012).
- 653 71. Tarasov, A., Vilella, A. J., Cuppen, E., Nijman, I. J. & Prins, P. Sambamba: fast processing
654 of NGS alignment formats. *Bioinformatics* **31**, 2032–2034 (2015).
- 655 72. Guo, W. *et al.* CGmapTools improves the precision of heterozygous SNV calls and
656 supports allele-specific methylation detection and visualization in bisulfite-sequencing data.
657 *Bioinformatics* **34**, 381–387 (2018).
- 658 73. Zemach, A., McDaniel, I. E., Silva, P. & Zilberman, D. Genome-wide evolutionary analysis
659 of eukaryotic DNA methylation. *Science* **328**, 916–919 (2010).
- 660 74. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic
661 features. *Bioinformatics* **26**, 841–842 (2010).
- 662 75. Landau, D. A. *et al.* Locally disordered methylation forms the basis of intratumor methylome
663 variation in chronic lymphocytic leukemia. *Cancer Cell* **26**, 813–825 (2014).
- 664 76. Hansen, K. D., Langmead, B. & Irizarry, R. A. BSmooth: from whole genome bisulfite
665 sequencing reads to differentially methylated regions. *Genome Biol.* **13**, R83 (2012).
- 666 77. Wagih, O. ggseqlogo: a versatile R package for drawing sequence logos. *Bioinformatics* **33**,
667 3645–3647 (2017).
- 668 78. Reimand, J., Kull, M., Peterson, H., Hansen, J. & Vilo, J. g:Profiler--a web-based toolset for

669 functional profiling of gene lists from large-scale experiments. *Nucleic Acids Res.* **35**,
670 W193–200 (2007).

671 79. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative
672 genomics. *Genome Biol.* **20**, 238 (2019).

673 80. Venkatesh, B. *et al.* Elephant shark genome provides unique insights into gnathostome
674 evolution. *Nature* **505**, 174–179 (2014).

675 81. Barbosa-Morais, N. L. *et al.* The evolutionary landscape of alternative splicing in vertebrate
676 species. *Science* **338**, 1587–1593 (2012).

677 82. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq
678 quantification. *Nat. Biotechnol.* **34**, 525–527 (2016).

679 83. Cardoso-Moreira, M. *et al.* Gene expression across mammalian organ development. *Nature*
680 **571**, 505–509 (2019).

681 84. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**, e1002195–
682 e1002195 (2011).

683 85. Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7:
684 Improvements in Performance and Usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).

685 86. Nguyen, L. T., Schmidt, H. A., Von Haeseler, A. & Minh, B. Q. IQ-TREE: A fast and
686 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol.*
687 *Evol.* **32**, 268–274 (2015).

688 87. Ross, S. E., Angeloni, A., Geng, F.S., de Mendoza, A. & Bogdanovic, O. Developmental
689 remodelling of non-CG methylation at satellite DNA repeats. *Nucleic Acids Res.* In press
690 (2020).

691 88. Baubec, T. *et al.* Genomic profiling of DNA methyltransferases reveals a role for DNMT3B
692 in genic methylation. *Nature* **520**, 243–247 (2015).

693 89. Dhayalan, A. *et al.* The Dnmt3a PWWP domain reads histone 3 lysine 36 trimethylation and
694 guides DNA methylation. *J. Biol. Chem.* **285**, 26114–26120 (2010).

695 90. Li, Y. *et al.* DNA methylation regulates transcriptional homeostasis of algal endosymbiosis
696 in the coral model *Aiptasia*. *Sci Adv* **4**, eaat2142 (2018).

697 91. Barau, J. *et al.* The novel DNA methyltransferase DNMT3C protects male germ cells from
698 transposon activity. *Science* **354**, 909–912 (2016).

