

1 **Annotation of Chromatin States in 66 Complete Mouse Epigenomes During Development**

2

3

4 Arjan van der Velde^{1,2}, Kaili Fan¹, Junko Tsuji¹, Jill Moore¹, Michael Purcaro¹, Henry Pratt¹, Zhiping
5 Weng¹

6

7

8 1 Program in Bioinformatics and Integrative Biology, University of Massachusetts Medical School,
9 Worcester, MA, USA

10 2 Bioinformatics Program, Boston University, Boston, MA 02215, USA

11 To whom correspondence should be addressed. Tel: +1-508-856-8866; Fax: +1-508-856-2392; Email:
12 zhiping.weng@umassmed.edu

13

14 **ABSTRACT**

15 The morphologically and functionally distinct cell types of a multicellular organism are maintained by
16 epigenomes and gene expression programs. Phase III of the ENCODE Project profiled 66 mouse
17 epigenomes across twelve tissues at daily intervals from embryonic day 10.5 to birth. Applying the
18 ChromHMM algorithm to these epigenomes, we annotated eighteen chromatin states with
19 characteristics of promoters, enhancers, transcribed regions, repressed regions, and quiescent
20 regions throughout the developmental time course. Our integrative analyses delineate the tissue
21 specificity and developmental trajectory of the loci in these chromatin states. Approximately 0.3% of
22 each epigenome is assigned to a bivalent chromatin state, which harbors both active marks and the
23 repressive mark H3K27me3. Highly evolutionarily conserved, these loci are enriched in silencers
24 bound by Polycomb Repressive Complex proteins and the transcription start sites of their silenced
25 target genes. This collection of chromatin state assignments provides a useful resource for studying
26 mammalian development.

27

28 **INTRODUCTION**

29 Multicellular organisms maintain a myriad of cell types along separate lineages to carry out the
30 cellular programs required for development and survival. These cell types all have the same genome
31 but different epigenomes, characterized by chromatin accessibility, histone modifications, and DNA
32 methylation, which cooperate with trans-factors to regulate gene expression and downstream
33 activities. Thus, systematic annotation of epigenomes is essential for understanding genome
34 functions. Experimental techniques such as chromatin immunoprecipitation followed by sequencing

35 (ChIP-seq)¹⁻³, transposase accessible chromatin with sequencing (ATAC-seq)⁴, and whole-genome
36 bisulfite sequencing (WGBS)⁵ enable genome-wide profiling of histone marks, chromatin accessibility,
37 and DNA methylation, respectively. When several of these epigenetic marks have been profiled for a
38 given cell type, the results can be integrated using computational algorithms such as ChromHMM⁶,
39 Segway⁷, and IDEAS⁸ to classify genomic loci into a small number of chromatin states, such that the
40 chromatin state of a locus is predictive of its function in the given cell type.

41

42 Coordinated efforts by the ENCODE and Roadmap Epigenomics Consortia provided tremendous
43 insights into gene regulation in a diverse array of human cell and tissue types^{9,10}. The
44 mouseENCODE project furthered our understanding of multiple adult mouse cell types¹¹. Phase III of
45 the ENCODE Consortium generated 66 complete mouse epigenomes across 12 fetal tissues at four
46 to seven developmental time-points with a daily interval, each investigated with ten assays¹²: ATAC-
47 seq¹³, WGBS¹⁴, and ChIP-seq of eight histone marks¹³. The histone marks included histone 3 lysine 4
48 trimethylation (H3K4me3) and histone 3 lysine 9 acetylation (H3K9ac), enriched at promoters and
49 present at enhancers^{1,15-17}; H3K27ac, H3K4me1, and H3K4me2, enriched at enhancers^{1,15,17,18};
50 H3K36me3, enriched within the bodies of actively transcribed genes¹⁹; H3K27me3, catalyzed by and
51 guiding the Polycomb Repressive Complexes (PRC) of proteins to repress gene expression²⁰; and
52 H3K9me3, enriched in heterochromatin to silence repeats and gene clusters¹⁹. All these 66
53 epigenomes were accompanied by transcriptome sequencing (RNA-seq) data²¹, and 20 of the
54 biosamples were assayed by DNase-seq, another technique for measuring chromatin accessibility²²
55 (**Fig. 1a** and **Supplementary Table 1**). This body of data was generated by four ENCODE labs, with
56 the same type of data generated by the same lab, representing the most complete epigenetic data on
57 fetal mouse tissues, ideal for characterizing the epigenomic landscape during mammalian
58 development.

59

60 We applied ChromHMM⁶ to these 66 mouse epigenomes and defined 18 chromatin states (**Fig. 1b**).
61 Most of these mouse chromatin states recapitulated the 15 human chromatin states defined by the
62 Roadmap Epigenomics Consortium using a subset of five histone marks in human biosamples¹⁰, and
63 our novel states corresponded to a refinement of previously defined enhancer, bivalent, and quiescent
64 states. We observed a substantially larger variation of chromatin state assignments among the mouse
65 tissue types at a given developmental time-point than we did across all developmental time-points for
66 a single tissue. We further investigated one chromatin state in detail—TssBiv, a bivalent state
67 enriched in the transcription start sites (TSS) which harbors both active marks (H3K4me3, H3K4me2,
68 H3K4me1, and H3K9ac) and the repressive mark H3K27me3. We found that genomic loci in TssBiv
69 were substantially more evolutionarily conserved than loci in any of the other 17 Chromatin states.
70 Genes with bivalent TSSs were first identified in embryonic stem cells and thought to be poised for

71 activation or repression in response to developmental or environmental cues²³. Subsequently, such
72 bivalent domains were reported in differentiated cell types^{24–27}, but they have not been studied during
73 fetal development. Each fetal tissue harbors approximately 3000 bivalent genes and they are
74 repressed in expression in that specific tissue. These bivalent genes are highly enriched in
75 transcription factors (TFs) differentially expressed among the fetal tissues. Comparison with recently
76 defined silencers bound by the Polycomb Repressive Complex 2 (PRC2) proteins²⁸ revealed that both
77 the PRC2-bound silencers and the TSSs of their silenced genes are highly enriched in the bivalent
78 regions. Thus, the bivalent regions support an evolutionarily conserved silencing mechanism for
79 lineage-specific genes, in particular the master TFs controlling tissue development. Our
80 comprehensive annotation of chromatin states provides a resource for studying mammalian
81 development.

82

83

84 RESULTS

85 Chromatin states were defined using ATAC-seq, WGBS, and the ChIP-seq data of eight histone 86 marks

87 The 66 mouse fetal epigenomes, all complete with ten chromatin marks, represent a comprehensive
88 collection for chromatin state assignment (**Fig. 1a**). We used ChromHMM to learn 18 states jointly
89 from this dataset (**Fig. 1b, c**). ChromHMM chunks the genome into non-overlapping 200 base-pair
90 (bp) bins and assigns each of these genomic bins to one of the 18 chromatin states in each
91 biosample. We named our chromatin states in a way to be consistent with earlier ChromHMM
92 publications^{6,10,29}. Two of our learned states are proximal to active TSSs (Tss and TssFlnk,
93 approximately 1.5% of the mouse genome); two states associate with actively transcribed genes (Tx
94 and TxWk, 8.5%); five states are enhancer-related (Enh, EnhLo, EnhPois, EnhPr, and EnhG; 4.5%);
95 one bivalent state often falls near inactive TSSs (TssBiv, 0.3%); three states are repressive (ReprPC
96 and ReprPCWk enriched in H3K27me3, 5.5%; and Het in H3K9me3, 2.5%); and five states are
97 quiescent (QuiesG, Quies, Quies2, Quies3, and Quies4; 75%). The remaining 2% or so of the
98 genome could not be confidently assigned to any one state.

99

100 The assignments of the 18 states are supported by comparison with gene expression and epigenomic
101 data available for a subset of biosamples (**Supplementary Table 1**). Although both the active-TSS
102 (Tss) and the bivalent-TSS states (TssBiv) are highly enriched in CpG islands, Tss (along with the
103 TSS-flanking state TssFlnk) is only enriched in the TSSs of expressed genes (determined using RNA-
104 seq data in the corresponding biosample) while TssBiv is only enriched in the TSSs of the repressed
105 genes (**Fig. 1c**). The transcription-related states (Tx and TxWk) are enriched in the exons and introns
106 of expressed genes but not those of the repressed genes (**Fig. 1c**). Enh (high-signal enhancer) is the

107 state most enriched in the ChIP-signal of EP300, a histone acetyltransferase that preferentially binds
108 active enhancers^{30,31} (**Fig. 1c**). The relative enrichment of the 18 states for the ATAC signal are highly
109 consistent with their enrichments in DNase hypersensitive sites (DHS), determined using DNase-seq
110 data in the corresponding biosample (**Fig. 1c**).

111

112 **Contributions of the chromatin marks to the assignments of chromatin states**

113 To assess the contribution made by each of the eight histone marks, ATAC, and DNA methylation, we
114 asked how accurately the ten-mark model would be able to annotate a new epigenome missing data
115 for one of the marks. We addressed this question by removing the data for each mark individually
116 from the midbrain E13.5 epigenome and computing the Jaccard similarity index between the
117 chromatin state assignments of all genomic bins (each 200 bp long, which is the resolution of
118 ChromHMM) with the data for the remaining nine marks. If a genomic bin has a posterior probability
119 less than 0.5, then it is classified as unassigned. In general, when a mark is removed, the states most
120 severely affected were among those states most enriched in this mark in the ten-mark model
121 (compare **Fig. 1d** and the chromatin-mark probabilities in **1c**). However, the converse is not
122 necessarily true, reflecting the redundancy between the marks. For example, the removal of H3K27ac
123 affects the low-signal enhancer state (EnhLo) although the high-signal enhancer (Enh) state is even
124 more enriched in H3K27ac than EnhLo (**Fig. 1c-d**). H3K4me3 and H3K9ac, when removed
125 individually, did not have a major impact on any of the states although promoter states are enriched in
126 H3K4me3 and both promoter and enhancer states are enriched in H3K9ac (**Fig. 1c-d**), indicating that
127 the information contained by each of these two marks is already accounted for by the other nine
128 marks. On the other hand, H3K36me3, H3K27me3, and H3K9me3 each brings non-redundant
129 information to the ten-mark model, as all the states enriched in each of these marks were affected
130 when the mark was removed (**Fig. 1c-d**).

131

132 **Chromatin states are conserved between human and mouse**

133 The Roadmap Epigenomics Consortium previously defined 15 human chromatin states using five
134 histone marks in 127 human biosamples¹⁰. To investigate the conservation of chromatin state types
135 between human and mouse, we built a 15 state model using the same set of five histone marks in the
136 66 mouse fetal biosamples. This five-mark 15-state mouse model recapitulated 13 of the 15 human
137 states identified by the Roadmap Epigenomics Consortium, with nearly identical emission probabilities
138 and similar genome coverages. The 13 reproduced states, including the promoter, enhancer,
139 transcribed, repressed, and bivalent states, were enriched in at least one of the five histone marks.

140

141 The remaining two mouse chromatin states had similar chromatin-mark probabilities to, but different
142 genome coverages from, the two remaining human states¹⁰ (**Supplementary Fig. 1a**). These human

143 states—the weak transcription state TxWk and the weak repressed polycomb state ReprPCWk
144 (11.6% and 8.3% of the human genome)—had low signals for all five marks, and their assignments
145 were based on their weak enrichments in expressed and repressed genes, respectively¹⁰. We
146 identified a similar state with low signals for all marks in mouse, but although it was enriched within
147 gene bodies in general, it was not enriched in either expressed or repressed genes in particular. We
148 thus denoted it as the quiescent gene state (QuiesG, 25.17% of the mouse genome). We also
149 identified a minor state (0.13% of the mouse genome) marked by both H3K36me3 and H3K27me3;
150 we denote this state TxWk because regions in this state were assigned to the transcription state (Tx),
151 the repressed polycomb state (ReprPC), or the weak repressed polycomb state (ReprPCWk) in our
152 complete ten-mark, 18-state model. In summary, our results indicate that the chromatin states are
153 highly conserved between human and mouse, and ChromHMM is able to identify these states reliably.

154

155 **Addition of three more histone marks, chromatin accessibility, and DNA methylation further** 156 **clarified enhancer, bivalent, and quiescent states**

157 To investigate the impact of incorporating additional data in the annotation of chromatin states, we
158 constructed a 15-state model using all eight available histone marks (**Supplementary Fig. 1a**). We
159 compared this model, and the five-mark 15-state model described above, with our 18-state model that
160 further incorporated chromatin accessibility and DNA methylation data (ten marks in total, **Fig. 1**, also
161 included in **Supplementary Fig. 1a** to facilitate comparison with the five-mark and eight-mark
162 models). Comparison of the three ChromHMM models built with increasing numbers of epigenetic
163 marks (five, eight, and ten marks) revealed that assignments differ predominantly for the enhancer,
164 bivalent, and quiescent states (**Supplementary Fig. 1b-e**).

165

166 The five-mark model specified one enhancer state (Enh; 3.7% of the mouse genome) with high
167 H3K4me1 levels (**Supplementary Fig. 1a**). Genomic regions in this state were assigned to five
168 distinct enhancer states in the eight-mark model, which reflected different levels of three additional
169 enhancer marks (H3K4me2, H3K9ac, and H3K27ac). Among these five states in the eight-mark
170 model, the high-signal enhancer state Enh, which showed high levels for all these four enhancer
171 marks, occupied only 0.2% of the genome (**Supplementary Fig. 1a, d**). The high-signal enhancer
172 state Enh defined by the ten-mark model further showed high chromatin accessibility (ATAC signal)
173 and low DNA methylation, occupying 0.64% of the genome (**Supplementary Fig. 1a, d**). The ten-
174 mark model defined three additional enhancer states, with two of the three (EnhLo and EnhPois)
175 being regroupings of the genomic regions assigned to the four enhancer states in the eight-mark
176 model. The other enhancer state defined by the ten-mark model (EnhPr) corresponded to a subset of
177 the regions assigned one of the enhancer states by the eight-mark model, showing high chromatin

178 accessibility but low levels of enhancer marks (**Supplementary Fig. 1d**). Thus, the additional marks
179 led to refined definitions of enhancer states.

180

181 One example of a tissue-specific enhancer is located inside the housekeeping gene *Metap1d*
182 (methionyl aminopeptidase Type 1D, which functions in the mitochondria) and 10 kb upstream of the
183 *Dlx1* gene, which encodes a brain-specific homeobox transcription factor. *Dlx1* is highly expressed in
184 the forebrain (~200 transcripts per million or TPM), but not expressed in most other tissues (e.g., < 3
185 TPM in the liver). This region is annotated as a high-signal enhancer (Enh) in the forebrain, showing
186 high ATAC and H3K27ac signals and low DNA methylation. It is annotated as a quiescent gene
187 (QuiesG) in the liver due to its low ATAC and histone mark signals and high DNA methylation (**Fig.**
188 **1e**). A VISTA enhancer (accession: hs553) overlaps this region, and it is active in the forebrain and
189 cranial nerve of mouse embryos³².

190

191 The five-mark model annotated three bivalent states with high levels of the active marks H3K4me1
192 and H3K4me3, as well as high levels of the repressive mark H3K27me3; however, both the eight-
193 mark and ten-mark models only annotated one bivalent state, which additionally showed high levels of
194 other active marks (H3K4me2, H3K9ac, and ATAC) and low levels of DNA methylation
195 (**Supplementary Fig. 1a**). Roughly the same set of genomic regions were assigned to these bivalent
196 states across the three models, suggesting that the state definition became more complete when
197 more marks were available (**Supplementary Fig. 1e**).

198

199 The five-mark and eight-mark models annotated one quiescent state (Quies), which had very low
200 signals for all available histone marks. The ten-mark model defined three additional quiescent states
201 besides Quies. These four quiescent states all showed very low levels of the eight histone marks and
202 ATAC, but they differed in DNA methylation, with the Quies state (49.0% of the genome) showing very
203 high methylation levels and the Quies2 state (9.9% of the genome) showing very low levels of DNA
204 methylation (**Supplementary Fig. 1a**). The quiescent states in the three models cover roughly the
205 same set of genomic regions (**Supplementary Fig. 1b**).

206

207 **Variation of state assignments across tissues and along developmental time-points**

208 After carefully analyzing the properties of the chromatin states in the ten-mark model, we assessed
209 how variable the assignments of these states were among the 66 mouse epigenomes. We computed
210 the Jaccard similarity index on the genomic regions assigned to each state between tissues or
211 between developmental time-points. The enhancer states exhibited the greatest variability among
212 tissues or across time-points, while the promoter, quiescent, and transcription states showed the least
213 variability (**Fig. 2a**). The repressive state Het, enriched in H3K9me3, was almost as variable as the

214 enhancer states (**Fig. 2a**). Moreover, all chromatin states were more similar across time-points in the
215 same tissue than across tissues at the same time-point (**Fig. 2a**), consistent with the notion that the
216 epigenome is inherited within the cell lineage.

217

218 Temporal chromatin state transitions for each tissue occurred mostly between related states, e.g.,
219 among the promoter states (Tss, TssFlnk, and TssBiv) or among the enhancer states (Enh, EnhLo,
220 EnhPois, and EnhPr). We also observed a preference for temporal transitions into or out of the
221 quiescent states (**Fig. 2b, c**).

222

223 To investigate whether the variations captured by the chromatin states could recapitulate the
224 embryonic developmental trajectory, we applied the UMAP dimension-reduction technique³³ to the 66
225 tissue biosamples using levels of chromatin marks at the genomic bins assigned to each chromatin
226 state. H3K27ac signal levels at genomic bins assigned high-signal enhancers (Enh) in any of the 66
227 biosamples (in total 5.4% of the genome) cleanly segregated the 66 biosamples by tissue (**Fig. 2d,**
228 **left panel**). The liver (with an endoderm origin) and heart (mesoderm) biosamples formed two
229 separate clusters. Tissues with similar developmental origins were positioned near each other, with
230 the four brain regions (ectoderm), the lung (endoderm) and the digestive organs stomach and
231 intestine (endoderm), and limb and facial prominence (with cells from both endoderm and ectoderm
232 origins) forming three clusters (**Fig. 2d**). The kidney (mesoderm) biosamples were positioned right
233 next to the stomach, intestine, and lung (endoderm) biosamples. Furthermore, the earlier time-points
234 (open symbols) are segregated from later time-points (filled symbols). A similar UMAP analysis on
235 genomic bins assigned to the bivalent state (TssBiv) in any of the biosamples (in total 1.2% of the
236 genome) by the levels of the ten chromatin marks also led to clear segregation of the biosamples by
237 tissue, although there was some mixing between the lung biosamples and the stomach and intestine
238 biosamples (**Fig. 2d, right panel**). Thus, the epigenomic landscapes captured by chromatin states
239 Enh and TssBiv can accurately recapitulate the tissue lineages during embryonic development.

240

241 **Genome regions transit among TssBiv, Tss, and ReprPC states**

242 Over developmental time, regions assigned to the bivalent promoter state (TssBiv), which has both
243 active marks and the repressive H3K27me3 mark (**Fig. 1c**), can either lose repressive H3K27me3
244 and become active TSSs (Tss) or lose the active marks and transition into the repressive polycomb
245 (ReprPC) state (**Supplementary Fig. 2**). Roughly 0.3% of any particular epigenome is assigned to
246 the TssBiv state; cumulatively 1.2% of the genome is assigned to TssBiv across all tissues and time-
247 points. TssBiv is less than half as prevalent as Tss and ReprPC, which constitute 0.8% and 0.8% of
248 each epigenome and 2.2% and 5.5% of the genome overall, respectively. Almost all stretches of
249 TssBiv genomic bins are flanked by ReprPC genomic bins. As an example, the promoter of the *Dlx1*

250 gene is annotated as Tss in the forebrain, where it is highly expressed, and as TssBiv in the liver,
251 where it is not expressed, and the bivalent promoter in the liver is surrounded by ReprPC regions
252 (**Fig. 1e**). Among the genomic bins that are assigned TssBiv in any of the epigenomes, 64.7% are
253 assigned ReprPC in at least one epigenome and 68.1% are assigned Tss in at least one epigenome
254 (**Supplementary Fig. 3a**), indicating that a particular region is TssBiv in some tissue but becomes
255 monovalent (Tss or ReprPC) in other tissues. Intriguingly, the overall fraction of TssBiv genomic bins
256 decreased over the course of the development in all five tissues with seven time-points, although due
257 to the small number of time-points this was statistically significant only in the three brain tissues
258 (**Supplementary Fig. 3b**). This suggests that the resolution of TssBiv regions into a monovalent state
259 is important for development, especially in the brain.

260

261 **Bivalent genes are involved in fundamental biological processes**

262 We identified 14,558 bivalent regions, defined as stretches of TssBiv genomic bins surrounded by
263 repressive chromatin states in any of the 66 biosamples (see **Methods**). These bivalent regions
264 overlapped 14,729 GENCODE-annotated TSSs (**Supplementary Table 2**), belonging to 6,797 genes
265 (**Supplementary Table 3**). There were 1,077 genes that were bivalent in all 12 tissues (i.e., having at
266 least one bivalent TSS at one or more time-points of every tissue), and these genes were highly
267 enriched in Gene Ontology (GO) terms related to embryonic development of myriad organs and
268 systems, regulation of fundamental cellular processes, and modulation of cell-cell communications
269 (**Supplementary Fig. 4a** and **Supplementary Table 4a, b**).

270

271 The liver had 5,482 bivalent genes (i.e., having at least one bivalent TSS at one or more time-points),
272 74% more than the other 11 tissues on average, and 1,291 of these 5,482 genes were not bivalent in
273 the other 11 tissues. GO analysis on the 1,291 liver-only bivalent genes revealed terms that were
274 involved in the development of a wide variety of organs other than the liver, such as heart, kidney,
275 smooth muscle, brain, and cytoskeleton (**Supplementary Fig. 4b** and **Supplementary Table 4c, d**).
276 We observed similar results for bivalent genes specific to other tissues. Thus, the bivalent genes in
277 each fetal tissue reflect the regulatory pathways that are unused by the developmental program of
278 that specific tissue.

279

280 **Bivalent genes exhibit repressed transcription**

281 We further analyzed the expression of the 25,215 genes that were expressed (≥ 1 TPM) in at least
282 one of the 66 biosamples, among which 6,324 were among our list of bivalent genes (**Methods**). We
283 found that the bivalent genes in a tissue had lower expression levels than non-bivalent genes

284 according to RNA-seq data in the same tissue. Across the 66 biosamples, the expression levels of
285 bivalent genes were 5.2 ± 1.7 TPM, much lower than the expression levels of non-bivalent genes
286 (39.8 ± 2.1 TPM; Wilcoxon rank-sum test P-value $< 2.2 \times 10^{-16}$). Furthermore, the genes that were not
287 bivalent in any of the time-points of a tissue were expressed 7.79-fold higher (Wilcoxon rank-sum test
288 P-values $\leq 2.2 \times 10^{-16}$) than the genes that were bivalent at all time-points of the tissue
289 (**Supplementary Fig. 5**). In a particular tissue, genes that were bivalent at different time-points were
290 largely consistent (forebrain in **Fig. 3a**; all tissues in **Supplementary Fig. 6**). For example, 1,830
291 genes were bivalent at all seven time-points of the liver; only 439 such genes would be expected if the
292 time-points were independent of one another (P-value $< 2.2 \times 10^{-16}$; Binomial test). Genes bivalent at
293 the earliest time-point but not the latest time-point were expressed at significantly lower levels earlier
294 in development; likewise, genes bivalent at the latest time-point but not at the earliest time-point were
295 expressed at lower levels later in development (midbrain in **Fig. 3b**; all tissues in **Supplementary Fig.**
296 **7**). Both of these two sets of genes were expressed at significantly higher levels than genes bivalent
297 at all time-points in the same tissue (**Fig. 3b**, **Supplementary Fig. 7**). Overall, the average expression
298 level of a TSS across the time-points in a tissue is anti-correlated with the number of time-points at
299 which the TSS is in a genomic bin assigned to the TssBiv chromatin state; in sharp contrast, a
300 positive correlation is observed between expression and the duration the TSS is in a genomic bin
301 assigned to the Tss chromatin state (**Fig. 3c**; **Supplementary Fig. 8**). Thus, the expression of
302 bivalent genes is repressed in a tissue- and time-point-specific manner.

303

304 **Bivalent genes are highly enriched in tissue-specific transcription factors**

305 We compared the 6,797 bivalent genes (6,324 expressed in at least one of the 66 biosamples) with a
306 curated list of 552 TFs with known DNA binding motifs in both mouse and human³⁴, of which 535 were
307 expressed in at least one of the 66 biosamples. A majority of the 535 TFs (338, 63.2%) were among
308 the 6,324 bivalent genes (Chi-square P-value $< 2.2 \times 10^{-16}$). For both TF and non-TF genes, those that
309 were bivalent were significantly more tissue-specific than those that were not bivalent (2.47-fold and
310 1.79-fold higher in median tissue specificity for TFs and non-TFs, respectively, Wilcoxon rank-sum
311 test P-values $< 2.2 \times 10^{-16}$; **Fig. 3d**).

312

313 Consistent with earlier findings in embryonic stem cells^{35,36}, a majority of the bivalent TSSs in our
314 mouse fetal biosamples (mean = 62.5% across the 66 biosamples) overlapped CpG islands, much
315 higher than non-bivalent TSSs (mean = 29.8%; Chi-square P-values in all 66 biosamples $< 2.2 \times 10^{-16}$).
316 The enrichment is highly significant for the TSSs of both the TF genes (mean = 64.4% for bivalent
317 TSSs vs. 43.5% for non-bivalent TSSs; P-values $< 2.2 \times 10^{-16}$) and the non-TF genes (62.3% vs.

318 29.5%, P-value < 2.2×10^{-16}). CpG promoters are known to be less tissue-specific than non-CpG
319 promoters³⁷, which may seem at odds with our above finding that bivalent genes were significantly
320 more tissue-specific than non-bivalent genes (**Fig. 3d**). To investigate the apparent contradiction, we
321 separated bivalent and non-bivalent TSSs into CpG and non-CpG sub-groups. Indeed, each CpG
322 sub-group is significantly less tissue-specific than the non-CpG subgroup with the same valency, yet
323 the bivalent group is significantly more tissue-specific than the non-bivalent group when CpG and
324 non-CpG promoters are combined (**Supplementary Fig. 9**).

325

326 We examined the TFs with the highest tissue-specificity scores, and a vast majority of these TFs were
327 bivalent. Seventy-five TFs had tissue-specificity scores higher than 6, meaning that the highest
328 expression level was at least as high as the expression levels in all other tissues combined
329 (**Methods**). Of these, 64 were bivalent and the other 11 were not; we illustrate their tissue-specific
330 gene expression (**Fig. 4a**) and the chromatin state assignments around eight example TFs (**Fig. 4b-i**).
331 Two paralogous TFs, *Gata4*, and *Gata1* (**Fig. 4d, e**), illustrate bivalent and non-bivalent genes. *Gata4*,
332 a bivalent gene, is predominantly expressed in the heart, consistent with its well-known role in
333 regulating cardiac development³⁸; it is also expressed at low levels in the stomach and intestine but
334 not in other tissues. Accordingly, its TSS shows broad regions of the Tss state in the heart and
335 narrower Tss regions surrounded by TssBiv and ReprPC regions in the stomach and intestine, while
336 the TSS is covered by only TssBiv and ReprPC regions in other tissues (**Fig. 4e**). In comparison,
337 *Gata1*, a non-bivalent gene, is a key regulator of erythrocyte development³⁹ and is predominantly
338 expressed in the liver. Consistently, the non-bivalent TSS of *Gata1* shows a broad Tss domain in the
339 liver and a narrow Tss domain during early time-points of heart, but it is labeled Quies in other tissues
340 (**Fig. 4d**). Thus, there are two distinct modes of gene repression: bivalent TSSs or quiescent TSSs.

