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

Bluma J Lesch[1], Sherman J Silber[2], John R McCarrey[3] & David C Page[1,4,5]

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

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

We focused our attention on the evolution of epigenetic 'poising' in germ cells. Epigenetic poising is defined by the simultaneous presence of two opposing histone modifications, the activating mark trimethylation of histone H3 at lysine 4 (H3K4me3) and the repressive mark trimethylation of histone H3 at lysine 27 (H3K27me3), as well as by transcriptional repression[12–14]. It has been best studied in embryonic stem cells (ESCs), where it is associated with genes involved in lineage specification. In ESCs, the presence of the activating H3K4me3 mark is thought to poise these otherwise silent genes for activation following receipt of a differentiation cue, whereas the H3K27me3 mark maintains repression in the pluripotent state[12]. Consistent with this model, in mouse embryos *in vivo*, Hox genes move sequentially from a poised state to an active, H3K4me3-only state, concurrent with their activation in an anterior-to-posterior direction[15,16]. Hox genes and other developmental regulators are poised in mouse and human germ cells but are not expressed in developing gametes; rather, genes poised in germ cells are expressed in somatic tissues during embryogenesis[17–20]. We hypothesized that evolution of epigenetic poising in mammalian germ cells reflects the evolution of a transcriptional program controlling somatic gene expression and morphogenesis in embryos.

Evaluating this hypothesis demanded that we carry out four tasks: characterization of the poised state in the germ cells of multiple mammalian species; examination of the relationship between conservation of germ cell poising and conservation of developmental function in mammals; definition of the relationship between differences in germ cell poising and differences in developmental function among specific evolutionary lineages; and reconstruction of the evolutionary origins of these relationships by comparison to non-mammalian taxa.

We used comparative epigenetic profiling in mammalian germ cells to perform these tasks. We examined genome-wide expression, H3K4me3, and H3K27me3 data in male germ cells from six species spanning 300 million years of evolution. We found that evolution of poising in male germ cells parallels evolution of somatic gene expression and development in the embryo. We propose an ancient evolutionary relationship between germline chromatin state and embryonic gene expression.
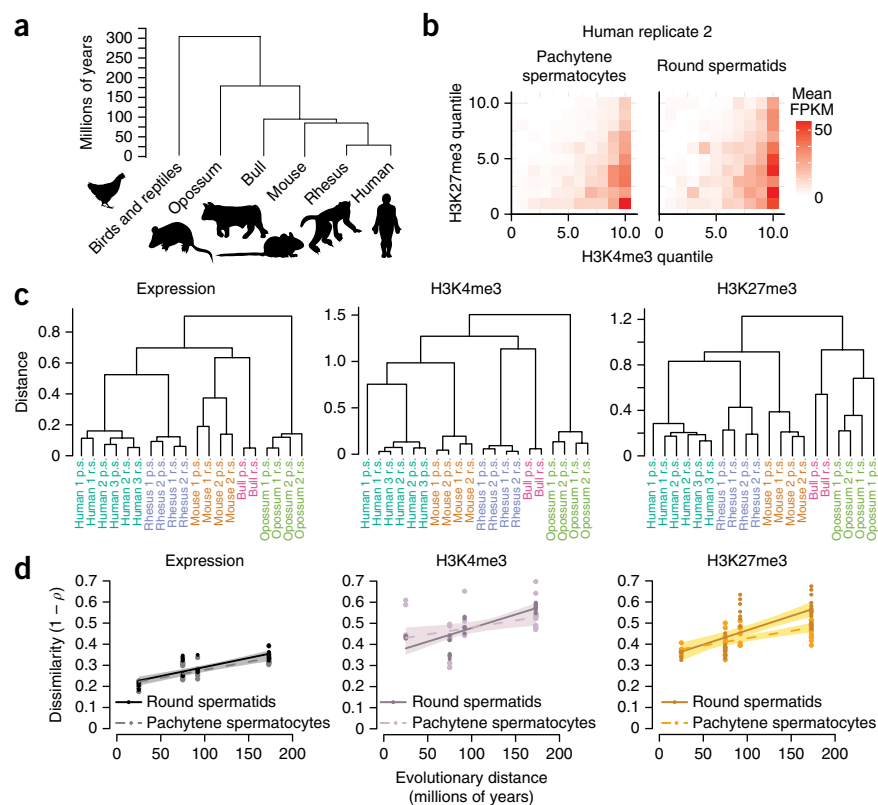
## RESULTS

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

We collected H3K4me3 and H3K27me3 chromatin immunoprecipitation and sequencing (ChIP-seq) data, as well as RNA-seq data, in sorted male germ cells from five species spanning 175 million years of evolutionary divergence in the mammalian lineage: human, rhesus macaque, mouse, bull, and opossum (**Fig. 1a**, **Supplementary Fig. 1**, and **Supplementary Table 1**). We obtained data from cells at two time

**Figure 1** Gene expression and chromatin state in the mammalian germ line. (**a**) Phylogeny of the mammalian species included in this study. (**b**) Heat maps showing mean gene expression level as a function of H3K4me3 and H3K27me3 quantile in human pachytene spermatocytes and round spermatids. Similar heat maps for all samples are shown in **Supplementary Figure 2**. (**c**) Hierarchical clustering of data sets by expression, H3K4me3, or H3K27me3 using $1 - \rho$ (where $\rho$ is Spearman's correlation coefficient) as a distance metric. p.s., pachytene spermatocytes; r.s., round spermatids. (**d**) Divergence in expression, H3K4me3, and H3K27me3 for pachytene spermatocytes and round spermatids as a function of evolutionary distance, using $1 - \rho$ as a dissimilarity metric. Lines represent best linear fit to pachytene spermatocyte (dashed) and round spermatid (solid) data. The shaded area surrounding each line represents the 95% confidence interval.



points during male germ cell development: prophase of meiosis I (pachytene spermatocytes) and after completion of meiosis (round spermatids). These cell types are unique in that they can be identified and reliably collected from whole testes in multiple mammalian species without the use of transgenes or genetic markers (Online Methods). In addition, they differ substantially from each other in their place in the cell cycle and in the physical state of their chromatin: pachytene spermatocytes are tetraploid cells in meiotic prophase, with large nuclei and synapsed pairs of homologous chromosomes, whereas round spermatids are haploid cells that have completed meiosis and have small, compact nuclei. Inclusion of both cell types in our study allowed us to control for the effects of chromatin compaction and physical state of the nucleus during spermatogenic development.

Before turning to analyses of the poised state, we first examined the relationships among the H3K4me3, H3K27me3, and expression data sets collected from different species. In all species, expression levels were positively correlated with H3K4me3 signal and negatively correlated with H3K27me3 signal, consistent with the known association of these marks with gene activity and repression, respectively (**Fig. 1b** and **Supplementary Fig. 2**). Our data included three biological replicates for human, two biological replicates for rhesus macaque, mouse, and opossum, and one biological replicate for bull (**Supplementary Fig. 3**); biological replicates were highly similar as evaluated by principal-component analysis (**Supplementary Fig. 4a**) or by hierarchical clustering (**Fig. 1c** and **Supplementary Fig. 4b**). For H3K4me3, clustering accurately separated pachytene spermatocytes from round spermatids in each species but did not fully recapitulate phylogenetic relationships, as has been previously reported for the H3K4me3 mark in other tissues[10]. In contrast, we derived correct or nearly correct phylogenies from both expression and H3K27me3 data. For expression, H3K4me3, and H3K27me3 data, greater dissimilarity between samples corresponded to greater evolutionary divergence between species (**Fig. 1d** and **Supplementary Note**).