699

700 **Data and materials availability**

701 Sequencing data have been deposited in the Gene Expression Omnibus (GEO) under the
702 accession number GSE141609.

703

704 **Code Availability**

705 The analysis code is available on <https://github.com/AlexdeMendoza/BrainZoo>.

706

707 **Acknowledgements**

708 We would like to dedicate this paper to the memory of Jose Luis Gomez-Skarmeta, a dear friend
709 and colleague who was instrumental in igniting this project and contributing to this work, but who
710 sadly passed away during the revision process. We would also like to thank Nacho Maeso for
711 critical reading of this manuscript and suggestions, and Manuel Irimia for advice on isoform
712 quantification. We thank Professor Norman Saunders (University of Melbourne) for sharing
713 opossum material. We would like to thank Juan Pascual-Anaya for granting access to the hagfish
714 genome assembly. We thank to “Semilleria las Ganchozas” for providing advice about material
715 required for this project. This work was supported by the Australian Research Council (ARC)
716 Centre of Excellence program in Plant Energy Biology (CE140100008). RL was supported by a
717 Sylvia and Charles Viertel Senior Medical Research Fellowship, ARC Future Fellowship
718 (FT120100862), and Howard Hughes Medical Institute International Research Scholarship. AdM
719 was funded by an EMBO long term fellowship (ALTF 144-2014). JLG-S was supported by the

720 Spanish government (grant no. BFU2016- 74961-P) and the institutional grant Unidad de
721 Excelencia María de Maeztu (no. MDM-2016-0687). BV was supported by the Biomedical
722 Research Council of the Agency for Science, Technology and Research of Singapore. FG is
723 supported by an ARC Future Fellowship (FT160100267). CR is supported by a NSF grant (IOS-
724 1354898).

725

726 **Author contributions**

727 OB, AdM and RL designed the study. AdM, OB, DP and RL prepared MethylC-seq libraries which
728 were sequenced by JP, JRN and DP. The data were analysed by AdM with help from SBu. JLG-
729 S, EdICM, CA, CWR, FG, TD, BV, JRE, BB, SBe and HE provided the biological samples. The
730 manuscript was written by AdM, OB and RL. All authors commented on the final manuscript.

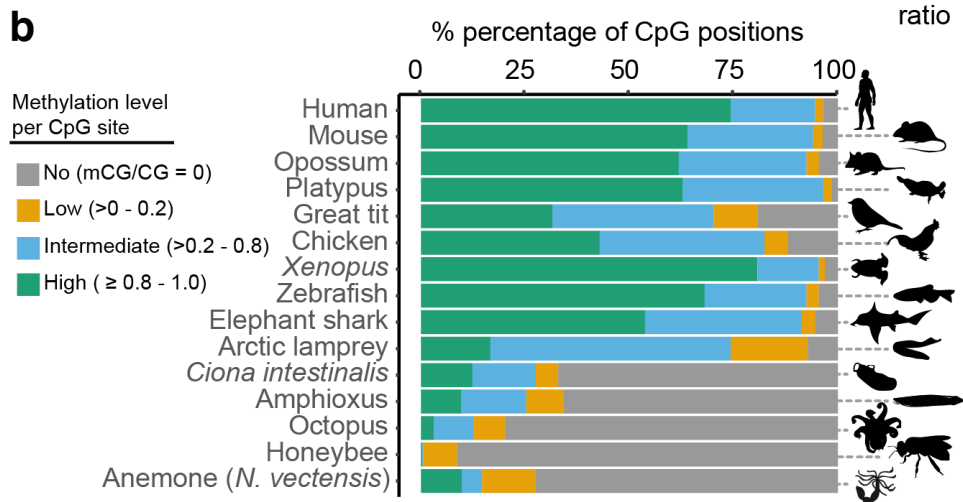
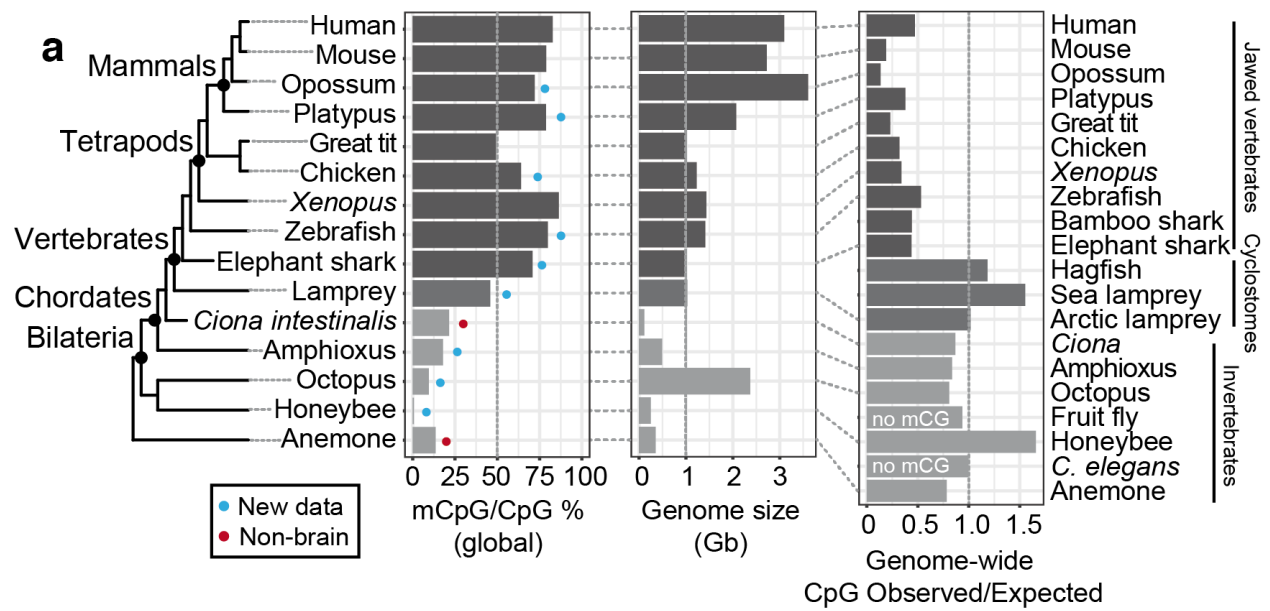
731