341

342 Other bivalent TFs show similar tissue specificity in their chromatin patterns—adopting the Tss state
343 in the tissues where they are expressed while being in the TssBiv state flanked by ReprPC regions in
344 the tissues that they are not expressed. The homeobox-containing transcription factor *Dlx1* is required
345 for the migration of progenitor cells from the subcortical telencephalon to the neocortex as well as the
346 differentiation of these progenitors into GABAergic neurons⁴⁰. It is expressed in the forebrain and
347 facial prominence; accordingly, its TSS adopts a highly active state in these tissues and the TssBiv-
348 ReprPC repressive states in other tissues (**Fig. 1e**). *Arx* is another homeobox-containing transcription
349 factor (**Fig. 4b**) important for the maturation and migration of GABAergic interneurons, and loss-of-
350 function mutations of *ARX* cause lissencephaly (smooth brain) in humans⁴¹. *En2* encodes a
351 homeobox transcription factor that is expressed at high levels in Purkinje cells and it functions as a
352 transcriptional repressor for neurodevelopment, and *En2* mutant mice display defective cerebellar

353 patterning and a reduction of Purkinje cell number⁴². *En2* is expressed only in the midbrain and
354 hindbrain and shows the corresponding tissue-specific chromatin patterns (**Fig. 4c**). Wilms' tumor-1
355 (*WT1*), which encodes a transcription factor and RNA-binding protein, is essential for kidney
356 development⁴³. It is predominantly expressed in the kidney and at lower levels in the heart, stomach,
357 and intestine. Its TSS is in the Tss state in the kidney and shows a broad TssBiv domain in the heart
358 while being TssBiv-ReprPC in other tissues (**Fig. 4f**). The forkhead transcription factor *Foxq1* is
359 required for the maturation of the abundant mucin-producing foveolar cells that line the mucosal
360 surface in the developing gastrointestinal tract⁴⁴. *Foxq1* is expressed in the gastrointestinal tissues
361 and in the Tss state in these tissues, but bivalent in other tissues (**Fig. 4g**). *Evx2* is required for the
362 morphogenesis of limbs⁴⁵, which is consistent with its expression and chromatin pattern (**Fig. 4h**).
363 Finally, the aristaless-like homeobox 1 transcription factor *Alx1* plays an important role in the
364 development of craniofacial mesenchyme, the first branchial arch, and the limb bud, and a complete
365 loss of function of ALX1 protein causes severe disruption of early craniofacial development in
366 humans⁴⁶. Consistent with its functions, *Alx1* is predominantly expressed in the embryonic facial
367 prominence and shows the corresponding chromatin profile (**Fig. 4i**).

368

369 **Genomic regions assigned to TssBiv are highly conserved evolutionarily**

370 Genomic bins assigned to the bivalent state (TssBiv) are much more evolutionarily conserved than
371 the genomic bins assigned to any of the other 17 chromatin states (**Fig. 5a**). In each biosample, we
372 calculated the mean PhyloP⁴⁷ score in each 200-bp genomic bin and then averaged these mean
373 PhyloP scores for the genomic bins assigned to each chromatin state (**Methods**). The TssBiv state
374 showed the highest PhyloP scores (0.51 averaged over the 66 biosamples), substantially higher
375 (Wilcoxon signed-rank test P-values $< 2.2 \times 10^{-16}$) than the transcription-related states Tx (0.41) and
376 EnhG (0.42), the active TSS state Tss (0.36), the high-signal enhancer state Enh (0.30), which were
377 in turn substantially higher than the remaining 13 states, with Quies2 (0.02) being the lowest (**Fig. 5a**).

378

379 For enhancer-related states (Enh, EnhLo, EnhPois, and EnhPr), the assigned regions in the four brain
380 tissues (forebrain, midbrain, hindbrain, and neural tube) had the highest PhyloP scores, the regions in
381 the liver had the lowest PhyloP scores, and the other seven tissues were in between (**Fig. 5a**). There
382 were some variations in the PhyloP scores over the time-points within each tissue (**Supplementary**
383 **Fig. 10**), but the four brain tissues were clearly the highest and the liver the lowest (**Fig. 5b**). For
384 example, the average PhyloP score of Enh genomic bins was 0.42 for midbrain, while it was 0.13 for
385 liver (Wilcoxon rank-sum test P-value = 5.8×10^{-4} for comparing the 7 midbrain time-points with the 7
386 liver time-points). We examined the transposon content in these Enh genomic bins and found that
387 40.6% of the Enh genomic bins in the liver overlapped annotated transposons, while only 14.1-17.5%

388 of those in the four brain tissues did (**Fig. 5c**), which explained their substantially different levels of
389 evolutionary conservation. These results suggest that the liver tissue has adopted some ancient
390 transposon sequences as enhancers.

391

392 We directly examined the evolutionary conservation of the TSSs of TFs, stratified by whether they
393 resided in a TssBiv genomic bin or not (the two bottom-right panels in **Supplementary Fig. 10**). The
394 average PhyloP score of the TF TSSs in TssBiv genomic bins was 0.82, substantially higher than that
395 of the TF TSSs not in TssBiv genomic bins (0.53, Wilcoxon rank-sum test P-value $< 2.2 \times 10^{-16}$ for
396 comparing the two groups in 66 biosamples). Combined with our aforementioned findings that TFs are
397 highly enriched in bivalent regions, these results indicate that TFs with bivalent TSSs play a key role
398 in evolutionarily conserved pathways driving tissue development.

399

400 **Genomic regions assigned to TssBiv are enriched in PRC2-bound silencers and their target** 401 **TSSs**

402 We used a set of 1800 silencers bound by Polycomb Group 2 proteins (PRC2), identified using ChIA-
403 PET assays targeting PRC2 component proteins in mouse embryonic stem cells²⁸, to further annotate
404 the chromatin states we defined in fetal mouse tissues. The PRC2-bound silencers overlapped
405 extensively with the 14,558 bivalent regions (defined as TssBiv genomic bins surrounded by
406 repressive bins; see **Methods**): 1069 out of 1800 silencers overlapped bivalent regions by at least
407 50% of the lengths of the silencers, while on average only 21 silencers overlapped with random
408 regions with matching sizes as the bivalent regions (Z-score = 140; P-value $< 2.2 \times 10^{-16}$). In individual
409 biosamples, the center locations of most silencers fall in the genomic bins assigned TssBiv or ReprPC
410 ($24 \pm 4\%$ and $28 \pm 6\%$ of the silencer centers, corresponding to 85.7- and 36.4-fold enrichment over
411 the genomic footprints of these states), consistent with the enrichment of these two states in
412 H3K27me3, the histone mark that PRC2 recognizes specifically.

413

414 The enrichment of the PRC2-bound silencers with chromatin states varied by silencer types. The
415 silencers were clustered into four groups according to their H3K27ac signal profiles across the fetal
416 mouse tissues²⁸, a subset of the data we used to define chromatin states (H3K27ac is one of the ten
417 marks used to train our ten-mark model). Group 1 silencers (N = 371) had the highest H3K27ac
418 signals in the fetal mouse tissues²⁸, and the centers of these silencers were in the Tss and Enh states
419 in some biosamples, especially in the brain, but not so much in the liver (**Supplementary Fig. 11**).
420 Group 2 silencers (N = 126) were depleted in H3K27ac in all fetal mouse tissues²⁸, and the centers of
421 most of these silencers were in quiescent states in all tissues (**Supplementary Fig. 11**). Group 3 and
422 4 silencers (N = 683 and 620) had intermediate levels of H3K27ac (higher in Group 3 than in Group

423 4)²⁸, and their centers mostly fell in TssBiv and ReprPC states (**Supplementary Fig. 11**). We included
424 in these alluvial plots the chromatin assignments in mouse embryonic stem cells (ES) using the ten-
425 mark model with ENCODE data on 7 histone marks (missing H3K4me2, ATAC, and WGBS), which
426 show similar chromatin state assignments as in the fetal tissues (**Supplementary Fig. 11**). To
427 normalize for the genomic footprint of each genomic state, we compared genomic bins assigned to
428 TssBiv (the least abundant state; **Fig. 1c**) with an equal number of genomic bins randomly drawn from
429 the other states in individual biosamples for their overlap with each group of PRC2-bound silencers.
430 TssBiv showed the highest enrichment for Group 1 and Group 3 silencers and moderate enrichment
431 for Group 4 silencers; ReprPC showed moderate enrichment for all groups of silencers; Tss showed
432 moderate enrichment for only Group 1 silencers; and none of the other states showed enrichment
433 (**Fig. 6a**).

434
435 The ChIA-PET data further provided the target TSSs for each PRC2-bound silencer²⁸, and these
436 TSSs also predominantly fell in the TssBiv, Tss, and ReprPC states, with the percentages of TSSs in
437 active chromatin states ranked in the descending order for Group 1, 3, 4, and 2 silencers
438 (**Supplementary Fig. 12**). Again, the brain regions showed higher percentages of TSSs in the Tss
439 state than the liver for Group 1 silencers (e.g., 57.9% for forebrain and 13.4% for the liver;
440 **Supplementary Fig. 12**). After normalizing for the genomic footprints of the chromatin states, TssBiv
441 showed a strong enrichment for the target TSSs of all four groups of silencers, while Tss and ReprPC
442 showed weak enrichment (**Fig. 6b**). Among the 75 tissue-specific TFs (**Fig. 4a**), 44 of the 62 bivalent
443 TFs but none of the 13 non-bivalent TFs were targeted by the PRC2-bound silencers (Fisher's exact
444 P-value = 1.6×10^{-6}). Among the seven example bivalent TFs (**Fig. 4b-i**), five were targeted by the
445 silencers (*En2*, *Gata4*, *Wt1*, *Foxq1*, and *Evx2*).

446
447

448 DISCUSSION

449 We defined 18 chromatin states by integrating data on eight histone marks (ChIP-seq), chromatin
450 accessibility (ATAC-seq), and DNA methylation (WGBS) in 66 biosamples across fetal mouse
451 development (**Fig. 1**). We recapitulated the human states previously defined using fewer marks¹⁰ and
452 refined enhancer, bivalent, and quiescent states. Regions annotated in these states showed higher
453 variations among tissues and lower developmental variations across time-points in the same tissue
454 (**Fig. 2a**), and the variations were specific enough to distinguish the tissue-of-origin for the 66
455 biosamples (**Fig. 2c**). Our chromatin state annotation should provide a useful resource for studying
456 mammalian development.

457

458 We define two types of repressive states: ReprPC and ReprPCWk, the two states highly enriched in
459 H3K27me3, jointly occupy 3.7% of the genome, and Het, the state highly enriched in H3K9me3,
460 occupies 1.8% of the genome. However, Zaret and colleagues reported much larger genomic
461 footprints for H3K27me3 domains (~10% of the human genome) and H3K9me3 domains (~20% of the
462 human genome)⁴⁸. They pointed out that if H3K9me3 and H3K27me3 ChIP-seq data were not
463 normalized to input chromatin from the same experiment, reads for those marks would be under-
464 represented, which could result in smaller H3K27me3 and H3K9me3 domains. We did normalize all
465 histone mark ChIP-seq data to the input chromatin from the same experiment, and our signal files for
466 H3K27me3 and H3K9me3 showed the same enriched regions as in the earlier work⁴⁸; thus, the
467 smaller genomic footprints of our ReprPC, ReprPCWk, and Het states were not a normalization
468 artifact. We also define five quiescent states (Quies, Quies2, Quies3, Quies4, QuiesG) collectively
469 occupy 80.5% of the mouse genome. These states show closed chromatin, very low levels of histone
470 marks, and varying levels of DNA methylation. Except for Quies2, the other four quiescent states
471 show low levels of H3K27me3 and H3K9me3 (**Fig. 1c**), the two repressive histone marks, and could
472 encompass some of the H3K27me3 and H3K9me3 domains. Thus, we directly compared ChromHMM
473 states with the H3K27me3 and H3K9me3 domains in the same IMR90 cell line as Becker et al., and
474 found that 17% of H3K27me3 domains and 60% of H3K9me3 domains were in quiescent states;
475 nevertheless, the ReprPC and ReprPCwk states were the most enriched in H3K27me3 domains and
476 the Het state was the most enriched in H3K9me3 domain. Thus, Quies, Quies3, Quies4, and QuiesG
477 states contain large portions of low-signal H3K27me3 and H3K9me3 domains.

478
479 Because enhancers and promoters have been examined extensively in previous ChromHMM
480 studies^{6,10,29}, we decided to focus on the TssBiv state in the current study. TssBiv has the smallest
481 genomic footprint (0.3% of the genome in a particular biosample) among the 18 states, yet TssBiv is
482 discovered consistently by the five-mark, eight-mark, and ten-mark models. TssBiv is particularly
483 conserved evolutionarily, on average more conserved than genomic regions assigned to any other
484 states (**Fig. 5**). We define 14,558 bivalent regions upon an integration of data in 66 biosamples, and
485 roughly half of these regions overlap GENCODE-defined TSSs and the other half are intergenic. The
486 bivalent TSSs show low mRNA levels in a tissue and developmental time-point specific manner (**Fig.**
487 **3**). These TSSs are highly enriched in tissue-specific TF genes (**Fig. 3, 4**). The TF TSSs in the TssBiv
488 state are much more evolutionarily conserved than the TF TSSs in other chromatin states (the two
489 bottom-right panels in **Supplementary Fig. 10**). Comparison with the recent ChIA-PET data²⁸
490 revealed that the bivalent regions are highly enriched in PRC2-bound silencers and their target TSSs.
491 Meanwhile, the TSSs of the target genes of PRC2-bound silencers are highly enriched in the TssBiv
492 state in individual biosamples (**Fig. 6**). Taken together, these results indicate that TssBiv is a
493 chromatin state that marks evolutionarily conserved PRC2-bound silencers and their target TSSs. It is

494 intriguing that both PRC2-bound silencers and their target TSSs possess the same epigenetic
495 signature and hence are assigned the same TssBiv state. This is perhaps not surprising because they
496 are recognized by the PRC2 protein complex. Along this line of reasoning, Enhancers and active
497 TSSs also share some epigenetic features (open chromatin, high levels of active marks such as
498 histone acetylation, and low DNA methylation; Enh and Tss states in **Fig. 1c**).

499

500 Our systematic analysis of bivalent regions in mouse fetal tissues complement earlier studies on
501 bivalent regions in other cell types and biological systems. Bivalent regions were first discovered in
502 embryonic stem cells²³, where their functions have been extensively studied. They have been shown
503 to repress their associated genes and yet allow them to be poised for quick responses to stimuli.
504 When embryonic stem cells differentiate, these bivalent genes become monovalent, retaining either
505 the active marks or the repressive mark, and accordingly be expressed or repressed¹⁹. Subsequent
506 studies reported bivalent domains in the differentiating CD4+ T cells²⁷, the multipotent cranial neural
507 crest cells²⁶, adult intestinal villi cells with regenerative potential²⁴, and terminally differentiated
508 medium spiny neurons in the striatum²⁵. In each of these studies, disruption of Polycomb group
509 proteins led to the activation of the bivalent genes but not genes marked by H3K27me3 only^{24,25},
510 suggesting that bivalency is a mechanism for persistent gene repression from embryonic stem cells to
511 terminally differentiated cells.

512

513 Our analysis of bivalent genes in mouse fetal tissues indicates that they have low expression levels in
514 the tissues where they are bivalent and are enriched for developmental transcription factors under
515 tissue- and time-point-specific repression. A repressed gene can be in a quiescent chromatin state,
516 which corresponds to low levels of all histone marks and high DNA methylation, such as GATA1 (**Fig.**
517 **4d**). Alternatively, it can be in an H3K9me3-enriched Het state accompanied by low levels of active
518 histone marks and high levels of DNA methylation (**Fig 1d**). However, a majority of the bivalent TSSs
519 in fetal tissues overlap CpG islands (mean = 62.5% across the 66 biosamples, vs. 29.8% for non-
520 bivalent TSSs). DNA-hypomethylated CpG islands recruit both Polycomb group and Trithorax group
521 proteins to lay down H3K27me3 and H3K4me3 marks respectively, and the expression level of the
522 gene reflects the competition between Polycomb-mediated repression and Trithorax-mediated
523 activation^{49,50}. As a result, the interplay between the TssBiv, Tss, and ReprPC chromatin states
524 (**Supplementary Fig. 3a**) reflects the main mechanism—distinct from quiescent or Het chromatin
525 states—for silencing genes with CpG-rich TSSs in a tissue-specific manner throughout fetal
526 development and possibly in adulthood.

527

528 In conclusion, we present genome-wide annotations of 18 chromatin states using ten chromatin marks
529 all assayed in a mouse developmental matrix—twelve fetal tissues across 4-7 developmental time-

530 points at daily intervals from E11.5 to birth. These comprehensive annotations enabled us to
531 investigate the changes of chromatin profiles across tissue and time-points and connect the changes
532 with gene expression. In particular, we analyzed bivalent regions in detail and found these
533 evolutionarily conserved regions to be highly enriched in master transcriptional factors important for
534 regulating tissue-specific developmental processes. More broadly, our results suggest that bivalent
535 regions represent a mechanism for silencing CpG-rich genes in a tissue- and time-point-specific
536 manner.

537

538

539 **METHODS**

540 **Experimental data processing for mouse epigenome construction and chromatin state**

541 **definition**

542 We downloaded datasets processed for the mouse genome (mm10) from the ENCODE Portal^{12,51}
543 (<http://encodeproject.org>) that corresponded to eight histone marks (H3K4me1, H3K4me2, H3K4me3,
544 H3K9ac, H3K27ac, H3K36me3, H3K9me3, H3K27me3), ATAC-seq, and WGBS for each of 66
545 epigenomes (**Supplementary Table 1**). All biosamples were from the C57BL/6 mouse strain. For
546 each histone mark, two biological replicates of the ChIP experiment were performed, and for each
547 epigenome, two replicates of the control (input) experiment were performed. We ran ChromHMM⁶ on
548 the 66 epigenomes at the default 200-bp resolution, using the histone ChIP-seq BAM files and the
549 relevant control files for each dataset. For ATAC-seq data, each BAM file was converted to a signal
550 track as follows. Reads were extended to their fragment size and counts-per-million were calculated
551 for all non-overlapping 200-bp genomic windows. Quantile normalization was then applied across the
552 entire data set and the normalized signal was binarized, using a threshold of 0.5. For WGBS data,
553 BED files containing CpG percentages were downloaded from the ENCODE portal (**Supplementary**
554 **table 1**), mean %CpG was calculated for all non-overlapping 200-bp genomic windows and after
555 combining the two replicates for each biosample, binarization was applied, at a cutoff of 50% CpG.

556

557 We defined 18 chromatin states using ChromHMM⁶ using the processed data described above on the
558 10 marks and assigned each 200-bp genomic bin (13,627,678 of them in total for the entire mouse
559 genome) to one of the 18 chromatin states in each biosample. We used the genomic bins with
560 posterior probability > 0.5 for the downstream analysis; these bins composed 99% of the genome on
561 average.

562

563 **Enrichment of chromatin states in other annotations (Fig. 1c)**

564 We assessed the chromatin states assignments in each of the 66 epigenomes for their enrichments in
565 three types of annotations (**Fig. 1c**, the right panel titled Enrichment): (1) for CpG islands, we

566 downloaded `cpGIslandExtUnmasked.txt` from the UCSC Genome Browser; (2) we used GENCODE
567 version M4 for gene-related annotations (transcription start sites or TSS, transcription end sites or
568 TES, gene, exon, and intron); and (3) we used epigenetic annotations (EP300 and CTCF ChIP-seq
569 peaks and DHS).

570

571 For every chromatin state, we computed its enrichment for each annotation, defined as the observed
572 joint probability (P) of a chromatin state and an annotation occurring together over the expected joint
573 probability (i.e., assuming the state and the annotation occur independently):

574

$$575 \quad \text{Enrichment} = P(\text{chromatin state } i, \text{ annotation } j) / P(\text{chromatin state } i) \times P(\text{annotation } j)$$

576

577 For visualization (the right panel of **Fig. 1c** titled Enrichment), the enrichments were scaled between 0
578 and 1:

$$579 \quad \text{Enrichment}_{\text{scaled}} = (\text{Enrichment} - \text{Enrichment}_{\text{min}}) / (\text{Enrichment}_{\text{max}} - \text{Enrichment}_{\text{min}})$$

580

581 We further integrated the RNA-seq data (**Supplementary Table 1**) processed with the ENCODE
582 uniform processing pipeline to compute the enrichment of the chromatin states in expressed or
583 repressed genes for each of the 66 epigenomes¹². For plotting the enrichment panels in **Fig. 1c**, we
584 clustered genes into either expressed or repressed groups in each biosample based on an
585 expression-level cutoff determined using a two-component Gaussian mixture model. The expression
586 levels (in TPM) for the two replicates of each biosample were averaged.

587

588 We calculated the enrichment of the chromatin states in EP300 and CTCF ChIP-seq peaks and
589 DNase hypersensitive sites (the right-most panel in **Fig. 1c**) for those epigenomes that had the EP300
590 and CTCF ChIP-seq or DNase-seq data available in the corresponding tissues and time-points
591 (**Supplementary Table 1**). For the EP300 ChIP-seq data, the BAM files from two biological replicates
592 were pooled, and peaks were called using MACS2⁵² with the q-value cutoff of 0.01. For the CTCF
593 ChIP-seq data, the optimal IDR thresholded peaks⁵³ defined by the ENCODE uniform ChIP-seq
594 pipeline were used¹². For the DNase-seq data, the hotspots defined by the ENCODE uniform DNase-
595 seq processing pipeline were used¹².

596

597 **Partial epigenome simulation and construction (Fig. 1d)**

598 To assess the reliability of chromatin state assignments on epigenomes that lacked the data for one of
599 the ten chromatin marks, for each biosample we simulated ten partial epigenomes, starting with the
600 ten-mark epigenome and omitting the data for each mark individually. We applied the ten-mark 18-
601 state ChromHMM model to the available data on the remaining nine marks and compared the

602 resulting chromatin states assignments with the chromatin state assignments of the ten-mark
603 epigenome by computing the Jaccard similarity between all genomic bins (**Fig. 1d**). The chromatin
604 states with Jaccard similarity less than 0.5 were labeled as misassigned in the missing-one-mark
605 epigenomes.

606
607 For the comparison with PRC2-bound silencers in embryonic stem cells, we also performed chromatin
608 state assignment on embryonic stem cells, with data on seven histone marks (**Supplementary Table**
609 **1**), missing H3K4me2, ATAC, and DNA methylation data. We simulated the effect of missing three
610 marks using midbrain and forebrain samples. These chromatin state assignments of the seven-mark
611 epigenomes were used to define bivalent genes and compared with the bivalent genes defined using
612 the chromatin state assignments of the ten-mark epigenomes (see below).

613
614 **Chromatin state variations across tissues and time-points (Fig. 2a)**

615 We computed Jaccard similarity between a pair of epigenomes by comparing the chromatin states at
616 the corresponding genomic bins between the two epigenomes.

617
618 **UMAP analysis of the epigenomes (Fig. 2d)**

619 We performed two-dimensional visualization of the 66 epigenomes using UMAP³³ analysis on two
620 sets of 200-bp genomic bins: those assigned to the Enh state or the TssBiv state in one or more
621 biosamples. For the Enh genomic bins, UMAP was provided with the H3K27ac signal levels across
622 the 66 biosamples and the following parameters were used: $n_neighbors = 7$, $min_dist = 0.5$, $seed =$
623 11 . For the TssBiv genomic bins, UMAP was provided with the signal levels of all ten marks across
624 the 66 biosamples and the following parameters were used: $n_neighbors = 10$, $min_dist = 0.04$, $seed$
625 $= 12$.

626
627 **Identification of bivalent TSSs and bivalent genes (Fig. 3, 4)**

628 We developed a method to identify bivalent TSSs and bivalent genes by their chromatin states in
629 each epigenome, described as follows. We first converted each epigenome to a character string using
630 an 18-letter alphabet (one symbol for each state). Regular expressions were then used to extract
631 punctate (median length 1800 bp) bivalent domains (stretches of contiguous genomic bins) in each
632 epigenome, defined as bivalent chromatin states flanked by quiescent or heterochromatin states
633 (ReprPC, ReprPCWk, Quies, Quies2, Quies3, Quies4, or QuiesG state). We used the union (14,558
634 regions across all tissue time-points, median 3,514 per tissue time-point, neighboring regions were
635 not merged) of the detected genomic regions matching our regular expression for downstream
636 analyses. Of the 14,558 regions detected in the 66 biosamples collectively, 14,729 regions
637 overlapped GENCODE-annotated TSSs; we denote these *bivalent TSSs*. We further define a *bivalent*

638 *gene* as having at least one bivalent TSS, yielding 6,797 genes that are bivalent in any of the 12
639 tissues.

640

641 We detected on average ~3,400 bivalent genes per tissue, defined as genes that are bivalent in any
642 of the time-points in the tissue. We performed Gene Ontology (GO) analysis on bivalent genes using
643 the PANTHER tool⁵⁴. The genes used in the Gene Ontology (GO) analysis, of which the results are
644 listed in **Supplementary Table 4** were obtained as follows: TSSs extracted from the M4 GENCODE
645 annotations were intersected with the bivalent regions detected in each tissue. For each tissue, genes
646 for which one or more TSSs intersected were retained. Then, the 1,077 genes that were found to
647 have TSSs overlapping bivalent regions in *all* tissues were used as input for the GO analysis
648 (**Supplementary Table 4a, b**). Another set of 1,291 genes was obtained using the same process,
649 except genes were collected that had TSSs in bivalent regions *only* in liver samples and *not* in any
650 other 11 tissues (**Supplementary table 4c, d**). Gene IDs were translated into gene names prior to
651 submission to PANTHER. For six gene IDs, no matching gene name was found, leaving 1,074 and
652 1,288 genes in the “all tissues” and the “liver-only” gene sets for submission. PANTHER was run on
653 the GO “Biological Process” ontology, using Fisher’s exact test and FDR for P-value calculations.

654

655 **Gene annotations and identification of transcription factors (Fig. 3d, Fig. 4, supplementary**
656 **tables 2-4)**

657 GENCODE M4 gene annotations were used to identify genes and transcription start sites (TSSs). To
658 avoid double-counting TSSs, coinciding TSSs were merged. To identify transcription factors, we used
659 the list of transcription factors and their homologs in mouse and human³⁴. Ensembl IDs were obtained
660 by mapping gene names to the GENCODE M4 annotations⁵⁵. 552 TFs matched IDs in the GENCODE
661 M4 mouse annotations.

662

663 **Evolutionary analysis (Fig. 5a-b, Supplementary Fig. 10)**

664 We averaged the mouse 60-way phyloP⁴⁷ score across the genomic positions in each 200-bp
665 genomic bin. We then average this per-bin score for all the genomic bins assigned to a particular
666 chromatin state in each biosample to obtain the average PhyloP score per state per biosample
667 (**Supplementary Fig. 10**, first 18 panels). For each tissue (**Fig. 5a**), the PhyloP scores from the
668 biosamples at different time-points were further averaged. For the TF TSSs (**Supplementary Fig. 10**,
669 the two bottom-right panels), we used the PhyloP score for genomic bins where each TF TSS resided
670 in, stratified by whether that bin was assigned to the TssBiv state or not.

671

672 **Overlap of Enh regions with annotated transposons (Fig. 5c)**

673 We used transposon annotations in the mouse genome from Repbase⁵⁶ to analyze the Enh state
674 across different tissues (**Fig. 5c**). We overlapped the genomic bins assigned to the Enh state in each
675 biosample with annotated transposons, requiring at least 1-bp overlap. The percentage of all genomic
676 bins that overlapped transposons was used as control (gray dashed line in **Fig. 5c**).

677 678 **Analysis of PRC2-bound silencers (Fig. 6, Supplementary Fig. 11, 12)**

679 We used the 18,000 PRC2-bound silencers classified into four groups based on their H3K27ac signal
680 in mouse fetal tissues²⁸. We overlapped the PRC2-bound silencers with our 14,558 bivalent regions,
681 requiring at least half of the length of a silencer length to overlap. We randomly selected genomic
682 regions with the same lengths as the bivalent regions to act as controls. Furthermore, we assigned
683 each silencer to a chromatin state in a particular biosample according to which chromatin state the
684 center of the silencer falls in.

685
686 We included embryonic stem cells in this analysis (ES-Bruce4). These cells were derived from
687 C57BL/6, the same strain of mice from which the tissues were harvested. We only had data on seven
688 histone marks on embryonic stem cells (**Supplementary Table 1**), and simulation of this partial
689 epigenome (see above Methods) showed no major impact on the assignment of the TssBiv state and
690 the resulting bivalent genes. Simulating using midbrain and forebrain samples, we found that most
691 bivalent genes were identified using the partial epigenome. For example, among the 2,250 bivalent
692 genes in the midbrain E11.5 sample, 2,014 (89.5%) were identified using the partial epigenome.