## Identification of poised chromatin in male germ cells
To identify genes associated with poised chromatin in germ cells from each species, we calculated read counts in 4-kb intervals surrounding
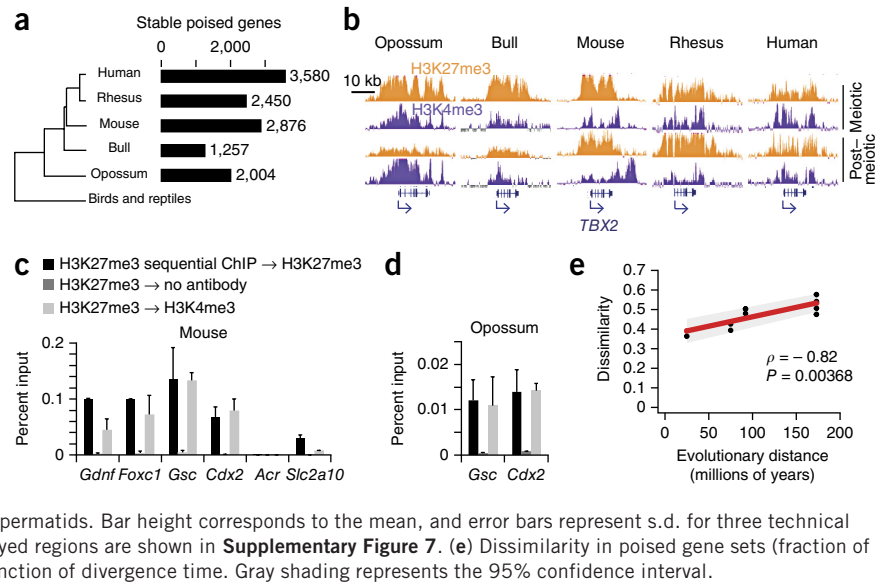
the transcriptional start sites of annotated genes, after normalizing for library size and subtracting input signal (**Supplementary Data**). Throughout our analysis, we considered only genes with orthologs in all five species (a total of 14,362 orthology groups). Genes with signal above a threshold of 0.5 reads per million for H3K4me3 and H3K27me3 and whose expression level was equal to or less than 5 FPKM were called as poised. In each species, we further filtered for genes at which the poised state was retained in both pachytene spermatocytes and round spermatids, implying that it is stable across much of spermatogenic development. Stably poised gene sets identified in this manner were robust to changes in ChIP and expression thresholds (**Supplementary Figs. 5** and **6**).

This approach identified 1,200–3,600 poised genes in each species (**Fig. 2a,b** and **Supplementary Table 2**). We verified that co-occurrence of high H3K4me3 and H3K27me3 signals at poised genes represented the simultaneous presence of the two marks on the same DNA molecule, not heterogeneity of chromatin state within our cell population, by performing sequential ChIP at the promoters of representative poised genes in both mouse (**Fig. 2c** and **Supplementary Fig. 7a**) and opossum (**Fig. 2d** and **Supplementary Fig. 7b**) round spermatids. We confirmed that four of four mouse and two of two opossum poised promoters were simultaneously marked by both H3K4me3 and H3K27me3, demonstrating that these genes are marked by a bona fide poised state in round spermatids. In general, between one-quarter and three-quarters of poised genes were shared by any two species (**Supplementary Table 3**), and dissimilarity in poising was positively correlated with evolutionary divergence time (**Fig. 2e**). Regardless of evolutionary distance, overlap in poised gene sets for each species pair was greater than expected by chance ($P < 1 \times 10^{-15}$ for all pairs, Fisher's exact test).

## Conserved poising at developmental regulators in mammals
Four hundred and five genes were poised in all five species ($P < 1 \times 10^{-280}$, compared to the number expected for five-way overlap; Online

**Figure 2** The poised chromatin state in the mammalian germ line. (**a**) Number of stable poised genes (called as poised in both pachytene spermatocytes and round spermatids) in each species. (**b**) Input-subtracted gene tracks showing H3K4me3 and H3K27me3 signal in all five species at one representative poised gene, *TBX2*. (**c**) qPCR data from sequential ChIP experiments at four representative poised promoters, one H3K4me3-only promoter (*Acr*), and one H3K27me3-only promoter (*Slc2a10*) in mouse round spermatids. Bar height corresponds to the mean, and error bars represent s.d. for three biological replicates (*Gdnf* and *Foxc1*) or three technical replicates (*Gsc*, *Cdx2*, *Acr*, and *Slc2a10*). Browser tracks corresponding to the assayed regions are shown in **Supplementary Figure 7**. (**d**) qPCR data from sequential ChIP experiments at two representative poised promoters in opossum round spermatids. Bar height corresponds to the mean, and error bars represent s.d. for three technical replicates. Browser tracks corresponding to the assayed regions are shown in **Supplementary Figure 7**. (**e**) Dissimilarity in poised gene sets (fraction of poised genes not shared) for pairs of species as a function of divergence time. Gray shading represents the 95% confidence interval.



Methods) (**Fig. 3a** and **Supplementary Table 4**). These genes were well distributed across chromosomes (**Supplementary Fig. 8a**). At the sequence level, the promoter regions of these genes were significantly better conserved than those of human-specific poised genes ($P < 1 \times 10^{-14}$, Welch $t$ test) or genes with conserved retention of H3K27me3 but not necessarily H3K4me3 ($P < 1 \times 10^{-4}$, Welch $t$ test) (**Supplementary Fig. 8b**). The set of genes poised in the germ lines of all five mammalian species, henceforth referred to as 'core' poised genes, was strongly enriched for genes encoding transcription factors, with a particularly striking enrichment for homeodomain-containing transcription factors (**Fig. 3b** and **Supplementary Fig. 8c**).

We used predefined Gene Ontology (GO) categories to confirm enrichment of genes encoding transcription factors in the core poised gene set. The 405 core poised genes were significantly enriched for genes belonging to the GO category 'sequence-specific DNA-binding transcription factor activity' (GO:0003700) when compared to all genes with five-way orthologs, to human- or mouse-specific poised genes, or to genes with conserved retention of H3K27me3 but not necessarily H3K4me3 (**Fig. 3c**).