732 **Competing interests**

733 The authors declare no competing interests.

734

735 **Figure legends**



736

737 **Fig. 1 | Brain methylomes reflect the vertebrate-invertebrate CG methylation boundary. a,**

738 Global brain CpG methylation, genome size, and CpG genome content across animal species.

739 Schematic representation of established animal phylogeny on the left-hand side. Newly generated

740 WGBS datasets marked with a blue circle, WGBS samples from non-neural tissue marked with a

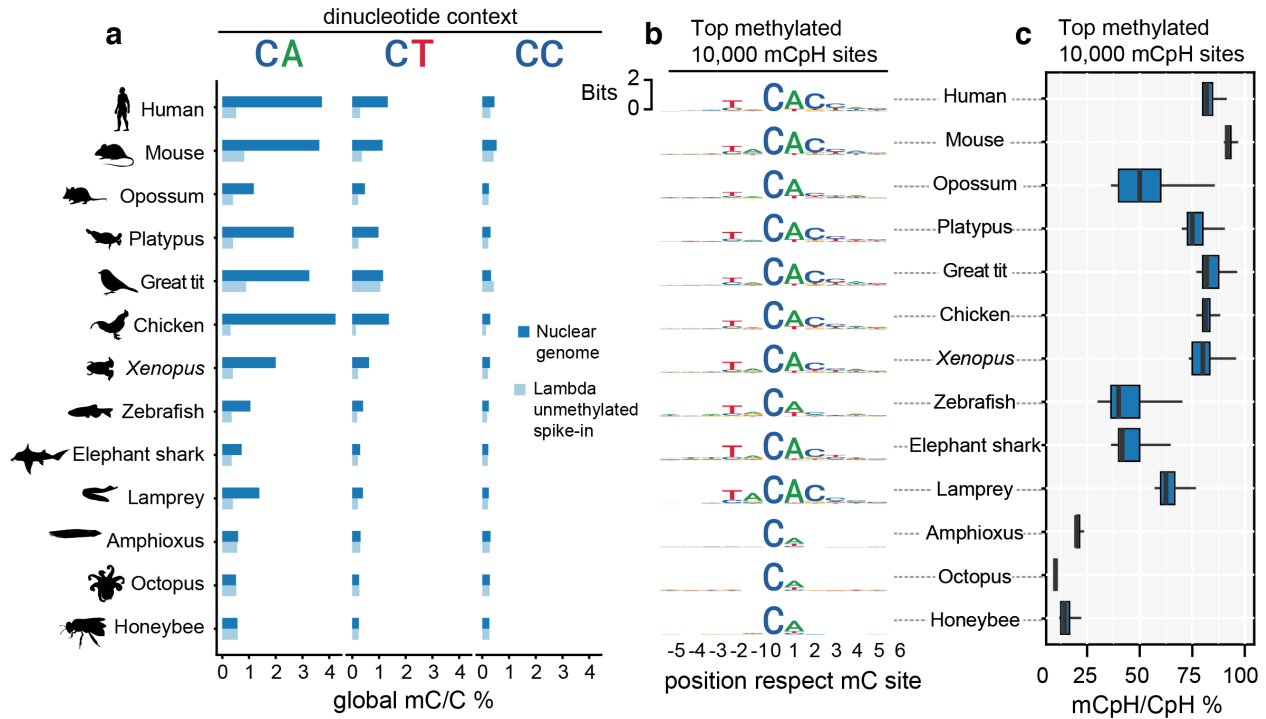
741 red circle. The *Ciona intestinalis* sample corresponds to muscle tissue⁷³, and sea anemone

742 *Nematostella vectensis* sample corresponds to a gastrula sample⁶⁴. Genome size represents the

743 genome assembly size. b, Proportion of CpG sites classified according to methylation levels

744 (mC/C). Only sites with coverage $\geq 10x$ were considered. Silhouettes of human, platypus, octopus

745 and honeybee obtained from phylopic.org.



746

747 **Fig. 2 | Neural CpH methylation is restricted to vertebrate brains. a**, Global methylation levels

748 in brain samples classified per dinucleotide context. Dark blue represents the global methylation

749 level on the nuclear chromosomes (excluding mitochondrial genome) and pale blue represents

750 the bisulfite reaction non-conversion rate for each library, calculated as the methylation levels on

751 an unmethylated lambda phage DNA spike-in. **b**, Sequence motifs found surrounding the most