693 694 **Data availability**

695 All experimental data used in this paper can be accessed at the encode Portal
696 (<http://www.encodeproject.org/>), using the accession IDs listed in **Supplementary Table 1**.

697 698 **Code Availability**

699 The code used to extract genomic regions based on regular expression can be found on GitHub, at
700 <https://github.com/weng-lab/stateregexp.git>.

701 702 **Data visualization via a UCSC track hub**

703 We made a track hub (https://users.wenglab.org/vanderva/trackhub/chromhmpaper/hub_0.txt) for
704 the UCSC genome browser⁵⁷ to visualize all the data and annotations used in this study listed below.

705 The trackhub can be accessed via a UCSC session:

706 https://genome.ucsc.edu/s/Kaili/ChromHMM_paper.

707
708 1. ten-mark, 18-state chromatin state assignments (in dense mode)

- 709 BigWig experimental data complete for 66 biosamples (in hide mode):
- 710 a. ChIP-seq of eight histone marks
- 711 b. ATAC-seq
- 712 c. WGBS
- 713 d. RNA-seq
- 714 e. DNase when available
- 715 f. EP300 ChIP-seq when available
- 716 g. CTCF ChIP-seq when available
- 717 2. ES-Bruce4 chromatin state assignments (in dense mode)
- 718 BigWig experimental data for ES-Bruce4 (in hide mode)
- 719 a. ChIP-seq of seven histone marks
- 720 b. RNA-seq
- 721 c. EP300 ChIP-seq
- 722 d. CTCF ChIP-seq
- 723 3. Turn on the GENCODE gene annotation (in pack mode)
- 724 4. Turn on the CpG island track from UCSC (in dense mode)
- 725 5. Bivalent regions (in dense mode)
- 726 6. PRC-bound silencers and their target TSSs in two tracks (in dense mode)
- 727 7. Turn on the PhyloP conservation track (in full mode)
- 728 8. Turn on VISTA enhancer track hub (in hide mode)
- 729 9. Mouse cCREs (in hide mode)

730
731

732 **Acknowledgments**

733 We thank ENCODE Consortium members for generating the ATAC-seq, ChIP-seq, RNA-seq, WGBS,
734 and DNase-seq data on the 66 mouse embryonic biosamples and making them freely available. This
735 work was supported in part by the National Institutes of Health grants HG009446 and HG007000 to
736 ZW.

737

738 **Author contributions**

739 AV: computational analysis, writing; KF: computational analysis; JT: computational analysis; JM:
740 computational analysis; MP: computational analysis; HP: computational analysis, writing; ZW, project
741 conception, design and management, writing.

742

743

744 **FIGURE CAPTIONS**

745 **Figure 1: Overview of the 66 epigenomes and 18 chromatin states during mouse**
746 **embryogenesis.**

- 747 a. Twelve tissues at 4-7 developmental time-points have ChIP-seq data for eight histone marks
748 (green boxes), ATAC-seq data, and DNA methylation (DNAm) data, totaling 66 complete
749 epigenomes. Twenty-one of these epigenomes also have DNase-seq data (red dots).
750 Embryonic stem cells (orange box) have ChIP-seq data for seven histone marks, and are
751 missing H3K4me2, ATAC-seq, and DNAm.
- 752 b. Eighteen chromatin states are defined by ChromHMM across the 66 complete epigenomes.
- 753 c. Histone-mark probabilities, genome coverage, and overlapping genomic features including
754 gene expression, regulatory features (P300 binding, CTCF binding, and DNase I
755 hypersensitive sites), and distances to the TSSs of expressed and repressed genes are shown
756 for each chromatin state. The enrichments for the categories are the averaged values across
757 tissues and time-points.
- 758 d. Jaccard similarities between the partial epigenomes with each mark omitted and the ten-mark
759 E13.5 midbrain epigenome.
- 760 e. The *Dlx1* locus is displayed with chromatin states (color-coded as in **a**) in the forebrain and the
761 liver for all seven time-points. Also shown are the signals of several histone marks (scale: 0–
762 50) that differ between forebrain and liver (for E11.5, E13.5, E15.5, and P0 only, due to space
763 constraints), along with ATAC and DNA methylation signals. A transgenic mouse embryo is
764 shown on top of the enhancer region, indicating the forebrain-specific activity of this enhancer.
765 A CpG island that overlaps with the bivalent region at the TSS of *Dlx1* is shown at the bottom
766 of the panel.

767

768 **Figure 2: Variations of the chromatin states across tissues and their transitions along the**
769 **developmental trajectory.**

- 770 a. Jaccard similarity between different time-points in the same tissue (y-axis) versus the similarity
771 between different tissues at the same time-point (x-axis). Error bars indicate the range
772 between the first and third quartiles.
- 773 b. Transitions between chromatin states along midbrain developmental time-points. For clarity,
774 only the genomic bins assigned TSS-related states (Tss, TssFlnk, and TssBiv) at one or more
775 time-points are included.
- 776 c. Same as b but for genomic bins assigned enhancer-related states (Enh, EnhLo, EnhPois, and
777 EnhPr) at one or more developmental time-points.
- 778 d. Visualization of the 66 epigenomes in two dimensions using the UMAP technique. (Left)
779 UMAP was given the H3K27ac signals in the Enh genomic bins across the 66 epigenomes.

780 There were 735,048 such genomic bins, which were assigned Enh in one or more
781 epigenomes. (Right) UMAP was given the signals of all ten marks in the TssBiv genomic bins
782 across the 66 epigenomes. There were 156,752 such genomic bins, which were assigned
783 TssBiv in one or more epigenomes.

784

785 **Figure 3: Count and expression of bivalent genes along developmental time-points.**

- 786 a. The number of bivalent genes at 1 to 7 time-points in the midbrain. Observed and expected
787 numbers of genes are in red and in gray respectively.
- 788 b. Median expression levels of three groups of genes: (green) bivalent at the earliest time-point
789 but not at the last time-point, (blue) bivalent at the last time-point but not at the first time-point,
790 and (pink) bivalent at all time-points.
- 791 c. Distribution of gene expression, with genes grouped by the total number of time-points at
792 which their TSSs are in the bivalent state TssBiv (left) or in the active state Tss (right) in the
793 forebrain. The total number of genes in each group is shown below each box plot in
794 parentheses. For all boxplots, whiskers show 95% confidence intervals, boxes represent the
795 first and third quartiles, the vertical midline is the median, and outliers are omitted. There is a
796 negative correlation between expression and the duration of the bivalent state and a positive
797 correlation between expression and the duration of the active state (P-values $< 2.2 \times 10^{-16}$).
- 798 d. Violin plots show the distributions of tissue-specificity scores for bivalent and non-bivalent
799 genes that encode transcription factors (TFs) and non-TFs. Medians are shown in black bars
800 with values indicated. P-values are shown for three comparisons as indicated.

801

802 **Figure 4: Expression profiles and chromatin states for the transcription factors with the**
803 **highest tissue-specificity scores.**

- 804 a. Hierarchical clustering of expression profiles for the TFs with tissue specificity scores greater
805 than 6, with 75 TFs in total. Rows on the top show the maximal expression level across all
806 biosamples (intensities of red), bivalency status (brown for 62 bivalent TFs, and yellow for 13
807 non-bivalent TFs), and tissue specificity score (intensities of green).
- 808 b-i. Example TFs and the chromatin state assignments near their loci. Among these, Gata1 (d) is a
809 non-bivalent TF and the rest are bivalent TFs: Arx (b), En2 (c), Gata4 (e), Wt1 (f), Foxq1 (g),
810 Evx2 (h), and Alx1 (i). Each gene name is near the 5'-end of the gene, and CpG islands are
811 indicated as green boxes beneath each gene. Chromatin states are colored as in Fig. 1.

812

813 **Figure 5: Evolutionary conservation of genomic regions by chromatin state.**

- 814 a. The PhyloP conservation score (phyloP60way for mm10) for genomic regions assigned to
815 each chromatin state. Colors correspond to tissues.

- 816 **b.** PhyloP score for genomic bins assigned to Enh in all 12 tissues.
817 **c.** Percentage of bins assigned to Enh that overlap with transposons, for all 12 tissues.

818

819 **Figure 6: PRC2-bound silencers and their target TSSs are enriched in the TssBiv and ReprPC**
820 **states.**

- 821 **a.** Percentage of PRC2-bound silencers whose centers overlap a genomic bin assigned to the
822 TssBiv, Tss, ReprPC, or other chromatin states. Silencers were divided into four groups by
823 Ngan et al.²⁸ according to H3K27ac signals in mouse fetal tissue biosamples. To normalize for
824 the differential genomic coverage of the chromatin states, the same numbers of genomic bins
825 were randomly drawn in the other states to match the number of genomic bins in TssBiv in
826 each biosample. States are colored as in Fig. 1b and the average of the other 15 states is
827 shown as a gray dashed line.
828 **b.** Same as **a** but for the TSSs targeted by the PRC2-bound silencers defined by Ngan et al.²⁸.

829

830

831 **SUPPLEMENTARY FIGURE CAPTIONS**

832 **Supplementary Fig. 1: Comparison of the five-mark, eight-mark, and ten-mark models.**

833 **a** Emission probabilities for the five-mark 15-state model, the eight-mark 15-state model, and the ten-
834 mark 18-state model, with the ten-mark model reproduced from **Fig. 1c** for easy comparison with the
835 other two models. **b-e** Alluvial plots illustrate the correspondence of chromatin states across the three
836 models in forebrain e13.5. **b** With genomic bins assigned to a quiescent state by all three models
837 omitted, 3,743,342 genomic bins are shown. **c** All 13,627,678 200-bp bins in the genome. **d** Genomic
838 bins assigned to Enh by one of the models. **e** Genomic bins assigned to TssBiv by one of the models.

839

840 **Supplementary Fig. 2: Chromatin state transition from early to late time-points.**

841 Genomic bins assigned to TssBiv in the forebrain at one or more time-points are included. States are
842 colored as in Fig. 1c.

843

844 **Supplementary Fig. 3: Comparison of genomic coverages of TssBiv, Tss and ReprPC, and the**
845 **decrease of TssBiv coverage over time.**

846 **a** A Venn diagram shows the overlap of genomic regions assigned to the TssBiv, Tss, or ReprPC
847 state in any of the 66 epigenomes. **b** red lines show the percentages of the genome in the TssBiv
848 state over the time course of development for each tissue. Only the five tissues with seven time-points
849 are included. P-values for linear fit (blue dashed line, with the 95% confidence interval in gray shaded
850 area) are provided.

851

852 **Supplementary Fig. 4: Word clouds for enriched GO terms in bivalent genes.**

853 Gene ontology (GO) enrichment analyses were performed using the PANTHER tool for two groups of
854 genes: genes with bivalent TSSs in **(a)** all 12 tissues; **(b)** in the liver and not in any other tissues. For
855 each analysis, a summary of significantly enriched GO terms is presented as a word cloud. See
856 Supplementary Table 4 for full PANTHER results.

857

858 **Supplementary Fig. 5: Expression levels of genes with or without bivalent TSSs.**

859 For each tissue, the expression levels of genes (in TPM) are plotted, stratified by whether it has a
860 bivalent TSS at all time-points. For each box plot, the total number of genes in each group is shown at
861 the bottom. Outliers are omitted for clarity. Wilcoxon P-values for comparing the two groups of genes
862 in each tissue are provided.

863

864 **Supplementary Fig. 6: Number of genes that are bivalent in a certain number of time-points.**

865 For each tissue, the total number of genes that are deemed bivalent in a certain number of time-points
866 is plotted in red, compared with the expected number (in grey) if genes were randomly assigned to be
867 bivalent at each time-point.

868

869 **Supplementary Fig. 7: Expression of genes with bivalent TSS at early vs. late time-points.**

870 Median expression levels of genes stratified into three distinct categories are plotted: genes deemed
871 bivalent at the first time-point but not at the last (early-bivalent genes; blue line); genes deemed
872 bivalent at the last time-point but not at the first (late-bivalent genes; green); and genes with bivalent
873 TSS at all time-points (all-bivalent genes; red dashed line). Wilcoxon rank-sum test P-values for the
874 comparisons between early- and late-bivalent genes for their expression levels at the first time-point
875 (green P-values); between the early- and late-bivalent genes at the last time-point (blue P-values);
876 and between all-bivalent genes vs. early- and late-bivalent genes (red P-values). n.s. stands for not
877 significant.

878

879 **Supplementary Fig. 8: Correlations between gene expression and the duration of bivalent and
880 active TSSs.**

881 The expression levels of genes with a certain duration (number of time-points) of being bivalent **(a)** or
882 active **(b)** in each tissue. Numbers in parentheses indicate the number of genes for each duration. P-
883 values were computed with ANOVA with multiple-testing correction.

884

885 **Supplementary Fig. 9: Tissue specificity for groups of genes classified by whether their TSSs
886 overlap CpG islands.**

887 Tissue-specificity scores are shown for all genes, and subsets of genes depending on whether they
888 encoded TFs, they have a bivalent TSS, and whether the TSSs overlap CpG islands. Wilcoxon P-
889 values for comparing the CpG and non-CpG groups of genes are provided.

890

891 **Supplementary Fig. 10: Evolutionary conservation for genomic bins assigned to each**
892 **chromatin state in each biosample.**

893 Average PhyloP scores are plotted for the 18 chromatin states. The last two panels (bottom right) are
894 TSSs of transcription factors stratified by whether they fall in a TssBiv genomic bin or not. The
895 thickness of a line corresponds to the standard error. Tissues are colored accordingly.

896

897 **Supplementary Fig. 11: Chromatin state assignments for the center positions of PRC2-bound**
898 **silencers.**

899 Four groups of PRC2-bound silencers correspond to those in Fig. 6a are plotted across time-points in
900 the forebrain (**a-d**) and liver (**e-h**). The state assignments for mouse embryonic stem cells (ES) are
901 included for comparison. States are colored as in Fig. 1c.

902

903 **Supplementary Fig. 12: Chromatin state assignments for the TSSs targeted by PRC2-bound**
904 **silencers.**

905 This figure corresponds to supplementary Fig. 12 but for the TSSs targeted by PRC2-bound silencers.

906

907

908 **SUPPLEMENTARY TABLES**

909 **Supplementary Table 1: Input datasets and their ENCODE accessions.** ENCODE file accession
910 IDs for all input files. **a.** BAM files for histone ChIP-seq datasets and controls. **b.** BED files with CpG
911 calls from WGBS. **c.** RNA-seq TPM matrices for the two replicates of each biosample. **d.** BAM files for
912 ATAC-seq. **e.** BAM files for DNase-seq.

913

914 **Supplementary Table 2: Bivalent TSSs in each biosample.** GENCODE M4 TSS annotations were
915 intersected with bivalent regions in each biosample. Sites occupying the same genomic position were
916 merged.

917

918 **Supplementary Table 3: Bivalent genes and their expression levels.** **a.** Expression levels are
919 reported in TPM, in each tissue and time-point. **b.** The number of bivalent genes shared between any
920 pair of biosamples. **c.** The number of bivalent genes shared between any pair of tissues. Diagonal
921 numbers indicate the total number of bivalent genes in each tissue. **d.** Bivalent state of the TSSs of

922 the genes in each biosample. **e.** Bivalent regions defined across all biosamples. **f.** Union of bivalent
923 regions detected in all biosamples, as determined by regular expression (see Methods).

924

925 **Supplementary Table 4: GO enrichment analysis using the PANTHER tool.** **a.** PANTHER output
926 for genes that are bivalent in all tissues. **b.** List of genes submitted for analysis in **a.** **c.** PANTHER
927 output for genes that are bivalent exclusively in the liver. **d.** List of genes submitted for analysis in **c.**

928

929

930 REFERENCE

- 931 1. Barski, A. *et al.* High-Resolution Profiling of Histone Methylations in the Human Genome. *Cell* vol.
932 129 823–837 (2007).
- 933 2. Robertson, G. *et al.* Genome-wide profiles of STAT1 DNA association using chromatin
934 immunoprecipitation and massively parallel sequencing. *Nat. Methods* **4**, 651–657 (2007).
- 935 3. Johnson, D. S., Mortazavi, A., Myers, R. M. & Wold, B. Genome-wide mapping of in vivo protein-
936 DNA interactions. *Science* **316**, 1497–1502 (2007).
- 937 4. Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J. Transposition of
938 native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding
939 proteins and nucleosome position. *Nat. Methods* **10**, 1213–1218 (2013).
- 940 5. Lister, R. *et al.* Human DNA methylomes at base resolution show widespread epigenomic
941 differences. *Nature* **462**, 315–322 (2009).
- 942 6. Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and characterization.
943 *Nat. Methods* **9**, 215–216 (2012).
- 944 7. Hoffman, M. M. *et al.* Unsupervised pattern discovery in human chromatin structure through
945 genomic segmentation. *Nature Methods* vol. 9 473–476 (2012).
- 946 8. Zhang, Y., An, L., Yue, F. & Hardison, R. C. Jointly characterizing epigenetic dynamics across
947 multiple human cell types. *Nucleic Acids Res.* **44**, 6721–6731 (2016).
- 948 9. Bernstein, B. E. *et al.* An integrated encyclopedia of DNA elements in the human genome. *Nature*
949 **489**, 57–74 (2012).
- 950 10. Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330
951 (2015).
- 952 11. Yue, F. *et al.* A comparative encyclopedia of DNA elements in the mouse genome. *Nature* **515**,
953 355–364 (2014).
- 954 12. The ENCODE Project Consortium, Jill E. Moore, *et al.* In press. Expanded Encyclopedias of DNA
955 Elements in the Human and Mouse Genomes. *Nature* (2020).
- 956 13. Gorkin, D. *et al.* Systematic mapping of chromatin state landscapes during mouse development.

- 957 *bioRxiv* 166652 (2017) doi:10.1101/166652.
- 958 14. He, Y. *et al.* Spatiotemporal DNA Methylation Dynamics of the Developing Mammalian Fetus.
959 *bioRxiv* 166744 (2017).
- 960 15. Heintzman, N. D. *et al.* Distinct and predictive chromatin signatures of transcriptional promoters
961 and enhancers in the human genome. *Nat. Genet.* **39**, 311–318 (2007).
- 962 16. Heintzman, N. D. *et al.* Histone modifications at human enhancers reflect global cell-type-specific
963 gene expression. *Nature* **459**, 108–112 (2009).
- 964 17. Wang, Z. *et al.* Combinatorial patterns of histone acetylations and methylations in the human
965 genome. *Nat. Genet.* **40**, 897–903 (2008).
- 966 18. Creyghton, M. P. *et al.* Histone H3K27ac separates active from poised enhancers and predicts
967 developmental state. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 21931–21936 (2010).
- 968 19. Guttman, M. *et al.* Chromatin signature reveals over a thousand highly conserved large non-
969 coding RNAs in mammals. *Nature* **458**, 223–227 (2009).
- 970 20. Margueron, R. & Reinberg, D. The Polycomb complex PRC2 and its mark in life. *Nature* **469**,
971 343–349 (2011).
- 972 21. He, P. *et al.* The changing mouse embryo transcriptome at whole tissue and single-cell
973 resolution. *Nature* **32**.
- 974 22. Hesselberth, J. R. *et al.* Global mapping of protein-DNA interactions in vivo by digital genomic
975 footprinting. *Nat. Methods* **6**, 283–289 (2009).
- 976 23. Bernstein, B. E. *et al.* A bivalent chromatin structure marks key developmental genes in
977 embryonic stem cells. *Cell* **125**, 315–326 (2006).
- 978 24. Jadhav, U. *et al.* Acquired Tissue-Specific Promoter Bivalency Is a Basis for PRC2 Necessity in
979 Adult Cells. *Cell* **165**, 1389–1400 (2016).
- 980 25. von Schimmelmann, M. *et al.* Polycomb repressive complex 2 (PRC2) silences genes
981 responsible for neurodegeneration. *Nat. Neurosci.* **19**, 1321–1330 (2016).
- 982 26. Minoux, M. *et al.* Gene bivalency at Polycomb domains regulates cranial neural crest positional
983 identity. *Science* **355**, (2017).
- 984 27. Wei, G. *et al.* Global mapping of H3K4me3 and H3K27me3 reveals specificity and plasticity in
985 lineage fate determination of differentiating CD4⁺ T cells. *Immunity* **30**, 155–167 (2009).
- 986 28. Ngan, C. Y. *et al.* Chromatin interaction analyses elucidate the roles of PRC2-bound silencers in
987 mouse development. *Nat. Genet.* (2020) doi:10.1038/s41588-020-0581-x.
- 988 29. Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human cell types.
989 *Nature* **473**, 43–49 (2011).
- 990 30. Eckner, R. *et al.* Molecular cloning and functional analysis of the adenovirus E1A-associated 300-
991 kD protein (p300) reveals a protein with properties of a transcriptional adaptor. *Genes Dev.* **8**,
992 869–884 (1994).

- 993 31. Yao, T.-P. *et al.* Gene Dosage–Dependent Embryonic Development and Proliferation Defects in
994 Mice Lacking the Transcriptional Integrator p300. *Cell* **93**, 361–372 (1998).
- 995 32. Visel, A., Minovitsky, S., Dubchak, I. & Pennacchio, L. A. VISTA Enhancer Browser--a database
996 of tissue-specific human enhancers. *Nucleic Acids Res.* **35**, D88–92 (2007).
- 997 33. Becht, E. *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nat.*
998 *Biotechnol.* (2018) doi:10.1038/nbt.4314.
- 999 34. Stergachis, A. B. *et al.* Conservation of trans-acting circuitry during mammalian regulatory
1000 evolution. *Nature* **515**, 365–370 (2014).
- 1001 35. Ku, M. *et al.* Genomewide analysis of PRC1 and PRC2 occupancy identifies two classes of
1002 bivalent domains. *PLoS Genet.* **4**, e1000242 (2008).
- 1003 36. Riising, E. M. *et al.* Gene silencing triggers polycomb repressive complex 2 recruitment to CpG
1004 islands genome wide. *Mol. Cell* **55**, 347–360 (2014).
- 1005 37. Saxonov, S., Berg, P. & Brutlag, D. L. A genome-wide analysis of CpG dinucleotides in the
1006 human genome distinguishes two distinct classes of promoters. *Proc. Natl. Acad. Sci. U. S. A.*
1007 **103**, 1412–1417 (2006).
- 1008 38. Kuo, C. T. *et al.* GATA4 transcription factor is required for ventral morphogenesis and heart tube
1009 formation. *Genes Dev.* **11**, 1048–1060 (1997).
- 1010 39. Pevny, L. *et al.* Development of hematopoietic cells lacking transcription factor GATA-1.
1011 *Development* **121**, 163–172 (1995).
- 1012 40. Anderson, S. A., Eisenstat, D. D., Shi, L. & Rubenstein, J. L. Interneuron migration from basal
1013 forebrain to neocortex: dependence on Dlx genes. *Science* **278**, 474–476 (1997).
- 1014 41. Kitamura, K. *et al.* Mutation of *ARX* causes abnormal development of forebrain and testes in mice
1015 and X-linked lissencephaly with abnormal genitalia in humans. *Nat. Genet.* **32**, 359–369 (2002).
- 1016 42. Jankowski, J., Holst, M. I., Liebig, C., Oberdick, J. & Baader, S. L. Engrailed-2 negatively
1017 regulates the onset of perinatal Purkinje cell differentiation. *J. Comp. Neurol.* **472**, 87–99 (2004).
- 1018 43. Kreidberg, J. A. WT1 and kidney progenitor cells. *Organogenesis* **6**, 61–70 (2010).
- 1019 44. Verzi, M. P., Khan, A. H., Ito, S. & Shivdasani, R. A. Transcription factor foxq1 controls mucin
1020 gene expression and granule content in mouse stomach surface mucous cells. *Gastroenterology*
1021 **135**, 591–600 (2008).
- 1022 45. Héroult, Y., Hraba-Renevey, S., van der Hoeven, F. & Duboule, D. Function of the *Evx-2* gene in
1023 the morphogenesis of vertebrate limbs. *EMBO J.* **15**, 6727–6738 (1996).
- 1024 46. Uz, E. *et al.* Disruption of *ALX1* causes extreme microphthalmia and severe facial clefting:
1025 expanding the spectrum of autosomal-recessive *ALX*-related frontonasal dysplasia. *Am. J. Hum.*
1026 *Genet.* **86**, 789–796 (2010).
- 1027 47. Pollard, K. S., Hubisz, M. J., Rosenbloom, K. R. & Siepel, A. Detection of nonneutral substitution
1028 rates on mammalian phylogenies. *Genome Res.* **20**, 110–121 (2010).

- 1029 48. Becker, J. S. *et al.* Genomic and Proteomic Resolution of Heterochromatin and Its Restriction of
1030 Alternate Fate Genes. *Mol. Cell* **68**, 1023–1037.e15 (2017).
- 1031 49. Schuettengruber, B., Bourbon, H.-M., Di Croce, L. & Cavalli, G. Genome Regulation by Polycomb
1032 and Trithorax: 70 Years and Counting. *Cell* **171**, 34–57 (2017).
- 1033 50. Holoch, D. & Margueron, R. Mechanisms Regulating PRC2 Recruitment and Enzymatic Activity.
1034 *Trends Biochem. Sci.* **42**, 531–542 (2017).
- 1035 51. Sloan, C. A. *et al.* ENCODE data at the ENCODE portal. *Nucleic Acids Res.* **44**, D726–32 (2016).
- 1036 52. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).
- 1037 53. Li, Q., Brown, J. B., Huang, H. & Bickel, P. J. Measuring reproducibility of high-throughput
1038 experiments. *Ann. Appl. Stat.* **5**, 1752–1779 (2011).
- 1039 54. Mi, H. *et al.* PANTHER version 11: expanded annotation data from Gene Ontology and
1040 Reactome pathways, and data analysis tool enhancements. *Nucleic Acids Res.* **45**, D183–D189
1041 (2017).
- 1042 55. Harrow, J. *et al.* GENCODE: The reference human genome annotation for The ENCODE Project.
1043 *Genome Res.* **22**, 1760–1774 (2012).
- 1044 56. Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive elements in
1045 eukaryotic genomes. *Mob. DNA* **6**, 11 (2015).
- 1046 57. Kent, W. J. *et al.* The human genome browser at UCSC. *Genome Res.* **12**, 996–1006 (2002).
- 1047

Figure 1: Overview of the 66 epigenomes and 18 chromatin states during mouse embryogenesis.

a. Twelve tissues at 4-7 developmental time-points have ChIP-seq data for eight histone marks (green boxes), ATAC-seq data, and DNA methylation (DNAm) data, totaling 66 complete epigenomes. Twenty-one of these epigenomes also have DNase-seq data (red dots). Embryonic stem cells (orange box) have ChIP-seq data for seven histone marks, and are missing H3K4me2, ATAC-seq, and DNAm.

b. Eighteen chromatin states are defined by ChromHMM across the 66 complete epigenomes.

c. Histone-mark probabilities, genome coverage, and overlapping genomic features including gene expression, regulatory features (P300 binding, CTCF binding, and DNase I hypersensitive sites), and distances to the TSSs of expressed and repressed genes are shown for each chromatin state. The enrichments for the categories are the averaged values across tissues and time-points.

d. Jaccard similarities between the partial epigenomes with each mark omitted and the ten-mark E13.5 midbrain epigenome. The *Dlx1* locus is displayed with chromatin states (color-coded as in **a**) in the forebrain and the liver for all seven time-points. Also shown are the signals of several histone marks (scale: 0–50) that differ between forebrain and liver (for E11.5, E13.5, E15.5, and P0 only, due to space constraints), along with ATAC and DNA methylation signals. A transgenic mouse embryo is shown on top of the enhancer region, indicating the forebrain-specific activity of this enhancer. A CpG island that overlaps with the bivalent region at the TSS of *Dlx1* is shown at the bottom of the panel.