We then asked whether, in addition to encoding proteins with a shared molecular function as sequence-specific transcription factors, the set of core poised genes had a unifying biological function during development. We examined enrichment of GO biological function categories in the set of core poised genes. We found that enriched GO categories described processes involved in patterning and organ formation, including 'embryonic organ morphogenesis', 'anterior/posterior pattern specification', 'limb development', and 'gastrulation' (**Fig. 3d** and **Supplementary Table 5**). Core poised genes were not enriched for germ-cell-related functions such as meiosis, spermatid development, or sperm maturation. These findings imply that the core poised genes have a shared biological role, specifically, transcriptional regulation of body patterning and somatic tissue specification during embryogenesis.
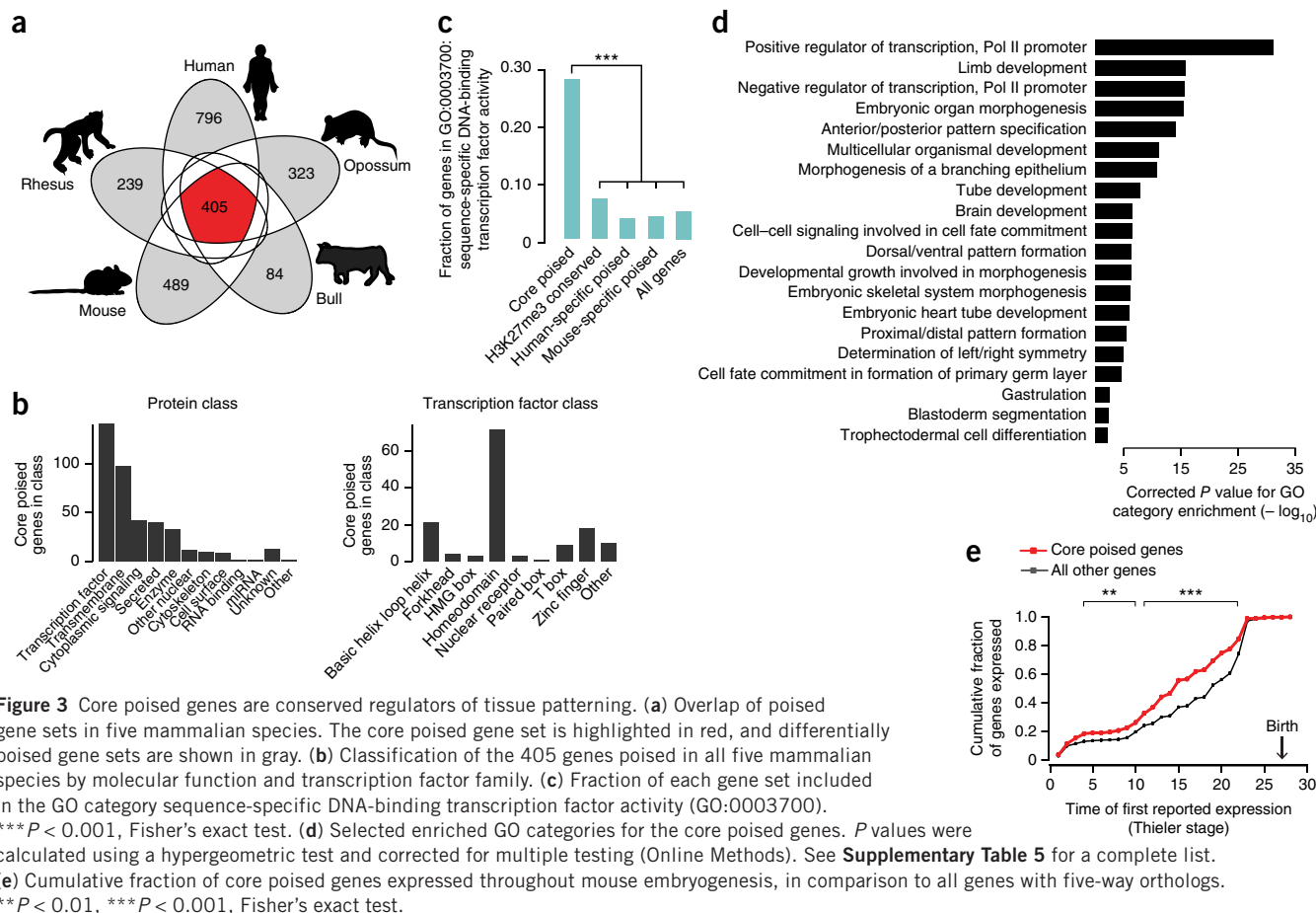
If the core poised genes are involved in patterning and tissue specification, they should be expressed during an interval in embryogenesis when these processes are occurring. We queried the Mouse Genome Informatics (MGI) Gene Expression Database[21], which contained expression data for 13,837 mouse genes at the time of our study, to determine the interval when each of these genes is first expressed. Core poised genes were enriched for expression during embryogenesis

when compared to other genes with orthologs in all five species; this difference was especially evident between gastrulation and the end of somite formation (Thieler stages (TS) 11–22; **Fig. 3e**). A subset of core poised genes involved in trophectoderm specification (for example, *Cdx2* (encoding caudal-type homeobox 2), *Hand1* (heart and neural crest derivatives expressed transcript 1), and *Tpbg* (trophoblast glycoprotein)) was also expressed early in embryogenesis (TS 2–4)[22].

We then examined specific cases where the regulatory hierarchies involved in body part and organ field specification are well defined and found that core poised genes are central to these processes. Core poised genes sit at the top of such specification hierarchies, including *PTF1A* (pancreas-specific transcription factor 1a) and *PDX1* (pancreatic and duodenal homeobox 1) in pancreatic development[23]; *EN1* and *EN2* (engrailed 1 and 2), *OTX1* (orthodenticle homeobox 1), *LMX1A* (LIM homeobox transcription factor α), and *GBX2* (gastrulation brain homeobox 2) in cerebellar development[24]; *NKX2-5* (NK homeobox 2-5) and *HAND1* and *HAND2* (heart and neural crest derivatives expressed transcript 1 and 2) in heart development[25,26]; and *MSX1* and *MSX2* (Msh homeobox 1 and 2) and *NKX2-2* (NK homeobox 2-2) in neural tube regionalization[27].

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

We conclude that the core poised genes constitute critical upstream regulators of gene expression during mammalian embryogenesis. Indeed, many of the core poised genes participate in developmental 'kernels', conserved genetic circuits controlling specification of body part progenitor fields[29,30]. Kernel architecture extends deep into the metazoan lineage and is highly conserved across metazoa, partly because extensive regulatory interactions among kernel constituents

**Figure 3** Core poised genes are conserved regulators of tissue patterning. (**a**) Overlap of poised gene sets in five mammalian species. The core poised gene set is highlighted in red, and differentially poised gene sets are shown in gray. (**b**) Classification of the 405 genes poised in all five mammalian species by molecular function and transcription factor family. (**c**) Fraction of each gene set included in the GO category sequence-specific DNA-binding transcription factor activity (GO:0003700). ***$P < 0.001$, Fisher's exact test. (**d**) Selected enriched GO categories for the core poised genes. $P$ values were calculated using a hypergeometric test and corrected for multiple testing (Online Methods). See **Supplementary Table 5** for a complete list. (**e**) Cumulative fraction of core poised genes expressed throughout mouse embryogenesis, in comparison to all genes with five-way orthologs. **$P < 0.01$, ***$P < 0.001$, Fisher's exact test.

mean that perturbation of any one component can have catastrophic effects on body part patterning[29,30]. In fact, we found that knockout alleles of 83 of the core poised genes (21%) resulted in embryonic lethality in mouse, whereas knockout alleles of only 10% of all genes with orthologs in all five mammalian species resulted in the same phenotype ($P = 3.77 \times 10^{-11}$, Fisher's exact test)[31].

## Differences in germline poising between species

Given that genes with strongly conserved roles in metazoan development exhibit conservation of poising in germ cells, we wondered whether genes that gain species-specific developmental functions might also acquire species-specific poising in germ cells. To address this question, we turned our attention to the five sets of genes poised specifically in only one of the five species evaluated ('differentially poised' genes), acknowledging that a subset of these genes may be misassigned as specific owing to false negative calls in one or more of the other four species (**Fig. 4a,b** and **Supplementary Table 6**). None of the five sets of differentially poised genes was strongly enriched for GO developmental functions. However, where comparative expression data in specific developmental structures was available, species differences
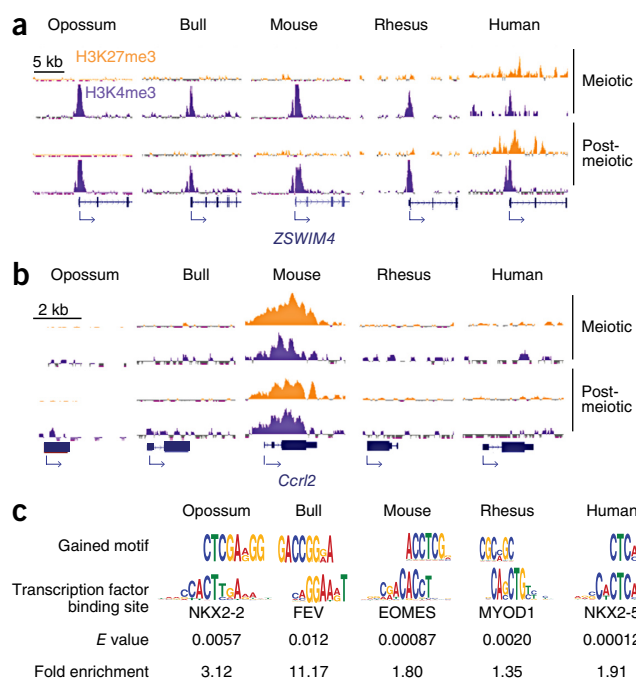


