

## MIT Open Access Articles

*Parallel evolution of male germline epigenetic poising and somatic development in animals*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

**Citation:** Lesch, Bluma J, Sherman J Silber, John R McCarrey, and David C Page. "Parallel Evolution of Male Germline Epigenetic Poising and Somatic Development in Animals." *Nature Genetics* 48, no. 8 (June 13, 2016): 888–894.

**As Published:** <http://dx.doi.org/10.1038/ng.3591>

**Publisher:** Springer Nature

**Persistent URL:** <http://hdl.handle.net/1721.1/107718>

**Version:** Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

**Terms of Use:** Article is made available in accordance with the publisher's policy and may be subject to US copyright law. Please refer to the publisher's site for terms of use.



## **Parallel evolution of male germline epigenetic poisoning and somatic development in animals**

Bluma J. Lesch<sup>1</sup>, Sherman J. Silber<sup>2</sup>, John R. McCarrey<sup>3</sup>, David C. Page<sup>1,4,5</sup>

<sup>1</sup>Whitehead Institute, 9 Cambridge Center, Cambridge, MA, USA.

<sup>2</sup>Infertility Center of St. Louis, St. Luke's Hospital, St. Louis, MO, USA.

<sup>3</sup>Department of Biology, University of Texas at San Antonio, San Antonio, TX, USA.

<sup>4</sup>Department of Biology, Massachusetts Institute of Technology, Cambridge, MA, USA.

<sup>5</sup>Howard Hughes Medical Institute, Whitehead Institute, Cambridge, MA, USA.

Correspondence to [dcpage@wi.mit.edu](mailto:dcpage@wi.mit.edu)

1 **Abstract**

2

3 Changes in gene regulation frequently underlie changes in morphology during evolution, and  
4 differences in chromatin state have been linked with changes in anatomical structure and gene  
5 expression across evolutionary time. Here, we assess the relationship between evolution of  
6 chromatin state in germ cells and evolution of gene regulatory programs governing somatic  
7 development. We examined the poised (H3K4me3/H3K27me3 bivalent) epigenetic state in male  
8 germ cells of five mammalian and one avian species. We find that core genes poised in germ  
9 cells of multiple amniote species are ancient regulators of morphogenesis that sit at the top of  
10 transcriptional hierarchies controlling somatic tissue development, while genes that gain poising  
11 in germ cells of individual species act downstream of core poised genes during development in a  
12 species-specific fashion. We propose that critical regulators of animal development gained an  
13 epigenetically privileged state in germ cells, manifested in amniotes by H3K4me3/H3K27me3  
14 poising, early in metazoan evolution.

15

16 **Introduction**

17

18 Together, maternal and paternal germ cells provide all the information needed to initiate  
19 formation of a new embryo at fertilization. Along with a haploid genome, germ cells carry gene  
20 regulatory information, which guides gene expression during sperm or egg development<sup>1,2</sup> and  
21 may impact development of the embryo in the next generation<sup>3-5</sup>. Changes in gene regulation  
22 contribute to evolution of morphology across diverse species<sup>6,7</sup>, and recent studies assessing  
23 chromatin state in specific tissues across multiple species have found that evolution of chromatin  
24 state is associated with evolution of gene expression and anatomical structure.<sup>8-11</sup> Building on  
25 this work, we reasoned that evolution of chromatin state in germ cells might be similarly  
26 associated with evolution of gene expression in somatic tissues of the embryo.

27

28 We focused our attention on evolution of epigenetic ‘poising’ in germ cells. Epigenetic poising  
29 is defined by the simultaneous presence of two opposing histone modifications, the activating  
30 mark H3K4me3 and the repressive mark H3K27me3, as well as by transcriptional repression<sup>12-14</sup>.  
31 It has been best studied in embryonic stem cells (ESCs), where it is associated with genes  
32 involved in lineage specification. In ESCs, the presence of the activating H3K4me3 mark is  
33 thought to ‘poise’ these otherwise silent genes for activation following receipt of a differentiation  
34 cue, while the H3K27me3 mark maintains repression in the pluripotent state<sup>12</sup>. Consistent with  
35 this model, in mouse embryos *in vivo*, *Hox* genes move sequentially from a poised state to an  
36 active H3K4me3-only state, concurrent with their activation in an anterior-to-posterior  
37 direction<sup>15,16</sup>. *Hox* genes and other developmental regulators are poised in mouse and human  
38 germ cells but are not expressed in developing gametes; rather, genes poised in germ cells are

39 expressed in somatic tissues during embryogenesis<sup>17-20</sup>. We hypothesized that evolution of  
40 epigenetic poising in mammalian germ cells reflects the evolution of a transcriptional program  
41 controlling somatic gene expression and morphogenesis in embryos.

42

43 Evaluating this hypothesis demanded we carry out four tasks: characterize the poised state in the  
44 germ cells of multiple mammalian species; relate conservation of germ cell poising to  
45 conservation of developmental function in mammals; define the relationship between differences  
46 in germ cell poising and differences in developmental function among specific evolutionary  
47 lineages; and reconstruct the evolutionary origins of these relationships by comparison to non-  
48 mammalian taxa.

49

50 We used comparative epigenetic profiling in mammalian germ cells to address these questions.

51 We examined genome-wide expression, H3K4me3, and H3K27me3 data in male germ cells of  
52 five mammalian species spanning 175 million years of evolution. We found that evolution of  
53 poising in male germ cells parallels evolution of somatic gene expression and development in the  
54 embryo. We propose an ancient evolutionary relationship between germline chromatin and  
55 embryonic gene expression.

56

## 57 **Results**

58

### 59 *Gene expression and chromatin state in male germ cells*

60 We collected H3K4me3 and H3K27me3 ChIP-seq data, as well as RNA-seq data, in sorted male  
61 germ cells from five species spanning 175 million years of evolutionary divergence in the

62 mammalian lineage: human, rhesus macaque, mouse, bull, and opossum (Fig. 1a, Supplementary  
63 Fig. 1, and Supplementary Table 1). We obtained data from cells at two time points during male  
64 germ cell development: prophase of meiosis I (pachytene spermatocytes), and following  
65 completion of meiosis (round spermatids). These cell types are unique in that they can be  
66 identified and reliably collected from whole testes in multiple mammalian species without the  
67 use of transgenes or genetic markers (see Methods). In addition, they differ substantially from  
68 each other in their place in the cell cycle and in the physical state of their chromatin: pachytenes  
69 are tetraploid cells in meiotic prophase, with large nuclei and synapsed pairs of homologous  
70 chromosomes, while round spermatids are haploid cells that have completed meiosis and have  
71 small, compact nuclei. Inclusion of both cell types in our study allowed us to control for the  
72 effects of chromatin compaction and physical state of the nucleus during spermatogenic  
73 development.

74  
75 Before turning to the poised state, we first examined the relationships between H3K4me3,  
76 H3K27me3, and expression datasets collected from different species. In all species, expression  
77 levels were positively correlated with H3K4me3 signal and negatively correlated with  
78 H3K27me3 signal, consistent with the known association of these marks with gene activity and  
79 repression, respectively (Fig. 1b, Supplementary Fig. 2). Our data included three biological  
80 replicates for human, two for rhesus, mouse, and opossum, and one for bull (Supplementary Fig.  
81 3); biological replicates were highly similar as evaluated by principal component analysis  
82 (Supplementary Fig. 4a) or by hierarchical clustering (Fig. 1c, Supplementary Fig. 4b). For  
83 H3K4me3, clustering accurately separated pachytene spermatocytes from round spermatids  
84 within each species, but did not fully recapitulate phylogenetic relationships, as has been

85 previously reported for the H3K4me3 mark in other tissues<sup>10</sup>. In contrast, we derived correct or  
86 nearly-correct phylogenies from both expression and H3K27me3 data. For expression,  
87 H3K4me3, and H3K27me3 data, greater dissimilarity between samples corresponded to greater  
88 evolutionary divergence between species (Fig. 1d, Supplementary Note 1).

89

### 90 *Identification of poised chromatin in male germ cells*

91 To identify genes associated with poised chromatin in germ cells of each species, we calculated  
92 read counts in four-kilobase intervals surrounding the transcription start sites of annotated genes,  
93 after normalizing for library size and subtracting input signal (Supplementary Files 1-10).

94 Throughout our analysis, we considered only genes with orthologs in all five species (a total of  
95 14,362 orthology groups). Genes above a threshold of 0.5 reads per million for H3K4me3 and  
96 H3K27me3 signal, and whose expression level was equal to or less than 5 FPKM, were called as  
97 poised. Within each species, we further filtered for genes at which the poised state was retained  
98 in both pachytene spermatocytes and round spermatids, implying that it is stable across much of  
99 spermatogenic development. Stably-poised gene sets identified in this manner were robust to  
100 changes in ChIP and expression threshold (Supplementary Figs. 5 and 6).

101

102 This approach identified 1,200-3,600 poised genes in each species (Fig. 2a,b and Supplementary  
103 Table 2). We verified that co-occurrence of high H3K4me3 and H3K27me3 signals at poised  
104 genes represented the simultaneous presence of the two marks on the same DNA molecule, not  
105 heterogeneity of chromatin state within our cell population, by performing sequential ChIP at the  
106 promoters of representative poised genes in both mouse (Fig. 2c and Supplementary Fig. 7a) and  
107 opossum (Fig. 2d and Supplementary Fig. 7b) round spermatids. We confirmed that four out of

108 four mouse and two out of two opossum poised promoters were simultaneously marked by both  
109 H3K4me3 and H3K27me3, demonstrating that these genes are marked by a *bona fide* poised  
110 state in round spermatids. In general, between one quarter and three quarters of poised genes  
111 were shared between any two species (Supplementary Table 3), and dissimilarity in poising was  
112 positively correlated with evolutionary divergence time (Fig. 2e). Regardless of evolutionary  
113 distance, overlap in poised gene sets between each species pair was greater than expected by  
114 chance ( $p < 10^{-15}$  for all pairs, Fisher's exact test).

115

116

### 117 ***Conserved poising at developmental regulators in mammals***

118 Four hundred and five genes were poised in all five species ( $p < 10^{-280}$  compared to expected for  
119 five-way overlap, see Methods) (Fig. 3a and Supplementary Table 4). These genes were well  
120 distributed across chromosomes (Supplementary Fig. 8a). At the sequence level, the promoter  
121 regions of these genes were significantly better conserved than those of human-specific poised  
122 genes ( $p < 10^{-14}$ , Welch t-test) or genes with conserved retention of H3K27me3 but not  
123 necessarily H3K4me3 ( $p < 10^{-4}$ , Welch t-test) (Supplementary Fig. 8b). The set of genes poised  
124 in germ lines of all five mammalian species, henceforth referred to as 'core' poised genes, was  
125 strongly enriched for genes encoding transcription factors, with a particularly striking enrichment  
126 for homeodomain-containing transcription factors (Fig. 3b and Supplementary Fig. 8c).

