

1 **Comprehensive identification of fetal cis-regulatory**  
2 **elements in the human genome by single-cell multi-omics**  
3 **analysis**

4 Hao Yu<sup>1,2,8</sup>, Na Ai<sup>1,2,8</sup>, Ping Peng<sup>6,8</sup>, Yuwen Ke<sup>1,2,3</sup>, Xuepeng Chen<sup>7</sup>, Yun Li<sup>1,2</sup>, Ting  
5 Zhao<sup>1,2,3</sup>, Shan Jiang<sup>1,2,3</sup>, Jiang Liu<sup>1,2,3,5\*</sup> & Lan Jiang<sup>1,2,3,4\*</sup>

6 1 China National Center for Bioinformation, Beijing 100101, China.

7 2 CAS Key Laboratory of Genome Sciences and Information, Beijing Institute of  
8 Genomics, Chinese Academy of Sciences, Beijing, China.

9 3 College of Future Technology College, University of Chinese Academy of Sciences,  
10 Beijing 100049, China.

11 4 Sino-Danish College, University of Chinese Academy of Sciences, Beijing 100049,  
12 China.

13 5 CAS Center for Excellence in Animal Evolution and Genetics, Chinese Academy of  
14 Sciences, Kunming, China.

15 6 Department of Obstetrics and Gynaecology, Peking Union Medical College  
16 Hospital, Beijing, China.

17 7 Bioland Laboratory (Guangzhou Regenerative Medicine and Health Guangdong  
18 Laboratory), Guangzhou, China.

19 8 These authors contributed equally.

20

21 \* Corresponding authors: Correspondence to Jiang Liu or Lan Jiang. Email:  
22 liuj@big.ac.cn(J.L.); jiangl@big.ac.cn(L.J).

23

24 **Abstract**

25 The regulatory programs driving early organogenesis in human is complex and still  
26 poorly understood. We performed parallel profiling of gene expression and chromatin  
27 accessibility to 28 human fetal tissue samples representing 14 organs in the first  
28 trimester. Collectively, we have generated 415,793 single-cell profiles. By integration  
29 analysis of transcriptome and chromatin accessibility, we detected 225 distinct cell  
30 types and 848,475 candidate accessible cis-regulatory elements (aCREs). By linking  
31 regulatory elements to their putative target genes, we identified not only 108,699  
32 enhancers, but also 23,392 silencers elements. We uncovered thousands of genes  
33 regulated by both enhancers and silencers in an organ or cell-type-specific manner.  
34 Furthermore, our unique approach revealed a substantial proportion of distal DNA  
35 elements are transcribed CREs (tCREs), which show both open chromatin signal and  
36 transcription initiation activity of non-coding transcript. The landscape of fetal  
37 cis-regulatory elements facilitates the interpretation of the genetic variant of complex  
38 disease and infer the cell type of origin for cancer. Overall, our data provide a  
39 comprehensive map of the fetal cis-regulatory elements at single-cell resolution and a  
40 valuable resource for future study of human development and disease.

## 41 **INTRODUCTION**

42 Developing and adult human tissues use different cis-regulatory elements but many  
43 adult chronic diseases including cancer may have a developmental origin<sup>1-3</sup>. Human  
44 fetal development is an exceedingly complex and fascinating process of transforming  
45 a single-cell zygote into a fully functioning organism within a mere span of 40 weeks<sup>4</sup>.

46 And the rudimentary formation of all organ systems raised from three primary germ  
47 layers (ectoderm, mesoderm, and endoderm) is completed by gestational week 16<sup>5-7</sup>.

48 A fundamental question is how the precursor cells with the same genetic material  
49 differentiate into diverse organs and cell types.

50 Leveraging single-cell molecular profiling techniques, many efforts have been carried  
51 out to explore cell heterogeneity and the development process in one or more  
52 organs<sup>8-11</sup>. But the majority of these were focused on transcriptome instead of  
53 chromatin states, which may prime to transcription or keep the epigenetic memory to  
54 adult cells<sup>12</sup>. Here, we performed massively parallel assays of 5' single-cell RNA  
55 sequencing (scRNA-seq) and single-cell assay for transposase-accessible chromatin  
56 and sequencing (scATAC-seq) for 14 human fetal organs. We characterize the  
57 chromatin accessibility, transcription initiation activity, interaction of target genes of  
58 cis-regulatory elements by integrative analysis of two assays to delineate the  
59 regulatory landscape of early organogenesis. Multiple modality rich information did  
60 uncover spatiotemporal dynamics of distal DNA elements driving human fetal  
61 development and help us further understand epigenomic change underlying disease  
62 pathogenesis.

## 63 **RESULTS**

64 We collected 1, 2, 13, and 12 fetal organ samples from four human donors ranging  
65 from gestational week 8 to gestational week 16 (Fig. 1a and Supplementary Fig. 1a, b).

66 For each sample, we parallelly generated matched 5' scRNA-seq and scATAC-seq

67 profiles by the droplet-based platform through the optimized protocol. All libraries  
68 were prepared with a capture target of 8000 cells.

69 After quality control, a total of ~3.1 billion read pairs were retained from  
70 scATAC-seq (Supplementary Table 1). These reads constitute 269,920 valid cells.

71 Taken system error into account, we merged multiplet cells about 8% of each library  
72 and removed doublet cells about 10% of each library (see Methods). Insert size  
73 distribution and TSS enrichment analysis confirms the high quality of our ATAC-seq  
74 data (Supplementary Fig. 1c, d). We observed an average level of 9,622 median  
75 fragments per cell among 28 samples. Finally, 230,732 high-quality cells with  
76 balanced sample sources are used for downstream analysis.

77 For the matched scRNA-seq for each sample, we applied stringent quality control for  
78 the number of detected genes and mitochondrial read counts. Doublets were removed  
79 by DoubletFinder (see Methods). In total, we profiled gene expression in 185,061  
80 individual cells, on average 2,150 genes per cell (Supplementary Fig. 1a and  
81 Supplementary Table 1).

## 82 **Annotating cell types**

83 Using SeuratV3<sup>13</sup>, we combined single-cell gene expression profiles from all samples  
84 and subjected them to batch effect removal and followed by Louvain clustering and  
85 UMAP visualization (Fig. 1b and Supplementary Fig. 1e). For the 42 major clusters  
86 identified, more than half of them are organ-specific, while others are derived from  
87 several organs. C10 (cluster 10) and C34 are mainly from the lung, while C8, C13,

88 C17, and C24 are a mixture of more than 7 organs. Surprisingly, the mixture clusters  
89 represent different common cell types and co-express specific marker genes. For  
90 example, C8 expresses endothelial cell markers PLVAP, as well as C13, which  
91 expresses enteric nervous system markers ELAVL4 (Supplementary Fig.1 f,g).  
92 Because large cell numbers and apparent heterogeneity exist in many of the 42 major  
93 clusters, we went into second round Louvain clustering. We identified sub-clusters  
94 within each major cluster and got 335 sub-clusters in total. We assign cell type labels  
95 to scRNA-seq major clusters and sub-clusters according to known marker genes from  
96 literature and HCL references<sup>11</sup> (Supplementary Table 2). Through 2 rounds of  
97 clustering, we were able to identify common cell types across samples while retaining  
98 organ-specific cell types.

99 Next, we transferred cell type labels from 5' scRNA-seq data to scATAC-seq data  
100 within each organ. We computed gene activity scores for scATAC-seq data, aligned  
101 cells from scATAC-seq to cells from scRNA-seq in low dimension space, and got a  
102 best-fitted label for each cell using ArchR<sup>14</sup>. As some labels have very few cells in  
103 scATAC-seq data, we set a cut-off removing transfer results with a low  
104 signal-to-noise ratio (Supplementary Fig. 1h) and finally got 225 reliable labels with  
105 paired pseudo-bulk profiles of gene expression and chromatin accessibility (Fig. 1c  
106 and Supplementary Fig. 1i). To facilitate the exploration of this dataset, we provide an  
107 online interface ([http://genome.ucsc.edu/cgi-bin/hgTracks?hgsid=1140461557\\_BMEZ54Vfu607BWs6t5LASyfZT5sj](http://genome.ucsc.edu/cgi-bin/hgTracks?hgsid=1140461557_BMEZ54Vfu607BWs6t5LASyfZT5sj)).

108 **Identify consensus accessible chromatin sites**

109 To construct a map of the cis-regulatory elements marked by chromatin accessibility,  
110 we called peaks for each cell type and took the iterative overlap peak merging  
111 procedure eliminating redundant peaks using ArchR (see Methods). As a result, the  
112 most significant signal in the form of 501bp peaks are caught and a master list of  
113 848,475 consensus accessible chromatin sites are constructed, spanning 14% area of  
114 the whole human genome (Supplementary Table 3).

115 Previous large-scale efforts such as ENCODE3 have mapped open chromatin regions  
116 for various tissues/organs and developmental stages mainly based on bulk  
117 DNase-seq<sup>15</sup> or bulk ATAC-seq<sup>16</sup>. However, to what extent, the list of cis-regulatory  
118 elements in the human genome is completed is still an open question. We calculated  
119 overlaps between peaks we identified and human DHSs of corresponding primary  
120 tissues from ENCODE3 (Fig. 2b and Supplementary Fig. 2a). As shown in Venn plots,  
121 more than half of the DHSs are detected in our data. And importantly, 153,496 novel  
122 peaks are uncovered in our data exclusively. Then, we probed into which tissues/cell  
123 types contributing most to the dataset's specific peaks. The majority of DHSs specific  
124 peaks are contributed from adult tissues (Supplementary Fig. 2b) while the majority  
125 of scATAC-seq specific peaks are contributed from common cell types such as  
126 neurons, macrophages, and endothelial cells, with limited overlap between sub cell  
127 types (Fig. 2c and Supplementary Fig. 2c-e). We proposed that common cell types  
128 distributing in various organs may be underrepresented in the bulk experiment, while  
129 clustering of single-cell data across organs can better capture cis-regulatory elements

130 of those cell types. A significance test in the box plot confirms the above viewpoint  
131 (Fig. 2d).