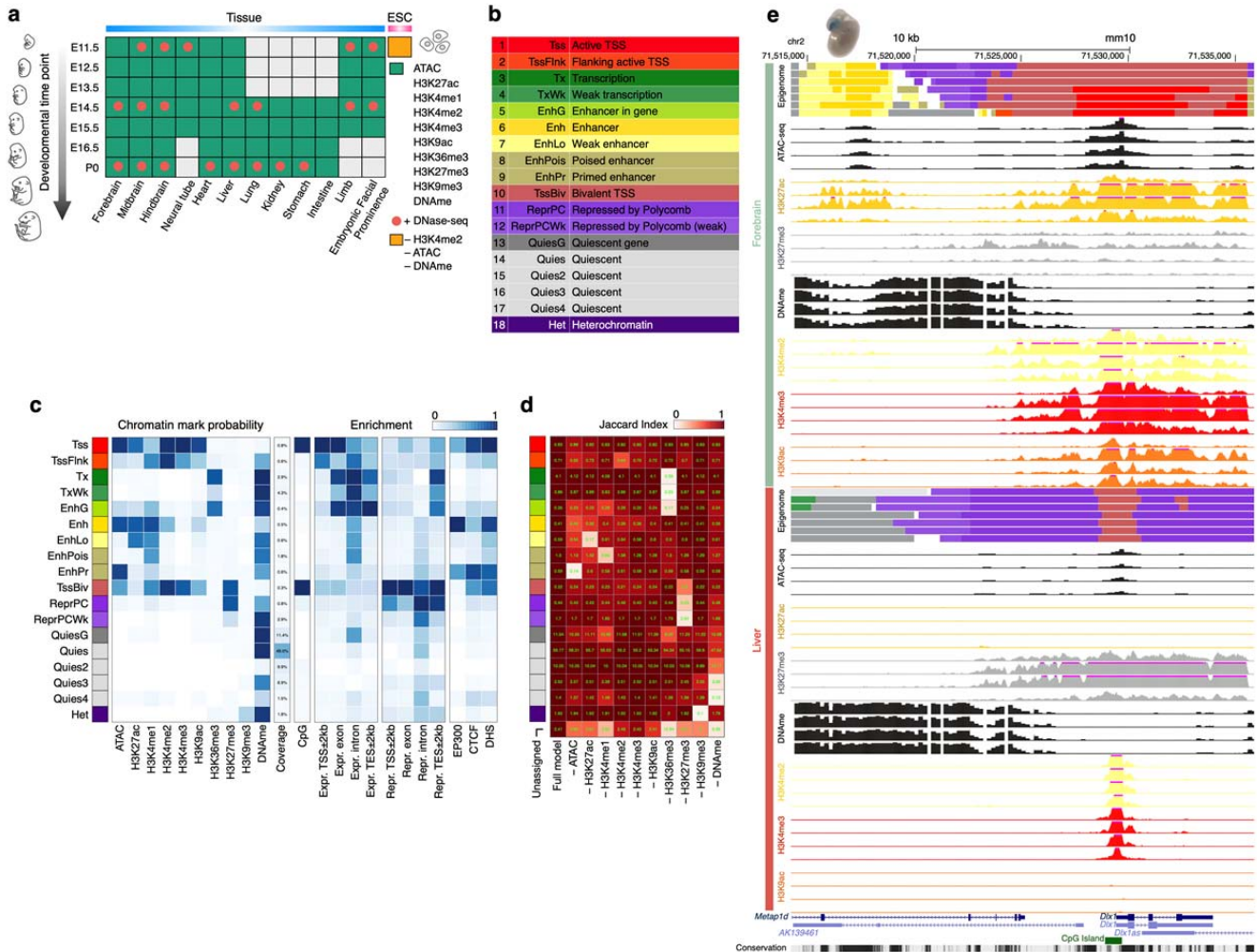


Figure 2: Variations of the chromatin states across tissues and their transitions along the developmental trajectory.

- a.** Jaccard similarity between different time-points in the same tissue (y-axis) versus the similarity between different tissues at the same time-point (x-axis). Error bars indicate the range between the first and third quartiles.
- b.** Transitions between chromatin states along midbrain developmental time-points. For clarity, only the genomic bins assigned TSS-related states (Tss, TssFlnk, and TssBiv) at one or more time-points are included.
- c.** Same as b but for genomic bins assigned enhancer-related states (Enh, EnhLo, EnhPois, and EnhPr) at one or more developmental time-points.

Visualization of the 66 epigenomes in two dimensions using the UMAP technique. (Left) UMAP was given the H3K27ac signals in the Enh genomic bins across the 66 epigenomes. There were 735,048 such genomic bins, which were assigned Enh in one or more epigenomes. (Right) UMAP was given the signals of all ten marks in the TssBiv genomic bins across the 66 epigenomes. There were 156,752 such genomic bins, which were assigned TssBiv in one or more epigenomes.

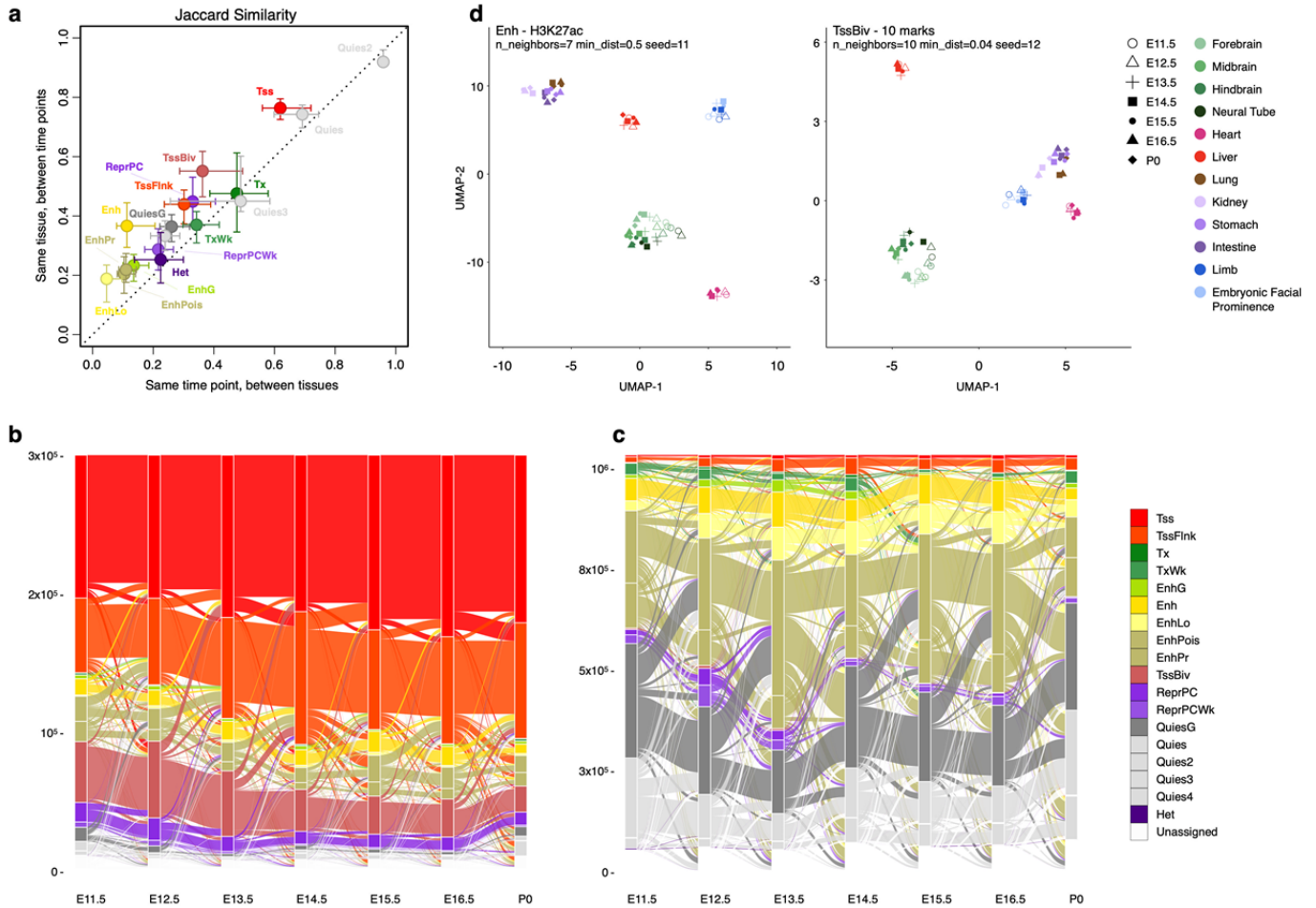


Figure 3: Count and expression of bivalent genes along developmental time-points.

a. The number of bivalent genes at 1 to 7 time-points in the midbrain. Observed and expected numbers of genes are in red and in gray respectively.

b. Median expression levels of three groups of genes: (green) bivalent at the earliest time-point but not at the last time-point, (blue) bivalent at the last time-point but not at the first time-point, and (pink) bivalent at all time-points.

c. Distribution of gene expression, with genes grouped by the total number of time-points at which their TSSs are in the bivalent state TssBiv (left) or in the active state Tss (right) in the forebrain. For all box plots, whiskers show 95% confidence intervals, boxes represent the first and third quartiles, and the vertical midline is the median, and outliers are omitted. The total number of genes in each group is shown below each box plot in parentheses. There is a negative correlation between expression and the duration of the bivalent state and a positive correlation between expression and the duration of the active state (P -values $< 2.2 \times 10^{-16}$).

Violin plots show the distributions of tissue-specificity scores for bivalent and non-bivalent genes that encode transcription factors (TFs) and non-TFs. Medians are shown in black bars with values indicated. P -values are shown for three comparisons as indicated.

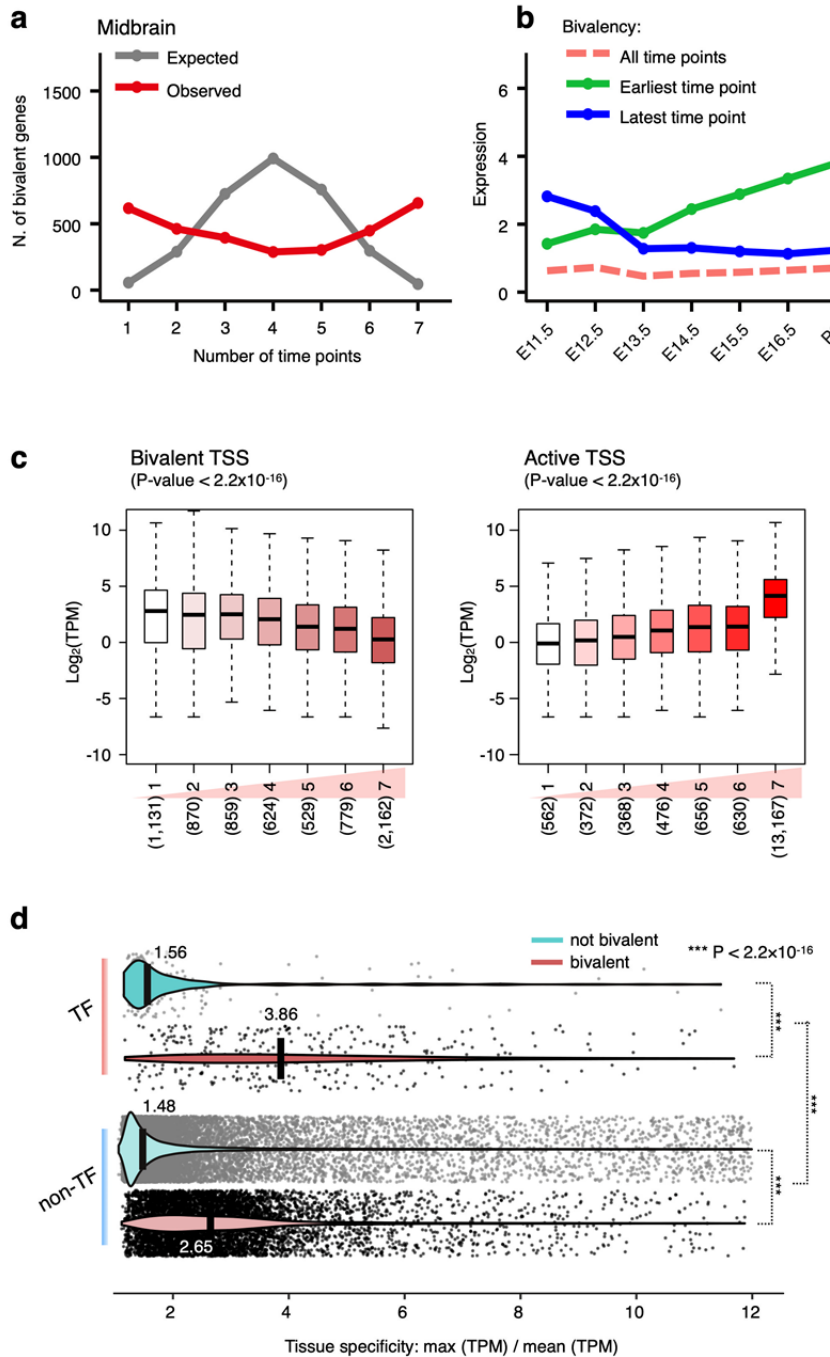


Figure 4: Expression profiles and chromatin states for the transcription factors with the highest tissue-specificity scores.

a. Hierarchical clustering of expression profiles for the TFs with tissue specificity scores greater than 6, with 75 TFs in total. Rows on the top show the maximal expression level across all biosamples (intensities of red), bivalency status (brown for 62 bivalent TFs, and yellow for 13 non-bivalent TFs), and tissue specificity score (intensities of green).

b-i. Example TFs and the chromatin state assignments near their loci. Among these, Gata1 (**d**) is a non-bivalent TF and the rest are bivalent TFs: Arx (**b**), En2 (**c**), Gata4 (**e**), Wt1 (**f**), Foxq1 (**g**), Evx2 (**h**), and Alx1 (**i**). Each gene name is near the 5'-end of the gene, and CpG islands are indicated as green boxes beneath each gene. Chromatin states are colored as in Fig. 1.

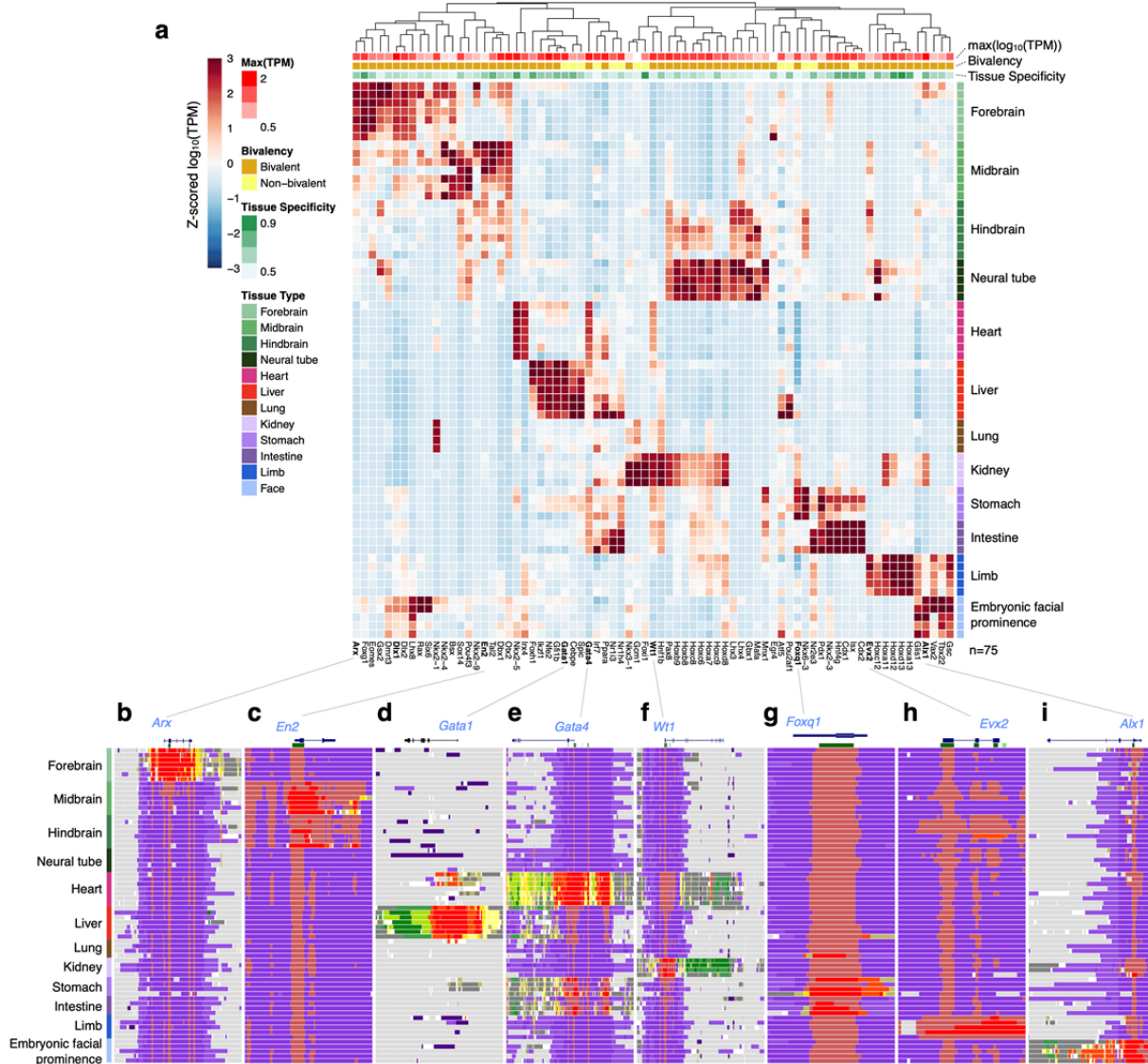


Figure 5: Evolutionary conservation of genomic regions by chromatin state.

a. The PhyloP conservation score (phyloP60way for mm10) for genomic regions assigned to each chromatin state. Colors correspond to tissues.

b. PhyloP score for genomic bins assigned to Enh in all 12 tissues.

c. Percentage of bins assigned to Enh that overlap with transposons, for all 12 tissues.

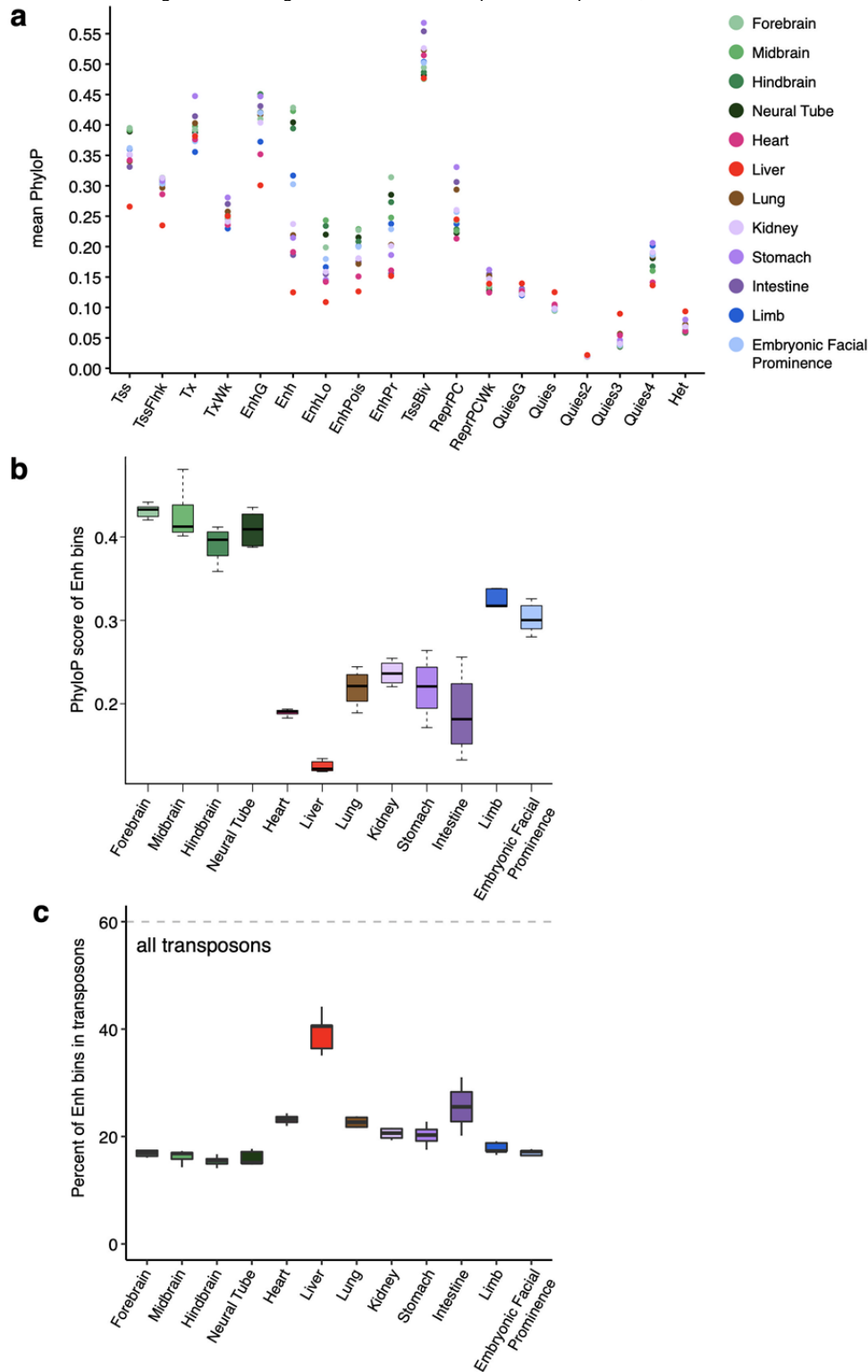
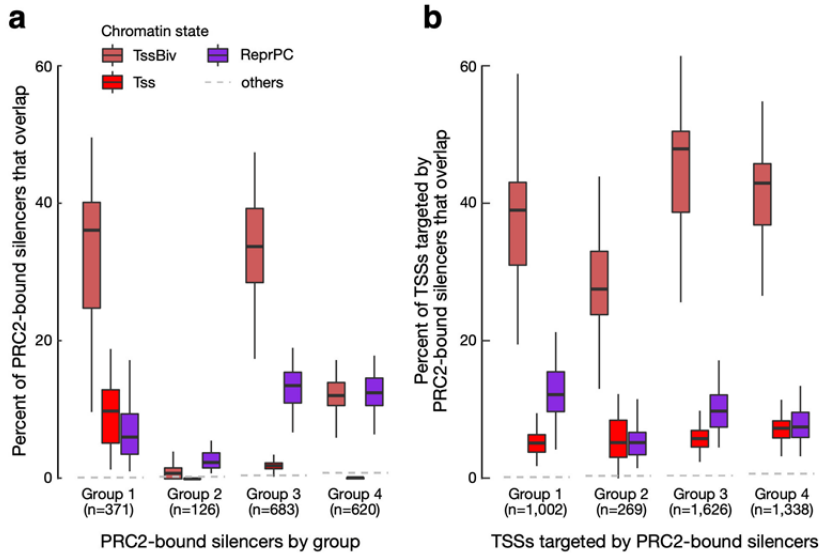
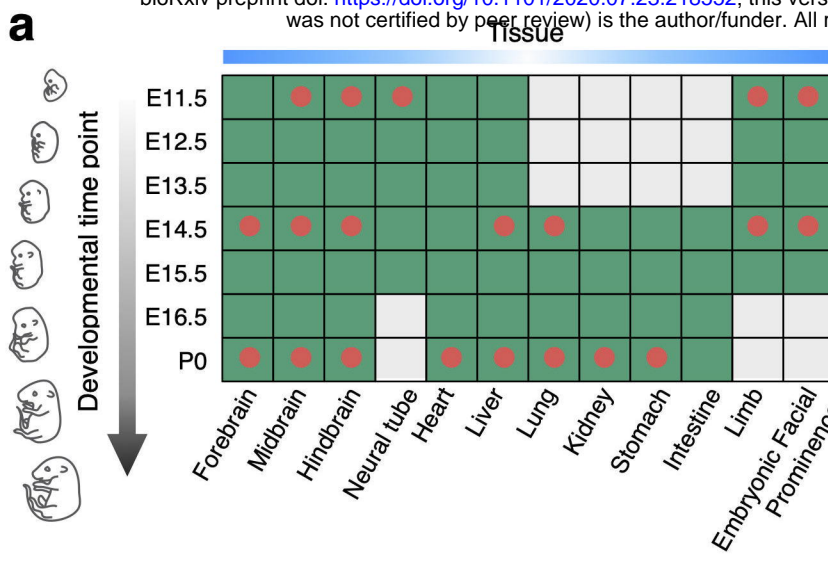


Figure 6: PRC2-bound silencers and their target TSSs are enriched in the TssBiv and ReprPC states.

a. Percentage of PRC2-bound silencers whose centers overlap a genomic bin assigned to the TssBiv, Tss, ReprPC, or other chromatin states. Silencers were divided into four groups by Ngan et al.²⁶ according to H3K27ac signals in mouse fetal tissue biosamples. To normalize for the differential genomic coverage of the chromatin states, the same numbers of genomic bins were randomly drawn in the other states to match the number of genomic bins in TssBiv in each biosample. States are colored as in Fig. 1b and the average of the other 15 states is shown as a gray dashed line.