**Figure 4** Gain of poising at genes with species-specific developmental roles. (**a,b**) Input-subtracted gene tracks showing H3K4me3 and H3K27me3 signal in all five species at the human-specific poised gene *ZSWIM4*, which is expressed in human but not mouse or bull placenta (**a**) and the mouse-specific poised gene *Ccrl2*, which is expressed in mouse but not human or bull placenta (**b**). (**c**) Representative gained motifs in the promoters of differentially poised genes in each species (top) aligned with binding motifs for transcription factors encoded by core poised genes. The core transcription factor corresponding to the binding site is indicated below each motif, along with the $E$ value (the $P$ value for motif enrichment multiplied by the number of motifs tested), and fold enrichment of the gained motif in comparison to orthologous sequences.
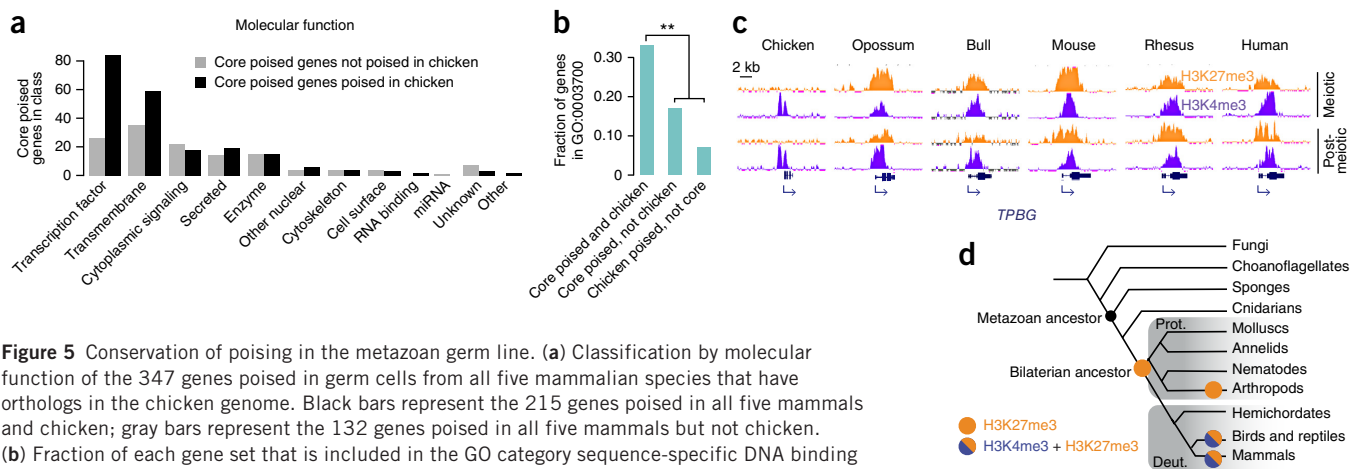
**Figure 5** Conservation of poising in the metazoan germ line. (**a**) Classification by molecular function of the 347 genes poised in germ cells from all five mammalian species that have orthologs in the chicken genome. Black bars represent the 215 genes poised in all five mammals and chicken; gray bars represent the 132 genes poised in all five mammals but not chicken. (**b**) Fraction of each gene set that is included in the GO category sequence-specific DNA binding transcription factor activity (GO:0003700). **$P < 0.01$, Fisher's exact test. (**c**) Input-subtracted gene tracks showing H3K4me3 and H3K27me3 signal in all five mammalian species and chicken at a gene poised specifically in mammals, *TPBG*. (**d**) Metazoan cladogram showing lineages with evidence for H3K4me3 and H3K27me3 (purple and orange circles) or H3K27me3 only (orange circles) at orthologs of the core poised genes in germ cells. When no circle is shown, appropriate data are not available. Shaded gray areas indicate the protostome (prot.) and deuterostome (deut.) lineages.

in germline poising were correlated with species differences in developmental expression[8,9,32,33]. For example, *ZSWIM4* and *LMF1*, which are poised specifically in human germ cells (**Fig. 4a** and **Supplementary Fig. 9a**), are expressed in human but not mouse or bovine placenta; likewise, *Ccrl2* and *Smug1* are poised specifically in mouse germ cells (**Fig. 4b** and **Supplementary Fig. 9b**) and expressed in mouse but not human or bovine placenta[32]. Similarly, *HIPK2* is poised specifically in human germ cells (**Supplementary Fig. 9a**) and has acquired a human-specific enhancer active during limb development[8].

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

We also identified cases of single-lineage loss, in which a gene was poised in all but one of the five species examined (**Supplementary Table 8**). As with single-lineage gains, some of these instances may be due to false negative poised gene calls in one species. However, some instances of single-lineage loss of poising are supported by previous reports of recent evolutionary divergence in expression or function of the associated gene. For example, of the 36 genes poised in four mammals but not in human, 3 are reported to have divergent expression patterns in human in comparison to other mammals (*AIM1*, *EPHA5*, and *THBS4*)[34–36], 1 is associated with differences in loss-of-function phenotype between human and mouse (*DOCK8*)[37], and 3 are associated with recent positive selection in the human lineage (*AIM1*,

*COL11A1*, and *LYPD1*)[38–40]. Like lineage-specific gain, lineage-specific loss of poising in the germ line may therefore reflect recent lineage-specific changes in developmental regulation and function.

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

First, we compared our mammalian data to a non-mammalian amniote, the chicken. Together with reptiles, birds constitute the closest living relatives of the mammalian clade (**Fig. 1a**)[41]. We collected ChIP-seq and RNA-seq data from chicken germ cells at time points matching those used for the five mammalian species (**Supplementary Data**). Using identical filtering conditions, we identified 1,716 genes poised in the chicken germ line (**Supplementary Table 9**). Of the 405 core poised genes we defined in mammals, 347 have orthologs in the chicken genome; of these, 215 (62%) are also poised in the chicken germ line. For the core poised genes in mammals, the set of genes whose orthologs were also poised in chicken was significantly enriched for sequence-specific transcription factors when compared to the core poised genes that were not poised in chicken (**Fig. 5a,b**). We conclude that epigenetic poising of developmental transcriptional regulators in germ cells is at least as old as the amniote common ancestor, placing its origin more than 300 million years ago. In at least five cases, core poised genes that were poised in mammalian but not chicken germ cells could also be correlated to differences in development between mammals and birds[42–45], supporting the hypothesis that acquisition of poising in the germ line is related to acquisition of somatic developmental function. For example, *TPBG* is poised in all five mammals but not in chicken (**Fig. 5c**), consistent with its early expression in trophectoderm, a mammal-specific structure[45,46]. Given previous reports of a multivalent chromatin state

comprising multiple repressive and activating histone marks at developmental genes in zebrafish sperm[47], it will be interesting to trace the evolutionary history of the poised state in additional non-amniote vertebrate species.