### 132 **Enhancer Validation by Comparison with VISTA database**

133 The VISTA database is a central resource for experimentally validated human and  
134 mouse noncoding fragments with gene enhancer activity as assessed in transgenic  
135 mice<sup>17</sup>. To know whether the validated enhancers are covered by our results, we drew  
136 comparisons on several levels. As showed in the bar plot, over 82% of VISTA  
137 enhancers are identified in corresponding organs in our dataset. Besides, about 96% of  
138 enhancers are covered without regard to organ sources (Supplementary Fig. 2f).  
139 VISTA enhancers are most enriched in the corresponding organ (Fig. 2e), which  
140 confirms the tissue specificity of enhancers. More importantly, we can go deep into  
141 the cell type level and explore which VISTA enhancers are open in each cell type,  
142 expanding our knowledge of enhancers' function (Supplementary Fig. 2g). For  
143 instance, most VISTA enhancers of the heart are open in cardiomyocytes,  
144 contributing to the expression of tissue-specific genes like FHL2 (Fig. 1d and Fig. 2f).

### 145 **Recognizing the pattern of accessible chromatin regions**

146 To connect accessible chromatin sites with biological cellular contexts, we  
147 constructed a binary matrix of 225 cell types  $\times$  848,475 peaks in which 1 denotes that  
148 the peak is open in the corresponding cell type. After characteristic clustering by rows  
149 and columns, we were able to visualize the binary matrix in a fashion of neatly  
150 arranged blocks on the diagonal (Fig. 2a). The hierarchical cluster by rows offers the

151 biological information on which cell types are most strongly associated with each  
152 peak group. The K-means cluster ( $K = 21$ ) by columns separates peaks into 21 groups.  
153 Based on lineage specificity for each group, we defined peaks in each group as a  
154 lineage specifier family (LSF). That is, LSF19 is mainly open in kidney epithelial  
155 cells, and may dominant cell differentiation and cell fate decision in the  
156 nephrogenesis.

157 We annotated each LSF with the best-fitted cell lineage based on associated cell types  
158 or regulators inferred by motif enrichment. For example, peaks in LSF2 are  
159 exclusively accessible among macrophages and are annotated as macrophage-related  
160 LSF. Peaks in LSF21 are universally open and over 64% of peaks are proximal to  
161 TSS ( $\pm 1\text{kb}$ ), which implies that promoter regions are less dynamic across all cell  
162 types. Peaks in LSF13 are open in about half of all cell types and we conjectured its  
163 universal function across organs. Motifs most enriched in LSF13 include  
164 *Atoch1/Tcf12/NeuroG2* and *Tcf21/MyoD/Twist2*, all of which are helix-loop-helix  
165 (HLH) transcription factors and act as key regulators of neurogenesis, myogenesis,  
166 and osteogenesis<sup>18</sup>.

### 167 **Developmental dynamics of chromatin accessibility**

168 To decipher molecular regulation mechanisms underlying LSFs, we sought to explore  
169 chromatin accessibility dynamics within each LSF and transcription regulators in  
170 lineage differentiation. Taken LSF19 as an example, we adopted an iterative strategy  
171 taking cell types of kidney epithelial and repeating K-means clustering ( $K = 10$ ) to



172 identify sub pattern of accessible chromatin states (Fig. 3a). And we denote subcluster  
173 3 of LSF19 as LSF19.3. This process produced informative substructures and  
174 uncovered a huge difference between progenitor cells (cap mesenchyme, CM) and  
175 differentiated cells (primitive vesicle, PV; proximal tubules, PT; Loop of Henle, LoH;  
176 distal tubules, DT).

177 The progenitor cells (CM) seem to have the most open chromatin states. Along with  
178 lineage differentiation, a lot of sites like LSF19.1/2 are turned off and other  
179 function-relevant regions like LSF19.6/7 are opened while some open states like  
180 LSF19.9/10 are maintained. The chromatin accessibility states are modified in a  
181 branch-determined way.

182 Furthermore, we found that motif enrichment was consistent with corresponding TF  
183 expression in each cell type (Fig. 3b). LSF19.1/2 were enriched with Six2/Six1 motifs  
184 and high expression of SIX2/SIX1 also appeared in progenitor cells like cap  
185 mesenchyme. This suggests that transcription factors are responsible for establishing  
186 and maintaining open chromatin states.

187 Next, we performed trajectory inference analysis to resolve lineage differentiation at  
188 the single-cell level using scATAC-seq data (Fig. 3c, d and Supplementary Fig. 3).

189 We were surprised to find that chromatin accessibility of LSF19 (only 46838 peaks)  
190 has sufficient information to distinguish different cell types and underlies  
191 differentiation order, which means PV emerges before other parts in the timeline. The  
192 DT and collecting duct (CD) are the final two segments of the kidney nephron with

193 the function of ions absorption and water reabsorption. However, the distal cells of  
194 the comma-shaped body (precursor of DT) invade the proximal tip of the UB  
195 (progenitor of CD) and fuse to form one continuous P/D axis at early stages. We  
196 captured a continuous reprogramming process along with the differentiation to DT at  
197 the single-cell level (Fig. 3e). The converge suggests that spatial organization or local  
198 function may be a more deterministic factor in chromatin accessibility states  
199 compared to cell origin. To understand how transcription factors, help to maintain cell  
200 states and play a role in lineage differentiation, we made an in-depth investigation on  
201 SIX2, which maintains cap mesenchyme in an undifferentiated state<sup>19</sup>. We found  
202 SIX2 as a transcription factor can also target the putative enhancer of SIX2 itself to  
203 positively regulate SIX2 expression. Then, we inferred the target genes of TF based  
204 on the association of TF target peaks. We found that the dynamics of chromatin  
205 accessibility of target peaks and expression of target genes of SIX2 have the same  
206 trend as SIX2 expression (Fig. 3f). This suggests that the dynamic of open chromatin  
207 states is driven by the expression and function of transcript factors, while  
208 cis-regulatory elements regulate gene expression in a forward way.

### 209 **Linking regulatory elements to cognate genes**

210 We next asked how distal regulatory elements regulate gene expression. Peak  
211 co-accessibility is often used to predict enhancer-promoter interactions<sup>20</sup>. However,  
212 the accessibility of ubiquitous opened promoters is usually moderately correlated with  
213 gene expression. Therefore, we leveraged the gene expression data and created a

214 correlation-based map between chromatin accessibility peaks and their cognate genes  
215 directly (see Methods).

216 Using correlation analysis, we identified 155,620 positive peak-to-gene links  
217 (associated with 108,699 peaks and 12,783 genes) and 34,287 negative peak-to-gene  
218 links (associated with 23,392 peaks and 7,628 genes) (Supplementary Table 4). Then  
219 we defined positive links as putative enhancer-gene pairs and negative links as  
220 putative silencer-gene pairs. For example, FHL2 plays an important role in  
221 cardiomyocyte differentiation by negatively regulating the calcineurin/NFAT  
222 signaling pathway. And we found the putative enhancers of FHL2 are exclusively  
223 open in two sub cell types of cardiomyocytes, confirming the accuracy of our results  
224 (Fig. 1d).

### 225 **Comparison with ReSE-identified silencers**

226 Pang and Snyder devised a lentiviral screening approach<sup>21</sup>, the repressive ability of  
227 silencer elements (ReSE), to systematically identify silencer regions in human cells.  
228 They assayed on K562, PMA-treated K562, and HepG2 cell lines, and identified a  
229 total of 5472 non-overlapping silencers. To validate our data, we compared our  
230 correlation-based silencers and ReSE-identified silencers and found an overlap of 174  
231 silencers. chr5:171602285-171602785 and chr19:48763298-48763798 are two  
232 examples with different distributions in 225 cell types (Fig. 4a-c). The former shows a  
233 sharp decline in expression when the accessibility of the silencer reaches a level of 0.2,  
234 and the latter is much milder with a downward tendency. Based on the sharp decline

235 or not, we can classify silencers into strong silencers or weak silencers (see Methods).

236 These two classes may underline two mechanisms: a switch way through repressed

237 epigenetic states to turn on or off target genes (strong silencers), and a competitive

238 way through transcriptional machinery interactions (weak silencers).

239 To take advantage of our large-scale data, we further predicted target genes for

240 ReSE-identified silencers. 2,113 silencers have at least one neighboring negative

241 correlated gene. Our data and analysis can add complementary information to

242 experimentally verified silencers in whole organism scales (Supplementary Table 5).

#### 243 **Adversarial regulation on the same gene**

244 To investigate the relationship between our classification of cis-elements and 21 LSFs,

245 we calculated enrichment for each category of cis-elements (Fig. 4d). Interestingly,

246 LSF1, LSF3, LSF10 are enriched with both silencer and enhancers, and they are all

247 related to the hemopoietic system, which underscores a complicated regulatory

248 fashion during hematogenesis, which is consistent with recent report<sup>22</sup>.

249 Although correlation analysis is based on one peak to one gene, the real situation is

250 that multiple cis-elements cooperatively or competitively regulate the same gene in a

251 cell-type-specific manner. We found a total of 6,091 genes which are the targets of

252 both putative enhancers and silencers (Supplementary Table 6) and focused on a set

253 of 94 genes identified at the whole organism level. Integrated genes expression and

254 open chromatin information allow us better resolve the complexity of regulation (Fig.

255 5a-c).

256 Of the 161 silencers, the majority are open in the hemopoietic system, which is  
257 consistent with the cis-regulatory elements & peak group enrichment analysis (Fig. 4d  
258 and Supplementary Fig. 4a, b). In line with expectations, the pattern of accessibility of  
259 the enhancer is almost the same as the gene expression, while the pattern of  
260 accessibility of the silencer is the opposite (Fig. 5a-c). The accessibility pattern of  
261 enhancer and silencer of the same gene are mutually exclusive and have a negative  
262 correlation. The underlying mechanism will require further investigation.