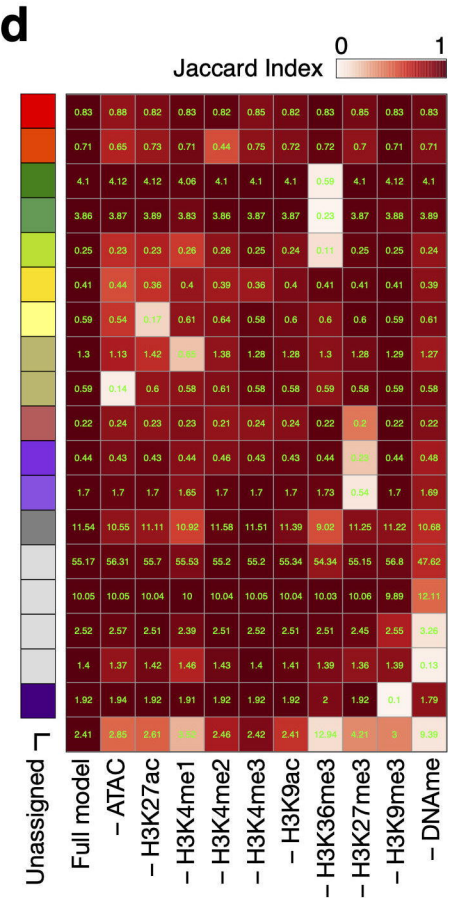
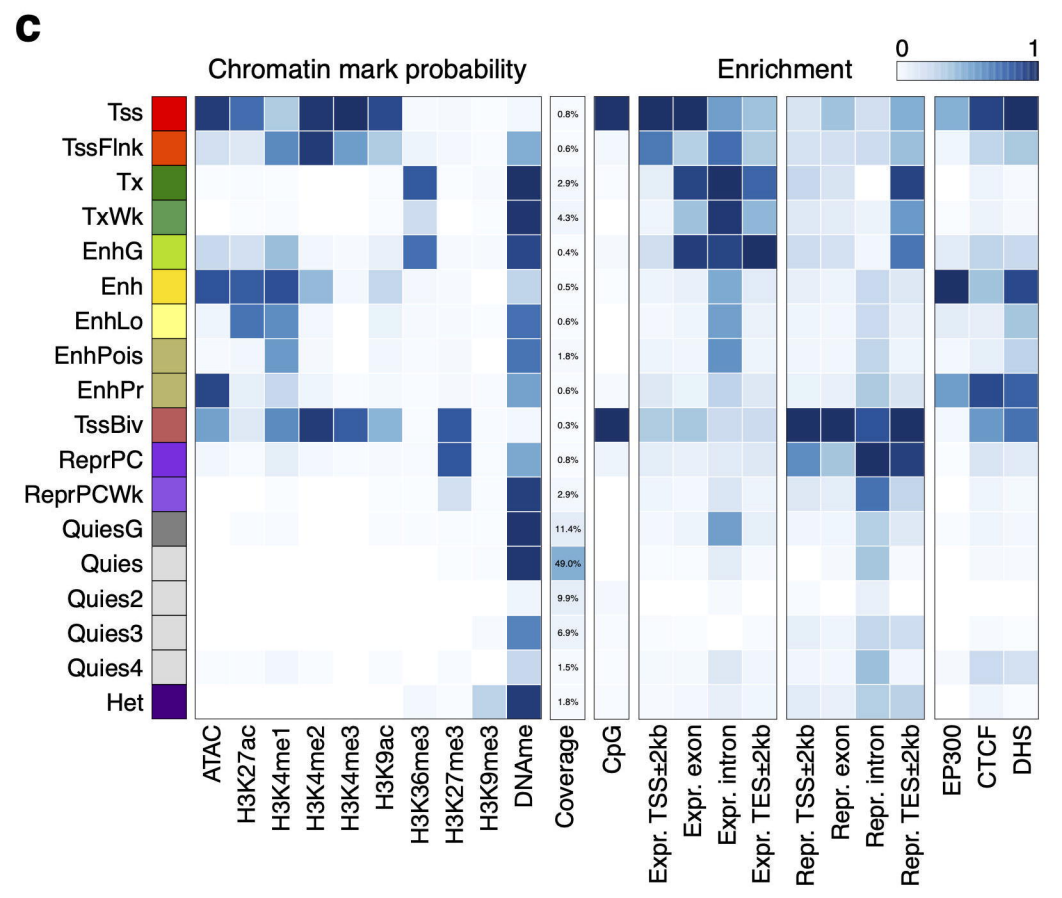
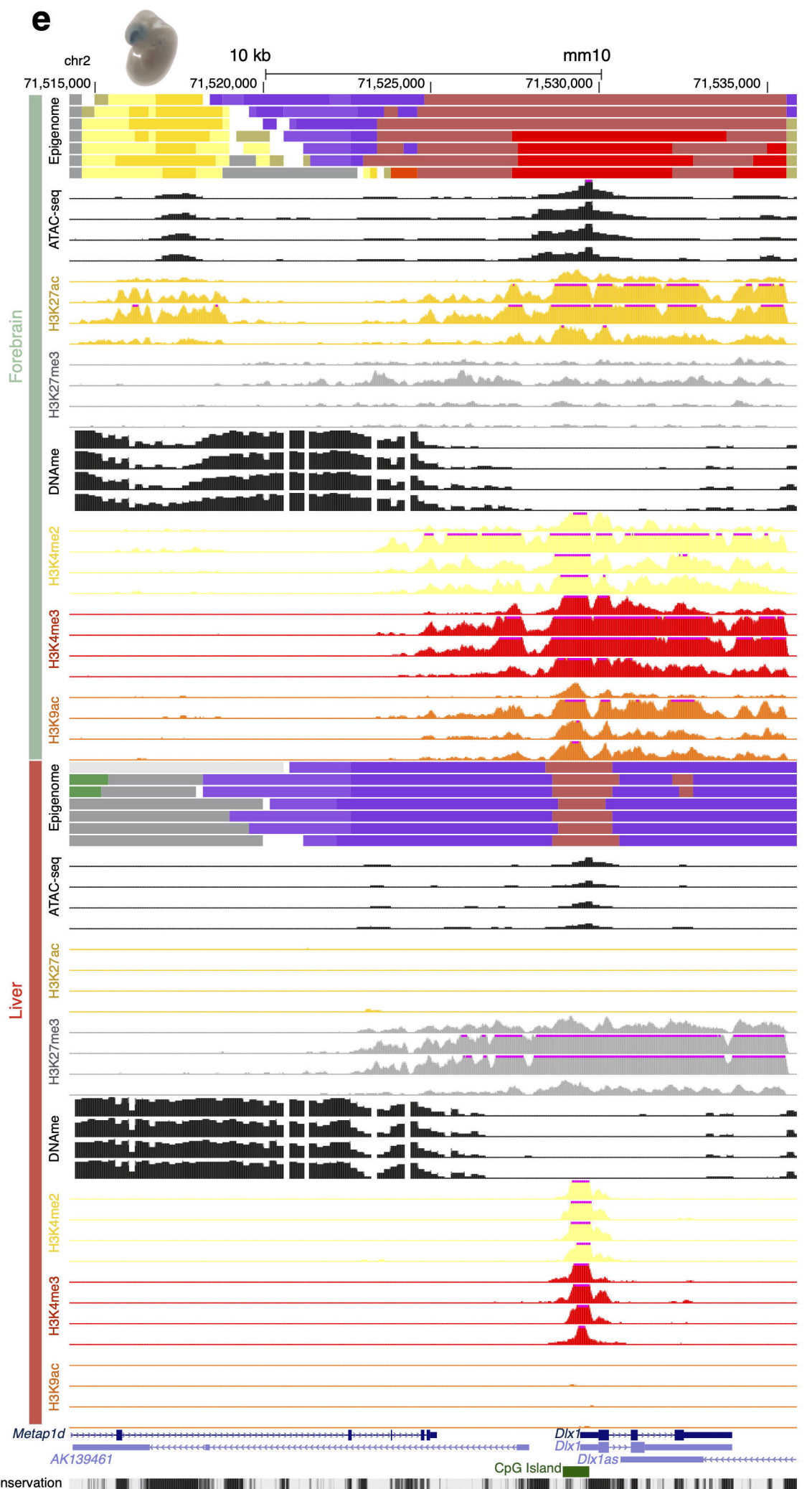
b. Same as **a** but for the TSSs targeted by the PRC2-bound silencers defined by Ngan et al.²⁶.





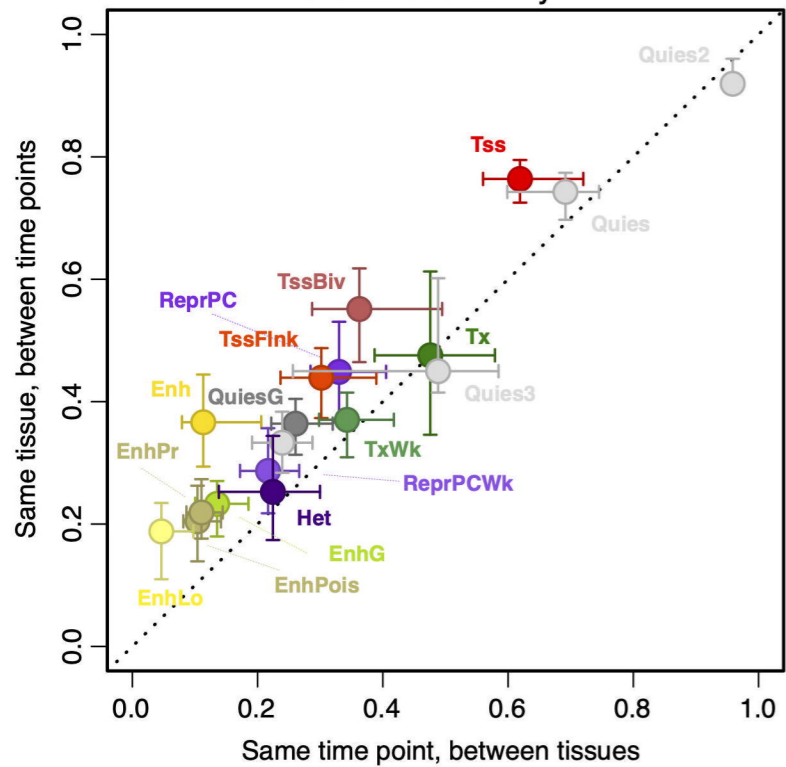
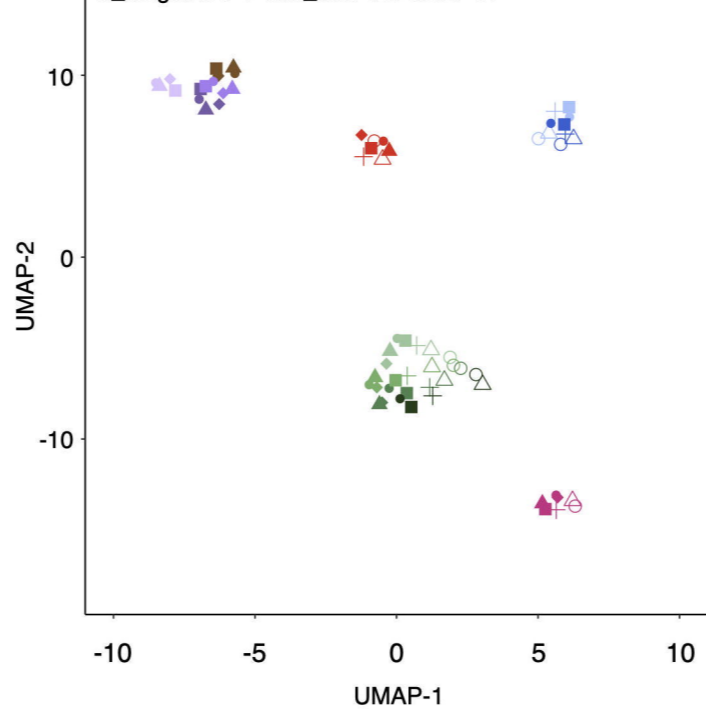
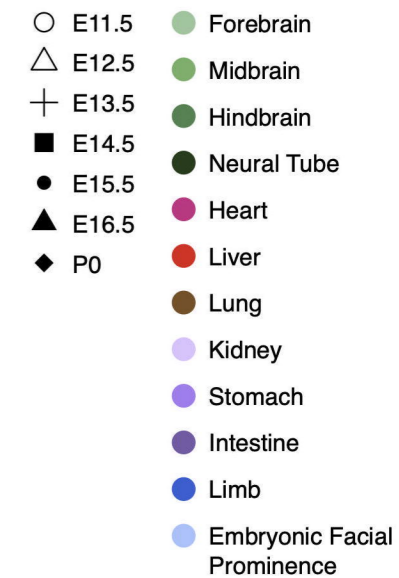
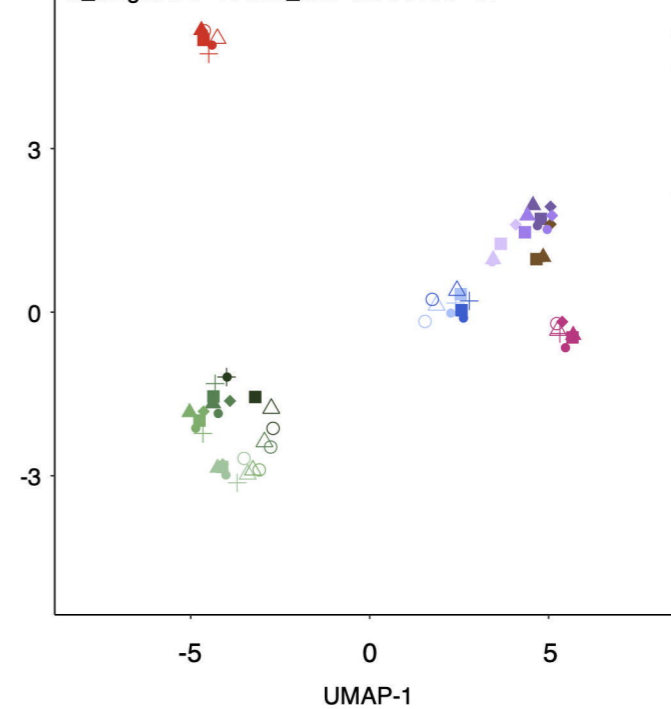
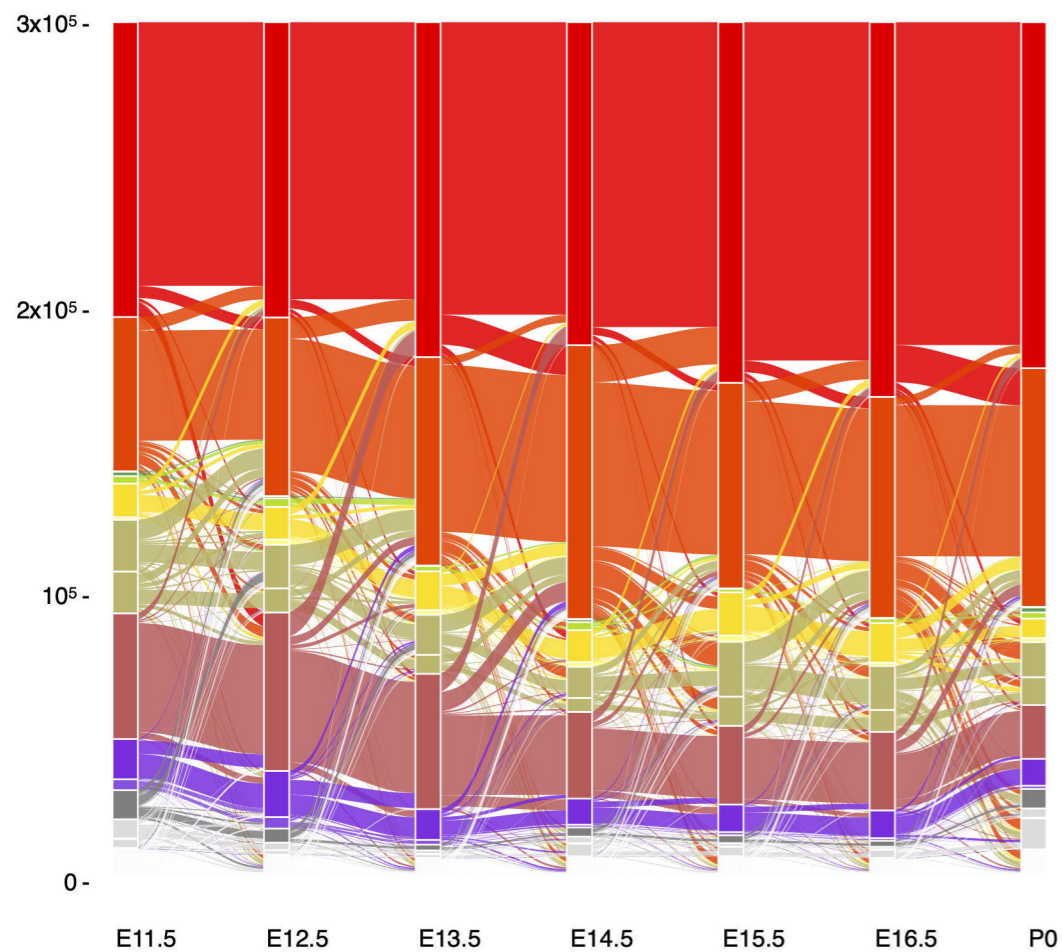
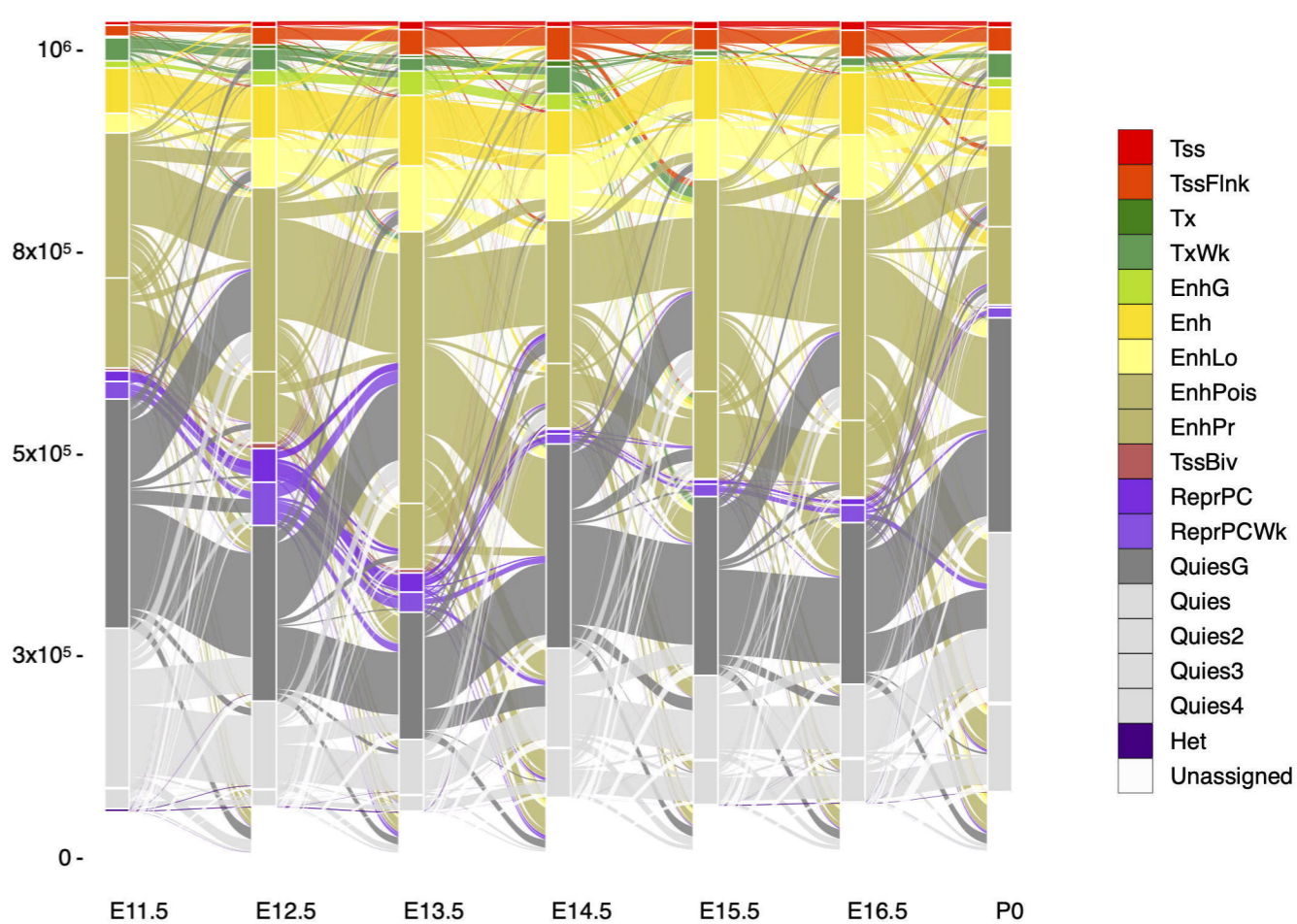
b

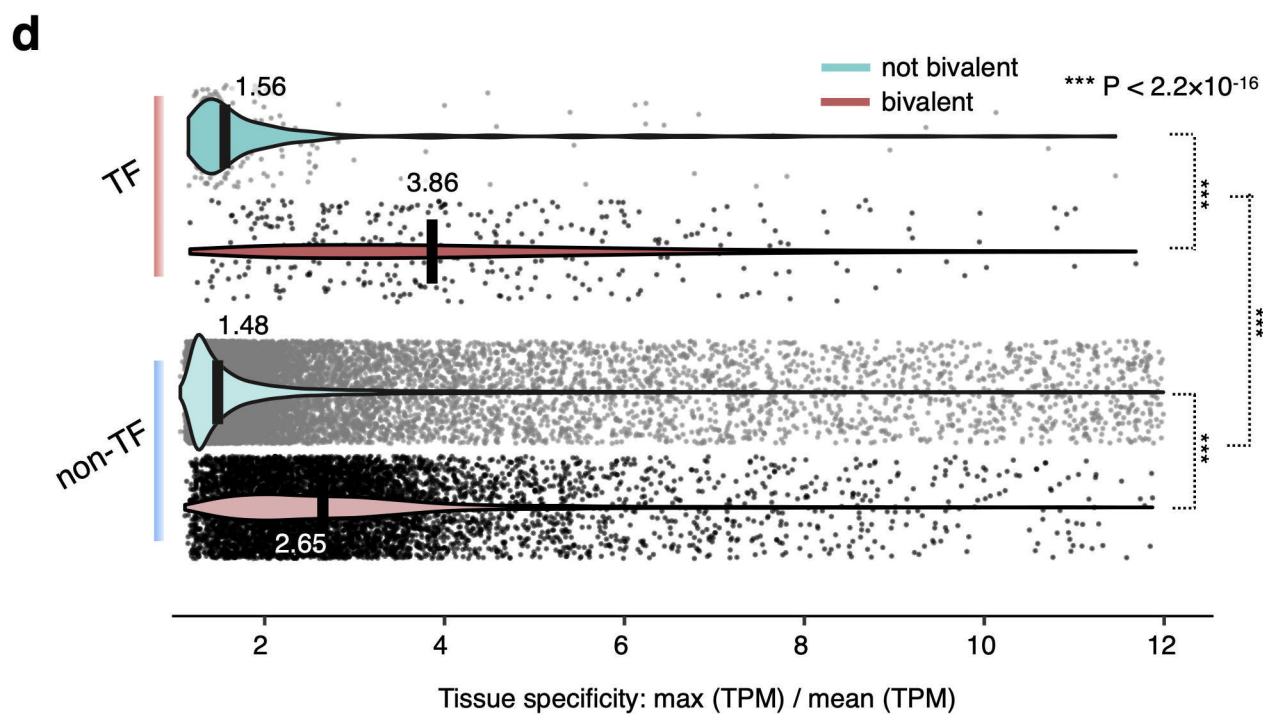
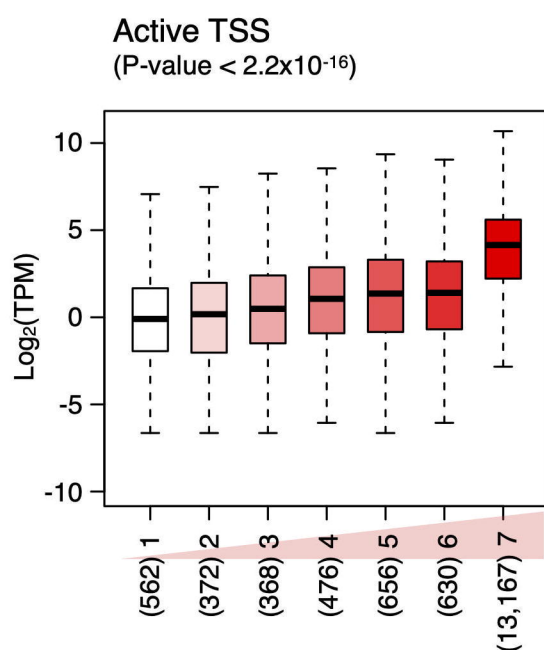
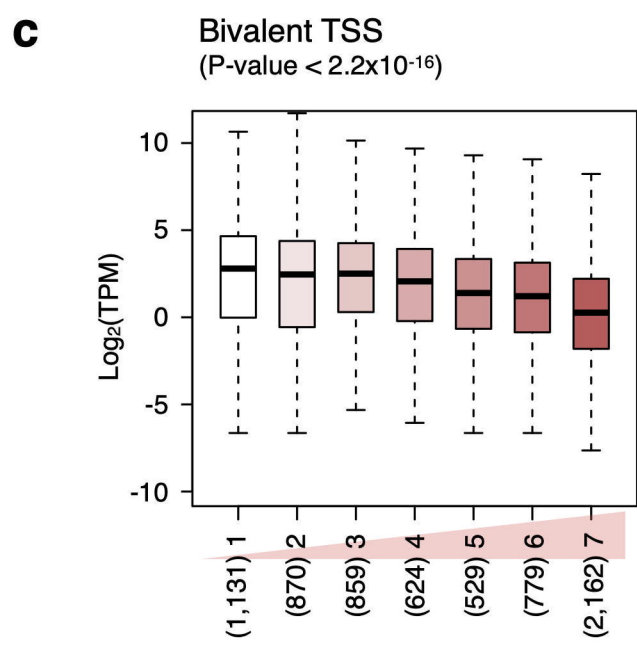
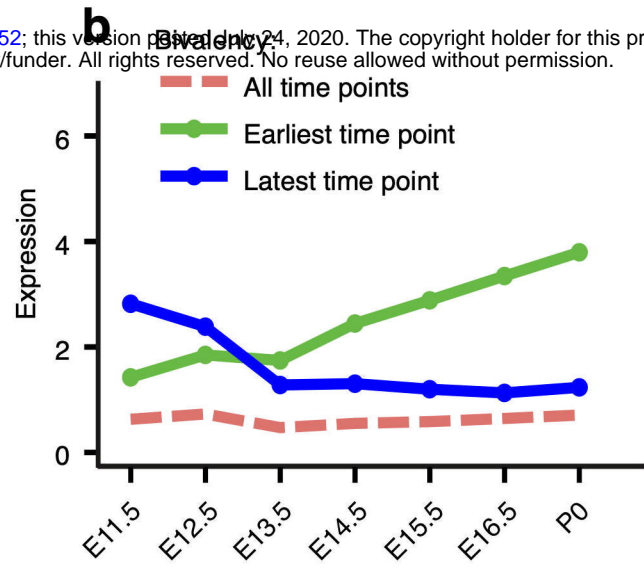
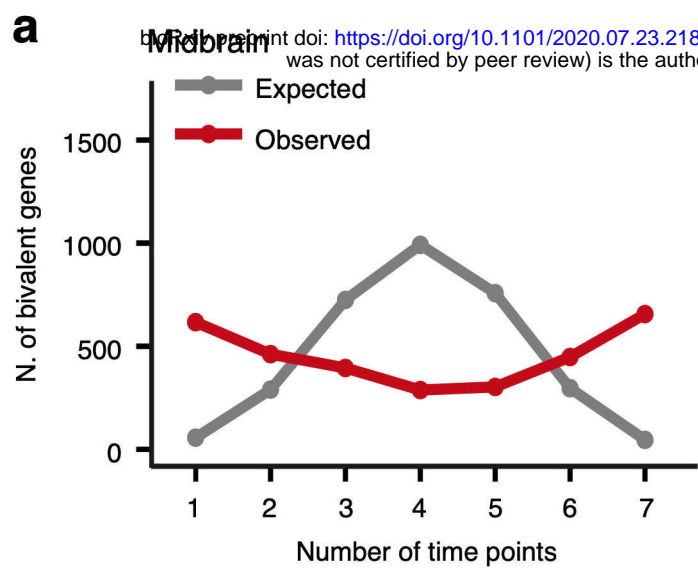
1	Tss	Active TSS
2	TssFlnk	Flanking active TSS
3	Tx	Transcription
4	TxWk	Weak transcription
5	EnhG	Enhancer in gene
6	Enh	Enhancer
7	EnhLo	Weak enhancer
8	EnhPois	Poised enhancer
9	EnhPr	Primed enhancer
10	TssBiv	Bivalent TSS
11	ReprPC	Repressed by Polycomb
12	ReprPCWk	Repressed by Polycomb (weak)
13	QuiesG	Quiescent gene
14	Quies	Quiescent
15	Quies2	Quiescent
16	Quies3	Quiescent
17	Quies4	Quiescent
18	Het	Heterochromatin