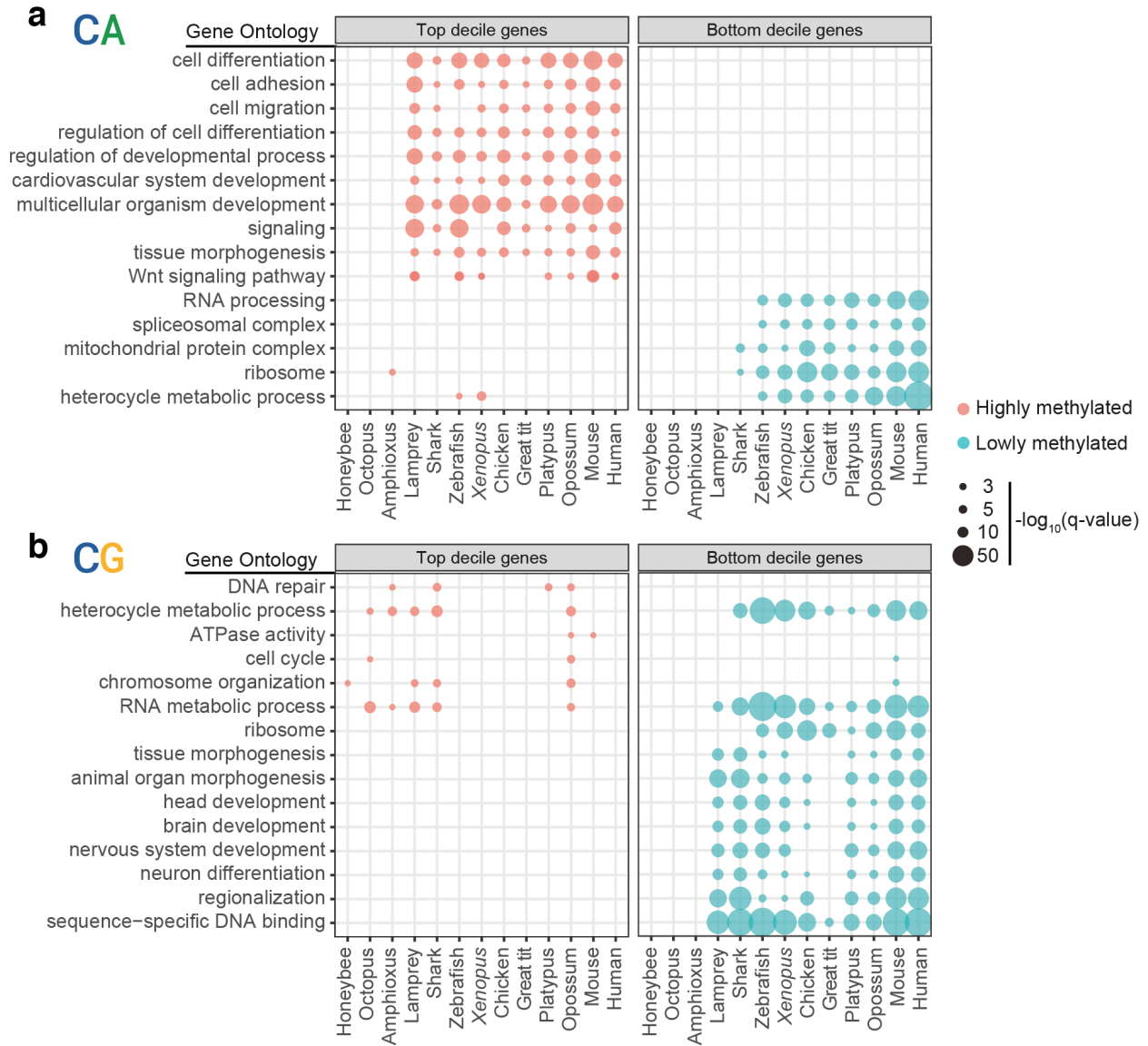
752 highly methylated CpH positions in each brain sample. Only CpH positions with coverage $\geq 10\times$

753 were considered. **c**, Methylation level (mC/C) for the top mCpH positions depicted in panel b.

754 Boxplot centre lines are medians, box limits are quartiles 1 (Q1) and 3 (Q3), whiskers are $1.5 \times$

755 interquartile range (IQR). Silhouettes of human, platypus, octopus and honeybee obtained from

756 phylopic.org.



757

758 **Fig. 3 | Conserved non-overlapping programs are associated with CpH and CpG**

759 **methylation. a**, Gene Ontology enrichments for genes showing the highest and lowest gene body