263 We next made an in-depth study on one gene, MMP14 (Fig. 5d), whose encoded  
264 protein are involved in the breakdown of extracellular matrix in normal physiological  
265 processes, such as embryonic development, reproduction, and tissue remodeling, as  
266 well as in disease processes, such as arthritis and metastasis. In our dataset, fibroblasts  
267 from different organs have high level expression while erythroid cells and immune  
268 cells have low level expression. In the track plot, silencers are from close to open and  
269 enhancers are from open to close with the decrease of expression level (Fig. 5d).  
270 There is a cliff-like change when the accessibility level of silencer 1 reaches the  
271 critical point of 0.3, which suggests a switch of regulatory modules (Fig. 5e). When  
272 under the critical point, the accessibility of enhancer 2, as well as the expression of  
273 MMP14, is highly variable, and enhancer 2 determines the expression level (Fig. 5e-g  
274 and Supplementary File 10). Once reaching the critical point, both enhancer and gene  
275 transcription is silenced. Mutually exclusiveness of chromatin accessibility between  
276 enhancers and silencers uncovers two regulatory modules, functioning in part of cells

277 antagonistically. We further probed into the silencer preference among different cell  
278 types. The ternary plot indicates that silencer 1 functions alone in erythroblast, while  
279 silencer 2/3 are co-accessible and functional in B cells and T cells (Fig. 5h). The  
280 cis-element selection may emerge along with the cell fate decision.

281 While the silencers in the above example are all strong silencers, we get quite curious  
282 about what if one gene is associated with a weak silencer. We took IFITM3 as an  
283 example and did the same analysis as MMP14 (Supplementary Fig. 4). Both the  
284 accessibility of the enhancer 3 and the expression of IFITM3 are mildly decreased as  
285 the silencer gets more accessible (Supplementary Fig. 4c-e). The antagonism between  
286 the enhancer and the silencer does make the expression of IFITM3 more variable  
287 (Supplementary Fig. 4f).

288 To compare these two different patterns, we would like to propose two models about  
289 adversarial regulation on the same gene: a switch model and a competitive model. The  
290 switch model is tightly associated with strong silencers, which turn off the enhancer  
291 and gene transcription simultaneously. As a result, the enhancers' function is  
292 restricted in a convergent triangular zone. That is, the enhancers only function in the  
293 absence of an active silencer, and the enhancers' activity converges to 0 with the  
294 activation of the silencer. The competitive model is involved with weak silencers and  
295 may have a relationship with competitive combination with the promoter. As a result,  
296 gene expression is highly variable and can be finely controlled in both positive and  
297 negative ways.

## 298 **Delineation of transcription initiation activity of distal regulatory elements**

299 Enhancer RNAs (eRNAs) are RNA molecules that are transcribed from genomic  
300 enhancer regions<sup>23</sup>. The previous study shows that the level of enhancer RNA  
301 expression positively correlates with the level of mRNA synthesis at nearby genes<sup>24</sup>.  
302 To decipher element functions in the transcription aspect, we quantitatively analyzed  
303 the transcription level of distal regulatory elements by leveraging 5' scRNA-seq. We  
304 observed strong enrichment of RNA signal at the center of distal ATAC peaks (Fig.  
305 6a). Transcription level and open chromatin states are positively correlated at sample  
306 level and cell type level with a large proportion of elements open but not transcribed  
307 (Fig. 6b and Supplementary Fig. 5a). To identify transcribed cis-regulatory elements  
308 (tCREs), in other words, open chromatin region with transcription initiation activity,  
309 at the whole organ scale, we applied a strict cut-off to each sample and merged tCREs  
310 lists into a master list of 190,356 regions (Supplementary Fig. 5b, c and  
311 Supplementary Table 7).

312 For each cell type, about 10% of open chromatin regions have non-coding  
313 transcription start site signal on average. Combining tCREs with peak-to-gene links,  
314 we found cell types with transcribed enhancers have significantly higher expression  
315 levels of target gene than cell types with an un-transcribed enhancer (Fig. 6c). We  
316 further identified 1361 peak-to-gene links in an eRNAs-dependent manner, 206 of  
317 which were associated with TF-encoding genes<sup>25</sup> (Supplementary Table 8). Open  
318 chromatin state is the necessary condition of transcription, and the level of eRNAs is a

319 determining factor in promoting target genes (Fig. 6d, e, and Supplementary Fig. 5f).  
320 To assess cell type specificity of the tCREs, we ordered tCREs according to their  
321 source peak groups and found a similar but more evident pattern with chromatin  
322 accessibility pattern, which may indicate higher specificity in cis-element  
323 transcription (Fig. 6f). We also note that universal open peaks have higher transcribed  
324 proportions and may have a specific function (Supplementary Fig. 5d). Multiple  
325 enhancers may be co-accessible and regulate the same gene. Based on this, we  
326 assumed that co-expressed cis-elements are likely to be functional elements instead of  
327 random non-coding transcription noise. We found about 54.2% of tCREs have a  
328 highly co-expressed patternner ( $cor>0.8$ ) (Fig. 6g). We also found more than half of  
329 our defined enhancer-to-gene pairs are associated with un-transcribed cis-elements,  
330 most of which cooperate with another transcribed enhancer to regulate the same target  
331 gene (Supplementary Fig. 5e). What's more, the remaining enhancers work alone  
332 without transcription signal, suggests that many enhancers function in a  
333 transcription-independent manner (Supplementary Fig. 5g). The precise molecular  
334 mechanism of different categories of enhancers needs further investigation.

### 335 **Enrichment analysis of GWAS signals in aCREs and tCREs LSFs**

336 To further our understanding of lineage specifier families, we applied stratified  
337 linkage disequilibrium score regression<sup>26,27</sup> and evaluated heritability enrichment in  
338 52 GWAS datasets (Supplementary Table 9) across these 20 LSFs. The spectrum of



339 traits evaluated covered blood cell physical traits, neurological, immunological,  
340 gastroenterology, metabolomic traits from UK Biobank data<sup>28</sup> and Broad LD Hub<sup>29</sup>.

341 We observed Immune-related LSF show similar heritability enrichment for immune  
342 traits (Supplementary Fig. 6a). Lupus, Crohn's disease, Rheumatoid Arthritis are  
343 significantly correlated with immune-related LSFs (T-cells, Immune system, and  
344 macrophage). The strongest enrichment of heritability for immunoglobulin A (IgA)  
345 deficiency is in T cells. Epithelial LSF dominated by different organs display specific  
346 enrichment features for organ-matched traits. Kidney epithelial are relevant with  
347 kidney-stone. The lung epithelial and gonad LSF both enrich in lung FEV1/FVC ratio.

348 Likewise, some blood cells' physiology traits and immune-related traits are  
349 significantly enriched in Erythroid LSF, T2D, and Fasting Glucose are highly  
350 correlated with Endocrine systems, which are consistent with prior knowledge.

351 Furthermore, we found that the enrichment tendency of heritability of two neuron  
352 LSFs is different. Neuron1 LSF, which is mainly contributed by the retina or neural  
353 portion of the eye, is part of the central nervous system. Neuron2 LSF, which is called  
354 enteric nervous system (ENS) LSF. The results of this GWAS heritability analysis  
355 showcase, several psychiatric traits, and major neurodegenerative disorders, like  
356 Schizophrenia, Neuroticism, highly correlate with Neuron1 LSF, in stark contrast  
357 with weak signal in Neuron2 LSF. It suggested that the eye is a 'window' into the  
358 brain, the accessibility and organization of the retina make it a convenient research  
359 tool with which to study processes in the CNS<sup>30</sup>. Unexpectedly, the Eye-related open

360 chromatin enriched variants of the muscle-skeletal system and connective diseases,  
361 which may suggest some unrecognized link.

362 The previous study suggests that some distal aCREs marked by ATAC-seq or  
363 DHS-seq signal don't have enhancer activity. Those regions maybe not binding by  
364 TFs or not interact with the promoter to drive gene expression, even they are open.  
365 Meanwhile, those open chromatin regions which have transcription initiation activity  
366 (tCREs) are more likely to be active enhancers, since the RNA signal suggests they  
367 are accessible by Pol II. Thus, we wonder whether tCREs are more enriched with  
368 GWAS signals and functionally relevant. For each open chromatin LSF, we identify  
369 the corresponding tCRE LSF (Supplementary Table 9). We calculate the GWAS  
370 signals enrichment similar to aCREs as described above. Interestingly, we found the  
371 enrichment of some traits and disease related SNPs are higher in tCREs than in  
372 aCREs LSF (Supplementary Fig. 6b-d). To avoid the trait heritability difference is  
373 caused by captured SNP number from aCREs LSF and tCREs LSF. We calculate

374  $\Pr(h_g^2)/\Pr(\text{SNPs})$  to measure LSF genetic associations and heritability. For Thyroid  
375 Disease, heritability was markedly enriched specifically within T cells associated  
376 tCREs LSF compared with aCREs, it indicated tCREs can capture trait heritability  
377 better than aCREs, it may cover more vital genetic signals (Supplementary Fig. 6b).

### 378 **Heritability enrichment identifies traits and disease-relevant fetal cell types**

379 Many common diseases have a developmental origin. Despite the remarkable success  
380 of genetic signal mapping in GWAS, the functional interpretation of GWAS remains

381 challenging. First, it is unclear in which tissues and cell types these variants are active,  
382 and how they disrupt specific biological networks to impact disease risk. Second,  
383 most disease-associated variants are located in non-protein-coding regions of the  
384 genome, and many are far away from the nearest known gene. We have evaluated the  
385 genetic risk of traits and disease for LSF, however, the most relevant cell types of  
386 certain diseases during organogenesis are poorly understood. CREs are bits of  
387 noncoding DNA that regulate the transcription of nearby genes. Here we can use each  
388 cell type top 10,000 specific CREs<sup>9</sup> to explore the cellular context in which  
389 disease-associated variants act.

390 The results revealed that risk variants for kidney stones and chronic kidney diseases  
391 were enriched in kidney tubule cells (Fig. 7a). For tubule cells, it comprised distinct  
392 subpopulations with differentially accessible chromatin regions. We further provide a  
393 finer genetic signal map of the tubule subpopulation. Distal tubule cell shows higher  
394 enrichment (q value <0.05) for kidney function-related traits (eGFR, BUNM, Urate)  
395 from the study by Wuttke et al. and Teumer et al<sup>31,32</sup>, and S-Shaped body cell type and  
396 LoH cells are both relevant to the kidney stone. Likewise, we find endocrine cells,  
397 which showed significant enrichment for fasting glucose (Fig. 7b).