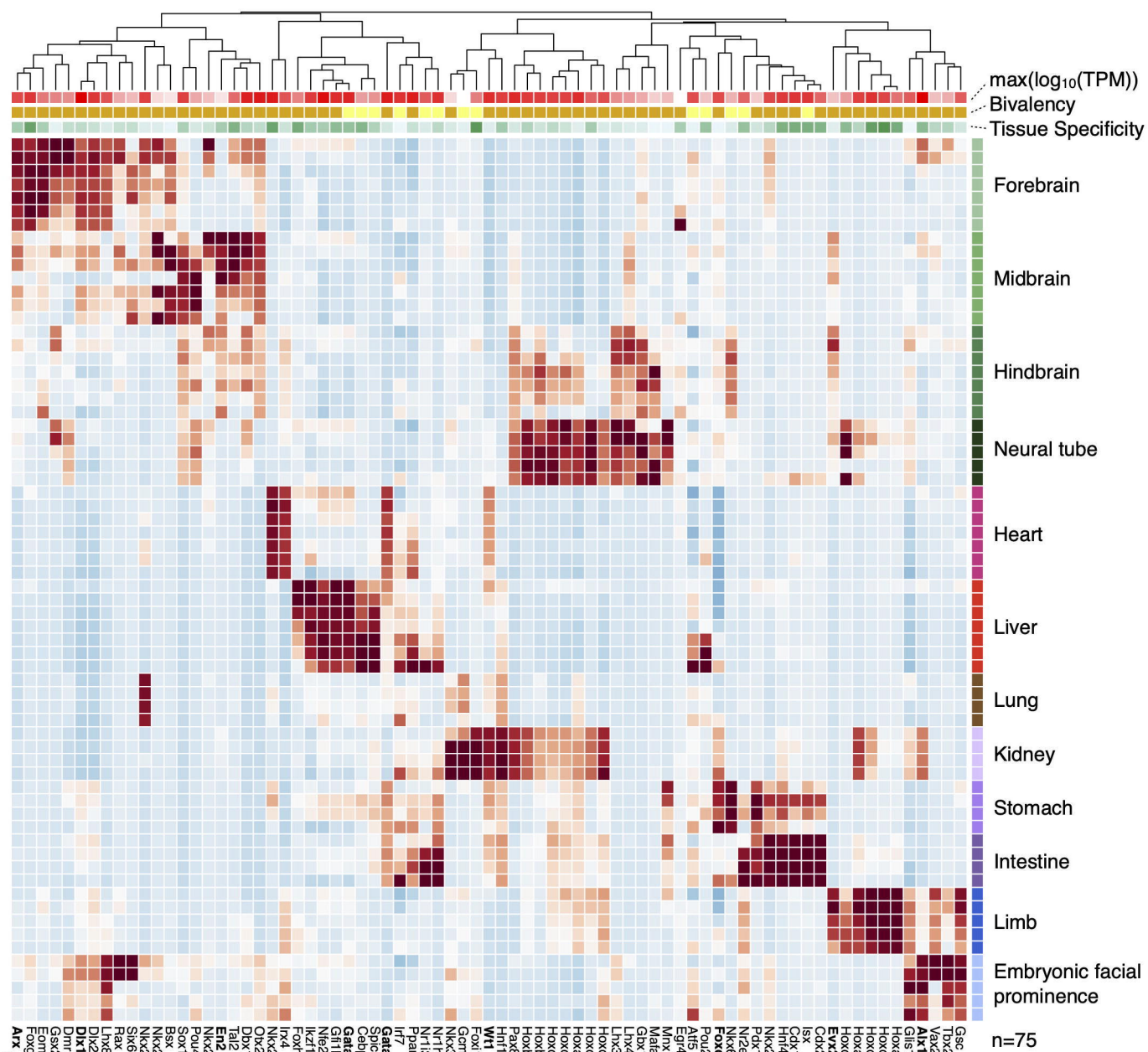
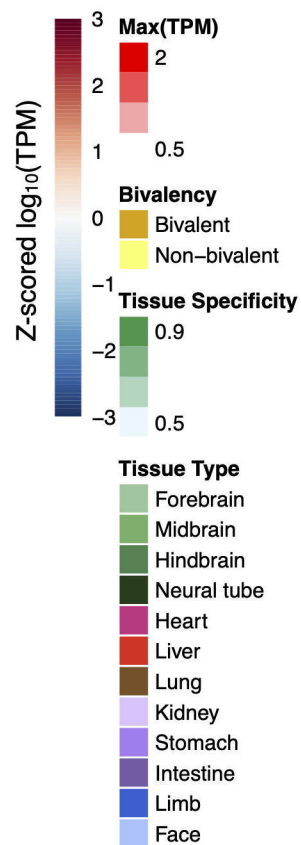
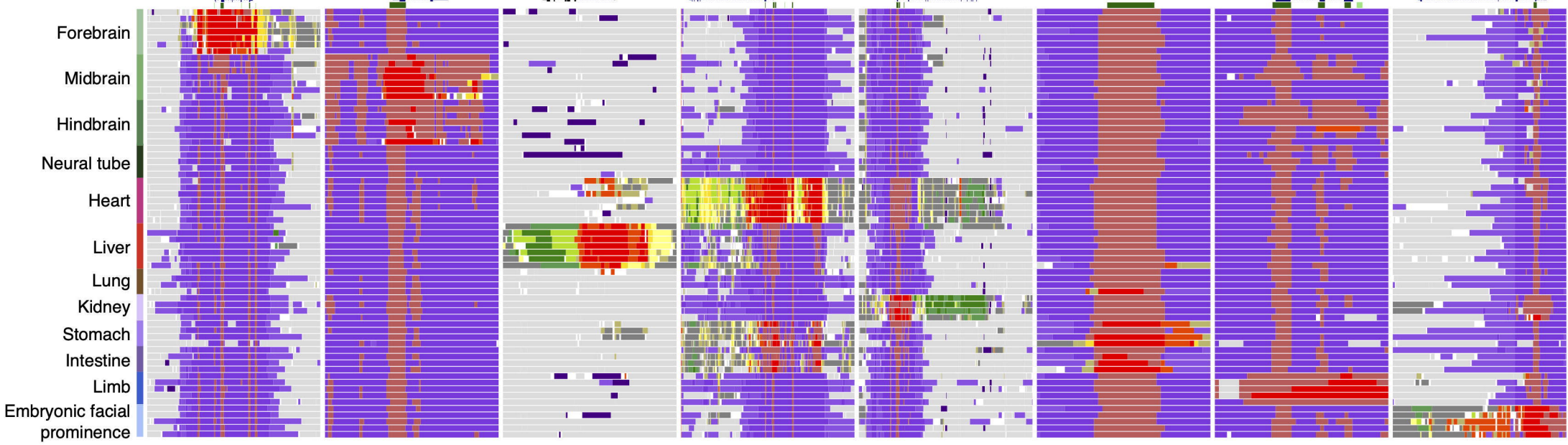


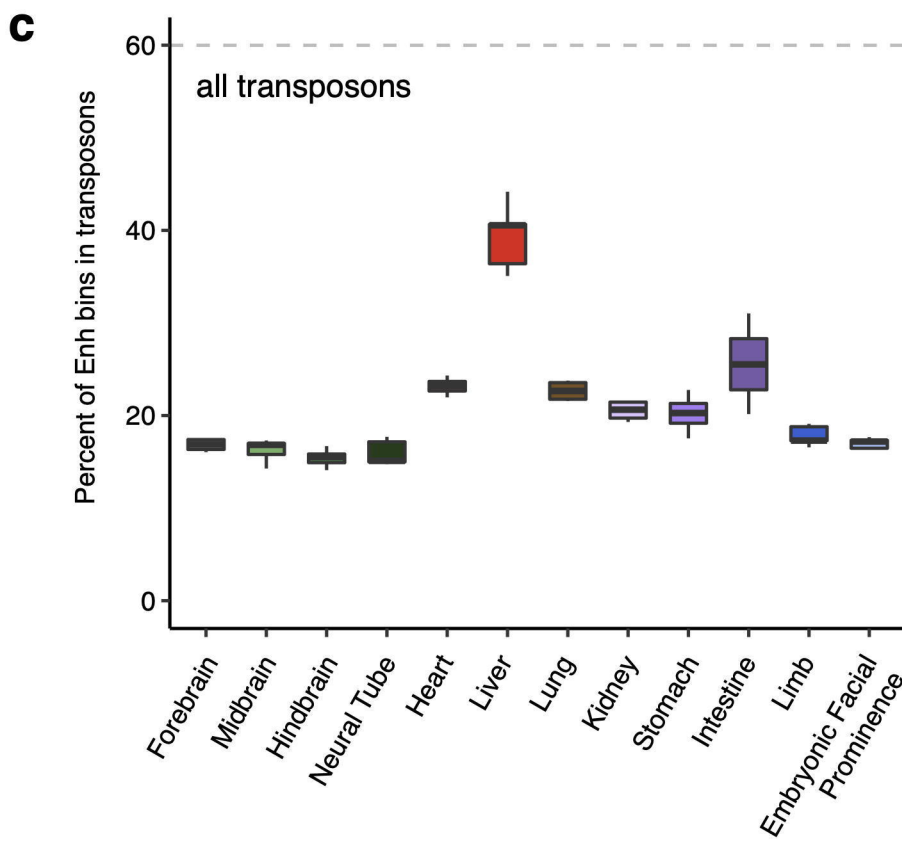
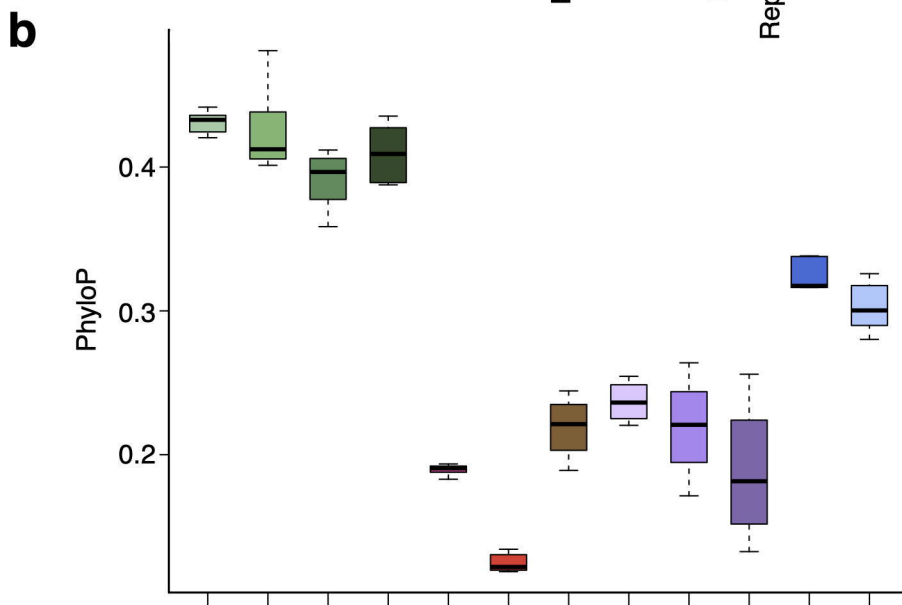
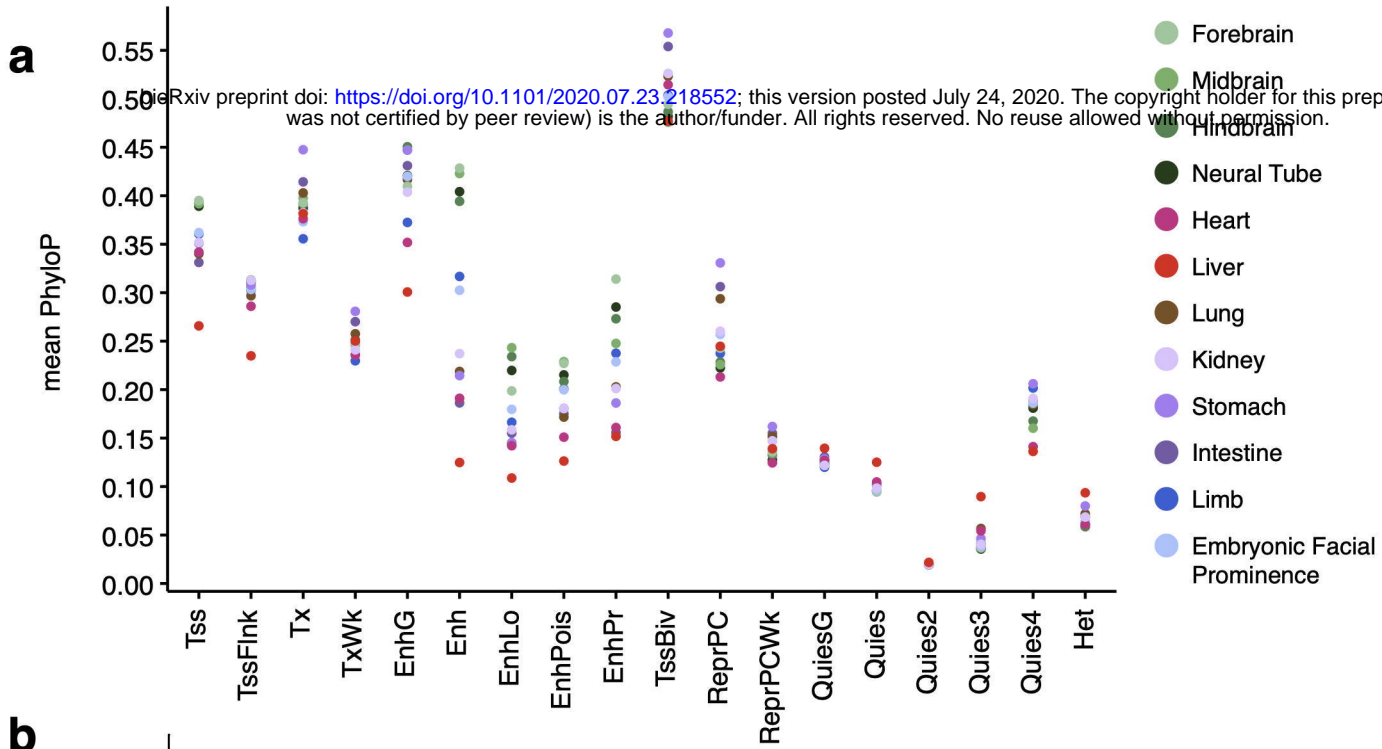
a

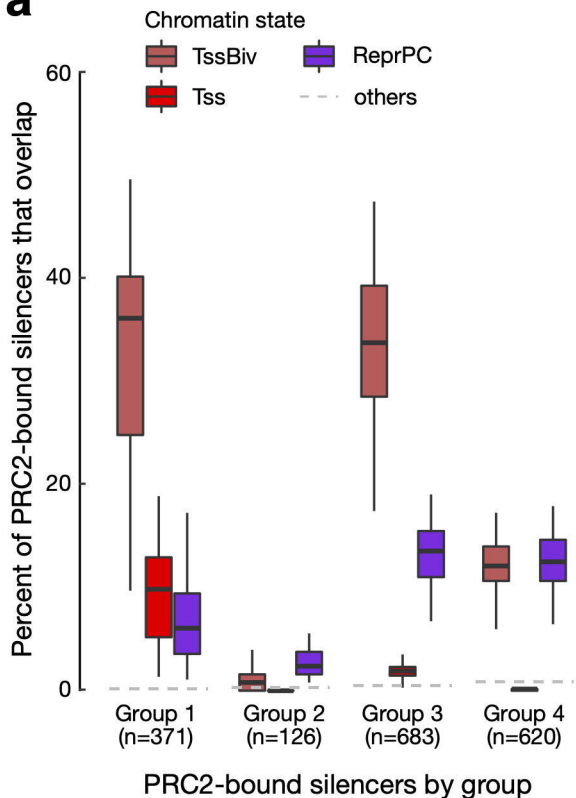
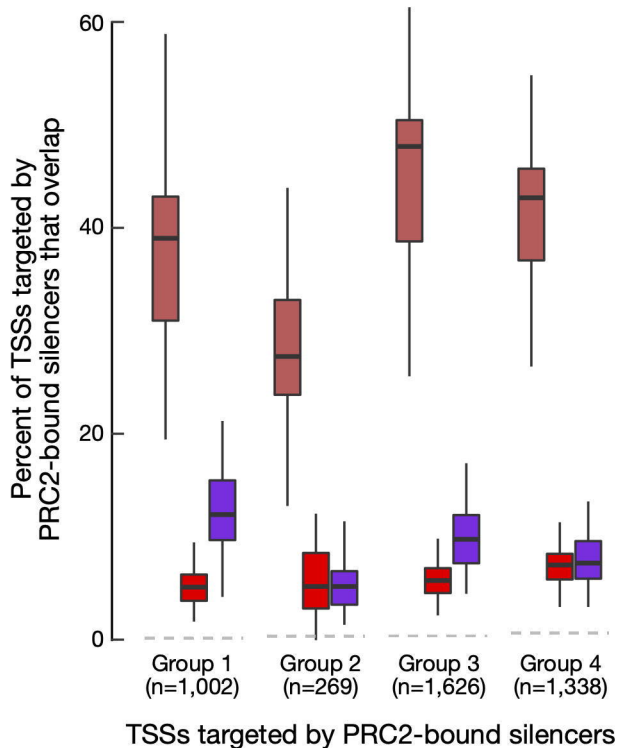
Jaccard Similarity

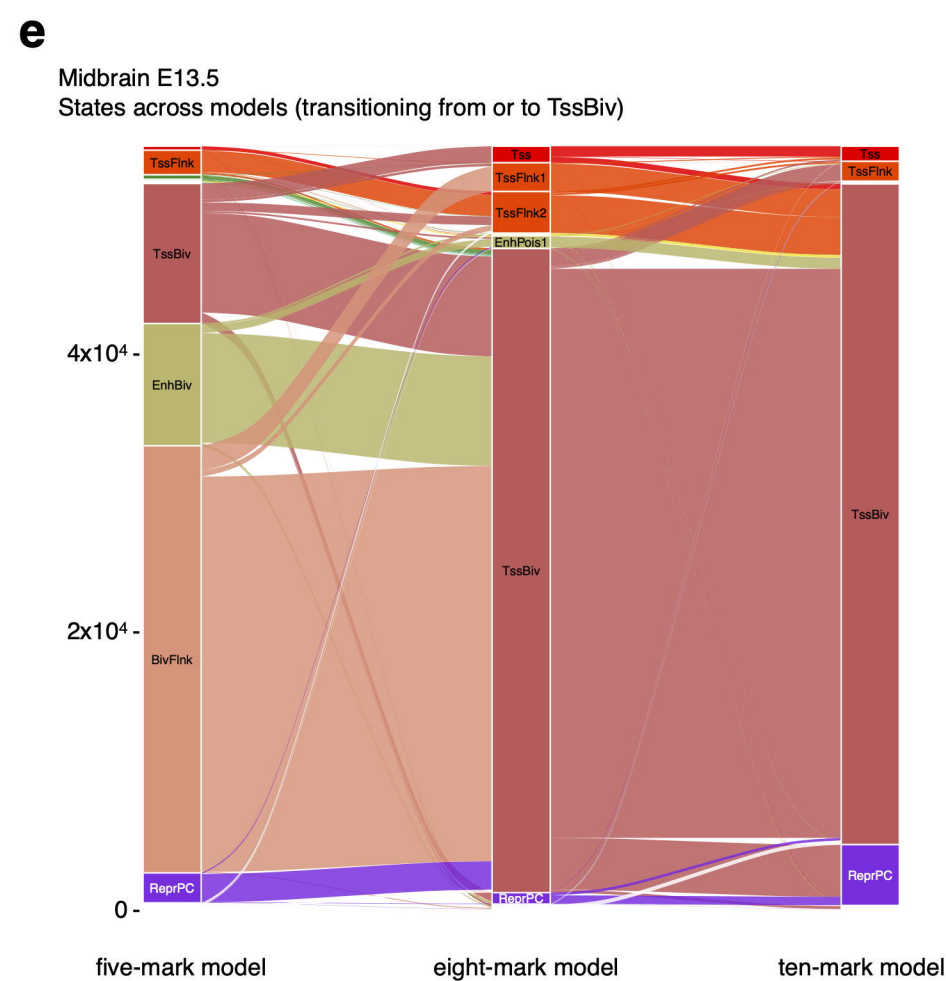
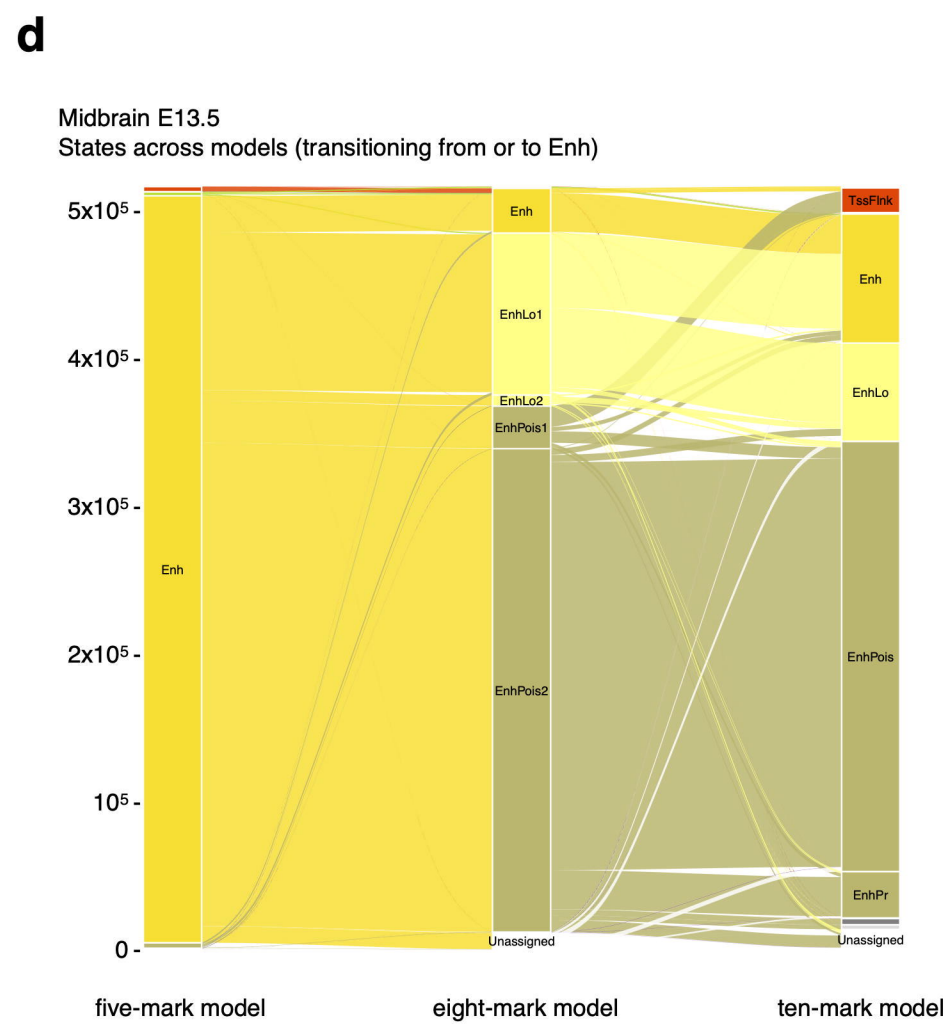
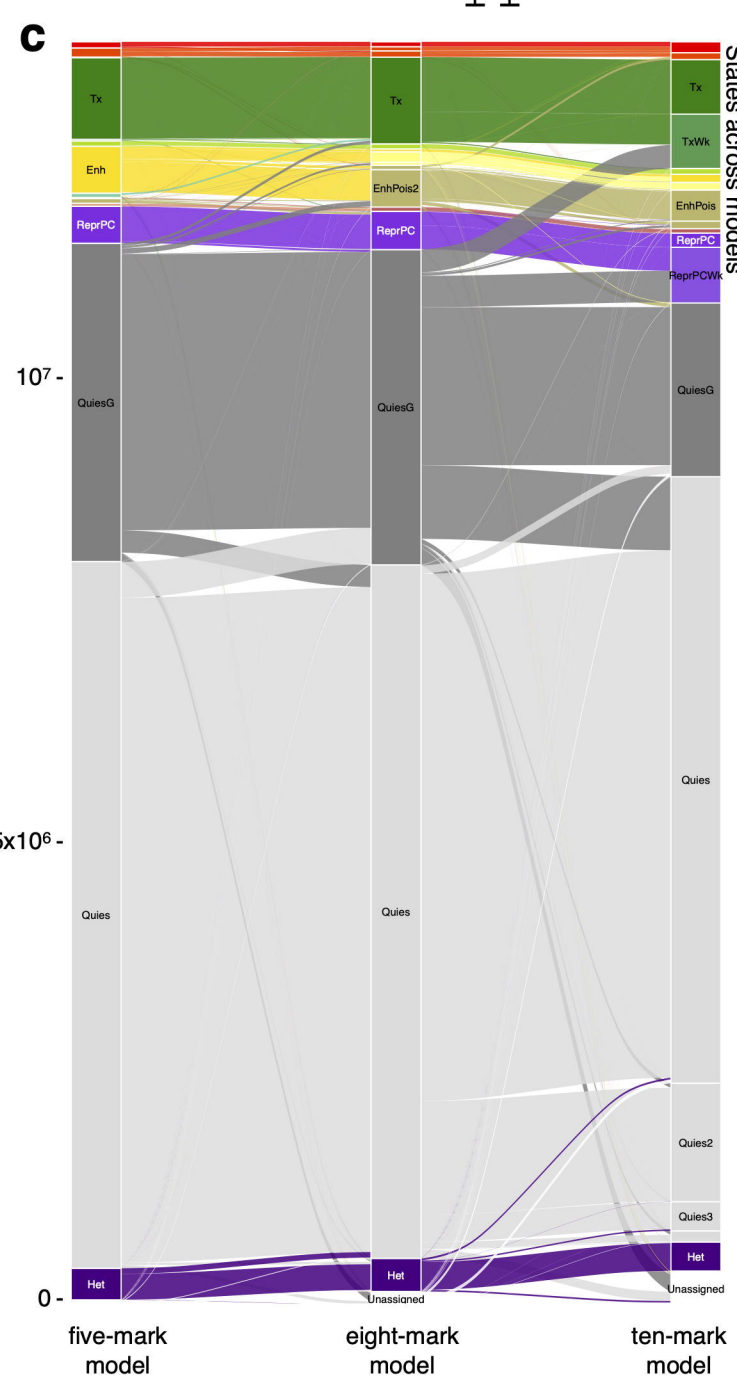
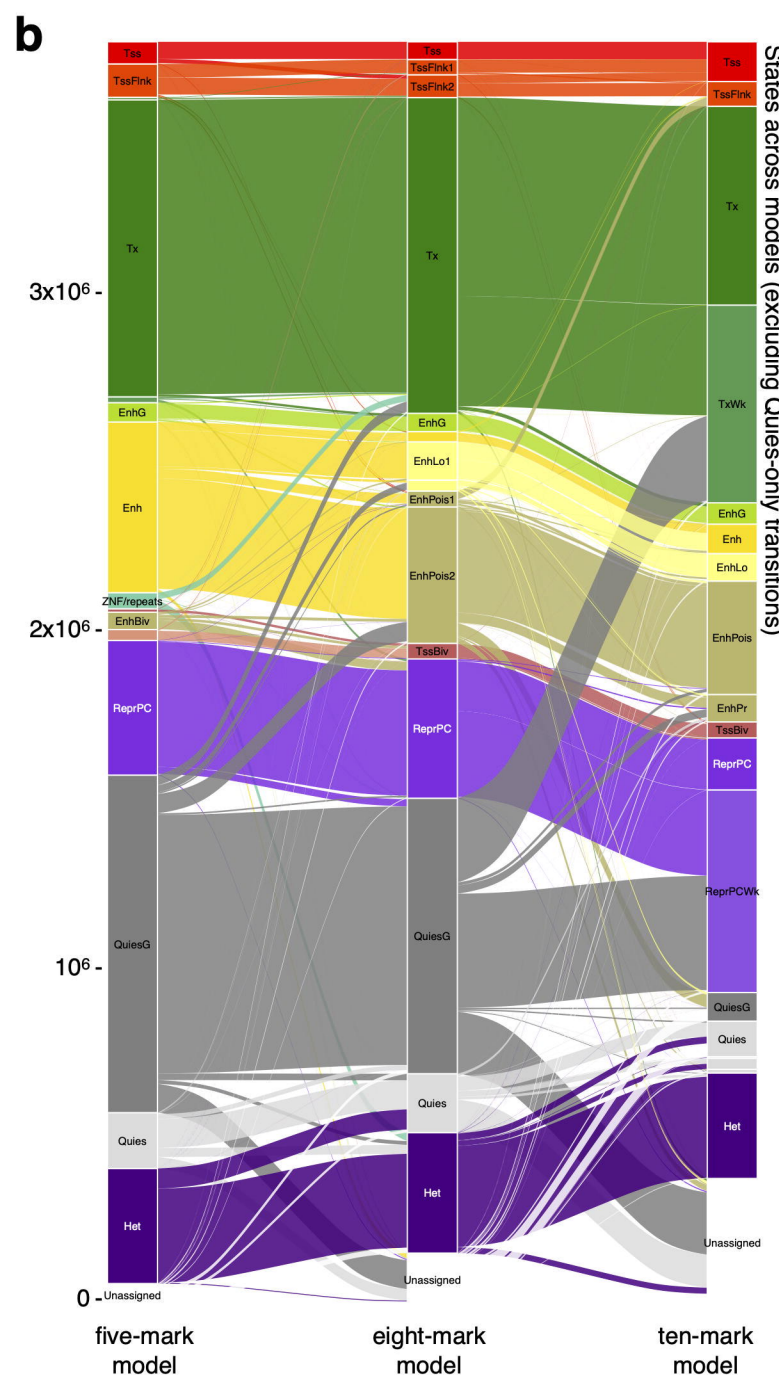
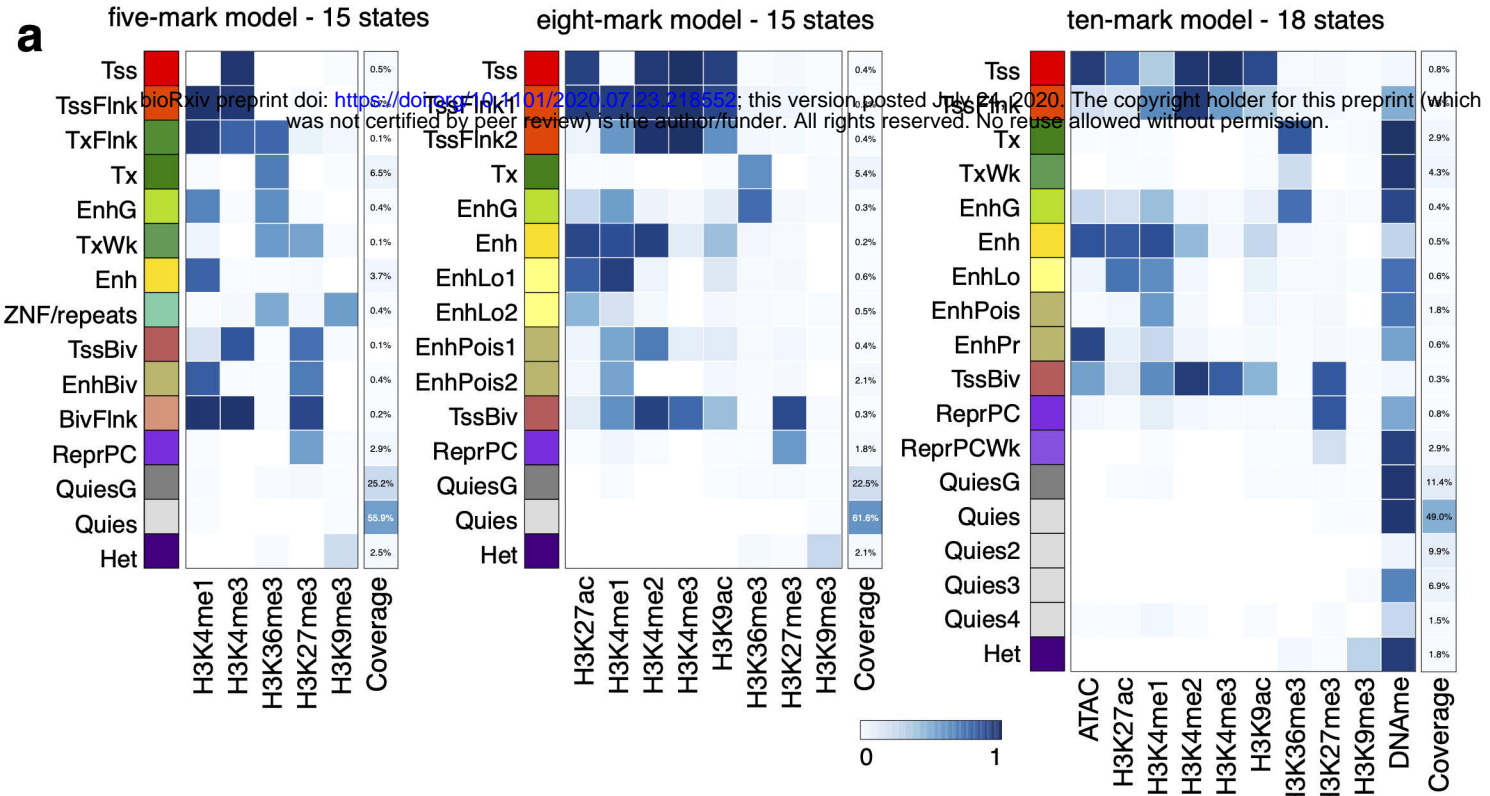
**d**Enh - H3K27ac
n_neighbors=7 min_dist=0.5 seed=11TssBiv - 10 marks
n_neighbors=10 min_dist=0.04 seed=12**b****c**