127

128 We used pre-defined Gene Ontology (GO) categories to confirm enrichment of transcription  
129 factor-encoding genes in the core poised gene set. The 405 core poised genes were significantly  
130 enriched for genes belonging to the GO category "sequence specific DNA binding transcription



131 factor activity” (GO:0003700) compared to all genes with five-way orthologs, to human- or  
132 mouse-specific poised genes, or to genes with conserved retention of H3K27me3 but not  
133 necessarily H3K4me3 (Fig. 3c).

134

135 We then asked whether, in addition to encoding proteins with a shared molecular function as  
136 sequence-specific transcription factors, the set of core poised genes had a unifying biological  
137 function during development. We examined enrichment of GO biological function categories in  
138 the set of core poised genes. We found that enriched GO categories described processes  
139 involved in patterning and organ formation, including “embryonic organ morphogenesis”,  
140 “anterior/posterior pattern specification”, “limb development”, and “gastrulation” (Fig. 3d and  
141 Supplementary Table 5). Core poised genes were not enriched for germ cell-related functions,  
142 such as meiosis, spermatid development, or sperm maturation. These findings imply that the  
143 core poised genes have a shared biological role, specifically, transcriptional regulation of body  
144 patterning and somatic tissue specification during embryogenesis.

145

146 If the core poised genes are involved in patterning and tissue specification, they should be  
147 expressed during an interval in embryogenesis when these processes are occurring. We queried  
148 the MGI Gene Expression Database<sup>21</sup>, which contained expression data for 13,837 mouse genes  
149 at the time of our study, to determine the interval when each of these genes is first expressed.  
150 Core poised genes were enriched for expression during embryogenesis compared to other genes  
151 with orthologs in all five species; this difference was especially evident between gastrulation and  
152 the end of somite formation (Thieler stage [TS] 11-22, Fig. 3e). A subset of core poised genes  
153 involved in trophectoderm specification (e.g. *Cdx2* [*Caudal type homeobox 2*], *Hand1* [*Heart*

154 *and neural crest derivatives expressed transcript 1*], and *Tpbg* [*Trophoblast glycoprotein*]) was  
155 also expressed early in embryogenesis (TS 2-4)<sup>22</sup>.

156

157 We then examined specific cases where the regulatory hierarchies involved in body part and  
158 organ field specification are well defined, and found that core poised genes are central to these  
159 processes. Core poised genes sit at the top of such specification hierarchies, including *PTF1A*  
160 (*Pancreas specific transcription factor 1a*) and *PDX1* (*Pancreatic and duodenal homeobox 1*) in  
161 pancreatic development<sup>23</sup>; *EN1/2* (*Engrailed 1 and 2*), *OTX1* (*Orthodenticle homeobox 1*),  
162 *LMX1A* (*LIM homeobox transcription factor alpha*) and *GBX2* (*Gastrulation brain homeobox 2*)  
163 in cerebellar development<sup>24</sup>; *NKX2-5* (*NK homeobox 2-5*), *HAND1* and *HAND2* (*Heart and*  
164 *neural crest derivatives expressed transcript 2*) in heart development<sup>25,26</sup>; and *MSX1/2* (*Msh*  
165 *homeobox 1 and 2*) and *NKX2-2* (*NK homeobox 2-2*) in neural tube regionalization<sup>27</sup>.

166

167 To obtain quantitative support for this finding, we examined the connectedness of core poised  
168 genes compared to other genes in the context of three experimentally-supported developmental  
169 regulatory networks: pancreas<sup>23</sup>, heart<sup>25</sup>, and cerebellum<sup>24</sup>. Considering all three of these  
170 networks together, core poised genes had more regulatory connections than other genes (mean  
171 4.50 compared to 2.27 connections,  $p=0.01765$  [one-sided Mann-Whitney U test]). Core poised  
172 genes also exhibited greater network centrality (betweenness centrality, the likelihood that a  
173 particular gene lies on the shortest path connecting two other genes) compared to other genes  
174 upregulated during differentiation and specification stages in an *in vitro* model of human cortical  
175 development<sup>28</sup> ( $p=1.43 \times 10^{-6}$ , one-sided Mann-Whitney U test).

176

177 We conclude that the core poised genes constitute critical upstream regulators of gene expression  
178 during mammalian embryogenesis. Indeed, many of the core poised genes participate in  
179 developmental ‘kernels’, conserved genetic circuits controlling specification of body part  
180 progenitor fields<sup>29,30</sup>. Kernel architecture extends deep into the metazoan lineage and is highly  
181 conserved across metazoa, partly because extensive regulatory interactions among kernel  
182 constituents mean that perturbation of any one component can have catastrophic effects on body  
183 part patterning<sup>29,30</sup>. Indeed, we found that knockout alleles of 83 of the core poised genes (21%)  
184 resulted in embryonic lethality in the mouse, whereas knockout alleles of only 10% of all genes  
185 with orthologs in all five mammalian species resulted in the same phenotype ( $p=3.77 \times 10^{-11}$ ,  
186 Fisher’s exact test)<sup>31</sup>.

187

### 188 ***Differences in germline poising between species***

189 Given that genes with strongly conserved roles in metazoan development exhibit conservation of  
190 poising in germ cells, we wondered if genes that gain species-specific developmental functions  
191 might also acquire species-specific poising in germ cells. To address this question, we turned  
192 our attention to the five sets of genes poised specifically in only one of the five species evaluated  
193 (‘differentially poised’ genes), acknowledging that a subset of these genes may be mis-assigned  
194 as ‘specific’ due to false negative calls in one or more of the other four species (Fig. 4a,b  
195 Supplementary Table 6). None of the five sets of differentially poised genes was strongly  
196 enriched for GO developmental functions. However, where comparative expression data in  
197 specific developmental structures was available, species differences in germline poising were  
198 correlated with species differences in developmental expression<sup>8,9,32,33</sup>. For example, *ZSWIM4*  
199 and *LMF1*, which are poised specifically in human germ cells (Fig. 4a and Supplementary Fig.

200 9a), are expressed in human but not mouse or bovine placenta; likewise, *Ccrl2* and *Smug1* are  
201 poised specifically in mouse germ cells (Fig. 4b and Supplementary Fig. 9b) and expressed in  
202 mouse but not human or bovine placenta<sup>32</sup>. Similarly, *HIPK2* is poised specifically in human  
203 germ cells (Supplementary Fig. 9a) and has acquired a human-specific enhancer active during  
204 limb development<sup>8</sup>.

205

206 We predicted that differential poising and differential expression during development in a given  
207 species would correspond to differences in regulatory sequence compared to orthologous genes  
208 in the other four species. To test this prediction, we searched for motifs enriched in the  
209 promoters of each of the five differentially poised gene sets relative to their non-poised orthologs  
210 (Supplementary Table 7). We found that motifs gained in poised promoters frequently  
211 corresponded to predicted binding motifs for transcription factors encoded by core poised genes  
212 (71% of human, 62% of rhesus, 51% of mouse, 33% of bull, and 60% of opossum gained motifs)  
213 (Fig. 4c, Supplementary Table 7). In general, motifs gained in differentially poised promoters  
214 were different in each species. Together with expression differences, these results imply that  
215 acquisition of epigenetic poising in germ cells may occur in parallel with gain of regulation by  
216 core poised genes during somatic development. Extension of germline poising to new  
217 developmental factors may facilitate their recruitment into ancient developmental circuits  
218 regulated by core poised genes.

219

220 We also identified cases of single-lineage losses, in which a gene was poised in all but one of the  
221 five species examined (Supplementary Table 8). As with single-lineage gains, some of these  
222 instances may be due to false-negative poised gene calls in one species. However, some

223 instances of single-lineage loss of poising are supported by previous reports of recent  
224 evolutionary divergence in expression or function of the associated gene. For example, of the 36  
225 genes poised in four mammals but not in human, three are reported to have divergent expression  
226 patterns in human compared to other mammals (*AIM1*, *EPHA5*, and *THBS4*)<sup>34-36</sup>, one is  
227 associated with differences in loss-of-function phenotype between human and mouse  
228 (*DOCK8*)<sup>37</sup>, and three are associated with recent positive selection in the human lineage (*AIM1*,  
229 *COL11A1*, and *LYPD1*)<sup>38-40</sup>. Like lineage-specific gains, lineage-specific loss of poising in the  
230 germ line may therefore reflect recent lineage-specific changes in developmental regulation and  
231 function.

232

### 233 ***Conservation of germline poising beyond mammals***

234 The set of core poised genes is notable both for its origins deep in the metazoan lineage and for  
235 its specificity to metazoa; 224 (55%) of the core poised genes have orthologs in the fly  
236 *Drosophila melanogaster* but not in the yeast *Saccharomyces cerevisiae*, compared to 36% of all  
237 genes with orthologs in all five mammals ( $p = 9.59 \times 10^{-16}$ , Fisher's exact test), implying that the  
238 majority of core poised genes arose before the divergence of protostomes and deuterostomes but  
239 after the divergence of animals from fungi. We asked when these genes might first have gained  
240 a specialized epigenetic state in the metazoan germ line.

241

242 First, we compared our mammalian data to a non-mammalian amniote, the chicken. Together  
243 with reptiles, birds constitute the closest living relatives of the mammalian clade (Fig. 1a)<sup>41</sup>. We  
244 collected ChIP- and RNA-seq data from chicken germ cells at time points matching those used  
245 for the five mammalian species (Supplementary File 11). Using identical filtering conditions, we

246 identified 1716 genes poised in the chicken germ line (Supplementary Table 9). Of the 405 core  
247 poised genes we defined in mammals, 347 have orthologs in the chicken genome; of these, 215  
248 (62%) are also poised in the chicken germ line. For the core poised genes in mammals, the set of  
249 genes whose orthologs were also poised in chicken was significantly enriched for sequence-  
250 specific transcription factors compared to the core poised genes that were not poised in chicken  
251 (Fig. 5a,b). We conclude that epigenetic poising of developmental transcriptional regulators in  
252 germ cells is at least as old as the amniote common ancestor, placing its origin more than 300  
253 million years ago. In at least five cases, core poised genes that were poised in mammalian but  
254 not chicken germ cells could also be correlated to differences in development between mammals  
255 and birds<sup>42-45</sup>, supporting the hypothesis that acquisition of poising in the germ line is related to  
256 acquisition of somatic developmental function. For example, *TPBG* is poised in all five  
257 mammals but not in chicken (Fig. 5c), consistent with its early expression in trophectoderm, a  
258 mammal-specific structure<sup>45,46</sup>. Given previous reports of a multivalent chromatin state  
259 comprised of multiple repressive and activating histone marks at developmental genes in  
260 zebrafish sperm<sup>47</sup>, it will be interesting to trace the evolutionary history of the poised state in  
261 additional non-amniote vertebrate species.