398 Dot plot shows the  $-\log_{10}(\text{q value of enrichment})$  for two chronic Inflammatory  
399 bowel diseases (IBD, Crohns' Disease) across all cell types in the large intestine (Fig.  
400 7c). Only one digestive-system sourced macrophage has significant enrichment. It

401 consisted of a recent study that reported a subtype of NOD2-driven Crohn's disease  
402 leads to dysregulated homeostasis of activated fibroblasts and macrophages<sup>33</sup>.

403 The most relevant cell type of heart traits cardiac arrhythmias and atrial fibrillation  
404 and flutter (AF) and Cardiac arrhythmias COPD comorbidities are cardiomyocytes  
405 (Fig. 7d, and Supplementary Fig. 6e). AF risk variant (rs7789585) is located in a  
406 cardiomyocyte's specific open chromatin region, which resides in the second intron of  
407 the KCNH2. Co-accessibility analysis suggests that KCNH2 is likely the target. This  
408 observation is consistent with a recent report that cardiomyocyte enhancers of  
409 potassium channel gene KCNH2 may be affected by noncoding risk variants  
410 associated with AF<sup>34</sup>. Collectively, we have assigned the most relevant fetal cell type  
411 for 10 traits or diseases (Supplementary Table 9).

#### 412 **Cell type of origin for cancer**

413 Cells from fetal tissue and tumor both grow and divide rapidly, and they share  
414 common cell surface markers and oncofetal antigens, include carcinoembryonic  
415 antigen (CEA), alpha-fetoprotein (AFP)<sup>35</sup>. To a certain extent, malignant tumor  
416 regulatory mechanisms resemble fetal cells, the fetal tissue in a single-cell resolution  
417 may provide the answer of the cell type of origin for the tumor. For example, a recent  
418 study has found that most adrenal NB tumor cells transcriptionally mirror early  
419 human embryos' noradrenergic chromaffin cells<sup>36</sup>. Moreover, another recent study  
420 reported a shared immunosuppressive oncofetal ecosystem in fetal liver and  
421 hepatocellular carcinoma<sup>37</sup>, suggesting fetal tissue may provide a better understanding

422 of the tumor ecosystem. The large-scale cross organ datasets generated in our study  
423 allow us to explore the similarity of fetal cell types with multiple cancer types. To  
424 ensure accuracy, we pay more attention to 9 tumor types from TCGA which have  
425 corresponding fetal tissue in our datasets, and their chromatin state was profiled by  
426 bulk ATAC-seq in a previous study<sup>38</sup>.

427 For each tumor sample, we inferred the putative cell type of origin based on the  
428 chromatin accessibility similarity with fetal cell types using Jaccard distance  
429 (Supplementary Table 10). We observed almost all patients show accordant  
430 preference on specific cell types based on chromatin accessibility and found  
431 cancer-associated cell types. Across 41 stomach adenocarcinoma (STAD) samples,  
432 the fetal stomach cell types which show the highest similarity score consistently to be  
433 Surface Mucous Pit Progenitor cells (Fig. 8b, d), which make mucus and stomach  
434 juices.

435 In a similar fashion to previous analysis of STAD-associated fetal cell types, we  
436 summary the top5 most similar fetal cell types for each tumor (Fig. 8b, and  
437 Supplementary Fig. 7a). Meanwhile, 54 cell types in our data have been annotated as  
438 a proliferative state based on CytoTRACE<sup>39</sup> inference and unique gene expression  
439 (Supplementary Table 10). To investigate whether the cancer-associated fetal cell  
440 types are enriched in proliferate or progenitor cell types, we use a hypergeometric test  
441 to compute the statistical significance of the intersection of cancer associated cell

442 types and proliferative state cells. We found colon adenocarcinoma (COAD) and  
443 STAD are clearly different from the other 7 cancer types, their associated fetal cell  
444 types are significantly enriched proliferate states (Fig. 8a). The COAD-associated cell  
445 type is enterocyte progenitor cell that sustains proliferating state in the large intestine,  
446 while STAD associated cell type is surface Mucous Pit progenitor in the stomach (Fig.  
447 8c, d). Moreover, both of them show similar chromatin states nearby CEA family  
448 genes with cancer (Supplementary Fig. 7b, c).

## 449 **DISCUSSION**

450 In this study, we leveraged single-cell profiling of RNA and chromatin to perform  
451 integration analysis and construct cis-regulatory elements atlas. The scale of the  
452 current analysis helped us to discern more details on the biological phenomenon and  
453 better understand transcription regulation. By comparing with the VISTA database,  
454 we got to know validated enhancers are open in which cell types. By integrating motif  
455 enrichment and gene expression, we confirmed transcription factors acting as key  
456 regulators of dynamics of open chromatin and lineage differentiation. By combining  
457 positive cis-elements with negative cis-elements, we found mutually exclusive  
458 modules regulating the same genes in a cell-type-specific manner, which may provide  
459 a potential way for disease treatment.

460 The cis-regulatory elements atlas of the current study provides a snapshot of fetal  
461 development. It would be more valuable to sample in continuous stages, offering a  
462 spatiotemporal perspective of lineage hierarchy and transcription regulation. More

463 advanced experimental technologies and algorithms will emerge and set a foundation  
464 for better resolving fetal development sometime in the future.

465 Fetal tissue with persistent differentiation potential finally developed functional  
466 mature normal adult tissue, whereas it also can switch to tumor disordered  
467 proliferation state in the oncogenic mutations stimulate. It looks like a one-direction  
468 irreversible event, whereas tumor tissue can break this order, it reactivates some  
469 cis-elements with normal fetal tissue which keep silent in adult tissue, switch cell  
470 status to benefit tumorigenesis. Our study builds a bridge between the two  
471 physiological states based on the similarity of the open state of chromatin and  
472 provides a new perspective for the exploration of the developmental origin of tumors.  
473 We systematically summarized fetal cell types have a similar regulatory mechanism  
474 with 9 primary tumors. In addition, for TCGA bulk level ATAC-Seq data of tumor  
475 tissues, it can observe cellular composition heterogeneity and complex  
476 microenvironment in tumor samples. And our findings show these oncofetal antigens  
477 are cell type-specific open in fetal tissues, which prefer proliferating state cell types  
478 with persistent multilineage differentiation potential, and these genes are also  
479 reactivated in tumor cells, which seems to support the previous hypothesis. However,  
480 we haven't detected adult tissue and cancer tumor chromatin state at a single cell level,  
481 so, we can't verify whether these cell types truly can happen oncofetal  
482 reprogramming.

483

484 **Acknowledgments**

485 Funding: This work was supported by the Strategic Priority Research Program  
486 XDB38020500 to L.J., National Key R&D Program of China 2019YFA0801702 to  
487 L.J., the Natural Science Foundation of China 31970760 to L.J.

488

489 **Author contributions**

490 L.J., J.L. and P.P. conceived the study. H.Y. and Y.K. facilitated its designs. P.P.,  
491 Y.K., S.J., and Y.L collected embryos sample. Y.K., S.J., T.Z., and Y.L performed  
492 scRNA-seq and scATAC-seq library construction. H.Y. and N.A. performed the  
493 bioinformatics analyses. N.A., H.Y., Y.K., X.C., L.J., and J.L. interpreted the data.  
494 H.Y., N.A., J.L. and L.J. wrote the paper with the assistance of the other authors.

495

496 **DATA AND CODE AVAILABILITY**

497 The accession number for the sequencing data reported in this paper is submitted to  
498 Genome Sequence Archive for Human (GSA-Human):

499 <https://ngdc.cncb.ac.cn/gsa-human/s/x70211Pp>. The processed files are uploaded to

500 Figshare: <https://figshare.com/projects/HumanProject/122983>. All codes are available  
501 upon reasonable request.