Using the chicken data, we further examined the scenario in which genes whose poised state was shared in distantly related species were not poised in the germ cells of intermediate lineages. Such a scenario implies either convergent evolution, requiring two independent epigenetic gains, or deep loss of poising followed by recent reacquisition, requiring a loss followed by a gain. Either explanation calls for at least two independent evolutionary events and is expected to be rarer than scenarios requiring either uninterrupted conservation (zero events) or single-lineage gain or loss (one event). Indeed, of the 11,188 genes with orthologs in all six species, 211 were poised in human, rhesus macaque, and mouse only, implying a single gain in the primate–rodent ancestor, in comparison to only 35 in human, opossum, and chicken and 14 in rhesus macaque, opossum, and chicken, each set requiring either convergence or deep loss in the placental lineage followed by a gain (**Supplementary Table 9**). For at least one gene among the 35 shared by human, opossum, and chicken (*NCS1*, encoding neuronal calcium sensor 1), the expression patterns in human and chicken were more similar than those in human and mouse, implying convergent evolution of expression[48,49].

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

## DISCUSSION

We show here that evolution of epigenetic poising in male germ cells is closely linked to evolution of somatic gene expression in developing mammalian embryos. Germline poising is conserved throughout the mammalian lineage at genes that are central to the transcriptional networks governing somatic development, and individual genes recruited to these networks in specific lineages also gain poising in germ cells. Poising of central developmental genes in male germ cells is at least as old as the amniote common ancestor, placing its origin at least 300 million years ago. Such deep conservation implies a functional role for the poised state in germ cells.

It is easy to envision a role for H3K27me3 at somatic genes in the germ line: this repressive mark reinforces silencing of genes whose expression in germ cells would disrupt their function and identity. However, we found that genes exhibiting conservation of H3K27me3 without H3K4me3 do not show the same functional enrichments as genes with conserved poising (**Fig. 3c** and **Supplementary Table 5**), indicating that H3K27me3 alone does not have the same biological role as H3K27me3/H3K4me3 bivalency. H3K4me3 may have a protective role at poised promoters as an antagonist of DNA methylation[5,52]. It is also possible that H3K4me3 helps to prepare poised genes for expression in somatic tissues following fertilization, similar to its proposed role in ESCs. Consistent with this hypothesis, altered

regulation of H3K4 methylation state in developing male germ cells in mouse perturbs somatic tissue development in embryos of the next generation[3]. The detailed mechanism by which this epigenetic information might be transmitted through fertilization remains unclear: modified histones may be carried in mature spermatozoa[19,20,53], or poised sites may be marked by an RNA or protein intermediate in sperm and reestablished in the early embryo. We found that published ChIP-seq data from mature mouse[20] and human[19] spermatozoa are consistent with retention of modified histones at core poised genes (**Supplementary Fig. 10b**), but this finding does not exclude the participation of additional factors in marking poised sites.

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

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

**URLs.** FASTX-Toolkit, http://hannonlab.cshl.edu/fastx_toolkit/index.html; FastQC, http://www.bioinformatics.babraham.ac.uk/projects/fastqc/; MGI Gene Expression Data, http://www.informatics.jax.org/gxd; MGI Phenotypes, Alleles and Disease Modules, http://www.informatics.jax.org/allele; R Project for Statistical Computing, http://www.R-project.org/.

## METHODS

Methods and any associated references are available in the online version of the paper.

1. Kimmins, S. & Sassone-Corsi, P. Chromatin remodelling and epigenetic features of germ cells. *Nature* **434**, 583–589 (2005).

2. Kurimoto, K. *et al.* Quantitative dynamics of chromatin remodeling during germ cell specification from mouse embryonic stem cells. *Cell Stem Cell* **16**, 517–532 (2015).

3. Siklenka, K. *et al.* Disruption of histone methylation in developing sperm impairs offspring health transgenerationally. *Science* **350**, aab2006 (2015).

4. Arico, J.K., Katz, D.J., van der Vlag, J. & Kelly, W.G. Epigenetic patterns maintained in early *Caenorhabditis elegans* embryos can be established by gene activity in the parental germ cells. *PLoS Genet.* **7**, e1001391 (2011).

5. Ihara, M. *et al.* Paternal poly (ADP-ribose) metabolism modulates retention of inheritable sperm histones and early embryonic gene expression. *PLoS Genet.* **10**, e1004317 (2014).

6. King, M.C. & Wilson, A.C. Evolution at two levels in humans and chimpanzees. *Science* **188**, 107–116 (1975).

7. Wray, G.A. The evolutionary significance of *cis*-regulatory mutations. *Nat. Rev. Genet.* **8**, 206–216 (2007).

8. Cotney, J. *et al.* The evolution of lineage-specific regulatory activities in the human embryonic limb. *Cell* **154**, 185–196 (2013).

9. Reilly, S.K. *et al.* Evolutionary changes in promoter and enhancer activity during human corticogenesis. *Science* **347**, 1155–1159 (2015).

10. Villar, D. *et al.* Enhancer evolution across 20 mammalian species. *Cell* **160**, 554–566 (2015).

11. Prescott, S.L. *et al.* Enhancer divergence and *cis*-regulatory evolution in the human and chimp neural crest. *Cell* **163**, 68–83 (2015).

12. Bernstein, B.E. *et al.* A bivalent chromatin structure marks key developmental genes in embryonic stem cells. *Cell* **125**, 315–326 (2006).

13. Mikkelsen, T.S. *et al.* Genome-wide maps of chromatin state in pluripotent and lineage-committed cells. *Nature* **448**, 553–560 (2007).

14. Azuara, V. *et al.* Chromatin signatures of pluripotent cell lines. *Nat. Cell Biol.* **8**, 532–538 (2006).

15. Soshnikova, N. & Duboule, D. Epigenetic temporal control of mouse Hox genes *in vivo*. *Science* **324**, 1320–1323 (2009).

16. Noordermeer, D. *et al.* Temporal dynamics and developmental memory of 3D chromatin architecture at Hox gene loci. *eLife* **3**, e02557 (2014).

17. Lesch, B.J., Dokshin, G.A., Young, R.A., McCarrey, J.R. & Page, D.C. A set of genes critical to development is epigenetically poised in mouse germ cells from fetal stages through completion of meiosis. *Proc. Natl. Acad. Sci. USA* **110**, 16061–16066 (2013).

18. Sachs, M. *et al.* Bivalent chromatin marks developmental regulatory genes in the mouse embryonic germline *in vivo*. *Cell Rep.* **3**, 1777–1784 (2013).19. Hammoud, S.S. *et al.* Distinctive chromatin in human sperm packages genes for embryo development. *Nature* **460**, 473–478 (2009).

20. Erkek, S. *et al.* Molecular determinants of nucleosome retention at CpG-rich sequences in mouse spermatozoa. *Nat. Struct. Mol. Biol.* **20**, 868–875 (2013).

21. Smith, C.M. *et al.* The mouse Gene Expression Database (GXD): 2014 update. *Nucleic Acids Res.* **42**, D818–D824 (2014).

22. Richardson, L. *et al.* EMAGE mouse embryo spatial gene expression database: 2014 update. *Nucleic Acids Res.* **42**, D835–D844 (2014).

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

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

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

26. Olson, E.N. Gene regulatory networks in the evolution and development of the heart. *Science* **313**, 1922–1927 (2006).

27. Arendt, D. & Nübler-Jung, K. Comparison of early nerve cord development in insects and vertebrates. *Development* **126**, 2309–2325 (1999).

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

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

30. Peter, I.S. & Davidson, E.H. Evolution of gene regulatory networks controlling body plan development. *Cell* **144**, 970–985 (2011).

31. Bult, C.J. *et al.* Mouse genome database 2016. *Nucleic Acids Res.* **44** D1, D840–D847 (2016).

32. Hou, Z.C. *et al.* Elephant transcriptome provides insights into the evolution of eutherian placentation. *Genome Biol. Evol.* **4**, 713–725 (2012).