262

263 Using the chicken data, we further examined the scenario in which genes whose poised state was  
264 shared in distantly related species were not poised in germ cells of intermediate lineages. Such a  
265 scenario implies either convergent evolution, requiring two independent epigenetic gains, or  
266 deep loss of poising followed by recent reacquisition, requiring a loss followed by a gain. Either  
267 explanation calls for at least two independent evolutionary events and is expected to be rarer than  
268 scenarios requiring either uninterrupted conservation (zero events) or single-lineage gains or

269 losses (one event). Indeed, among the 11,188 genes with orthologs in all six species, 211 were  
270 poised in human, rhesus, and mouse only, implying a single gain in the primate-rodent ancestor,  
271 compared to only 35 in human, opossum, and chicken or 14 in rhesus, opossum, and chicken,  
272 each requiring either convergence or deep loss in the placental lineage followed by a gain  
273 (Supplementary Table 9). For at least one gene among the 35 shared by human, opossum, and  
274 chicken (*NCSI*, *Neuronal calcium sensor 1*), expression patterns in human and chicken were  
275 more similar than in human and mouse, implying convergent evolution of expression.<sup>48,49</sup>

276

277 To examine the possibility that the origins of germline poisoning lie deeper in the metazoan  
278 lineage, we compared our set of core poised genes to Polycomb ChIP-microarray data in sorted  
279 *Drosophila* male germ cells<sup>50</sup>. We found that 5.2% of all genes whose promoters were marked  
280 by high levels of Polycomb in the *Drosophila* germ line were orthologs of the core poised genes,  
281 compared to 2.5% of genes with low Polycomb levels ( $p=1.059 \times 10^{-5}$ , Fisher's exact test).  
282 Overall, orthologs of core poised genes were enriched for Polycomb signal in the *Drosophila*  
283 germ line (Supplementary Fig. 10a). This effect is modest, but suggests that orthologs of some  
284 mammalian core poised genes may have acquired a specialized epigenetic state, characterized by  
285 Polycomb binding and H3K27me3, in germ cells before the emergence of the bilaterian  
286 ancestor<sup>51</sup>, and retained it independently in protostomes and deuterostomes (Fig. 5d).

287

## 288 **Discussion**

289

290 We show here that evolution of epigenetic poisoning in male germ cells is closely linked to  
291 evolution of somatic gene expression in developing mammalian embryos. Germline poisoning is

292 conserved throughout the mammalian lineage at genes that are central to the transcriptional  
293 networks governing somatic development, and individual genes recruited to these networks in  
294 specific lineages also gain poising in germ cells. Poising of central developmental genes in male  
295 germ cells is at least as old as the amniote common ancestor, placing its origin at least 300  
296 million years ago. Such deep conservation implies a functional role for the poised state in germ  
297 cells.

298  
299 It is easy to envision a role for H3K27me3 at somatic genes in the germ line: this repressive  
300 mark reinforces silencing of genes whose expression in germ cells would disrupt their function  
301 and identity. However, we found that genes exhibiting conservation of H3K27me3 without  
302 H3K4me3 do not show the same functional enrichments as genes with conserved poising (Fig.  
303 3c and Supplementary Table 5), indicating that H3K27me3 alone does not play the same  
304 biological role as H3K27me3/H3K4me3 bivalency. H3K4me3 may play a protective role at  
305 poised promoters as an antagonist of DNA methylation<sup>5,52</sup>. It is also possible that H3K4me3  
306 helps to prepare poised genes for expression in somatic tissues following fertilization, similar to  
307 its proposed role in ESCs. Consistent with this hypothesis, altered regulation of H3K4  
308 methylation state in developing male germ cells in the mouse perturbs somatic tissue  
309 development in embryos of the next generation<sup>3</sup>. The detailed mechanism by which this  
310 epigenetic information might be transmitted through fertilization remains unclear: modified  
311 histones may be carried in mature spermatozoa<sup>19,20,53</sup>, or poised sites may be marked by an RNA  
312 or protein intermediate in sperm and re-established in the early embryo. We found that  
313 published ChIP-seq data from mature mouse<sup>20</sup> and human<sup>19</sup> spermatozoa is consistent with



314 retention of modified histones at core poised genes (Supplementary Fig. 10b), but this finding  
315 does not exclude the participation of additional factors in marking poised sites.

316

317 Our study leverages comparative analysis of *in vivo* epigenomic data across multiple species to  
318 identify a set of genes that is epigenetically privileged in the mammalian germ line. This  
319 privileged state is manifested by H3K4me3/H3K27me3 bivalency in amniotes, and association  
320 of H3K27me3 with core members of this gene set extends deep in animal evolution to the  
321 common bilaterian ancestor. Future work in additional animal species will be important to better  
322 define the evolutionary history of this privileged epigenetic state in the metazoan germ line.

323

324 Together with existing studies from non-amniote species<sup>27,50,54-58</sup>, our data implicate core poised  
325 genes as ancient regulators of metazoan development that sit at the heart of somatic  
326 developmental networks, and differentially poised genes as agents of lineage-specific change. In  
327 mammalian germ cells, the poised state thus represents a memory of ancient developmental  
328 regulatory hierarchies and a device for understanding their evolution.

329

330

331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354

**URLs**

[http://hannonlab.cshl.edu/fastx\\_toolkit/index.html](http://hannonlab.cshl.edu/fastx_toolkit/index.html)  
<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>  
<http://www.informatics.jax.org/gxd>  
<http://www.informatics.jax.org/allele>  
<http://www.R-project.org/>

**Accession codes**

ChIP-seq and RNA-seq data are available at the Sequence Read Archive (SRA) under accession code SRP057141 and at the Gene Expression Omnibus (GEO) under accession code GSE68507.

**Acknowledgments**

This project was funded by an HHMI award to DCP, by a Hope Funds for Cancer Research postdoctoral fellowship to BJL, and by a Burroughs-Wellcome Career Award to BJL. We thank H. Skaletsky for statistical advice and analysis; R. Young for advice on ChIP-seq analysis and critical reading of the manuscript; and P. Reddien for critical reading of the manuscript.

**Author contributions**

B.J.L. designed the project, conducted experiments, analyzed data, and wrote the paper. D.C.P. designed the project and wrote the paper. S.J.S. provided human testis samples and contributed to writing the paper. J.R.M. isolated germ cells for all samples and contributed to writing the paper.

355 **Main text references**

356

- 357 1. Kimmins, S. & Sassone-Corsi, P. Chromatin remodelling and epigenetic features of germ  
358 cells. *Nature* **434**, 583-9 (2005).
- 359 2. Kurimoto, K. *et al.* Quantitative dynamics of chromatin remodeling during germ cell  
360 specification from mouse embryonic stem cells. *Cell Stem Cell* **16**, 517-32 (2015).
- 361 3. Siklenka, K. *et al.* Disruption of histone methylation in developing sperm impairs  
362 offspring health transgenerationally. *Science* **350** (2015).
- 363 4. Arico, J.K., Katz, D.J., van der Vlag, J. & Kelly, W.G. Epigenetic patterns maintained in  
364 early *Caenorhabditis elegans* embryos can be established by gene activity in the parental  
365 germ cells. *PLoS Genet* **7**, e1001391 (2011).
- 366 5. Ihara, M. *et al.* Paternal poly (ADP-ribose) metabolism modulates retention of inheritable  
367 sperm histones and early embryonic gene expression. *PLoS Genet* **10**, e1004317 (2014).
- 368 6. King, M.C. & Wilson, A.C. Evolution at two levels in humans and chimpanzees. *Science*  
369 **188**, 107-16 (1975).
- 370 7. Wray, G.A. The evolutionary significance of cis-regulatory mutations. *Nat Rev Genet* **8**,  
371 206-16 (2007).
- 372 8. Cotney, J. *et al.* The evolution of lineage-specific regulatory activities in the human  
373 embryonic limb. *Cell* **154**, 185-96 (2013).
- 374 9. Reilly, S.K. *et al.* Evolutionary changes in promoter and enhancer activity during human  
375 corticogenesis. *Science* **347**, 1155-9 (2015).
- 376 10. Villar, D. *et al.* Enhancer evolution across 20 mammalian species. *Cell* **160**, 554-66  
377 (2015).
- 378 11. Prescott, S.L. *et al.* Enhancer divergence and cis-regulatory evolution in the human and  
379 chimp neural crest. *Cell* **163**, 68-83 (2015).
- 380 12. Bernstein, B.E. *et al.* A bivalent chromatin structure marks key developmental genes in  
381 embryonic stem cells. *Cell* **125**, 315-26 (2006).
- 382 13. Mikkelsen, T.S. *et al.* Genome-wide maps of chromatin state in pluripotent and lineage-  
383 committed cells. *Nature* **448**, 553-60 (2007).
- 384 14. Azuara, V. *et al.* Chromatin signatures of pluripotent cell lines. *Nat Cell Biol* **8**, 532-8  
385 (2006).
- 386 15. Soshnikova, N. & Duboule, D. Epigenetic temporal control of mouse *Hox* genes in vivo.  
387 *Science* **324**, 1320-3 (2009).
- 388 16. Noordermeer, D. *et al.* Temporal dynamics and developmental memory of 3D chromatin  
389 architecture at *Hox* gene loci. *Elife* **3**, e02557 (2014).
- 390 17. Lesch, B.J., Dokshin, G.A., Young, R.A., McCarrey, J.R. & Page, D.C. A set of genes  
391 critical to development is epigenetically poised in mouse germ cells from fetal stages  
392 through completion of meiosis. *Proc Natl Acad Sci U S A* **110**, 16061-6 (2013).
- 393 18. Sachs, M. *et al.* Bivalent chromatin marks developmental regulatory genes in the mouse  
394 embryonic germline in vivo. *Cell Rep* **3**, 1777-84 (2013).
- 395 19. Hammoud, S.S. *et al.* Distinctive chromatin in human sperm packages genes for embryo  
396 development. *Nature* **460**, 473-8 (2009).
- 397 20. Erkek, S. *et al.* Molecular determinants of nucleosome retention at CpG-rich sequences in  
398 mouse spermatozoa. *Nat Struct Mol Biol* **20**, 868-75 (2013).
- 399 21. Smith, C.M. *et al.* The mouse Gene Expression Database (GXD): 2014 update. *Nucleic*  
400 *Acids Res* **42**, D818-24 (2014).

401 22. Richardson, L. *et al.* EMAGE mouse embryo spatial gene expression database: 2014  
402 update. *Nucleic Acids Res* **42**, D835-44 (2014).

403 23. Arda, H.E., Benitez, C.M. & Kim, S.K. Gene regulatory networks governing pancreas  
404 development. *Dev Cell* **25**, 5-13 (2013).

405 24. Oberdick, J. in *Handbook of the Cerebellum and Cerebellar Disorders* (eds. Manto, M.,  
406 Schmahmann, J., Rossi, F., Gruol, D. & Koibuchi, N.) 127-145 (Springer Netherlands,  
407 2013).

408 25. Cripps, R.M. & Olson, E.N. Control of cardiac development by an evolutionarily  
409 conserved transcriptional network. *Dev Biol* **246**, 14-28 (2002).

410 26. Olson, E.N. Gene regulatory networks in the evolution and development of the heart.  
411 *Science* **313**, 1922-7 (2006).

412 27. Arendt, D. & Nubler-Jung, K. Comparison of early nerve cord development in insects  
413 and vertebrates. *Development* **126**, 2309-25 (1999).