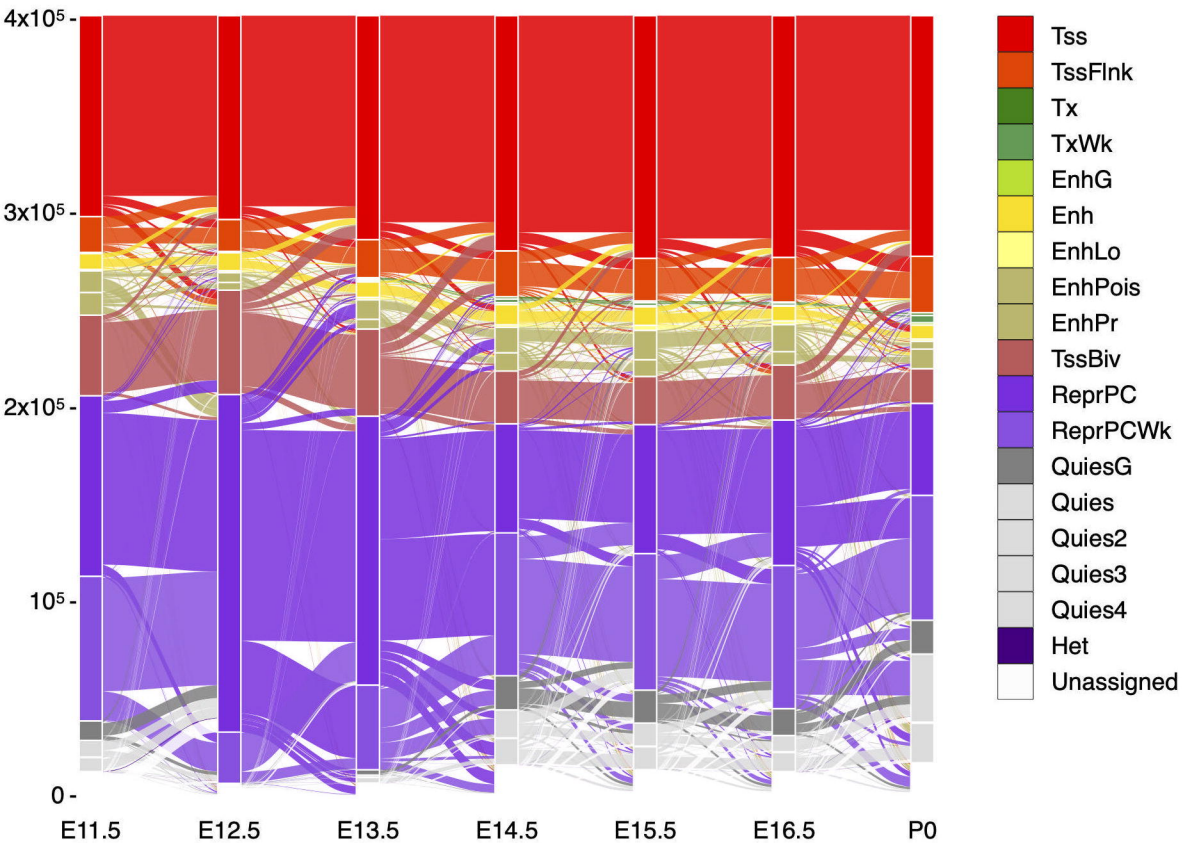
a**b***Arx***c***En2***d***Gata1***e***Gata4***f***Wt1***g***Foxq1***h***Evx2***i***Alx1*



a**b**

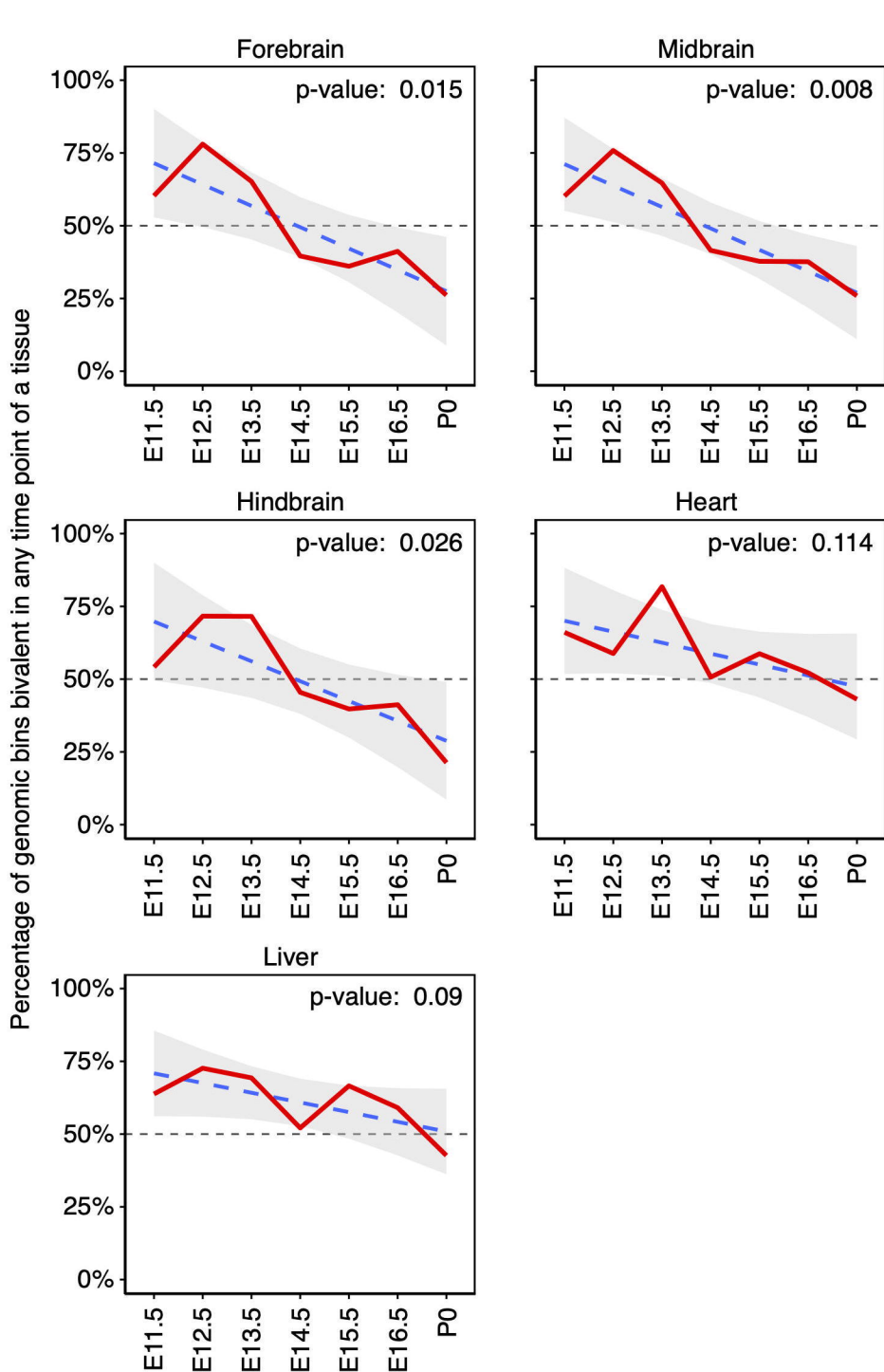
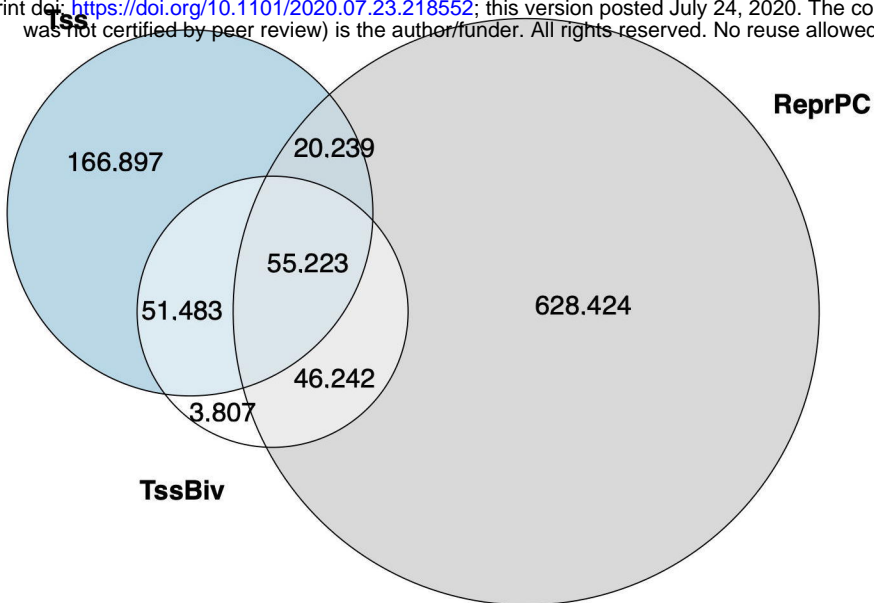


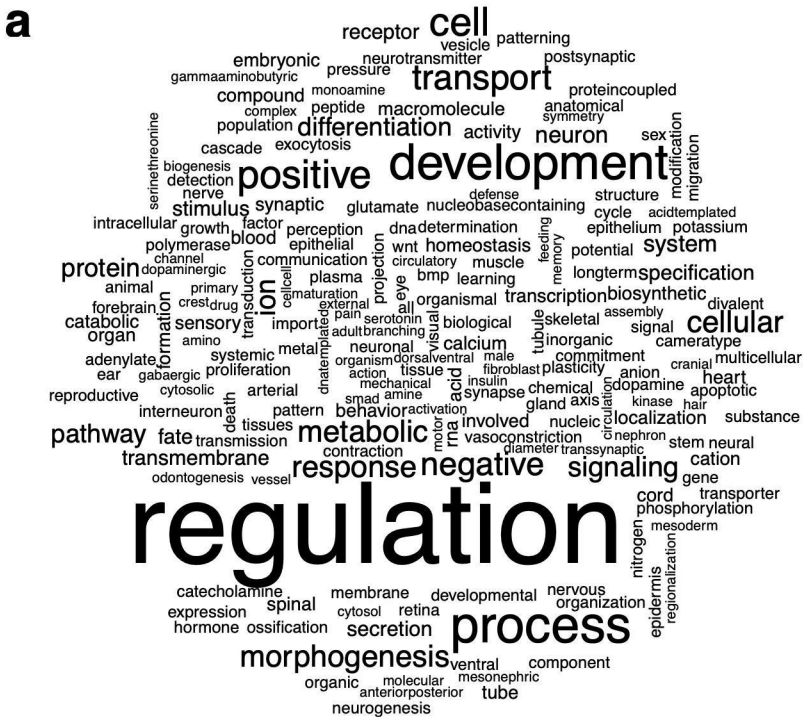
Forebrain
States across time points (transitioning from or to Tss, TssBiv, ReprPC)



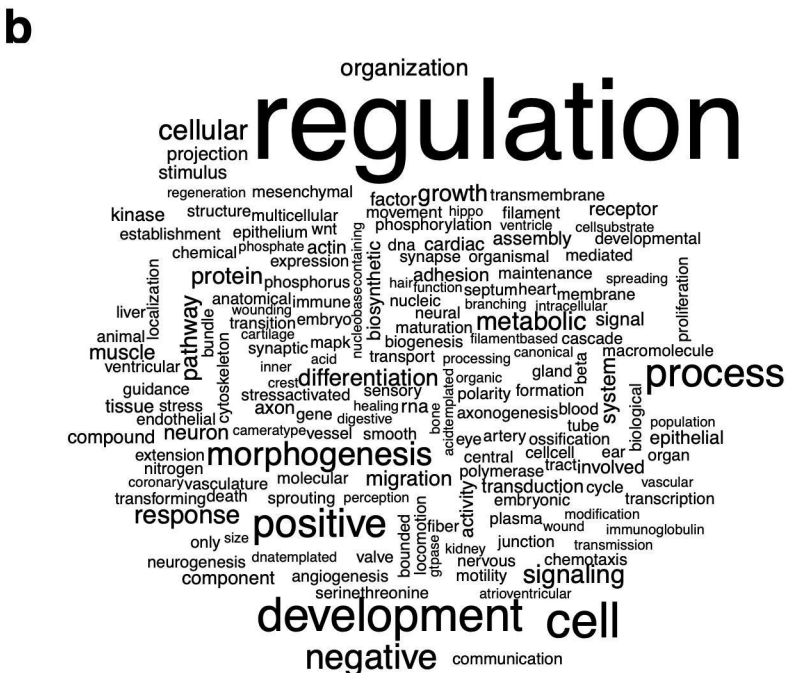
Total # genomic bins: 13,627,678

bioRxiv preprint doi: <https://doi.org/10.1101/2020.07.23.218552>; this version posted July 24, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