502



503 **Competing Interests statement**

504 The authors declare no competing interests.

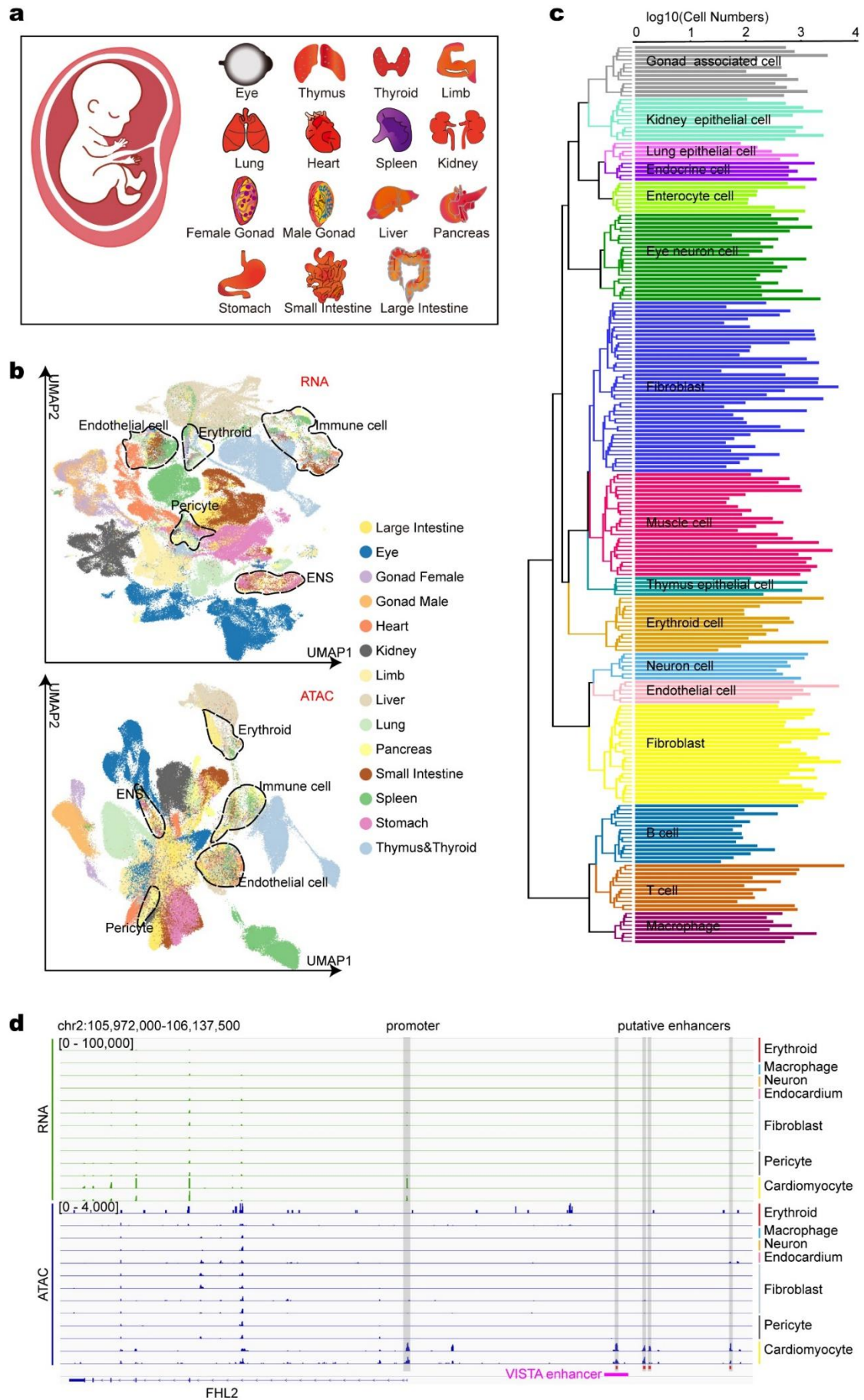
505

506 **Reference**

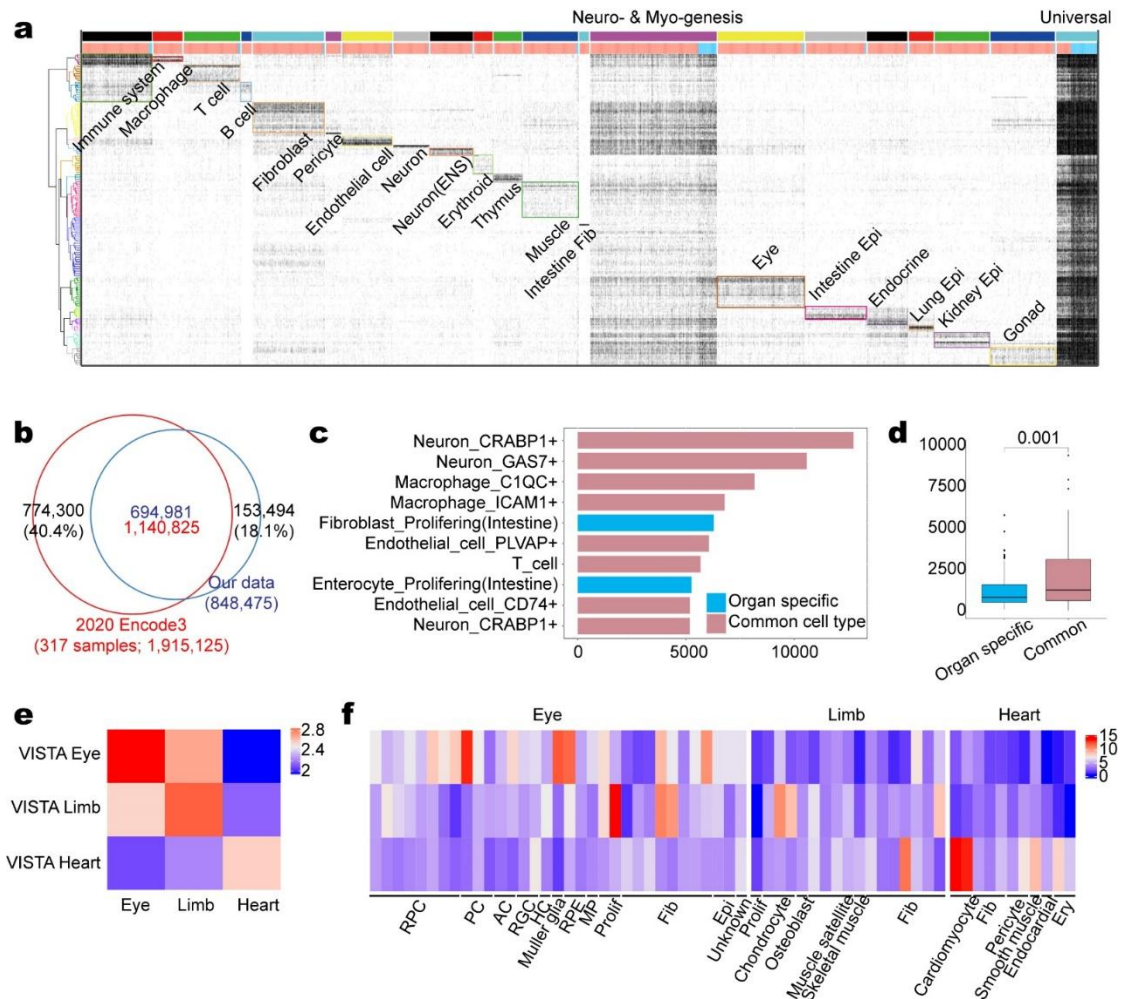
- 507 1. Barker, D.J. The origins of the developmental origins theory. *J Intern Med*  
508 **261**, 412-7 (2007).
- 509 2. Gluckman, P.D., Hanson, M.A., Cooper, C. & Thornburg, K.L. Effect of in  
510 utero and early-life conditions on adult health and disease. *N Engl J Med*  
511 **359**, 61-73 (2008).
- 512 3. Haniffa, M. *et al.* A roadmap for the Human Developmental Cell Atlas.  
513 *Nature* **597**, 196-205 (2021).
- 514 4. Wadhwa, P.D. *et al.* Behavioral perinatology: biobehavioral processes in  
515 human fetal development. *Regul Pept* **108**, 149-57 (2002).
- 516 5. Gerrard, D.T. *et al.* An integrative transcriptomic atlas of organogenesis in  
517 human embryos. *Elife* **5**(2016).
- 518 6. Hochane, M. *et al.* Single-cell transcriptomics reveals gene expression  
519 dynamics of human fetal kidney development. *PLoS Biol* **17**, e3000152  
520 (2019).
- 521 7. Hu, Y. *et al.* Dissecting the transcriptome landscape of the human fetal  
522 neural retina and retinal pigment epithelium by single-cell RNA-seq  
523 analysis. *PLoS Biol* **17**, e3000365 (2019).
- 524 8. Cao, J. *et al.* The single-cell transcriptional landscape of mammalian  
525 organogenesis. *Nature* **566**, 496-502 (2019).
- 526 9. Domcke, S. *et al.* A human cell atlas of fetal chromatin accessibility.  
527 *Science* **370**(2020).
- 528 10. Drokhlyansky, E. *et al.* The Human and Mouse Enteric Nervous System at  
529 Single-Cell Resolution. *Cell* **182**, 1606-1622.e23 (2020).
- 530 11. Han, X. *et al.* Construction of a human cell landscape at single-cell level.  
531 *Nature* **581**, 303-309 (2020).
- 532 12. Corso-Díaz, X., Jaeger, C., Chaitankar, V. & Swaroop, A. Epigenetic control  
533 of gene regulation during development and disease: A view from the  
534 retina. *Prog Retin Eye Res* **65**, 1-27 (2018).
- 535 13. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**,  
536 1888-1902.e21 (2019).
- 537 14. Granja, J.M. *et al.* Author Correction: ArchR is a scalable software package  
538 for integrative single-cell chromatin accessibility analysis. *Nat Genet* **53**,  
539 935 (2021).
- 540 15. Meuleman, W. *et al.* Index and biological spectrum of human DNase I  
541 hypersensitive sites. *Nature* **584**, 244-251 (2020).

- 542 16. Gorkin, D.U. *et al.* An atlas of dynamic chromatin landscapes in mouse  
543 fetal development. *Nature* **583**, 744-751 (2020).
- 544 17. Visel, A., Minovitsky, S., Dubchak, I. & Pennacchio, L.A. VISTA Enhancer  
545 Browser--a database of tissue-specific human enhancers. *Nucleic Acids*  
546 *Res* **35**, D88-92 (2007).
- 547 18. Zhang, Y. *et al.* Intricate gene regulatory networks of helix-loop-helix  
548 (HLH) proteins support regulation of bone-tissue related genes during  
549 osteoblast differentiation. *J Cell Biochem* **105**, 487-96 (2008).
- 550 19. Self, M. *et al.* Six2 is required for suppression of nephrogenesis and  
551 progenitor renewal in the developing kidney. *Embo j* **25**, 5214-28 (2006).
- 552 20. Cusanovich, D.A. *et al.* A Single-Cell Atlas of In Vivo Mammalian  
553 Chromatin Accessibility. *Cell* **174**, 1309-1324.e18 (2018).
- 554 21. Pang, B. & Snyder, M.P. Systematic identification of silencers in human  
555 cells. *Nat Genet* **52**, 254-263 (2020).
- 556 22. Liu, N. *et al.* Transcription factor competition at the  $\gamma$ -globin promoters  
557 controls hemoglobin switching. *Nat Genet* **53**, 511-520 (2021).
- 558 23. De Santa, F. *et al.* A large fraction of extragenic RNA pol II transcription  
559 sites overlap enhancers. *PLoS Biol* **8**, e1000384 (2010).
- 560 24. Kim, T.K. *et al.* Widespread transcription at neuronal activity-regulated  
561 enhancers. *Nature* **465**, 182-7 (2010).
- 562 25. Lambert, S.A. *et al.* The Human Transcription Factors. *Cell* **172**, 650-665  
563 (2018).
- 564 26. Bulik-Sullivan, B.K. *et al.* LD Score regression distinguishes confounding  
565 from polygenicity in genome-wide association studies. *Nat Genet* **47**,  
566 291-5 (2015).
- 567 27. Finucane, H.K. *et al.* Heritability enrichment of specifically expressed genes  
568 identifies disease-relevant tissues and cell types. *Nat Genet* **50**, 621-629  
569 (2018).
- 570 28. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and  
571 genomic data. *Nature* **562**, 203-209 (2018).
- 572 29. Anttila, V. *et al.* Analysis of shared heritability in common disorders of the  
573 brain. **360**(2018).
- 574 30. London, A., Benhar, I. & Schwartz, M. The retina as a window to the  
575 brain--from eye research to CNS disorders. *Nat Rev Neurol* **9**, 44-53  
576 (2013).
- 577 31. Teumer, A. *et al.* Genome-wide association meta-analyses, and  
578 fine-mapping elucidate pathways influencing albuminuria. *Nat Commun*  
579 **10**, 4130 (2019).
- 580 32. Wuttke, M. *et al.* A catalog of genetic loci associated with kidney function  
581 from analyses of a million individuals. *Nat Genet* **51**, 957-972 (2019).
- 582 33. Nayar, S. *et al.* A myeloid-stromal niche and gp130 rescue in  
583 NOD2-driven Crohn's disease. *Nature* **593**, 275-281 (2021).