414 28. van de Leemput, J. *et al.* CORTECON: A temporal transcriptome analysis of in vitro  
415 human cerebral cortex development from human embryonic stem cells. *Neuron* **83**, 51-68  
416 (2014).

417 29. Davidson, E.H. & Erwin, D.H. Gene regulatory networks and the evolution of animal  
418 body plans. *Science* **311**, 796-800 (2006).

419 30. Peter, I.S. & Davidson, E.H. Evolution of gene regulatory networks controlling body plan  
420 development. *Cell* **144**, 970-85 (2011).

421 31. Bult, C.J. *et al.* Mouse genome database 2016. *Nucleic Acids Res* **44**, D840-7 (2016).

422 32. Hou, Z.C. *et al.* Elephant transcriptome provides insights into the evolution of eutherian  
423 placentation. *Genome Biol Evol* **4**, 713-25 (2012).

424 33. Ozawa, M. *et al.* Global gene expression of the inner cell mass and trophectoderm of the  
425 bovine blastocyst. *BMC Dev Biol* **12**, 33 (2012).

426 34. Das, R. *et al.* DNMT1 and AIM1 Imprinting in human placenta revealed through a  
427 genome-wide screen for allele-specific DNA methylation. *BMC Genomics* **14**, 685  
428 (2013).

429 35. Caceres, M., Suwyn, C., Maddox, M., Thomas, J.W. & Preuss, T.M. Increased cortical  
430 expression of two synaptogenic thrombospondins in human brain evolution. *Cereb*  
431 *Cortex* **17**, 2312-21 (2007).

432 36. Olivieri, G. & Miescher, G.C. Immunohistochemical localization of EphA5 in the adult  
433 human central nervous system. *J Histochem Cytochem* **47**, 855-61 (1999).

434 37. McGhee, S.A. & Chatila, T.A. DOCK8 immune deficiency as a model for primary  
435 cytoskeletal dysfunction. *Dis Markers* **29**, 151-6 (2010).

436 38. Soejima, M., Tachida, H., Ishida, T., Sano, A. & Koda, Y. Evidence for recent positive  
437 selection at the human AIM1 locus in a European population. *Mol Biol Evol* **23**, 179-88  
438 (2006).

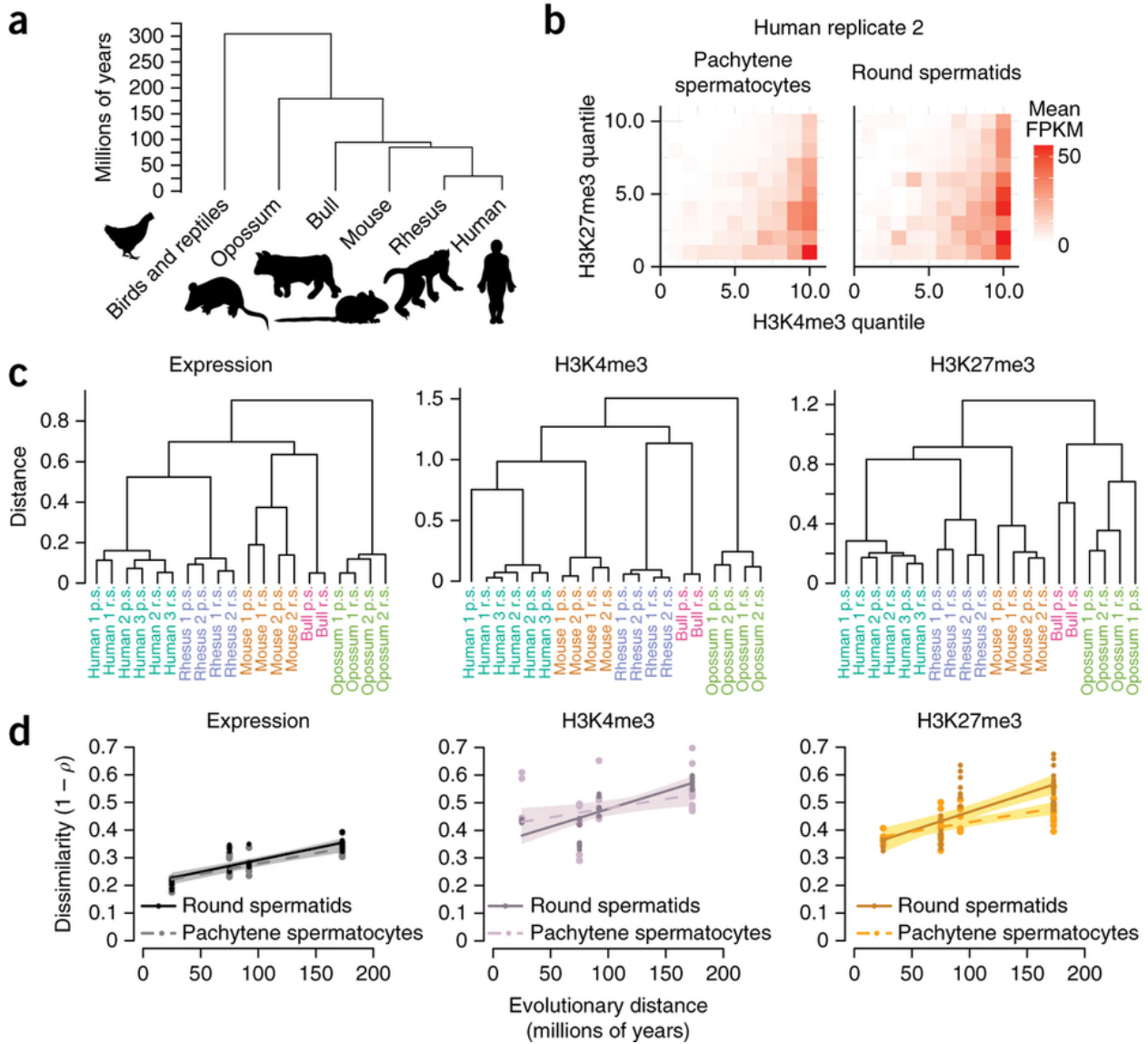
439 39. Liu, X. *et al.* Detecting signatures of positive selection associated with musical aptitude  
440 in the human genome. *Sci Rep* **6**, 21198 (2016).

441 40. Capra, J.A., Erwin, G.D., McKinsey, G., Rubenstein, J.L. & Pollard, K.S. Many human  
442 accelerated regions are developmental enhancers. *Philos Trans R Soc Lond B Biol Sci*  
443 **368**, 20130025 (2013).

444 41. Hedges, S.B., Marin, J., Suleski, M., Paymer, M. & Kumar, S. Tree of life reveals clock-  
445 like speciation and diversification. *Mol Biol Evol* **32**, 835-45 (2015).

- 446 42. Lynch, V.J. *et al.* Adaptive changes in the transcription factor HoxA-11 are essential for  
447 the evolution of pregnancy in mammals. *Proc Natl Acad Sci U S A* **105**, 14928-33 (2008).
- 448 43. Krol, A.J. *et al.* Evolutionary plasticity of segmentation clock networks. *Development*  
449 **138**, 2783-92 (2011).
- 450 44. Zhang, X.M., Ramalho-Santos, M. & McMahon, A.P. *Smoothened* mutants reveal  
451 redundant roles for Shh and Ihh signaling including regulation of L/R asymmetry by the  
452 mouse node. *Cell* **105**, 781-92 (2001).
- 453 45. Barrow, K.M., Ward, C.M., Rutter, J., Ali, S. & Stern, P.L. Embryonic expression of  
454 murine 5T4 oncofoetal antigen is associated with morphogenetic events at implantation  
455 and in developing epithelia. *Dev Dyn* **233**, 1535-45 (2005).
- 456 46. Sheng, G. & Foley, A.C. Diversification and conservation of the extraembryonic tissues  
457 in mediating nutrient uptake during amniote development. *Ann N Y Acad Sci* **1271**, 97-  
458 103 (2012).
- 459 47. Wu, S.F., Zhang, H. & Cairns, B.R. Genes for embryo development are packaged in  
460 blocks of multivalent chromatin in zebrafish sperm. *Genome Res* **21**, 578-89 (2011).
- 461 48. Chen, C. *et al.* Human neuronal calcium sensor-1 shows the highest expression level in  
462 cerebral cortex. *Neurosci Lett* **319**, 67-70 (2002).
- 463 49. Brawand, D. *et al.* The evolution of gene expression levels in mammalian organs. *Nature*  
464 **478**, 343-8 (2011).
- 465 50. El-Sharnouby, S., Redhouse, J. & White, R.A. Genome-wide and cell-specific epigenetic  
466 analysis challenges the role of Polycomb in *Drosophila* spermatogenesis. *PLoS Genet* **9**,  
467 e1003842 (2013).
- 468 51. Arthur, R.K. *et al.* Evolution of H3K27me3-marked chromatin is linked to gene  
469 expression evolution and to patterns of gene duplication and diversification. *Genome Res*  
470 **24**, 1115-24 (2014).
- 471 52. Lesch, B.J. & Page, D.C. Poised chromatin in the mammalian germ line. *Development*  
472 **141**, 3619-26 (2014).
- 473 53. Brykczynska, U. *et al.* Repressive and active histone methylation mark distinct promoters  
474 in human and mouse spermatozoa. *Nat Struct Mol Biol* **17**, 679-87 (2010).
- 475 54. Fortunato, S. *et al.* Genome-wide analysis of the sox family in the calcareous sponge  
476 *Sycon ciliatum*: multiple genes with unique expression patterns. *EvoDevo* **3**, 14 (2012).
- 477 55. Fortunato, S.A. *et al.* Calcisponges have a ParaHox gene and dynamic expression of  
478 dispersed NK homeobox genes. *Nature* **514**, 620-3 (2014).
- 479 56. Saudemont, A. *et al.* Complementary striped expression patterns of NK homeobox genes  
480 during segment formation in the annelid *Platynereis*. *Dev Biol* **317**, 430-43 (2008).
- 481 57. Larroux, C. *et al.* Developmental expression of transcription factor genes in a  
482 demosponge: insights into the origin of metazoan multicellularity. *Evol Dev* **8**, 150-73  
483 (2006).
- 484 58. Larroux, C. *et al.* Genesis and expansion of metazoan transcription factor gene classes.  
485 *Mol Biol Evol* **25**, 980-96 (2008).
- 486
- 487
- 488
- 489

490 **Figures**



491

492

493 **Figure 1. Gene expression and chromatin state in the mammalian germ line. a,** Phylogeny