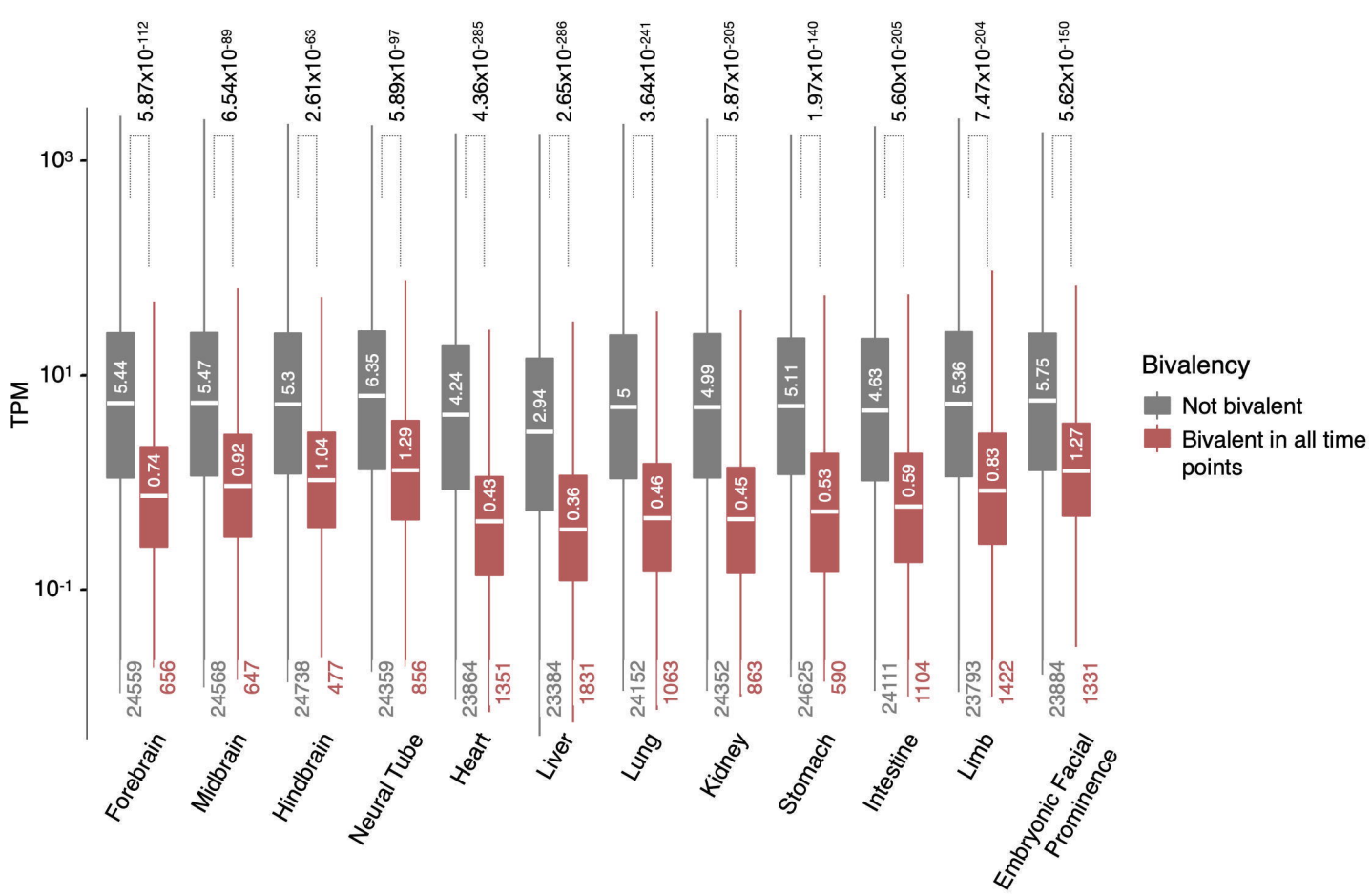


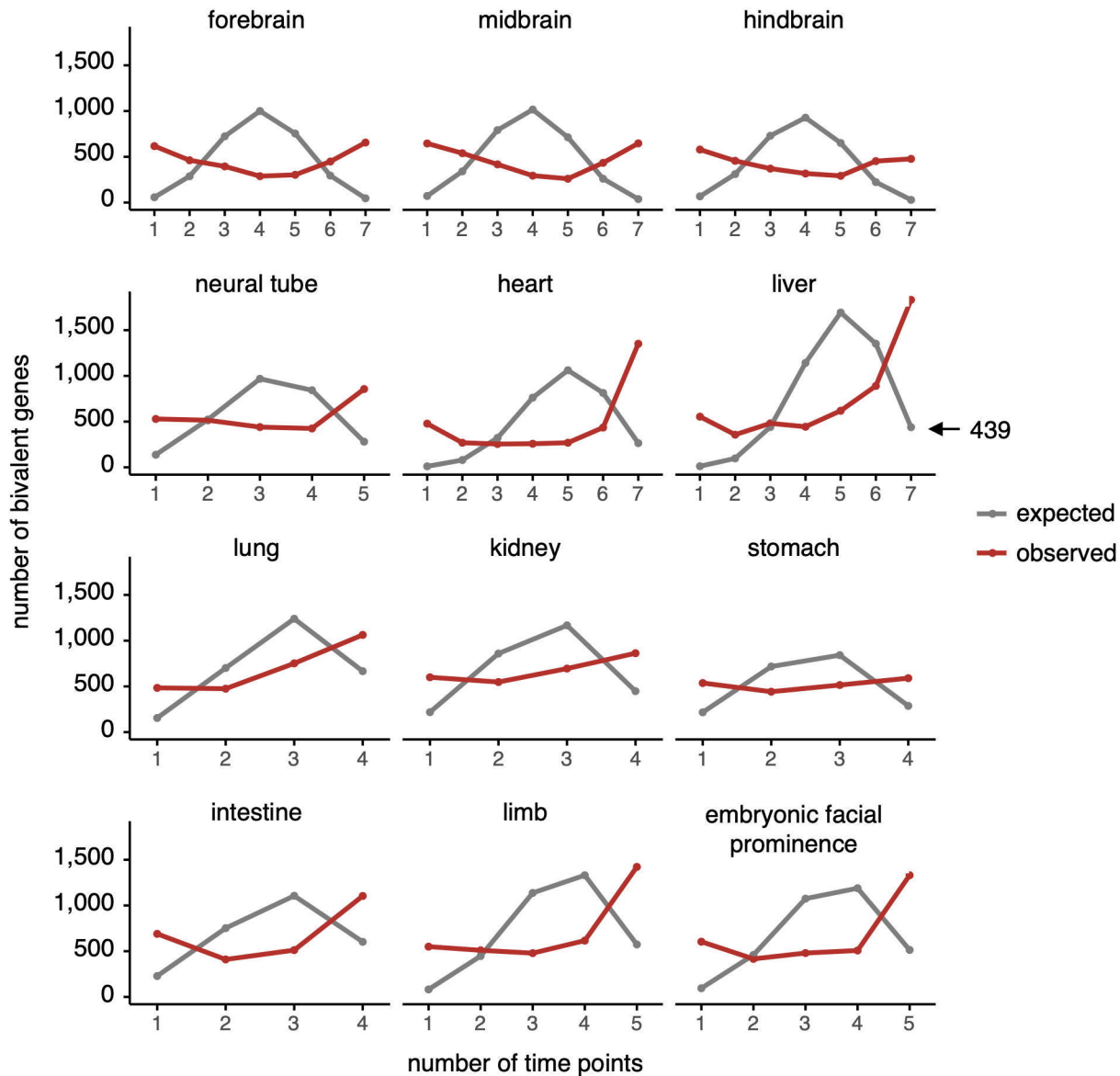


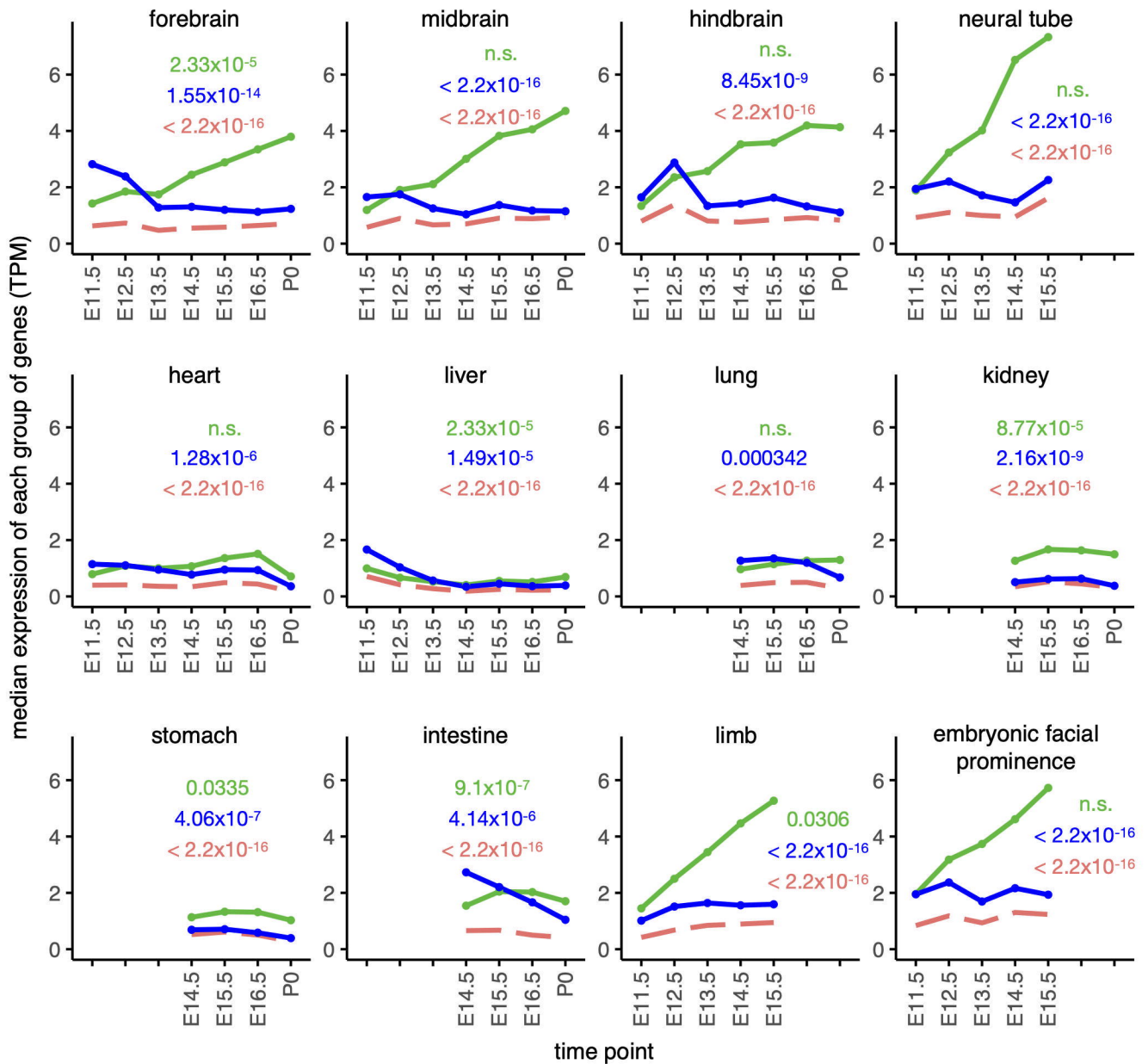
Bivalent TSS in all tissues

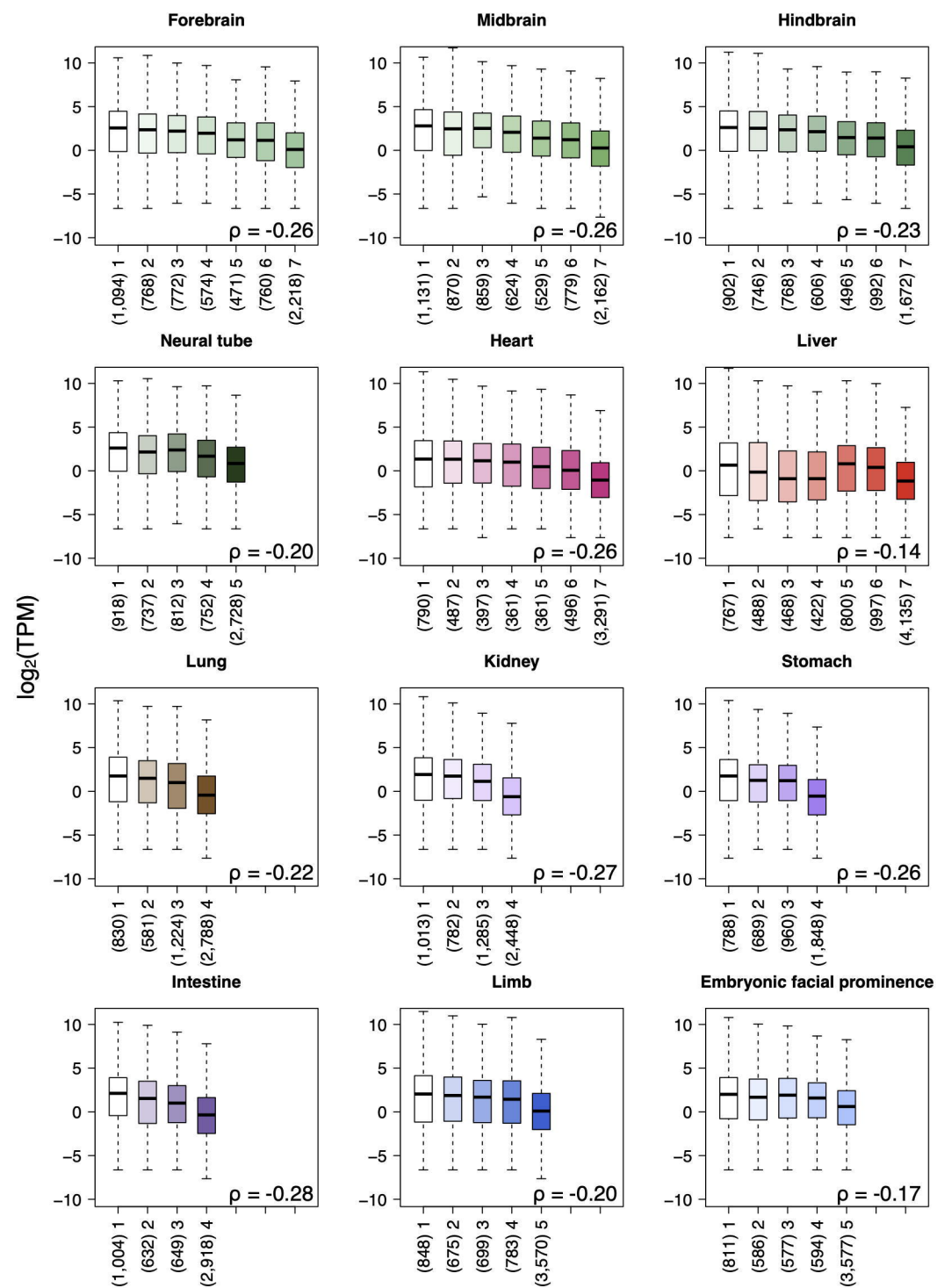
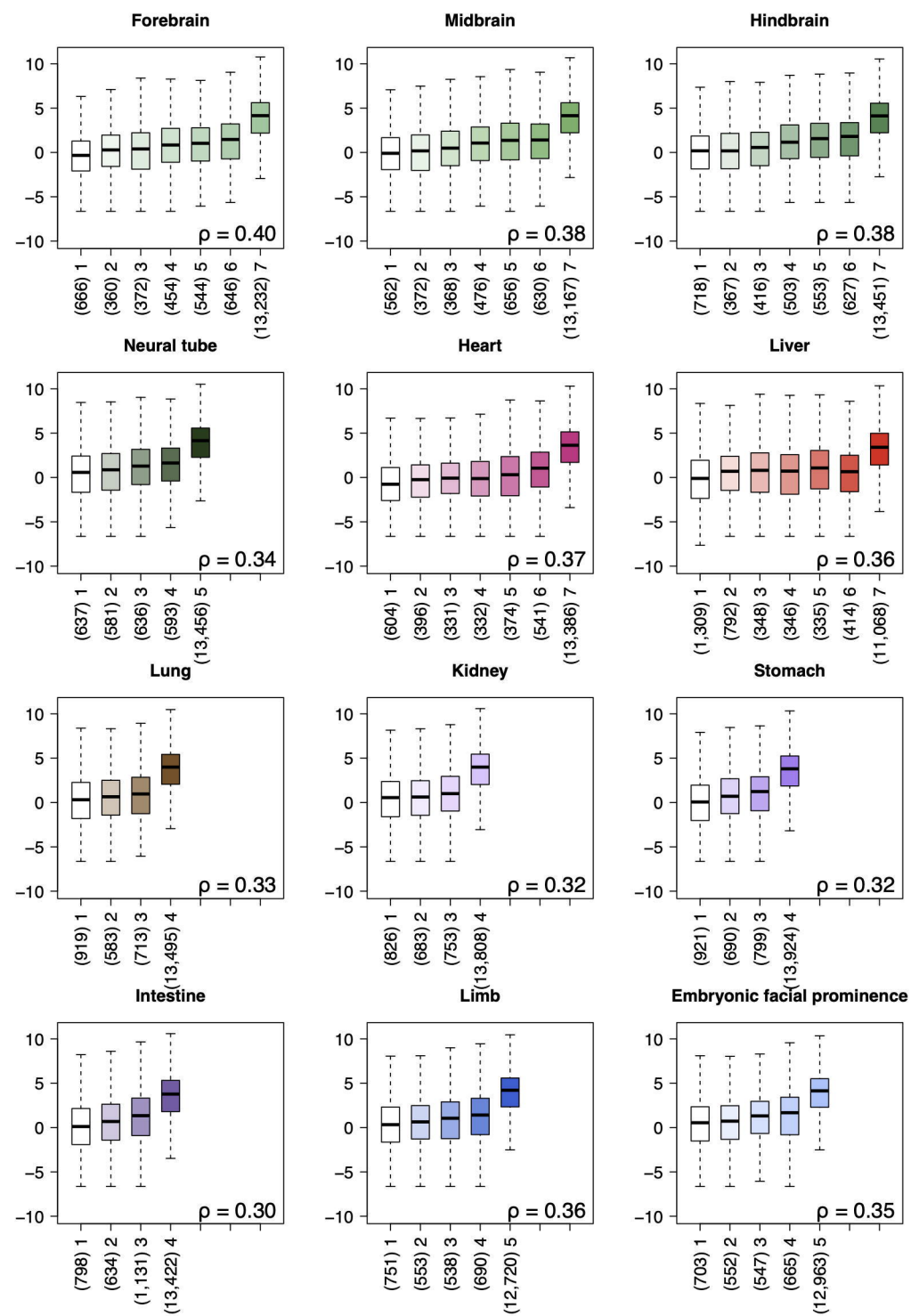


Bivalent TSS only in liver

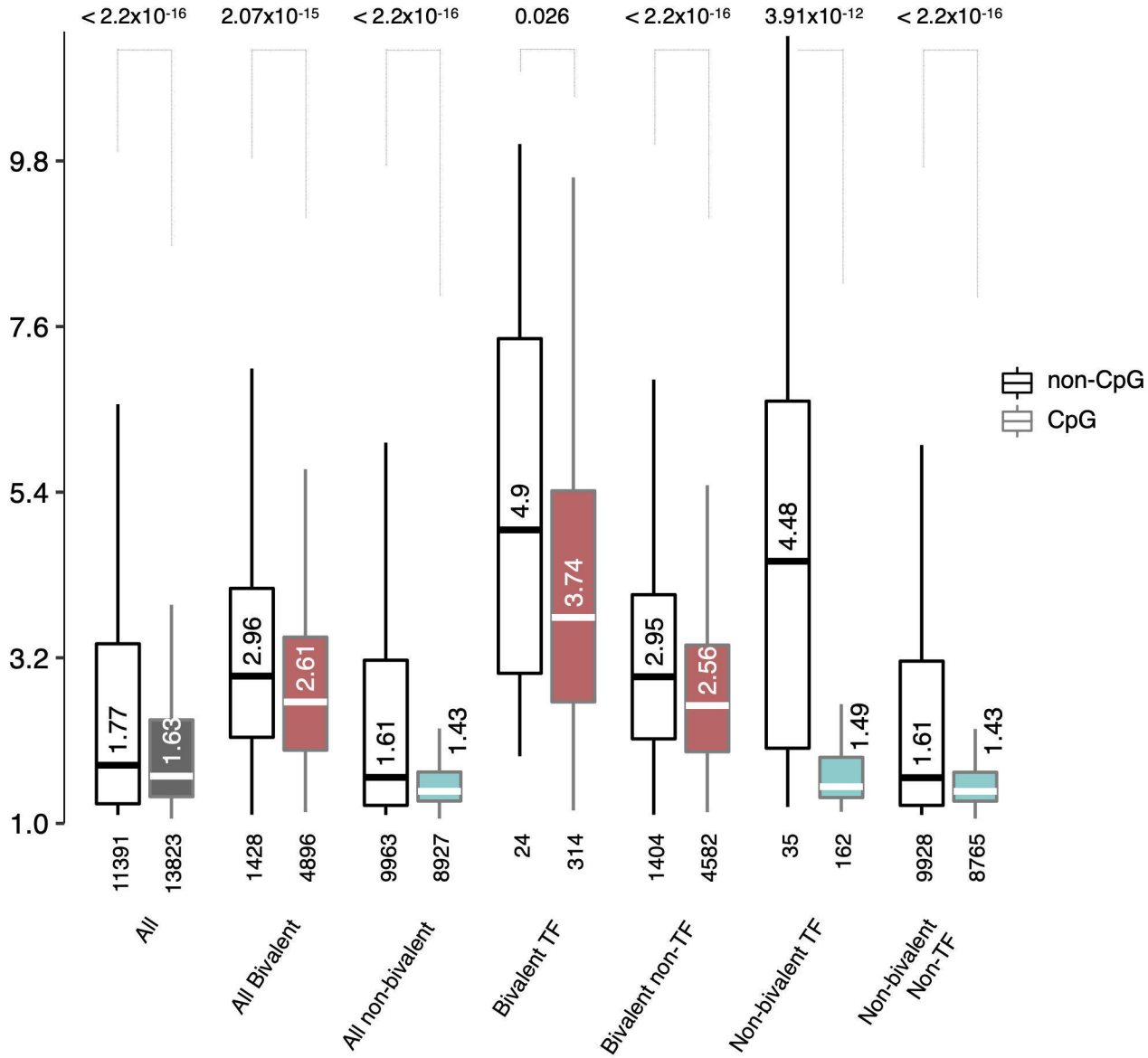


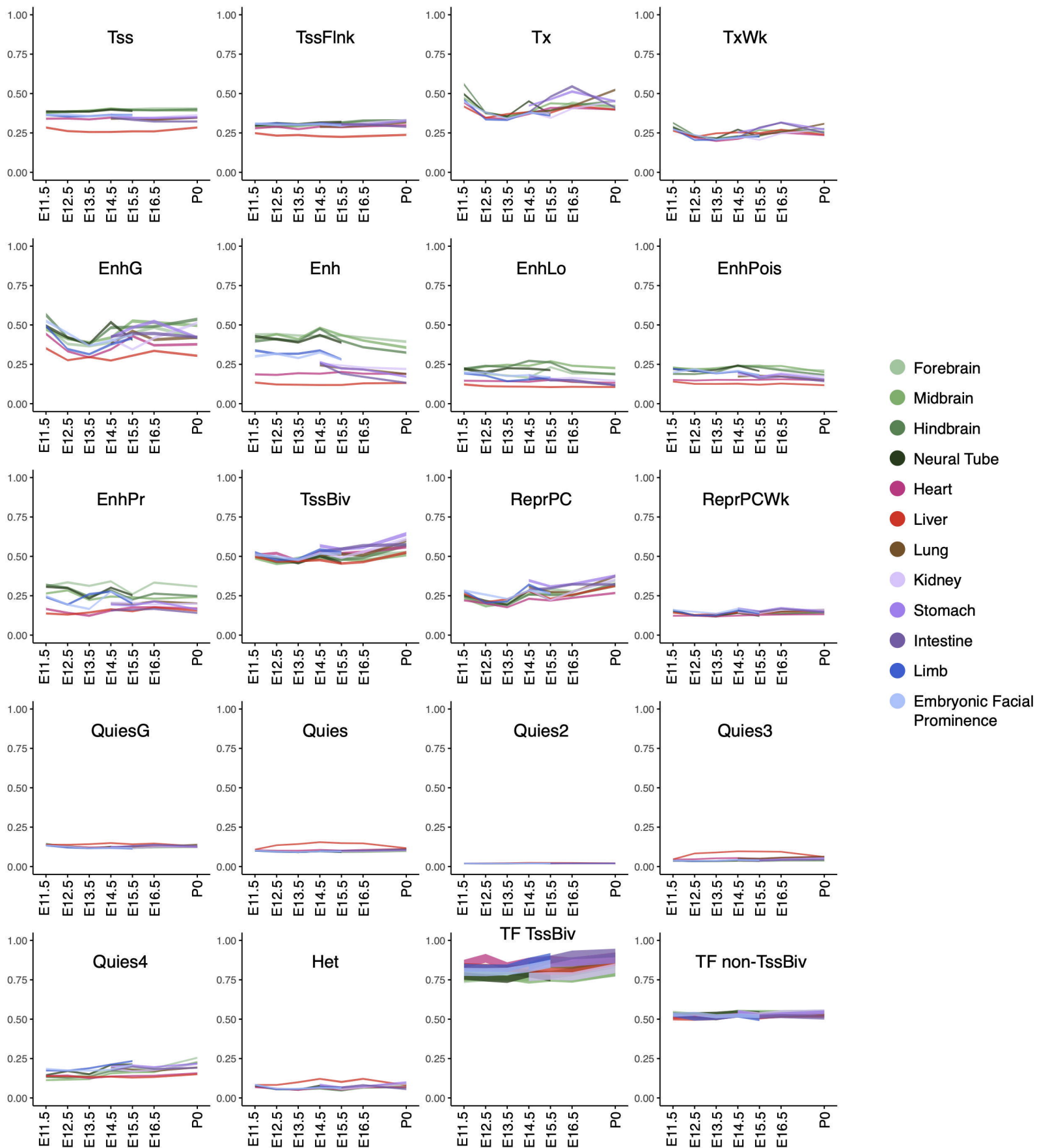




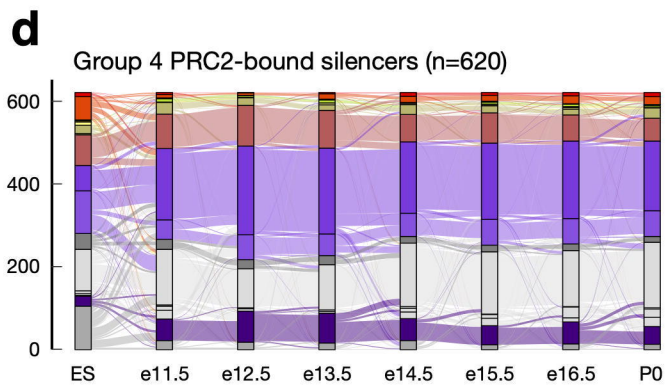
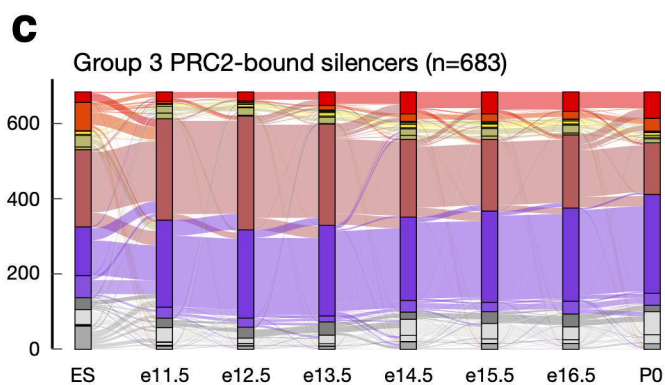
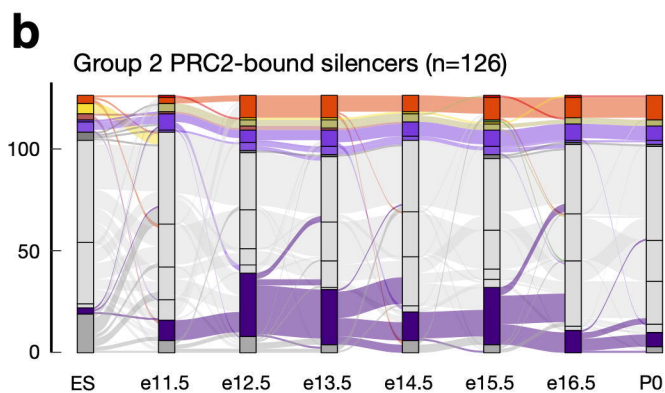
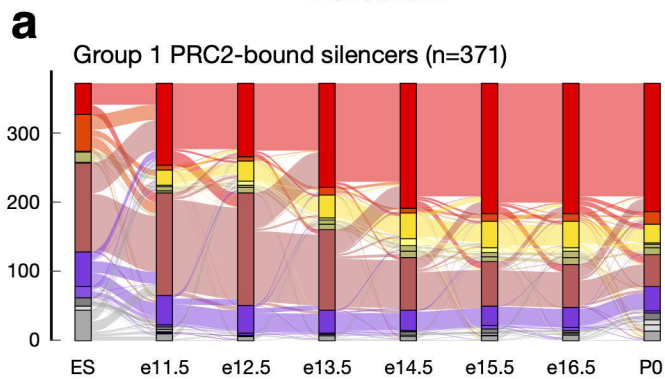
a**Bivalent TSS****b****Active TSS*** all P-values < 2.2x10⁻¹⁶

Tissue Specificity

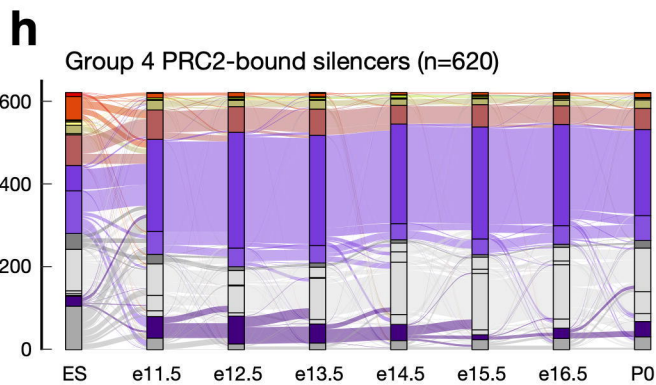
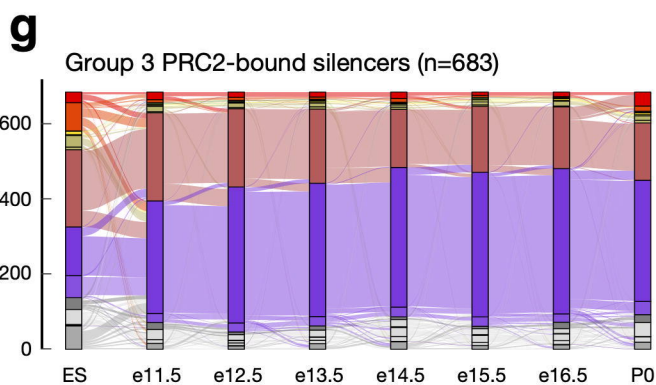
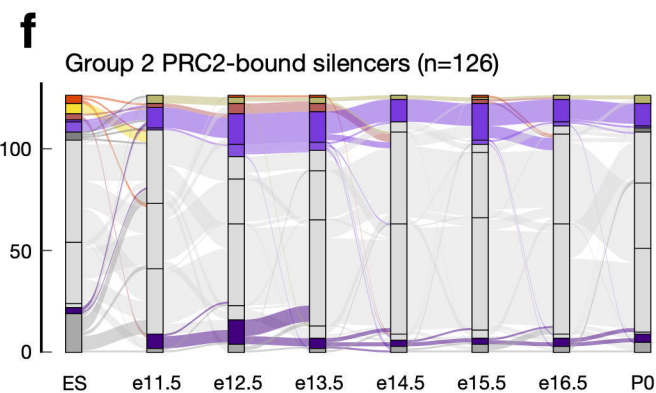
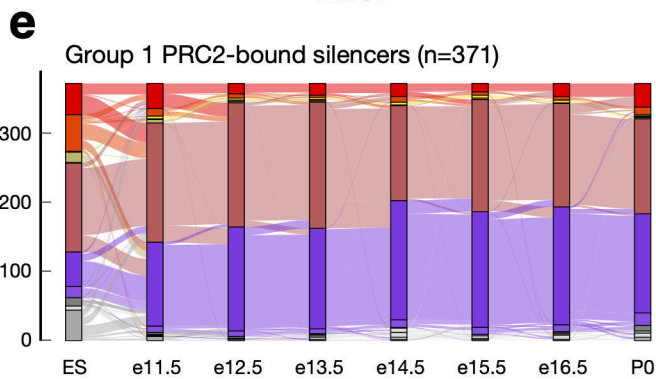




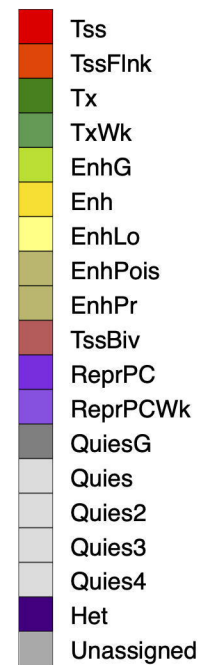
Forebrain



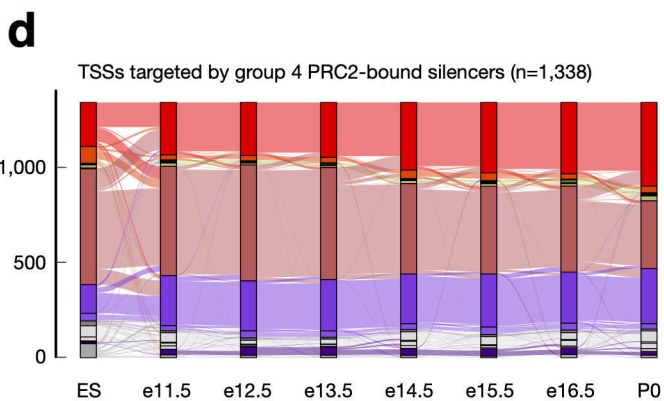
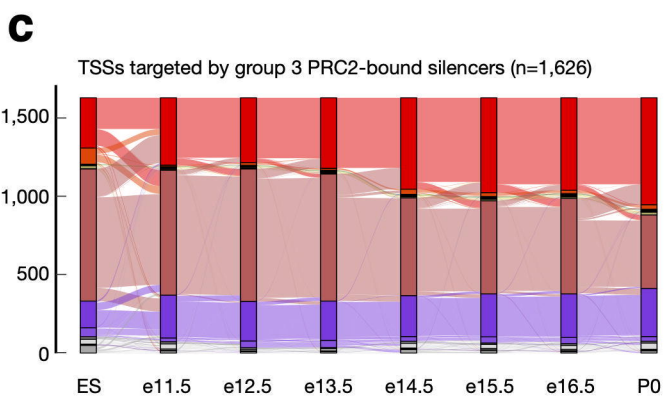
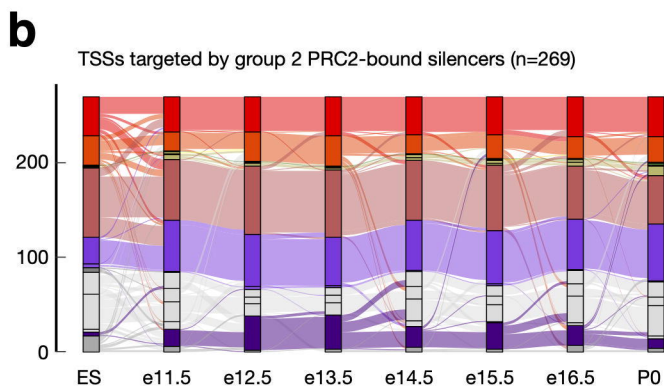
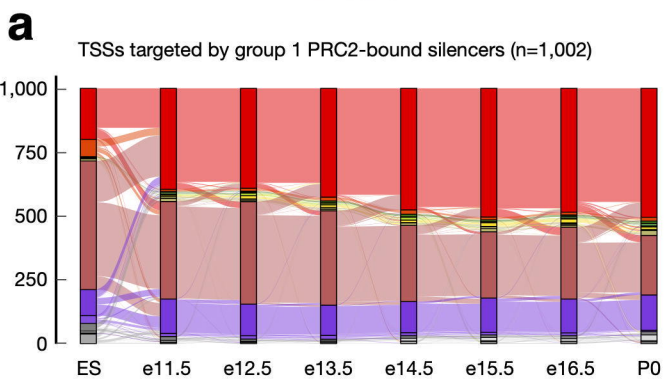
Liver



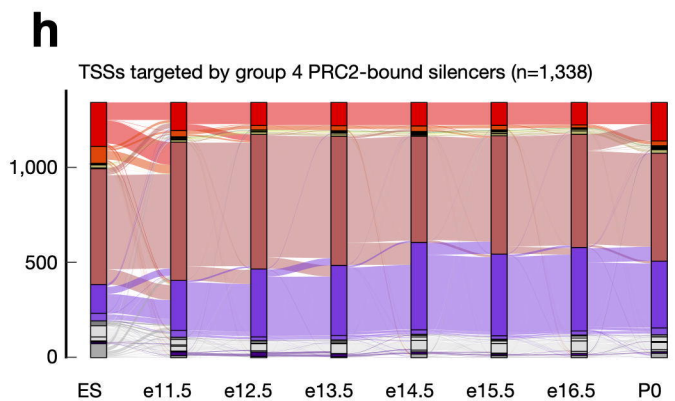
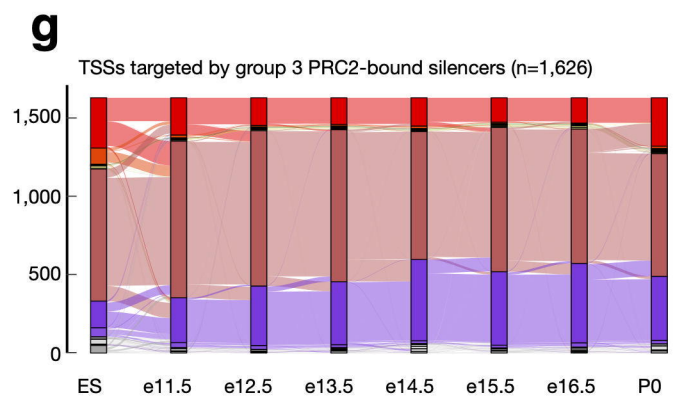
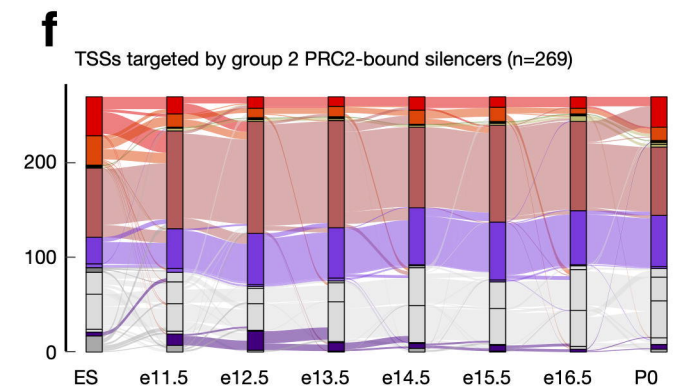
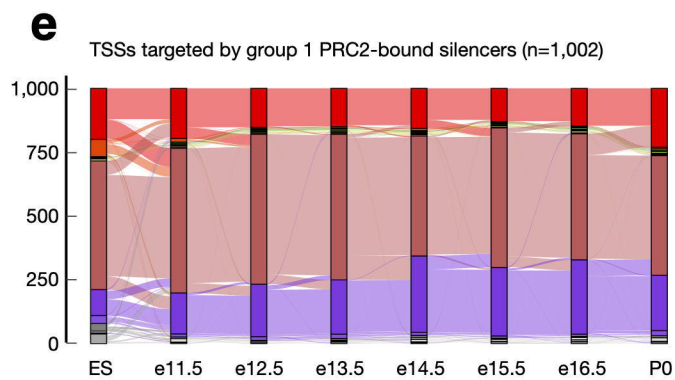
Chromatin State



Forebrain



Liver



Chromatin State