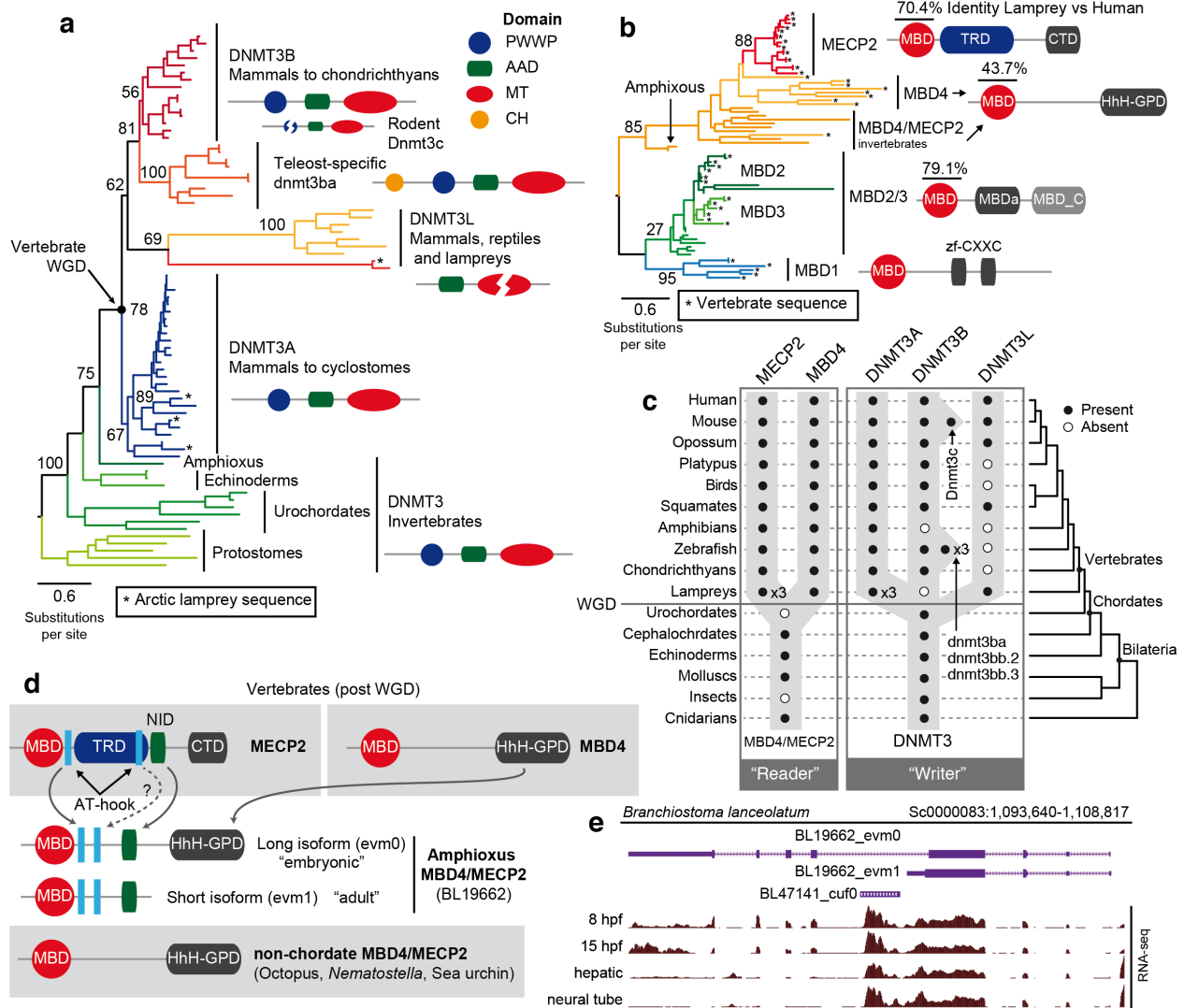
760 methylation levels in the CpA context, as defined by belonging to the top and bottom deciles in

761 each species. **b**, Gene Ontology enrichments for genes showing the highest and lowest

762 methylated levels in the CpG context. Q-values were obtained using the g:SCS algorithm

763 implemented in the gProfiler2 R package.

764



765

766 **Fig. 4 | Vertebrate origins of MECP2 and DNMT3A.** **a**, Maximum likelihood phylogenetic tree of

767 DNMT3 genes in animals. Nodal supports represent 100 bootstrap nonparametric replications.

768 Schematic protein domain configurations shown for each clade. PWWP, Pro-Trp-Trp-Pro motif

769 domain (PF00855). AAD ATRX, DNMT3, DNMT3L domain. MT, cytosine Methyltransferase

770 domain (PF00145). CH, Calponin Homology domain (PF00307). Asterisk highlights arctic