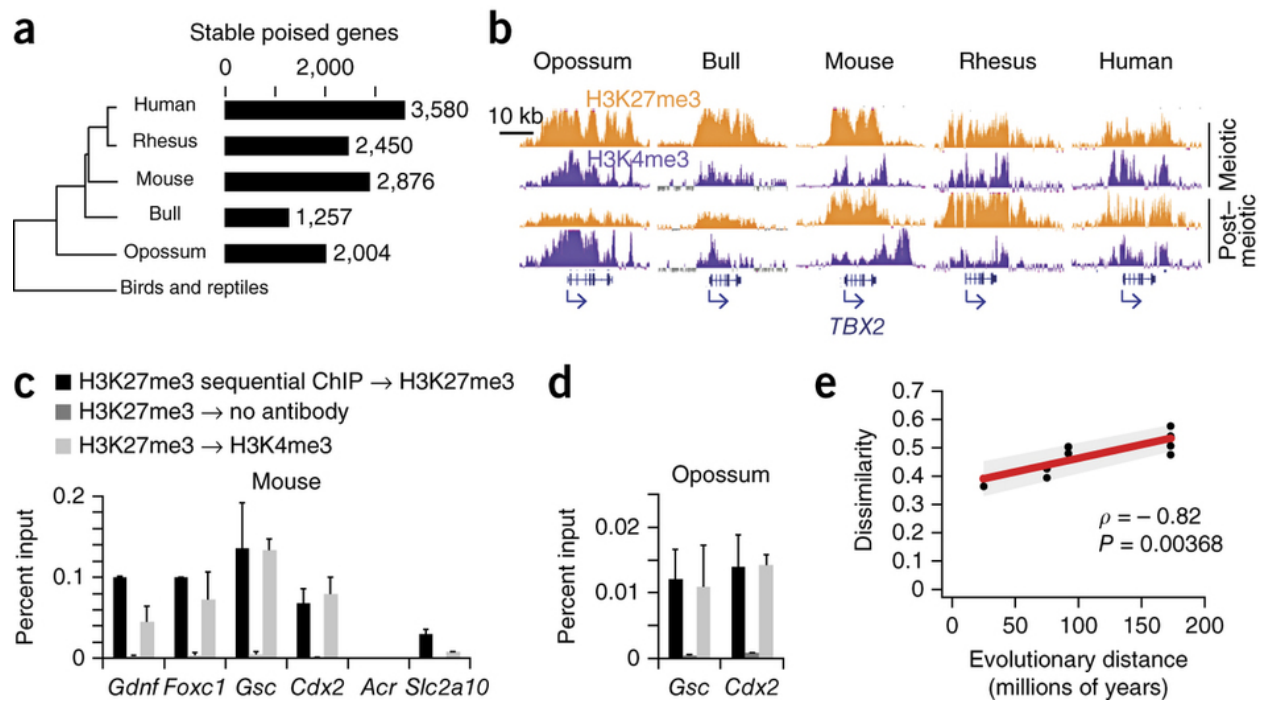
494 of the mammalian species included in this study. **b,** Heat maps showing mean gene expression

495 level as a function of H3K4me3 and H3K27me3 quantile in human pachytene spermatocytes and

496 round spermatids. Similar heat maps for all samples are shown in Supplementary Figure 2. **c,**

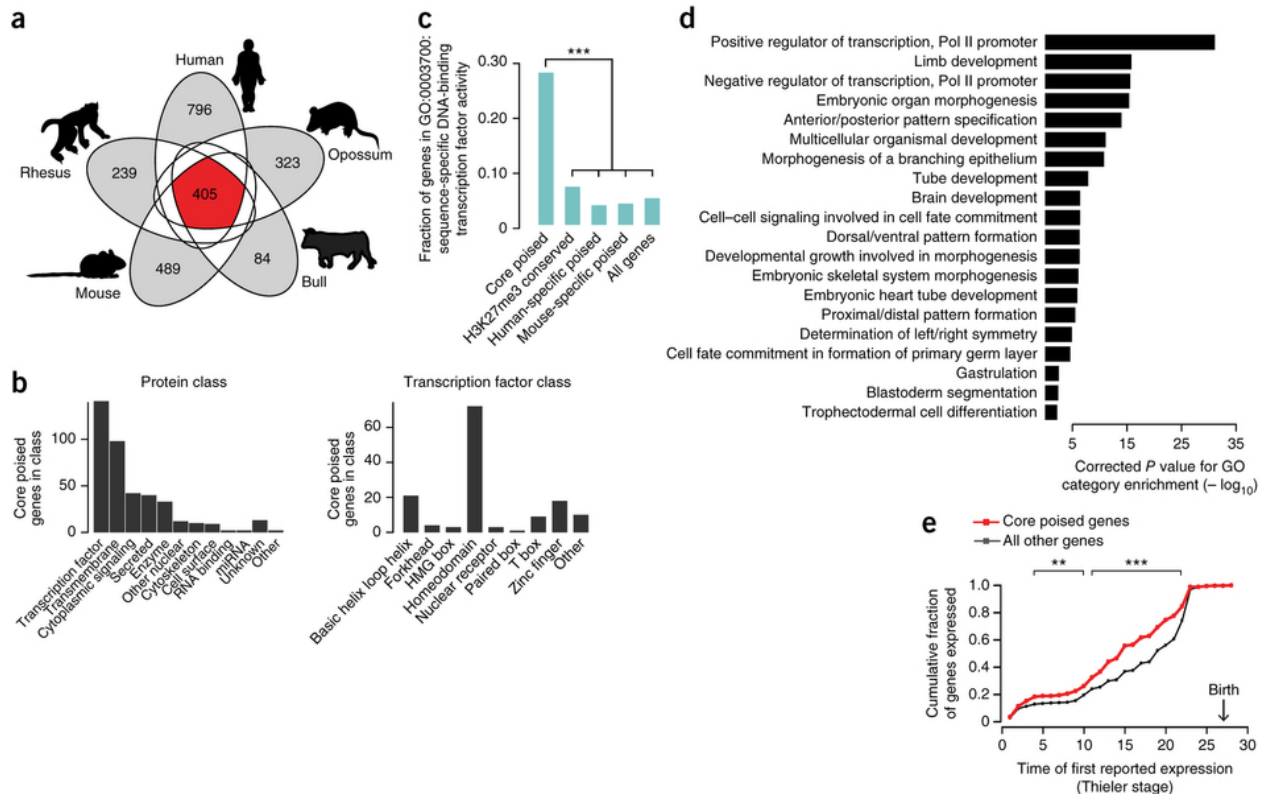
497 Hierarchical clustering of datasets by expression, H3K4me3, or H3K27me3 using  $1 - \rho$  ( $\rho$ ,

498 Spearman's correlation coefficient) as a distance metric. p.s., pachytene spermatocytes; r.s.,  
 499 round spermatids. **d**, Divergence in expression, H3K4me3, and H3K27me3 for pachytene  
 500 spermatocytes and round spermatids as a function of evolutionary distance, using  $1-\rho$  as a  
 501 dissimilarity metric. Lines represent best linear fit to pachytene spermatocyte (dashed) and  
 502 round spermatid (solid) data. The shaded area surrounding each line indicates 95% confidence  
 503 interval.



504  
 505 **Figure 2. The poised chromatin state in the mammalian germ line.** **a**, Numbers of stable  
 506 poised genes (called as poised in both pachytene spermatocytes and round spermatids) in each  
 507 species. **b**, Input-subtracted gene tracks showing H3K4me3 and H3K27me3 signal in all five  
 508 species at one representative poised gene, *TBX2*. kb, kilobases. **c**, Quantitative PCR (qPCR)  
 509 data from sequential ChIP experiments at four representative poised promoters, one H3K4me3-  
 510 only promoter (*Acr*), and one H3K27me3-only promoter (*Slc2a10*) in mouse round spermatids.  
 511 Bar height shows the mean and error bars represent s.d. for three biological replicates (*Gdnf* and

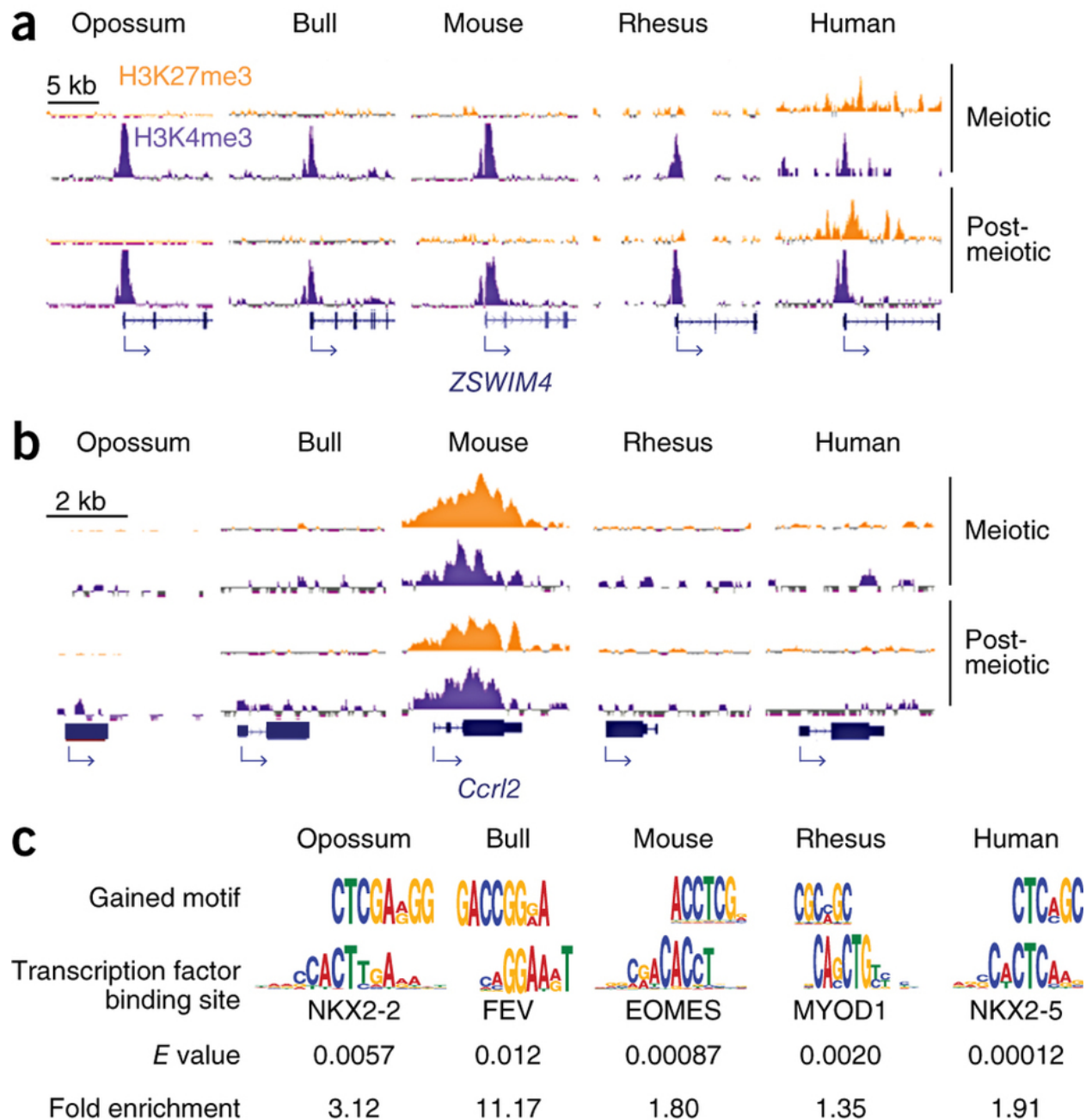
512 *Foxc1*) or three technical replicates (*Gsc*, *Cdx2*, *Acr*, *Slc2a10*). Browser tracks corresponding to  
 513 the assayed regions are shown in Supplementary Figure 7. **d**, qPCR data from sequential ChIP  
 514 experiments at two representative poised promoters in opossum round spermatids. Bar height  
 515 shows the mean and error bars represent s.d. for three technical replicates. Browser tracks  
 516 corresponding to the assayed regions are shown in Supplementary Figure 7. **e**, Dissimilarity in  
 517 poised gene sets (fraction of poised genes not shared) between pairs of species as a function of  
 518 divergence time. Grey shading indicates 95% confidence interval.  $\rho$ , Spearman's correlation  
 519 coefficient.