- 584 34. Hocker, J.D. *et al.* Cardiac cell type-specific gene regulatory programs and  
585 disease risk association. *Sci Adv* **7**(2021).
- 586 35. Wepsic, H.T. Overview of oncofetal antigens in cancer. *Ann Clin Lab Sci* **13**,  
587 261-6 (1983).
- 588 36. Dong, R. *et al.* Single-Cell Characterization of Malignant Phenotypes and  
589 Developmental Trajectories of Adrenal Neuroblastoma. *Cancer Cell* **38**,  
590 716-733.e6 (2020).
- 591 37. Sharma, A. *et al.* Onco-fetal Reprogramming of Endothelial Cells Drives  
592 Immunosuppressive Macrophages in Hepatocellular Carcinoma. *Cell* **183**,  
593 377-394.e21 (2020).
- 594 38. Corces, M.R. *et al.* The chromatin accessibility landscape of primary  
595 human cancers. *Science* **362**(2018).
- 596 39. Gulati, G.S. *et al.* Single-cell transcriptional diversity is a hallmark of  
597 developmental potential. *Science* **367**, 405-411 (2020).
- 598
- 599



601 Fig. 1 | Single-cell transcriptome and chromatin accessibility maps of human early fetus.  
602 a, Schematic of collected tissues.  
603 b, Upper panel: UMAP embedding of all 185,061 cells from the scRNA-seq data. Lower panel:  
604 UMAP embedding of all 212,776 cells from the scATAC-seq data. Each point represents a cell,  
605 colored by organ. Some common cell types across organs are outlined.  
606 c, Dendrogram showing relationships among 225 cell types. The bar chart on the right represents  
607 the number of cells in each cell types in the scATAC-seq data.  
608 d, Example locus around FHL2 with differential expression and accessibility across heart-related  
609 populations. Shadowed regions highlight the identified cis-regulatory elements.  
610



611

612 Fig. 2 | Identifying chromatin accessible sites and patterns in all cell types.

613 a, Chromatin accessibility at 848,475 peaks (x axis) across 225 cell types (y axis). The color code  
614 on top represents 21 LSFs. Orangered/deepskyblue color code represnets TSS distal/proximal  
615 peaks.

616 b, The overlap between DHSs from ENCODE3 paper and our ATAC peaks. DHSs from  
617 corresponding organs/tissues are used for comparison.

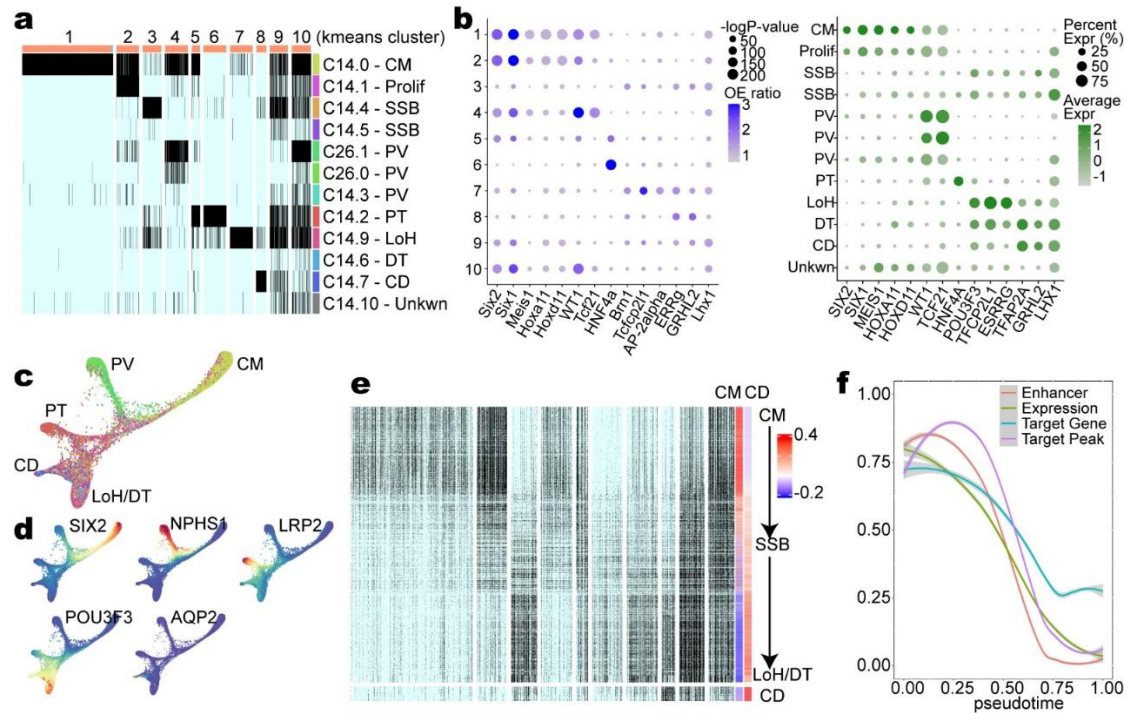
618 c, Top 10 cell types that contribute most to ATAC specific peaks (153,494 in Fig 2B).

619 d, Contribution to ATAC specific peaks stratified by two classes of cell types. Boxes denote  
620 medians and interquartile ranges (IQRs, 25–75%), whiskers represent 1.5 x IQRs.

621 e, Enrichment for VISTA enhancers within ATAC peaks in the corresponding organ.

622 f, Same as Fig. 2e, but in the cell type level. RPC, retinal progenitor cell; PC, photoreceptor cell;  
623 AC, amacrine cell; RGC, retinal ganglion cell; HC, horizontal cell; RPE, retinal pigment  
624 epithelium; MP, fetal mesenchymal progenitor cell; Prolif, proliferating cell; Fib, fibroblast; Epi,  
625 epithelial cell; Ery, erythroblast.

626



627

628 Fig.3 | Dynamics of open chromatin and driving transcription factors in nephrogenesis.

629 a, Sub-patterns of chromatin accessible states in G19 from Fig. 2a. All cell types are kidney  
 630 epithelial cells. CM, cap mesenchyme; Prolif, proliferating cells; SSB, S-shaped body; PV,  
 631 primitive vesicle; PT, proximal tubules; LoH, Loop of Henle; DT, distal tubules; CD, collecting  
 632 duct.

633 b, Left panel: Motif enrichment among 10 K-means clusters. Right panel: Expression level of  
 634 transcription factors among different cell types. The motifs and transcription factors are  
 635 corresponding in position.

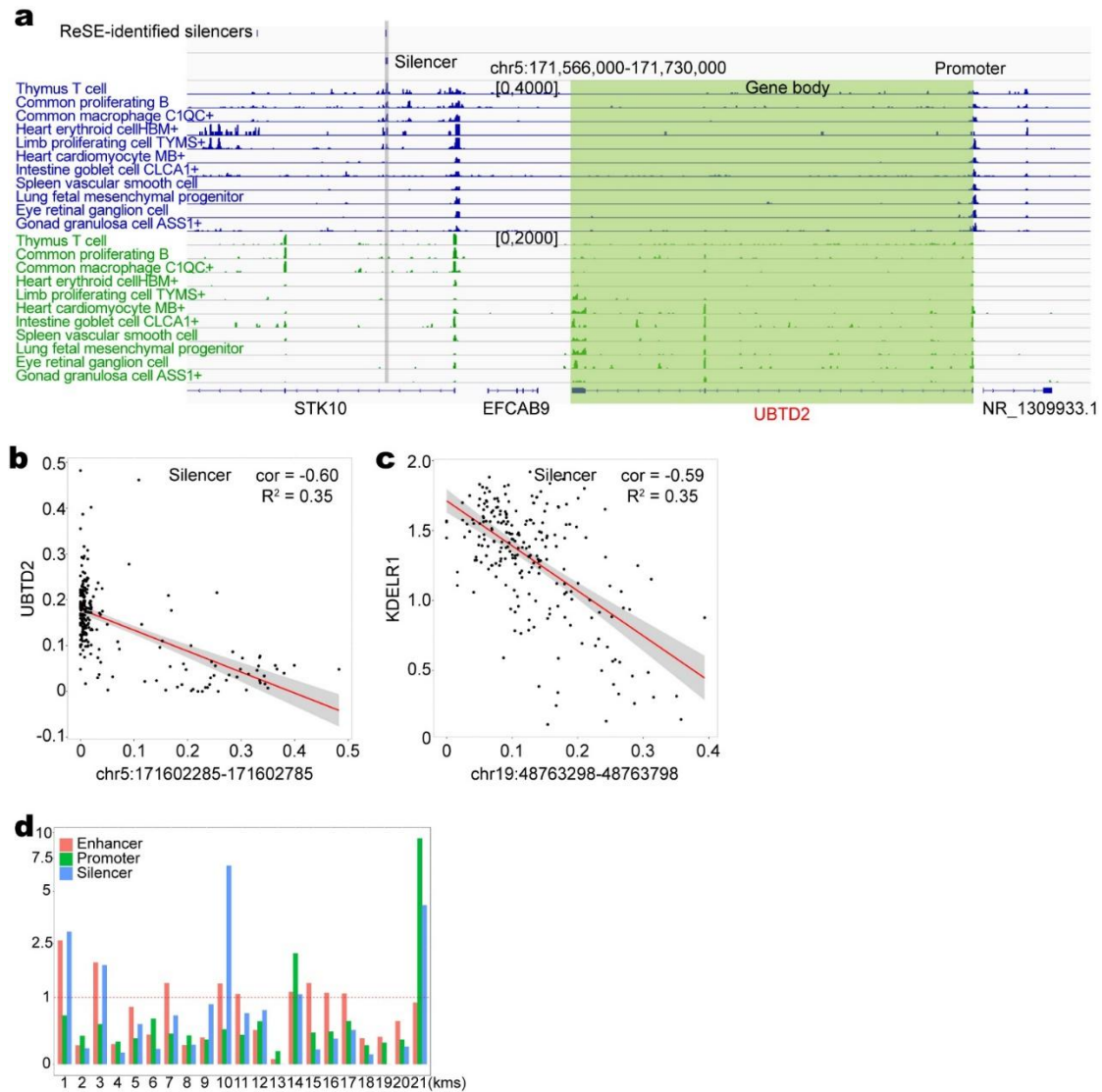
636 c, UMAP embedding of all 12,652 cells from the scATAC-seq data, colored by cell type in Fig.  
 637 3a.