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

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

35. Cáceres, M., Suwyn, C., Maddox, M., Thomas, J.W. & Preuss, T.M. Increased cortical expression of two synaptogenic thrombospondins in human brain evolution. *Cereb. Cortex* **17**, 2312–2321 (2007).

36. Olivieri, G. & Miescher, G.C. Immunohistochemical localization of EphA5 in the adult human central nervous system. *J. Histochem. Cytochem.* **47**, 855–861 (1999).

37. McGhee, S.A. & Chatila, T.A. *DOCK8* immune deficiency as a model for primary cytoskeletal dysfunction. *Dis. Markers* **29**, 151–156 (2010).

38. Soejima, M., Tachida, H., Ishida, T., Sano, A. & Koda, Y. Evidence for recent positive selection at the human *AIM1* locus in a European population. *Mol. Biol. Evol.* **23**, 179–188 (2006).

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

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

41. Hedges, S.B., Marin, J., Suleski, M., Paymer, M. & Kumar, S. Tree of life reveals clock-like speciation and diversification. *Mol. Biol. Evol.* **32**, 835–845 (2015).

42. Lynch, V.J. *et al.* Adaptive changes in the transcription factor HoxA-11 are essential for the evolution of pregnancy in mammals. *Proc. Natl. Acad. Sci. USA* **105**, 14928–14933 (2008).

43. Krol, A.J. *et al.* Evolutionary plasticity of segmentation clock networks. *Development* **138**, 2783–2792 (2011).

44. Zhang, X.M., Ramalho-Santos, M. & McMahon, A.P. *Smoothened* mutants reveal redundant roles for Shh and Ihh signaling including regulation of L/R asymmetry by the mouse node. *Cell* **105**, 781–792 (2001).

45. Barrow, K.M., Ward, C.M., Rutter, J., Ali, S. & Stern, P.L. Embryonic expression of murine 5T4 oncofoetal antigen is associated with morphogenetic events at implantation and in developing epithelia. *Dev. Dyn.* **233**, 1535–1545 (2005).

46. Sheng, G. & Foley, A.C. Diversification and conservation of the extraembryonic tissues in mediating nutrient uptake during amniote development. *Ann. NY Acad. Sci.* **1271**, 97–103 (2012).

47. Wu, S.F., Zhang, H. & Cairns, B.R. Genes for embryo development are packaged in blocks of multivalent chromatin in zebrafish sperm. *Genome Res.* **21**, 578–589 (2011).

48. Chen, C. *et al.* Human neuronal calcium sensor-1 shows the highest expression level in cerebral cortex. *Neurosci. Lett.* **319**, 67–70 (2002).

49. Brawand, D. *et al.* The evolution of gene expression levels in mammalian organs. *Nature* **478**, 343–348 (2011).

50. El-Sharnouby, S., Redhouse, J. & White, R.A. Genome-wide and cell-specific epigenetic analysis challenges the role of polycomb in *Drosophila* spermatogenesis. *PLoS Genet.* **9**, e1003842 (2013).

51. Arthur, R.K. *et al.* Evolution of H3K27me3-marked chromatin is linked to gene expression evolution and to patterns of gene duplication and diversification. *Genome Res.* **24**, 1115–1124 (2014).

52. Lesch, B.J. & Page, D.C. Poised chromatin in the mammalian germ line. *Development* **141**, 3619–3626 (2014).

53. Brykczynska, U. *et al.* Repressive and active histone methylation mark distinct promoters in human and mouse spermatozoa. *Nat. Struct. Mol. Biol.* **17**, 679–687 (2010).

54. Fortunato, S. *et al.* Genome-wide analysis of the sox family in the calcareous sponge *Sycon ciliatum*: multiple genes with unique expression patterns. *Evodevo* **3**, 14 (2012).

55. Fortunato, S.A. *et al.* Calcisponges have a ParaHox gene and dynamic expression of dispersed NK homeobox genes. *Nature* **514**, 620–623 (2014).

56. Saudemont, A. *et al.* Complementary striped expression patterns of NK homeobox genes during segment formation in the annelid *Platynereis*. *Dev. Biol.* **317**, 430–443 (2008).

57. Larroux, C. *et al.* Developmental expression of transcription factor genes in a demosponge: insights into the origin of metazoan multicellularity. *Evol. Dev.* **8**, 150–173 (2006).

58. Larroux, C. *et al.* Genesis and expansion of metazoan transcription factor gene classes. *Mol. Biol. Evol.* **25**, 980–996 (2008).

## ONLINE METHODS

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

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

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

**RNA isolation.** Unfixed aliquots of sorted cells were thawed on ice, washed once in cold PBS, resuspended in 350 μl of RLT Plus buffer from the RNeasy Mini kit (Qiagen, 74134), and then disrupted by drawing up and down five times through a 26-gauge insulin needle and syringe. Genomic DNA was removed using the genomic DNA eliminator columns supplied with the kit. The remainder of RNA isolation was performed using the RNeasy Mini kit according to the manufacturer's instructions. Samples were processed in batches of 2–6 in order of collection with no blinding. All biological replicates were processed in separate batches from each other.

**Chromatin immunoprecipitation.** For ChIP-seq analysis, between $5 \times 10^4$ and $5 \times 10^6$ cells were used as starting material, depending on the number obtained from sample isolation and sorting. Pachytene spermatocytes and round spermatids were treated identically. For human samples, fixed cells frozen in lysis buffer (1% SDS, 10 mM EDTA, and 50 mM Tris-HCl (pH 8)) were thawed on ice. For non-human samples, pellets of fixed cells were thawed on ice and then washed once in cold PBS and resuspended in 100 μl of lysis buffer. Once in lysis buffer, cells were incubated on ice for 5 min. Two hundred microliters of ChIP dilution buffer (0.01% SDS, 1.1% Triton X-100, 1.2 mM EDTA, 16.7 mM Tris-HCl (pH 8), and 167 mM NaCl) was added to each sample. Samples were sonicated in aliquots of 150 μl in 0.5-ml Eppendorf tubes at 4 °C using a Bioruptor (Diagenode) for 35 cycles on the high setting, with 30 s on and 30 s off. Aliquots of the same sample were again pooled and spun down at 12,000*g* for 5 min, and the chromatin-containing supernatant was moved to a fresh tube. Chromatin from each sample was split into two separate tubes

(150 μl in each), and 700 μl of dilution buffer, 50 μl of lysis buffer, and 100 μl proteinase inhibitor cocktail (Complete Mini tablets, Roche, 11836153001) were added to each tube. Fifty microliters of each sample was set aside as input. The remainder of the ChIP was performed as previously described[17], using 0.5 μg of antibody to H3K4me3 (Abcam, ab8580) or 1.0 μg of antibody to H3K27me3 (Abcam, ab6002). Samples were processed in batches of 2–4 in order of collection with no blinding; pachytene spermatocytes and round spermatids from a given sample were processed side by side, and all biological replicates were processed in separate batches from each other.

**Sequential chromatin immunoprecipitation.** Sequential ChIP was performed using the Re-ChIP-IT kit from Active Motif (53016) according to the manufacturer's instructions. $5 \times 10^7$ mouse round spermatids or $2 \times 10^7$ opossum round spermatids were used as starting material for each experiment. The first ChIP was performed with 2 μg of antibody to H3K27me3 (Abcam, ab6002). The second ChIP was performed with (i) 1.5 μg of antibody to H3K4me3 (Abcam, ab8580), (ii) 2.5 μl of antibody to H3K27me3 (Millipore, 07449) as a positive control, or (iii) 2 μl of nuclease-free water as a negative control.