520 **Figure 3. Core poised genes are conserved regulators of tissue patterning.** **a**, Overlap  
 521 between poised gene sets in five mammalian species. The core poised gene set is highlighted in  
 522 red and differentially poised gene sets in grey. **b**, Classification of molecular function and  
 523 transcription factor family for the 405 genes poised in all five mammalian species. **c**, Fraction of  
 524 each gene set included in the GO category “sequence specific DNA binding transcription factor  
 525

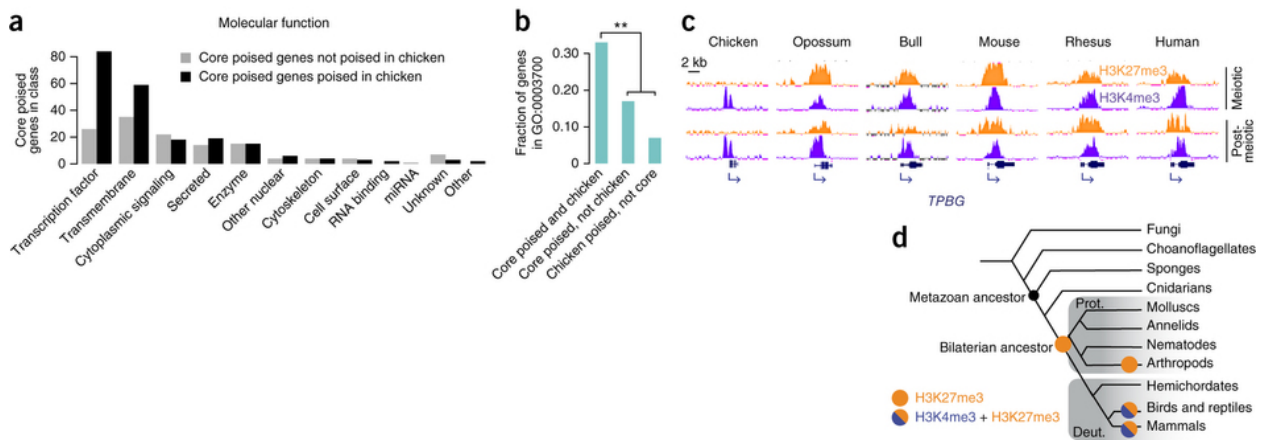


526 activity” (GO:0003700). \*\*\*p < 0.001 (Fisher’s exact test). **d**, Selected enriched GO categories  
 527 for the core poised genes. p-values were calculated using a hypergeometric test and corrected for  
 528 multiple testing (Methods). See Supplementary Table 5 for a complete list. **e**, Cumulative  
 529 fraction of core poised genes expressed throughout mouse embryogenesis, compared to all genes  
 530 with five-way orthologs. \*\*p<0.01, \*\*\*p < 0.001 (Fisher’s exact test).



531

532 **Figure 4. Gain of poising at genes with species-specific developmental roles.** **a, b,** Input-  
 533 subtracted gene tracks showing H3K4me3 and H3K27me3 signal in all five species at **(a)** the  
 534 human-specific poised gene *ZSWIM4*, which is expressed in human but not mouse or bull  
 535 placenta and **(b)** the mouse-specific poised gene *Ccrl2*, which is expressed in mouse but not  
 536 human or bull placenta. kb, kilobases. **c,** Representative gained motifs in the promoters of  
 537 differentially poised genes in each species (top) aligned with binding motifs for transcription  
 538 factors encoded by core poised genes. The core transcription factor corresponding to the binding  
 539 site is indicated below each motif, along with E-values (p-value for motif enrichment times  
 540 number of motifs tested) and fold enrichment of the gained motif compared to orthologous  
 541 sequences.



542

543 **Figure 5. Conservation of poising in the metazoan germ line.** **a,** Classification of molecular  
 544 function for the 347 genes poised in germ cells of all five mammalian species that have orthologs  
 545 in the chicken genome. Black bars represent the 215 genes poised in all five mammals and  
 546 chicken; grey bars represent the 132 genes poised in all five mammals but not chicken. **b,**  
 547 Fraction of each gene set that is included in the GO category “sequence specific DNA binding  
 548 transcription factor activity” (GO:0003700). \*\* $p < 0.01$  (Fisher’s exact test). **c,** Input-subtracted  
 549 gene tracks showing H3K4me3 and H3K27me3 signal in all five mammalian species and

550 chicken at a gene poised specifically in mammals, *TPBG*. kb, kilobases. **d**, Metazoan cladogram  
551 showing lineages with evidence for H3K4me3 and H3K27me3 (purple and orange circles) or  
552 H3K27me3 only (orange circles) at orthologs of the core poised genes in germ cells. When no  
553 circle is shown, appropriate data is not available. Shaded grey boxes indicate protostome (prot.)  
554 and deuterostome (deut.) lineages.  
555

556 **Online Methods**

557

558 **Human subjects**

559 These studies were approved by the Massachusetts Institute of Technology's Committee on the  
560 Use of Humans as Experimental Subjects. Informed consent was obtained from all subjects.

561

562 **Human sample collection and sorting**

563 Human testis samples were obtained from adult male patients undergoing vasectomy reversals at  
564 the Infertility Clinic of St. Louis. All men whose tissue was used in this study had a prior history  
565 of fertility demonstrated by at least one living child. Epididymal sperm quality and abundance  
566 proximal to the vasectomy site was assessed at the time of biopsy, and abundant, motile,  
567 morphologically normal sperm were confirmed for each patient. Testis biopsy samples were  
568 minced, dissociated using collagenase and trypsin, and then filtered to obtain a single-cell  
569 suspension as described<sup>59</sup>. Pachytene spermatocyte and round spermatid fractions were collected  
570 by StaPut<sup>59-61</sup>, and pooled fractions were counted on a hemocytometer. Purity was >95% for  
571 each sample, as assessed by counts of 100 cells from each fraction under phase optics. Cells  
572 were washed once in PBS and then split into two aliquots. One aliquot (for ChIP) was fixed in  
573 1% formaldehyde for 8 minutes at room temperature and then quenched with 2.5M glycine for 5  
574 minutes at room temperature, while the second (for RNA) was kept on ice during this time. Both  
575 fixed and unfixed aliquots were snap frozen in liquid nitrogen, then stored at -80C.

576

577

578 **Non-human sample collection and sorting**

579 Testes from rhesus monkeys were obtained from adult male animals undergoing necropsy for  
580 other purposes at the Texas Biomedical Research Institute (TBRI). The necropsy procedure was  
581 approved in advance by the TBRI Institutional Animal Care and Use Committee (IACUC).  
582 Procedures involving mice were approved in advance by the IACUC of the University of Texas  
583 at San Antonio. Testes were isolated from adult male CD1 mice (Charles River Laboratories),  
584 and tissue from several mice was pooled before cell separation. Testes from grey short-tailed  
585 opossums (*Monodelphis domestica*) were obtained from adult male animals culled from a colony  
586 maintained at the TBRI. Euthanasia of these animals was also approved by the TBRI IACUC.  
587 Testes from a bull and three roosters were obtained as abattoir material that would otherwise  
588 have been discarded. Tissue from the three roosters was pooled before cell separation. In each  
589 case, populations of pachytene spermatocytes and round spermatids were recovered using a  
590 StaPut gradient as described<sup>17,59,62-65</sup> and prepared for ChIP or RNA-Seq as described above and  
591 elsewhere<sup>17</sup>. Purity was assessed by counting 100 cells from each fraction under phase optics.  
592 Purity was 89-90% for *Monodelphis* samples, and >90% for samples from the other species.

593

594 **RNA isolation**

595 Unfixed aliquots of sorted cells were thawed on ice, washed once in cold PBS, resuspended in  
596 350 ul RLT Plus buffer from the RNEasy Mini kit (Qiagen #74134) and then disrupted by  
597 drawing up and down five times in a 26G insulin needle and syringe. Genomic DNA was  
598 removed using gDNA eliminator columns supplied with the kit. The remainder of the RNA  
599 isolation was performed using the RNEasy Mini kit according to the manufacturer's instructions.

600 Samples were processed in batches of 2-6 in the order of collection with no blinding. All  
601 biological replicates were processed in separate batches from each other.

602

### 603 **Chromatin immunoprecipitation**

604 For ChIP-seq, between  $5 \times 10^4$  and  $5 \times 10^6$  cells were used as starting material, depending on the  
605 number obtained from sample isolation and sorting. Pachytene spermatocytes and round  
606 spermatids were treated identically. For human samples, fixed cells frozen in lysis buffer (1%  
607 SDS, 10 mM EDTA, 50 mM Tris-HCl [pH 8]) were thawed on ice. For non-human samples,  
608 fixed cell pellets were thawed on ice, then washed once in cold PBS and resuspended in 100  $\mu$ l  
609 lysis buffer. Once in lysis buffer, cells were incubated on ice 5 minutes. 200  $\mu$ l ChIP dilution  
610 buffer (0.01% SDS, 1.1% Triton X-100, 1.2 mM EDTA, 16.7 mM Tris-HCl [pH 8], 167 mM  
611 NaCl) was then added to each sample. Samples were then sonicated in aliquots of 150  $\mu$ l in 0.5  
612 ml Eppendorf tubes at 4C using a BioRuptor (Diagenode) for 35 cycles on High setting, 30  
613 seconds on/30 seconds off. Aliquots of the same sample were then re-pooled and spun down at  
614 12,000  $\times$ g for 5 minutes, and the chromatin supernatant moved to a fresh tube. Chromatin from  
615 each sample was then split into two separate tubes (150  $\mu$ l in each), and 700  $\mu$ l dilution buffer, 50  
616  $\mu$ l lysis buffer, and 100 proteinase inhibitor cocktail (Complete Mini tablets, Roche  
617 #11836153001) were added to each tube. 50  $\mu$ l of each sample was set aside as input. The  
618 remainder of the ChIP was performed as previously described<sup>17</sup>, using 0.5  $\mu$ g of antibody to  
619 H3K4me3 (abcam #ab8580) or 1.0  $\mu$ g of antibody to H3K27me3 (abcam #ab6002). Samples  
620 were processed in batches of 2-4 in the order of collection with no blinding; pachytene  
621 spermatocytes and round spermatids from a given sample were processed side by side, and all  
622 biological replicates were processed in separate batches from each other.