638 d, Normalized gene activity score level of 5 marker genes.

639 e, Continuous change of chromatin accessibility states along differentiation of loop of Helen/  
 640 tubule. Each row represents a cell, which is ordered by pseudo-time. The bottom part is from  
 641 collect ducts as a reference.

642 f, Dynamics of SIX2 expression, chromatin accessibility of its upstream and downstream peaks  
 643 and downstream gene expression.

644



645

646 Fig. 4 | Comparison with ReSE-identified silencers.

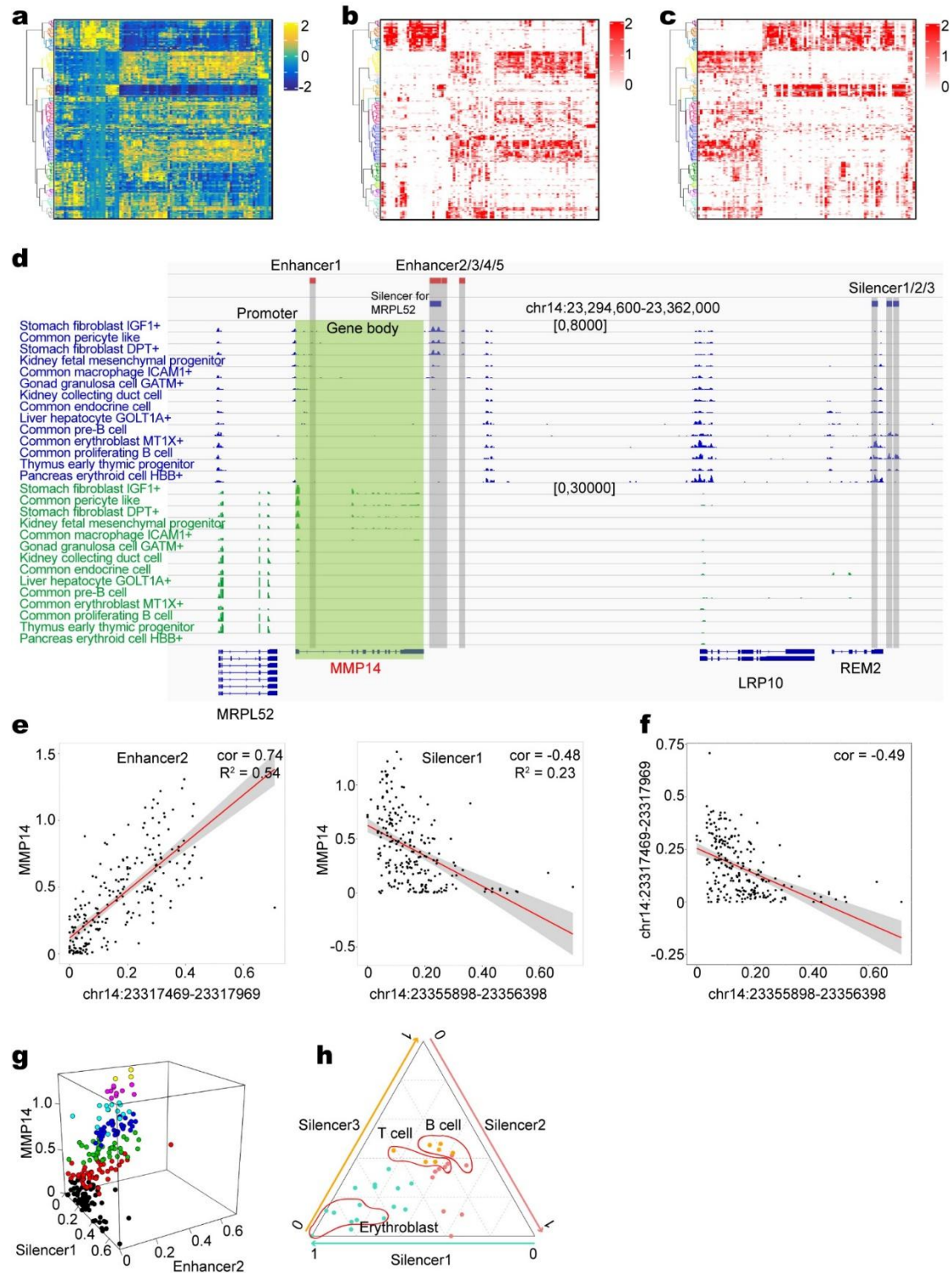
647 a, Example locus around UBTD2 with annotated cis-regulatory elements on the top. Cell types are  
 648 ordered according to the accessibility level of the silencer identified in both study.

649 b, Scatter plot demonstrates the silencer's accessibility level (x axis), along with UBTD2  
 650 expression level (y axis) of each cell type, related to Fig. 4a.

651 c, Scatter plot demonstrates the accessibility level of another overlapped silencer (x axis), along with  
 652 UBTD2 expression level (y axis) of each cell type.

653 d, Enrichment for all annotated cis-regulatory elements in different peak groups from Fig. 2a.



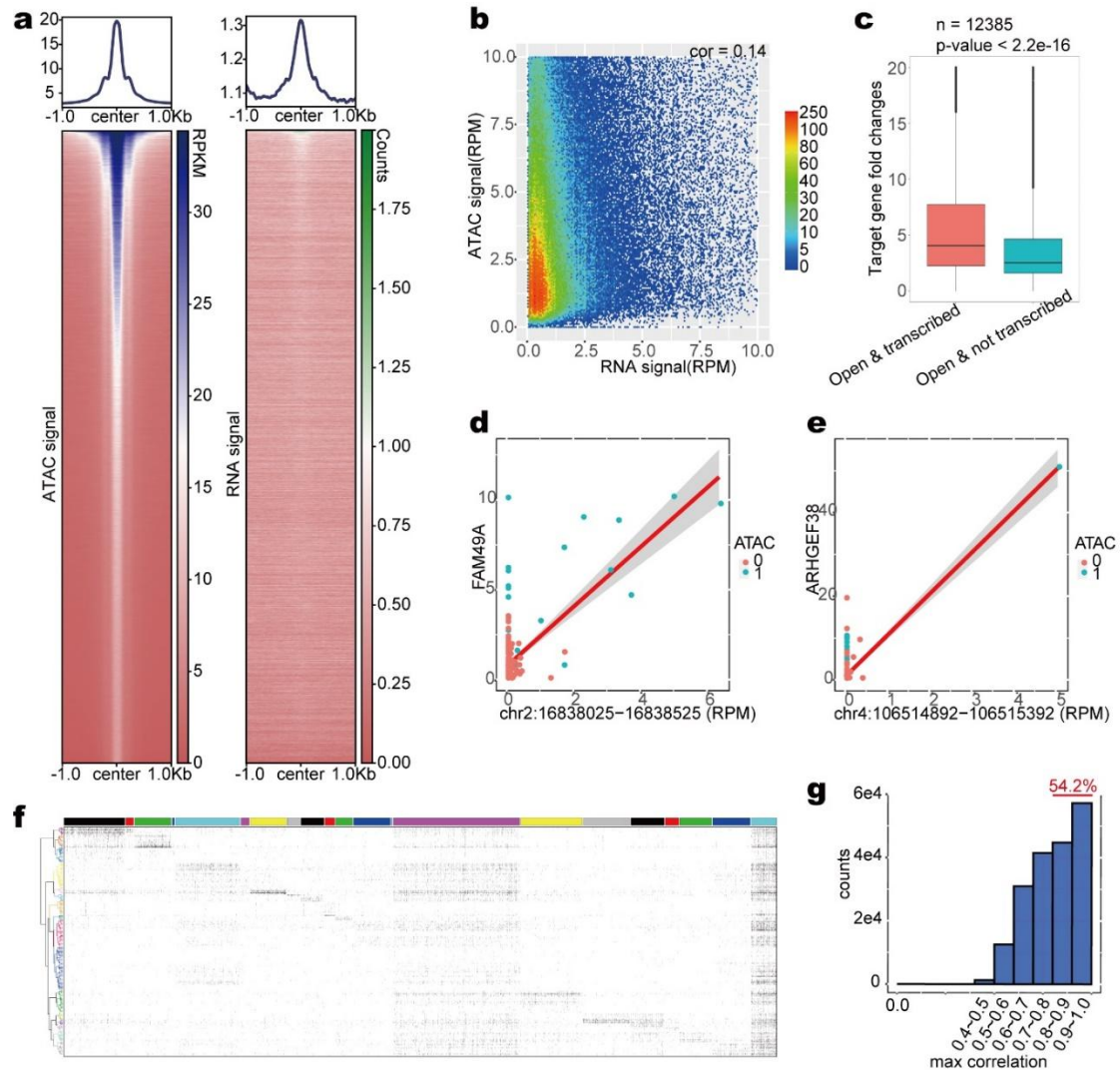


654

655 Fig. 5 | Combinational regulation of positive and negative cis-regulatory elements on the same  
 656 gene.

657 a-c, Heatmaps showing the combinational regulation on the same gene. (a) Relative gene  
 658 expression levels. Each column represents a gene. Each row represents a cell type. (b) Relative  
 659 chromatin accessibility levels of enhancers. (c) Relative chromatin accessibility levels of silencers.  
 660 Each column represents a peak associated with corresponding gene in (a). Median value is used to  
 661 represent multiple peaks linked to the same gene.