**qPCR.** All primer sequences are listed in **Supplementary Table 10**. qPCR was performed on an Applied Biosystems 7500 Fast Real-Time PCR instrument using Applied Biosystems Power SYBR Green PCR Master Mix with the following cycling conditions: 50 °C for 20 s, 95 °C for 10 min, 95 °C for 15 s, 60 °C for 1 min, and 70 °C for 30 s (go to the third step 39 times). For targets with biological replicates (*Gdnf* and *Foxc1*), we observed variable ChIP efficiency between experiments performed on different days. To compare the relative percent input values for experimental and control conditions across biological replicates, all values from a given experiment were proportionately scaled such that the value of the positive control condition (H3K27me3 → H3K27me3) was 0.1% of input.

**Antibodies.** Antibody to H3K4me3 (rabbit polyclonal; Abcam, ab8580) was used for ChIP-seq in all species and for sequential ChIP in mouse and opossum. This antibody has been validated for ChIP-seq applications in human (Histone Modification Antibody Validation Database)[66], rhesus macaque[67], and mouse[68], as well as non-mammalian species including *Drosophila*[66]. Antibody to H3K27me3 (mouse monoclonal; Abcam, ab6002) was used for ChIP-seq in all species and for sequential ChIP in mouse and opossum. This antibody has been validated for ChIP-seq applications in human[69], mouse[70], and chicken[71], as well as non-mammalian species including *Drosophila*[66] and *Caenorhabditis elegans*[66]. Antibody to H3K27me3 (rabbit polyclonal; Millipore, 07449) was used for sequential ChIP in mouse and opossum. This antibody has been validated for ChIP applications in mouse[68], human[66], and *Drosophila*[72].

**Illumina library preparation and sequencing.** RNA libraries were prepared using an Apollo 324 library prep instrument with supplied reagents (Integenx) for non-human samples and using a SMARTer stranded RNA prep kit (Clontech) for human samples, according to the manufacturer's instructions. ChIP libraries were prepared using a TruSeq ChIP sample prep kit (Illumina), according to the manufacturer's instructions, except that size selection was performed after (instead of before) PCR amplification. Data from mouse replicate 1 have previously been published[17]; for this sample, ChIP-seq and RNA-seq libraries were sequenced on an Illumina Genome Analyzer IIx with 36-bp single-end reads. All other libraries were sequenced on an Illumina HiSeq 2500, with 40-bp single-end reads for ChIP libraries and 100-bp or 40-bp paired-end reads for RNA-seq libraries (**Supplementary Table 1**).

**Sequence alignment.** We filtered all data sets for read quality using FASTX-Toolkit and assessed library quality using FastQC. We aligned ChIP-seq libraries to a species-appropriate genome build (hg19, rheMac2, mm10, bosTau7, monDom5, or galGal4) using Bowtie v1.1.1 (ref. 73) (**Supplementary Table 1**). For ChIP-seq data, we called peaks at a threshold of $P < 1 \times 10^{-6}$ using MACS v1.4 (ref. 74); the number and locations of peaks were used to evaluate the quality of the data set but were not used in our analysis. For all data sets, peak numbers were within the expected range for the histone modification (H3K4me3 or H3K27me3); variation in peak numbers within this range did

not strongly affect poised gene calls. For RNA-seq data, we aligned libraries using TopHat v2.0.11 (ref. 75) with Ensembl[76] (release 75) transcripts as a reference (-G flag). The default genome assemblies included in Ensembl release 75 matched those used for alignment for all species except bull. For bull, Ensembl coordinates (for bosTau6) were mapped to the bosTau7 assembly using CrossMap[77].

Because of the small number of cells for many of the sorted cell populations both ChIP and RNA-seq libraries tended to have high duplication levels (**Supplementary Table 1**). We treated duplicate reads conservatively for both ChIP-seq and RNA-seq data. For ChIP-seq data, where a minimum count threshold was used to call poised genes, we retained only one duplicate for analysis. For RNA-seq data, where a maximum threshold was used, we retained a maximum of 20 duplicate reads (TopHat default).

**Poised gene calls.** For ChIP data, we counted total reads in the 4-kb interval surrounding each transcriptional start site (Ensembl build 75) using htseq-count[78] with the intersection-nonempty option. The 4-kb interval (2 kb upstream and 2 kb downstream of the transcriptional start site) is standard for analysis of promoter-associated histone modifications. For each data set, total ChIP or input reads in each interval were normalized to reads per million, and the normalized input count was subtracted from the normalized ChIP count in each interval to obtain a final ChIP signal. For RNA, we obtained FPKM values using Cufflinks v2.2.1 (ref. 79), with Ensembl release 75 transcripts as a reference (-G option). Transcript values for both ChIP and RNA-seq data were summed to obtain a single H3K4me3, H3K27me3, and expression value for each gene (**Supplementary Data**). We conducted simulations in which we varied ChIP and expression thresholds and evaluated the numbers of poised genes called in each species, and we selected thresholds that included the maximum number of poised genes while remaining robust to small changes in threshold value (**Supplementary Figs. 5** and **6**, and **Supplementary Data**). We set thresholds of ≥0.5 input-subtracted reads per million for H3K4me3 signal, ≥0.5 input-subtracted reads per million for H3K27me3 signal, and ≤5 FPKM for expression. For each sample, a gene had to meet the thresholds for H3K4me3, H3K27me3, and expression in both pachytene spermatocytes and round spermatids (six data points in total) to be considered stably poised. For species with two biological replicates (rhesus macaque, mouse, and opossum), we used mean ChIP signal and FPKM values to call poised genes for that species; this approach yielded similar gene lists to either the union or intersection of the two replicates but was more robust to changes in threshold. For species with three biological replicates (human), we included genes called as poised in at least two of the three individual replicates. We note that these criteria are expected to result in greater sensitivity for poised gene calls in species with more replicates, as use of more replicates allows inclusion of genes that may fail to meet one of the six thresholds in a single replicate. We did observe the fewest poised gene calls in bull (one replicate) and the most in human (three replicates). These differences in sensitivity may result in a subset of false positives and negatives in lists of genes called as differentially poised between species. The list of conserved H3K27me3-only genes (**Fig. 3c**, **Supplementary Fig. 8b**, and **Supplementary Table 5**) was defined as the set of genes that met the H3K27me3 and expression thresholds in all five mammalian species but met the H3K4me3 threshold in fewer than four of the five mammalian species.

**Orthologous gene sets.** We required that a gene have orthologs in all five mammalian species to be included in our analysis. Because there is no strong a priori expectation that gene duplication would have a specific effect (loss, gain, or retention) on the chromatin state surrounding the transcriptional start site, we reasoned that exclusion of genes with one-to-many relationships could result in loss of biologically meaningful information and might introduce bias by excluding a non-random set of genes from the analysis. We therefore included genes with either one-to-one or one-to-many orthology relationships among the five species. We used the BioMart database[80] with Ensembl release 75 to find one-to-one and one-to-many orthologs. In total, we identified 14,362 orthology groups containing at least one gene from each species. 12,104 of these involved only one-to-one orthology relationships across all five species; we used this number to calculate *P* values for the significance of five-way overlaps. The 14,362 orthology groups included a total

of 15,492 human, 15,904 rhesus macaque, 16,253 mouse, 15,650 bovine, and 15,966 opossum genes, which together comprised the total gene set considered in our downstream analysis. When determining sets of overlapping poised genes, an orthology group was counted as overlapping if at least one gene belonging to the group was poised in each species (**Supplementary Data** and **Supplementary Code**).