623

624 **Sequential chromatin immunoprecipitation**

625 Sequential ChIP was performed using the Re-ChIP-IT kit from Active Motif (#53016) according  
626 to the manufacturer's instructions.  $5 \times 10^7$  mouse round spermatids or  $2 \times 10^7$  opossum round  
627 spermatids were used as starting material for each experiment. The first ChIP was performed  
628 with 2 ug of antibody to H3K27me3 (abcam, #ab6002). The second ChIP was performed with  
629 either (1) 1.5 ug of antibody to H3K4me3 (abcam, #ab8580), (2) 2.5 ul of antibody to  
630 H3K27me3 (Millipore #07449) as a positive control, or (3) 2 ul of water as a negative control.

631

632 **Quantitative PCR**

633 All primer sequences are listed in Supplementary Table 10. Quantitative PCR (qPCR) was  
634 performed on an Applied Biosystems 7500 Fast Real-Time PCR instrument using Applied  
635 Biosystems Power SYBR Green PCR Master Mix with the following cycling conditions:

636 1. 50° 0:20

637 2. 95° 10:00

638 3. 95° 0:15

639 4. 60° 1:00

640 5. 70° 0:30

641 6. Go to step 3 x39

642 For targets with biological replicates (*Gdnf* and *Foxc1*), we observed variable ChIP efficiency  
643 between experiments done on different days. To compare the relative percent input values  
644 between experimental and control conditions across biological replicates, all values from a given

645 experiment were proportionately scaled such that the value of the positive control condition  
646 (H3K27me3 → H3K27me3) was 0.1% input.

647

#### 648 **Antibodies**

649 Anti-H3K4me3 (rabbit polyclonal, Abcam #8580) was used for ChIP-seq in all species and for  
650 sequential ChIP in mouse and opossum. This antibody has been validated for ChIP-seq  
651 applications in human (Histone Modification Antibody Validation Database)<sup>66</sup>, rhesus  
652 macaque<sup>67</sup>, and mouse<sup>68</sup>, as well as non-mammalian species including *Drosophila*<sup>66</sup>. Anti-  
653 H3K27me3 (mouse monoclonal, Abcam #6002) was used for ChIP-seq in all species and for  
654 sequential ChIP in mouse and opossum. This antibody has been validated for ChIP-seq  
655 applications in human<sup>69</sup>, mouse<sup>70</sup>, and chicken<sup>71</sup>, as well as non-mammalian species including  
656 *Drosophila*<sup>66</sup> and *C. elegans*<sup>66</sup>. Anti-H3K27me3 (rabbit polyclonal, Millipore #07449) was used  
657 for sequential ChIP in mouse and opossum. This antibody has been validated for ChIP  
658 applications in mouse<sup>68</sup>, human<sup>66</sup>, and *Drosophila*<sup>72</sup>.

659

#### 660 **Illumina library preparation and sequencing**

661 RNA libraries were prepared using an Apollo 324 library prep instrument with supplied reagents  
662 (Integenx) for non-human samples, and using a SMARTer stranded RNA prep kit (Clontech) for  
663 human samples, according to the manufacturer's instructions. ChIP libraries were prepared  
664 using a TruSeq ChIP sample prep kit (Illumina), according to the manufacturer's instructions,  
665 except that size selection was performed after (instead of before) PCR amplification. Data from  
666 mouse replicate #1 has been previously published<sup>17</sup>; for this sample, ChIP- and RNA-seq  
667 libraries were sequenced on an Illumina GAII with 36-base-pair single-end reads. All other



668 libraries were sequenced on an Illumina HiSeq2500, with 40-base-pair single-end reads for ChIP  
669 libraries and 100-base-pair or 40-base-pair paired end reads for RNA-seq libraries (Table S1).

670

### 671 **Sequence alignment**

672 We filtered all datasets for read quality using FASTX-toolkit and assessed library quality using  
673 FastQC. We aligned ChIP-seq libraries to a species-appropriate genome build (hg19, rheMac2,  
674 mm10, bosTau7, monDom5, or galGal4) using Bowtie v1.1.1<sup>73</sup> (Table S1). For ChIP-seq data,  
675 we called peaks at a threshold of  $p < 10^{-6}$  using MACS v1.4<sup>74</sup>; the number and locations of peaks  
676 were used to evaluate the quality of the dataset, but were not used in our analysis. For all  
677 datasets, peak numbers were within the expected range for the appropriate histone modification  
678 (H3K4me3 or H3K27me3); variation in peak numbers within this range did not strongly affect  
679 poised gene calls. For RNA-seq data, we aligned libraries using Tophat v2.0.11<sup>75</sup> with  
680 Ensembl<sup>76</sup> (release 75) transcripts as a reference (-G flag). The default genome assemblies  
681 included in Ensembl release 75 matched those used for alignment for all species except bull. For  
682 bull, Ensembl coordinates (for bosTau6) were mapped to the bosTau7 assembly using  
683 CrossMap<sup>77</sup>.

684

685 Due to the small cell numbers of many of the sorted cell populations, both ChIP and RNA-seq  
686 libraries tended to have high duplication levels (Table S1). We treated duplicate reads  
687 conservatively for both ChIP- and RNA-seq data. For ChIP-seq data, where a minimum count  
688 threshold was used for poised gene calls, we retained only one duplicate for analysis. For RNA-  
689 seq data, where a maximum threshold was used, we retained a maximum of 20 (Tophat default)  
690 duplicate reads.

691

692 **Poised gene calls**

693 For ChIP data, we counted total reads in the 4-kb interval surrounding each transcript TSS  
694 (Ensembl build 75) using htseq-count<sup>78</sup> with the intersection-nonempty option. The 4-kb interval  
695 (2-kb upstream and 2-kb downstream of the TSS) is standard for analysis of promoter-associated  
696 histone modifications. For each dataset, total ChIP or input reads in each interval were  
697 normalized to reads per million, and the normalized input count was subtracted from the  
698 normalized ChIP count in each interval to get a final ChIP signal. For RNA, we obtained FPKM  
699 values using Cufflinks v2.2.1<sup>79</sup>, with Ensembl release 75 transcripts as a reference (-G option).  
700 Transcript values for both ChIP and RNA-seq data were summed to get a single H3K4me3,  
701 H3K27me3, and expression value for each gene. We conducted simulations in which we varied  
702 ChIP and expression thresholds and evaluated the numbers of poised genes called in each  
703 species, and selected thresholds that included the maximum number of poised genes while  
704 remaining robust to small changes in threshold value (Supplementary Figs. 5 and 6,  
705 Supplementary Data). We set thresholds of  $\geq 0.5$  input-subtracted reads per million for  
706 H3K4me3 signal,  $\geq 0.5$  input-subtracted reads per million for H3K27me3 signal, and  $\leq 5$  FPKM  
707 for expression. For each sample, a gene had to meet thresholds for H3K4me3, H3K27me3, and  
708 expression in both pachytene spermatocytes and round spermatids (six data points total) to be  
709 considered stably “poised”. For species with two biological replicates (rhesus, mouse, and  
710 opossum), we used mean ChIP signal and FPKM values to call poised genes for that species; this  
711 approach yielded similar gene lists to either the union or intersection of the two replicates, but  
712 was more robust to changes in threshold. For species with three biological replicates (human),  
713 we included genes called as poised in at least two out of the three individual replicates. We note

714 that these criteria are expected to result in greater sensitivity of poised gene calls for species with  
715 more replicates, since use of more replicates allows inclusion of genes that may fail to meet one  
716 of the six thresholds in a single replicate. We do observe the fewest poised gene calls in bull  
717 (one replicate) and the most in human (three replicates). These differences in sensitivity may  
718 result in a subset of false positives and negatives in lists of genes called as differentially poised  
719 between species. The list of conserved H3K27me3-only genes (Figures 3c, S8b, Supplementary  
720 Table 5) was defined as the set of genes that met H3K27me3 and expression thresholds in all  
721 five mammalian species, but met the H3K4me3 threshold in fewer than four out of five.

722

### 723 **Orthologous gene sets**

724 We required that a gene have orthologs in all five mammalian species in order to be included in  
725 our analysis. Since there is no strong *a priori* expectation that gene duplication would have a  
726 specific effect (loss, gain, or retention) on chromatin state surrounding the TSS, we reasoned that  
727 exclusion of genes with 1-to-many relationships could result in loss of biologically meaningful  
728 information and might introduce bias by excluding a non-random set of genes from the analysis.  
729 We therefore included genes with either 1-to-1 or 1-to-many orthology relationships among the  
730 five species. We used the BioMart database<sup>80</sup> with Ensembl release 75 to find 1-to-1 and 1-to-  
731 many orthologs among the five species. In total, we identified 14362 orthology groups  
732 containing at least one gene from each species. 12104 of these involved only 1-to-1 orthology  
733 relationships across all five species; we used this number to calculate p-values for the  
734 significance of five-way overlaps (see below). The 14362 orthology groups included a total of  
735 15492 human, 15904 rhesus, 16253 mouse, 15650 bovine, and 15966 opossum genes, which  
736 together comprised the total gene set considered in our downstream analysis. When determining

737 sets of overlapping poised genes, an orthology group was counted as overlapping if at least one  
738 gene belonging to the group was poised in each species (Supplementary Data and Supplementary  
739 Code).

740

## 741 **Statistics**

742 **Sample inclusion criteria.** Testis samples were excluded from the study if any morphological  
743 abnormality was observed in the intact tissue. For any ChIP-seq or RNA-seq dataset with  $<10^6$   
744 unique (nonduplicate) reads aligning uniquely to the genome, the associated biological sample  
745 and all datasets derived from it were excluded from analysis. These criteria were established  
746 prior to beginning the study. When possible, at least two biological replicates were obtained to  
747 allow for individual variation. For human data, three biological replicates were used in the final  
748 analysis to account for greater variability in genetic background compared to non-human species.

749 **Statistical tests.** Categorical data comparisons were evaluated using hypergeometric tests  
750 (Fisher's exact test). For comparisons of continuously distributed data, we used a two-sided  
751 Welch t-test for statistical comparison, which is robust to non-normal distributions at large  
752 sample sizes and also accounts for unequal variance between groups. We assessed variance  
753 using the Brown-Forsythe test. For ranked-list comparisons, we used a one-sided Mann-Whitney  
754 U test.