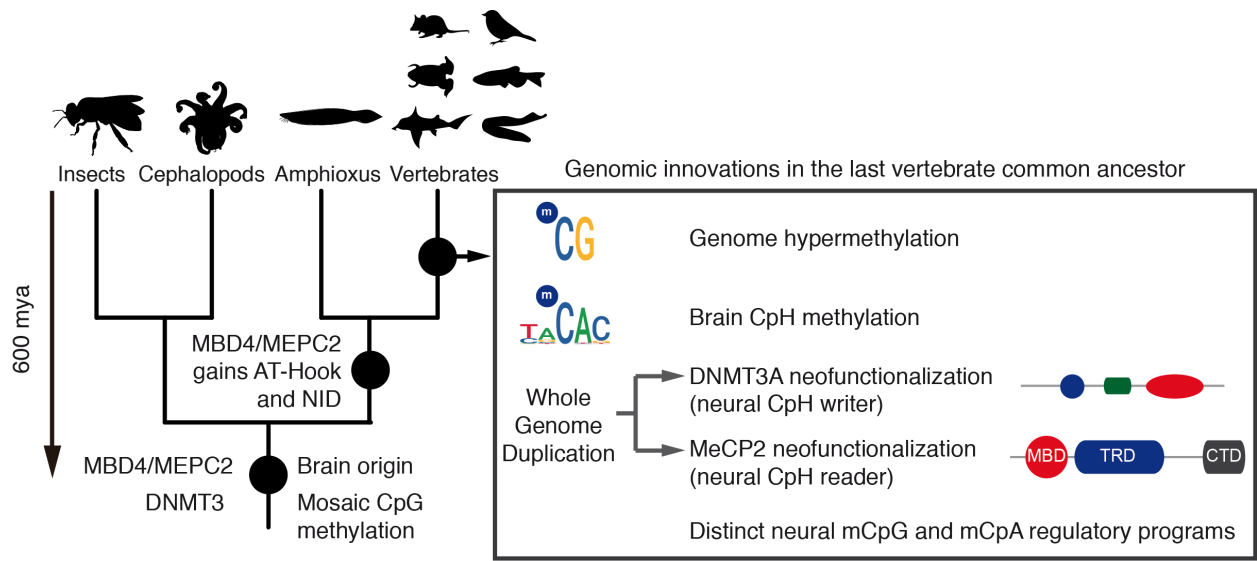
771 lamprey sequences. Broken domains indicate that the domain has large deletions in the given

772 clade. **b**, Maximum likelihood phylogenetic tree of the Methyl-CpG Binding Domain family in

773 animals. Nodal supports represent 100 bootstrap nonparametric replications. On the right, protein

774 domain structure of each clade, as defined by Pfam domains. MBD, Methyl Binding Domain

775 (PF01429). HhH-GPD, Thymine glycosylase (PF00730). MBDA, p55-binding region of MBD2/3
776 (PF16564). MBD_C, MBD2/3 C-terminal domain (PF14048). zf-CXXC, zinc finger (PF02008).
777 CTD, MECP2 C-Terminal Domain. TRD, MECP2 Transcriptional Repression Domain. Asterisks
778 highlight vertebrate sequences, percentages are shown for amino acid MBD identity between
779 lamprey and human orthologues. **c**, Distribution of MECP2/MBD4 and DNMT3 genes across
780 animal lineages. Absence of a dot indicates gene absence. Numbers indicate those
781 species/lineages that have multiple copies of a given gene. Dnmt3c in rodents and
782 dnmt3ba/bb.1/bb.2 are lineage-specific duplications of DNMT3B that have diverged in their
783 function or domain architecture. “x3” indicates lineage-specific duplications. On the right, the
784 phylogenetic relationships among animal lineages. **d**, Stepwise evolution of the MeCP2 and
785 MBD4 protein domains in vertebrates, amphioxus, and non-chordates. NID stands for the N-
786 CoR/SMRT interacting amino acids. **e**, Genome browser snapshot of amphioxus MBD4 locus.
787 The longer isoform with the capacity to repair DNA has higher expression in embryonic samples,
788 see further detail in Extended Data 10.
789
790
791
792



793

794 **Fig. 5 | The assembly of neural-CpH methylation.** Cladogram representing the evolutionary
 795 scenario of neural CpH methylation acquisition in vertebrates. Silhouettes of octopus and
 796 honeybee obtained from phylopic.org.