**Statistics.** *Sample inclusion criteria.* Testis samples were excluded from the study if any morphological abnormality was observed in the intact tissue. For any ChIP-seq or RNA-seq data set with $<1 \times 10^6$ unique (non-duplicate) reads aligning uniquely to the genome, the associated biological sample and all data sets derived from it were excluded from analysis. These criteria were established before beginning the study. When possible, at least two biological replicates were obtained to allow for individual variation. For human data, three biological replicates were used in the final analysis to account for greater variability in genetic background in comparison to non-human species.

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

*Gene set overlaps.* Overlaps between multiple (>2) gene sets were computed using the overLapper function from the systemPipeR package in R. The statistical significance of five-way overlap was derived using the formula

$$P < \binom{N}{m} \left[ \frac{\binom{N-m}{n-m}}{\binom{N}{n}} \right]^5$$

where *N* is the total number of genes with orthologs in all five species, *n* is the largest number of poised genes called in any single species (within the set of genes with orthologs in all five species), and *m* is the number of genes called as poised in all five species. Using our gene set, $N = 14{,}362$, $n = 3{,}580$, $m = 405$, and $P < 1 \times 10^{-300}$ in this calculation. However, we note that core poised genes are enriched for genes with true one-to-one orthologs, meaning that the groups being compared are not completely independent as the formula assumes. To account for this bias, we recalculated the overlap between gene sets, including only one-to-one orthologs in the analysis. With only one-to-one orthologs, $N = 12{,}104$, $n = 3{,}361$, $m = 401$, and $P < 1 \times 10^{-280}$. We report this *P* value as a conservative estimate of significance for five-way overlap.

**Gene ontology enrichment.** GO enrichment was evaluated using the GOStats package[81] in R. *P* values were adjusted both by conditioning out child categories and by subsequent correction for multiple testing using the Benjamini–Hochberg method.

**Clustering and divergence estimates.** We generated a distance matrix for expression data sets based on FPKM values and for ChIP data sets based on normalized counts around promoter regions (**Supplementary Data**), using $1 - \rho$ (Spearman's correlation) as a dissimilarity metric[49]. Clustering was performed and dendrograms were generated using the cluster package in R. Analysis of gene expression and promoter chromatin divergence was carried out using dissimilarity scores for each species pair, and data from each cell type were fit to a linear model (**Supplementary Code**).

**Principal-component analysis.** We used the PCA function from the FactoMineR package in R for principal-component analysis, with data scaled to unit value. Input data were the same processed data (normalized H3K4me3 signal, normalized H3K27me3 signal, or FPKM values) as were used for calling poised genes (**Supplementary Data**).

**Somatic tissue expression.** The MGI Gene Expression Database[21] was used to determine stages of gene expression for mouse poised genes. Only wild-type samples from the database were used in the analysis. At the time of our study,

the database included a total of 13,837 genes; 9,884 genes in the database had orthologs in all five mammalian species. The numbers of genes in the database with orthologs in all five species (359 core poised genes and 9,525 other genes) were used as denominators in calculating the fraction of genes expressed at each stage.

**Embryonic lethality.** We identified alleles associated with embryonic lethality by searching the MGI database using the Phenotypes, Alleles and Disease Models query[31] for "embryonic lethality" (phenotype ID MP:0008762) and filtering for null/knockout alleles. This search identified 2,812 alleles, corresponding to 1,881 genes; 1,570 of these genes had orthologs in all five mammalian species.

**Transcription factor class enrichment.** Genes were assigned to transcription factor classes (**Fig. 3b** and **Supplementary Fig. 8c**) according to Wingender et al.[82].

**Motif analysis.** We identified motifs enriched in the promoters of species-specific poised genes in two steps. In the first step, we used DREME[83] with default settings to detect motifs enriched in each set of differentially poised promoters (±1 kb from the transcriptional start site) in comparison to orthologous promoter regions from the other four species, with a threshold of $E < 0.05$. In the second step, to control for biases introduced by comparing different species, we used AME[84] to scan 100 random, equally sized sets of orthologous promoters for enrichment of the motifs detected in the first step in the species in which they were first detected in comparison to the other four species. The fraction of random promoter sets demonstrating enrichment at $E < 0.05$ constituted a raw $P$ value; these values were adjusted for multiple comparisons using the Benjamini–Hochberg approach to obtain a false discovery rate (FDR) for each enriched motif. We considered only motifs with FDR <0.10 in our subsequent analysis. To match enriched motifs with binding sites for known transcription factors, we used Tomtom[85] (with conditions -thresh 10 -evalue -dist ed) and pulled motifs from the JASPAR Core vertebrates[86] (205 total motifs) and Uniprobe mouse[87–89] (386 total motifs) databases.

***Drosophila* Polycomb data.** *Drosophila* ChIP-chip data were taken from El-Sharnouby et al.[50]. Using tiled enrichment values generated by the authors, we calculated average Polycomb enrichment within the regions 1 kb upstream or downstream of each *Drosophila* transcriptional start site and assigned these values to the associated gene. We designated genes with the top 25% of Polycomb signal as 'high Polycomb' and genes with the bottom 25% of Polycomb signal as 'low Polycomb'.

**Code availability.** Custom R scripts used in our analyses are included as **Supplementary Code**.

59. Bellvé, A.R. Purification, culture, and fractionation of spermatogenic cells. *Methods 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. USA* **65**, 192–199 (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**, 553–562 (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**, 332–340 (1988).
66. Egelhofer, T.A. et al. An assessment of histone-modification antibody quality. *Nat. Struct. Mol. Biol.* **18**, 91–93 (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**, 1408–1418 (2011).
68. Goldberg, A.D. et al. Distinct factors control histone variant H3.3 localization at specific genomic regions. *Cell* **140**, 678–691 (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**, 249–257 (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**, 1105–1111 (2009).
76. Yates, A. et al. Ensembl 2016. *Nucleic Acids Res.* **44** D1, D710–D716 (2016).
77. Zhao, H. et al. CrossMap: a versatile tool for coordinate conversion between genome assemblies. *Bioinformatics* **30**, 1006–1007 (2014).
78. Anders, S., Pyl, P.T. & Huber, W. HTSeq—a Python framework to work with high-throughput sequencing data. *Bioinformatics* **31**, 166–169 (2015).
79. Trapnell, C. et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat. Biotechnol.* **28**, 511–515 (2010).
80. Durinck, S. et al. BioMart and Bioconductor: a powerful link between biological databases and microarray data analysis. *Bioinformatics* **21**, 3439–3440 (2005).
81. Falcon, S. & Gentleman, R. Using GOstats to test gene lists for GO term association. *Bioinformatics* **23**, 257–258 (2007).
82. Wingender, E., Schoeps, T. & Dönitz, J. TFClass: an expandable hierarchical classification of human transcription factors. *Nucleic Acids Res.* **41**, D165–D170 (2013).
83. Bailey, T.L. DREME: motif discovery in transcription factor ChIP-seq data. *Bioinformatics* **27**, 1653–1659 (2011).
84. McLeay, R.C. & Bailey, T.L. Motif Enrichment Analysis: a unified framework and an evaluation on ChIP data. *BMC Bioinformatics* **11**, 165 (2010).
85. Gupta, S., Stamatoyannopoulos, J.A., Bailey, T.L. & Noble, W.S. Quantifying similarity between motifs. *Genome Biol.* **8**, R24 (2007).
86. 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** D1, D110–D115 (2016).
87. Badis, G. et al. Diversity and complexity in DNA recognition by transcription factors. *Science* **324**, 1720–1723 (2009).
88. Berger, M.F. et al. Variation in homeodomain DNA binding revealed by high-resolution analysis of sequence preferences. *Cell* **133**, 1266–1276 (2008).
89. 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–D122 (2015).