755 **Gene set overlaps.** Overlaps between multiple ( $>2$ ) gene sets were computed using the  
756 `overLapper` function from the `systemPipeR` package in R<sup>81</sup>. Statistical significance of five-way  
757 overlap was derived using the formula  $p < \binom{N}{m} \left[ \frac{\binom{N-m}{n-m}}{\binom{N}{n}} \right]^5$ , where  $N$  = total number of genes  
758 with orthologs in all five species,  $n$  = largest number of poised genes called in any single species  
759 (within the set of genes with orthologs in all five species), and  $m$  = number of genes called as

760 poised in all five species. Using our gene set,  $N=14362$ ,  $n=3580$ ,  $m=405$ , and  $p<10^{-300}$  in this  
761 calculation. However, we note that core poised genes are enriched for genes with true one-to-  
762 one orthologs, meaning that the groups being compared are not completely independent as the  
763 formula assumes. To account for this bias, we re-calculated the overlap between gene sets,  
764 including only 1-to-1 orthologs in the analysis. With 1-to-1 orthologs only,  $N=12104$ ,  $n=3361$ ,  
765  $m=401$ , and  $p<10^{-280}$ . We report this p-value as a conservative estimate of significance for the  
766 five-way overlap.

767

### 768 **Gene ontology enrichment**

769 GO enrichments were evaluated using the GOSTats package<sup>82</sup> in R. p-values were adjusted both  
770 by conditioning out child categories and by subsequent correction for multiple testing using the  
771 Benjamini-Hochberg method.

772

### 773 **Clustering and divergence estimates**

774 We generated a distance matrix for expression datasets based on FPKM, and for ChIP datasets  
775 based on normalized counts around the promoter regions (Supplementary Data), using 1- $\rho$   
776 (Spearman's correlation) as a dissimilarity metric<sup>49</sup>. Clustering was performed and dendrograms  
777 generated using the 'cluster' package<sup>83</sup> in R. Analysis of gene expression and promoter  
778 chromatin divergence was carried out using dissimilarity scores between each species pair, and  
779 data from each cell type was fit to a linear model (Supplementary Code).

780

781

### 782 **Principal component analysis**

783 We used the ‘PCA’ function from the FactoMineR package<sup>84</sup> in R for principal component  
784 analysis, with data scaled to unit value. Input data was the same processed data (normalized  
785 H3K4me3 signal, normalized H3K27me3 signal, or FPKM) as was used for calling poised genes  
786 (Supplementary Data).

787

### 788 **Somatic tissue expression**

789 The Mouse Genome Informatics (MGI) Gene Expression Database<sup>21</sup> was used to determine  
790 stages of gene expression for mouse poised genes. Only wild type samples from the database  
791 were used in the analysis. At the time of our study, the database included a total of 13837 genes;  
792 9884 genes in the database had orthologs in all five mammalian species. The numbers of genes  
793 in the database with orthologs in all five species (359 core poised genes and 9525 other genes)  
794 were used as denominators in calculating the fraction of genes expressed at each stage.

795

### 796 **Embryonic lethality**

797 We identified alleles associated with embryonic lethality by searching the MGI database using  
798 the Phenotypes, Alleles, and Disease Models query<sup>31</sup> for “embryonic lethality” (phenotype ID  
799 MP:0008762) and filtering for null/knockout alleles. This search identified 2812 alleles,  
800 corresponding to 1881 genes; 1570 of these genes had orthologs in all five mammalian species.

801

### 802 **Transcription factor class enrichment**

803 Transcription factor classes (Fig. 3b, Supplementary Fig. 8c) were assigned according to  
804 Wingender et al., 2013<sup>85</sup>.

805

## 806 **Motif analysis**

807 We identified motifs enriched in promoters of species-specific poised genes in two steps. In the  
808 first step, we used DREME<sup>86</sup> with default settings to detect motifs enriched in each set of  
809 differentially poised promoters (+/- 1kb from the transcription start site) compared to  
810 orthologous promoter regions of the other four species, with a threshold of  $E < 0.05$ . In the  
811 second step, to control for the biases introduced from comparing different species, we used  
812 AME<sup>87</sup> to scan 100 random, equally-sized sets of orthologous promoters for enrichment of the  
813 motifs detected in step (1) in the species in which they were first detected compared to the other  
814 four species. The fraction of random promoter sets demonstrating enrichment at  $E < 0.05$   
815 constituted a raw p-value; these values were adjusted for multiple comparisons using the  
816 Benjamini-Hochberg approach to get a false discovery rate for each enriched motif. We  
817 considered only motifs with  $FDR < 0.10$  in our subsequent analysis. To match enriched motifs  
818 with binding sites for known transcription factors, we used Tomtom<sup>88</sup> (conditions: -thresh 10 -  
819 evaluate -dist ed), and pulled motifs from the JASPAR Core vertebrates<sup>89</sup> (205 total motifs) and  
820 Uniprobe mouse<sup>90-92</sup> (386 total motifs) databases.

821

## 822 ***Drosophila* Polycomb data**

823 *Drosophila* CHIP-chip data was taken from El-Sharnouby et al, 2013<sup>50</sup>. Using tiled enrichment  
824 values generated by the authors, we calculated average Polycomb enrichment within 1 kilobase  
825 upstream or downstream of each *Drosophila* TSS and assigned these values to the associated  
826 gene. We designated genes in the top 25% of Polycomb signal as “high Polycomb”, and genes  
827 in the bottom 25% as “low Polycomb”.

828

829 **Code availability**

830 Custom R scripts used in our analyses are included as Supplementary material

831 (SupplementaryDataSet1.zip).



## Supplementary references

59. Bellve, A.R. Purification, culture, and fractionation of spermatogenic cells. *Methods in Enzymol* **225**, 84-113 (1993).
60. Shepherd, R.W., Millette, C.F. & DeWolf, W.C. Enrichment of primary pachytene spermatocytes from the human testes. *Mol Reprod Dev* **4**, 487-498 (1981).
61. Liu, Y. *et al.* Fractionation of human spermatogenic cells using STA-PUT gravity sedimentation and their miRNA profiling. *Sci Rep* **5**, 8084 (2015).
62. Lam, D.M., Furrer, R. & Bruce, W.R. The separation, physical characterization, and differentiation kinetics of spermatogonial cells of the mouse. *Proc Natl Acad Sci U S A* **65**, 192-9 (1970).
63. Longo, F.J., Cook, S. & Baillie, R. Characterization of an acrosomal matrix protein in hamster and bovine spermatids and spermatozoa. *Biol Reprod* **42**, 553-62 (1990).
64. Chan, J. *et al.* Characterization of the *CDKN2A* and *ARF* genes in UV-induced melanocytic hyperplasias and melanomas of an opossum (*Monodelphis domestica*). *Mol Carcinog* **31**, 16-26 (2001).
65. Oliva, R., Mezquita, J., Mezquita, C. & Dixon, G.H. Haploid expression of the rooster protamine mRNA in the postmeiotic stages of spermatogenesis. *Dev Biol* **125**, 332-40 (1988).
66. Egelhofer, T.A. *et al.* An assessment of histone-modification antibody quality. *Nat Struct Mol Biol* **18**, 91-3 (2011).
67. Liu, Y. *et al.* *Ab initio* identification of transcription start sites in the Rhesus macaque genome by histone modification and RNA-Seq. *Nucleic Acids Res* **39**, 1408-18 (2011).
68. Goldberg, A.D. *et al.* Distinct factors control histone variant H3.3 localization at specific genomic regions. *Cell* **140**, 678-91 (2010).
69. Guenther, M.G. *et al.* Chromatin structure and gene expression programs of human embryonic and induced pluripotent stem cells. *Cell Stem Cell* **7**, 249-57 (2010).
70. Shpargel, K.B., Starmer, J., Yee, D., Pohlers, M. & Magnuson, T. KDM6 demethylase independent loss of histone H3 lysine 27 trimethylation during early embryonic development. *PLoS Genet* **10**, e1004507 (2014).
71. Mitra, A. *et al.* Marek's disease virus infection induces widespread differential chromatin marks in inbred chicken lines. *BMC Genomics* **13**, 557 (2012).
72. Rebollo, R. *et al.* A snapshot of histone modifications within transposable elements in *Drosophila* wild type strains. *PLoS One* **7**, e44253 (2012).
73. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* **10**, R25 (2009).
74. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137 (2008).
75. Trapnell, C., Pachter, L. & Salzberg, S.L. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* **25**, 1105-11 (2009).
76. Yates, A. *et al.* Ensembl 2016. *Nucleic Acids Res* **44**, D710-6 (2016).
77. Zhao, H. *et al.* CrossMap: a versatile tool for coordinate conversion between genome assemblies. *Bioinformatics* **30**, 1006-7 (2014).
78. Anders, S., Pyl, P.T. & Huber, W. HTSeq—a Python framework to work with high-throughput sequencing data. *Bioinformatics* (2014).
79. Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* **28**, 511-5 (2010).
80. Durinck, S. *et al.* BioMart and Bioconductor: a powerful link between biological databases and microarray data analysis. *Bioinformatics* **21**, 3439-40 (2005).

81. R Core Team. R: A language and environment for statistical computing. (2015).
82. Falcon, S. & Gentleman, R. Using GOSTATS to test gene lists for GO term association. *Bioinformatics* **23**, 257-8 (2007).
83. Maechler, M., Rousseeuw, P., Struyf, A., Hubert, M. & Hornik, K. cluster: cluster analysis basics and extensions. version 2.0.3 edn (2015).
84. Husson, F., Josse, J., Le, S. & Mazet, J. FactoMineR: multivariate exploratory data analysis and data mining. version 1.32 edn (2016).
85. Wingender, E., Schoeps, T. & Donitz, J. TFClass: an expandable hierarchical classification of human transcription factors. *Nucleic Acids Res* **41**, D165-70 (2013).
86. Bailey, T.L. DREME: motif discovery in transcription factor ChIP-seq data. *Bioinformatics* **27**, 1653-9 (2011).
87. McLeay, R.C. & Bailey, T.L. Motif enrichment analysis: a unified framework and an evaluation on ChIP data. *BMC Bioinformatics* **11**, 165 (2010).
88. Gupta, S., Stamatoyannopoulos, J.A., Bailey, T.L. & Noble, W.S. Quantifying similarity between motifs. *Genome Biol* **8**, R24 (2007).
89. Mathelier, A. *et al.* JASPAR 2016: a major expansion and update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res* **44**, D110-5 (2016).
90. Badis, G. *et al.* Diversity and complexity in DNA recognition by transcription factors. *Science* **324**, 1720-3 (2009).
91. Berger, M.F. *et al.* Variation in homeodomain DNA binding revealed by high-resolution analysis of sequence preferences. *Cell* **133**, 1266-76 (2008).
92. Hume, M.A., Barrera, L.A., Gisselbrecht, S.S. & Bulyk, M.L. UniPROBE, update 2015: new tools and content for the online database of protein-binding microarray data on protein-DNA interactions. *Nucleic Acids Res* **43**, D117-22 (2015).
93. Wyckoff, G.J., Wang, W. & Wu, C.I. Rapid evolution of male reproductive genes in the descent of man. *Nature* **403**, 304-9 (2000).
94. Good, J.M. & Nachman, M.W. Rates of protein evolution are positively correlated with developmental timing of expression during mouse spermatogenesis. *Mol Biol Evol* **22**, 1044-52 (2005).
95. Soumillon, M. *et al.* Cellular source and mechanisms of high transcriptome complexity in the mammalian testis. *Cell Rep* **3**, 2179-90 (2013).

### Competing financial interests

The authors have no competing financial interests.