662 d, Example locus around MMP14 with annotated cis-regulatory elements on the top. Cell types are  
663 ordered according to expression level of MMP14.  
664 e, Scatter plot demonstrates the peak accessibility level (x axis), along with MMP14 expression  
665 level (y axis) of each cell type. Left is enhancer 2 and right is silencer 1 from Fig. 5d.  
666 f, Scatter plot demonstrates the accessibility level of silencer 1 (x axis), along with the  
667 accessibility level of enhancer 2 (y axis) of each cell type.  
668 g, 3D scatter plot showing the relationship among the accessibility level of enhancer 2, silencer 1  
669 and gene expression of MMP14.  
670 h, Ternary plot showing the silencer preference among different cell types. Only cell types with  
671 normalized expression level less than 0.20 are plotted.  
672



673

674 Fig. 6 | Transcription analysis uncover transcription-dependent and transcription-independent  
 675 enhancers

676 a, Heatmaps showing the ATAC/RNA signal around distal ATAC peaks. (Left) ATAC signal of  
 677 GW10 limb using RPM; (Right) RNA signal of GW10 Limb using number of read counts.

678 b, Smooth scatter plot demonstrates the peak transcription level (x axis), along with ATAC signal  
 679 intensity (y axis) of transcribed cis elements in each cell type. Only cis elements with open  
 680 chromatin state are shown.

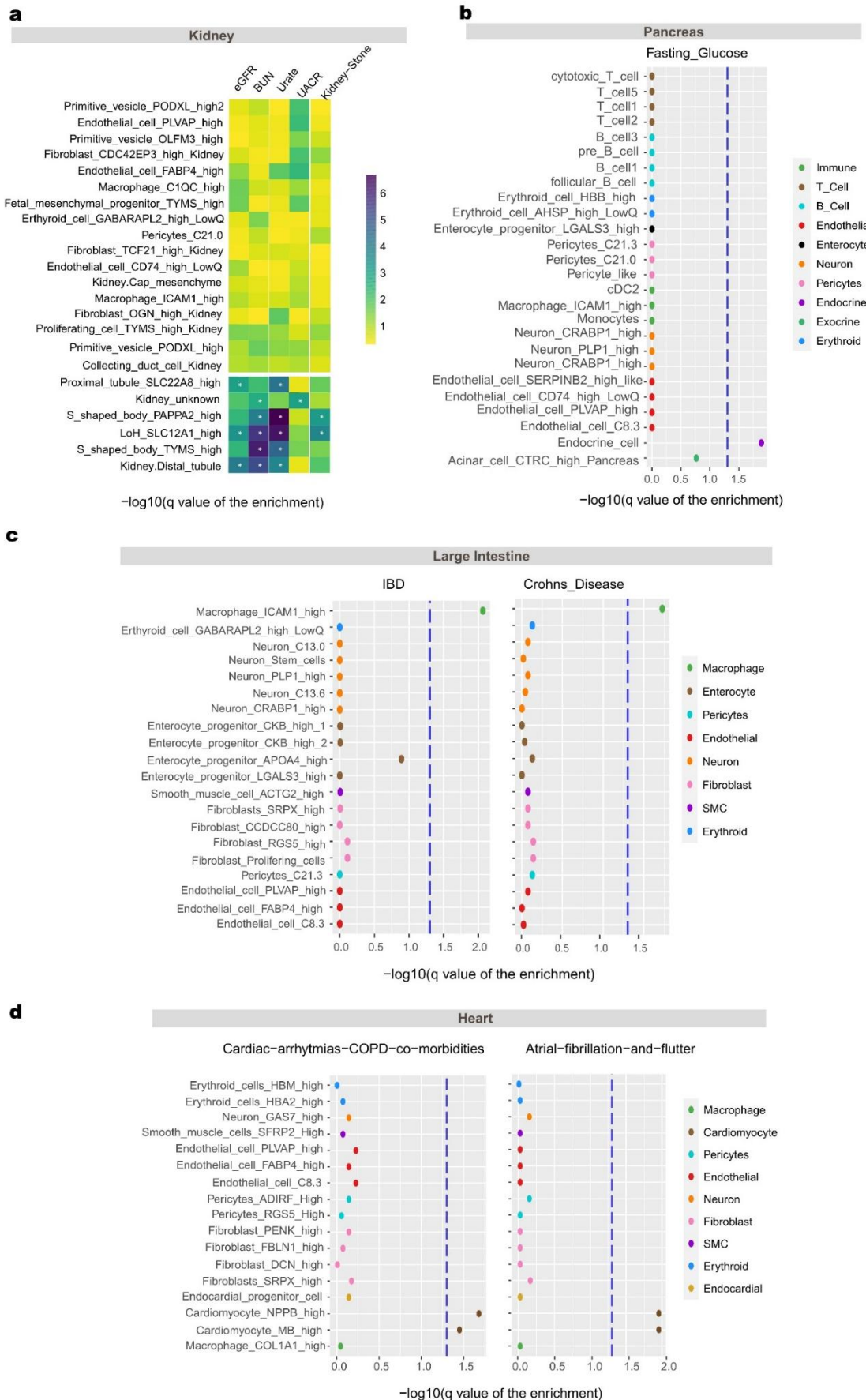
681 c, Relationship between enhancer transcription and target gene expression. Boxes denote medians  
 682 and interquartile ranges (IQRs, 25–75%), whiskers represent 1.5 x IQRs.

683 d,e, Scatter plot demonstrates transcription level of the peak (x axis), along with transcription level  
 684 of target gene (y axis) in each cell type.

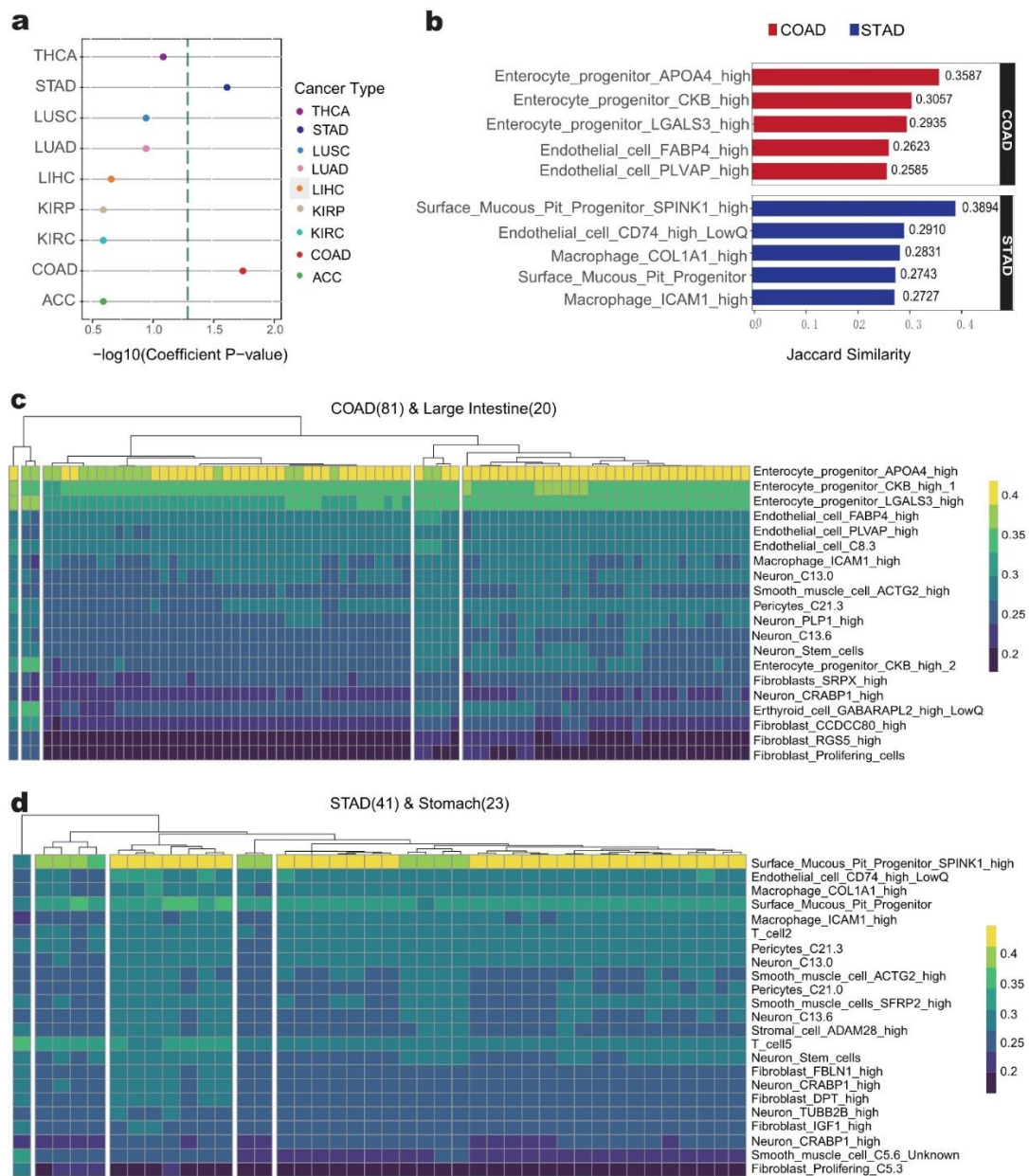
685 f, Transcription/not transcription at 190,356 transcribed cis elements (x axis) across 225 cell types  
 686 (y axis). The color code on top represents 21 accessibility patterns.

687 g, Frequency distribution of max correlation of co-expressed cis elements for each transcribed cis  
 688 element.

689



691 Fig. 7 | Enrichment analysis of GWAS signals in cell-type-specific chromatin regions.  
692 a, S-LDSC results suggests these disease and traits' susceptibility and heritability are cell type  
693 specificity.  
694 Heatmap show cell-type-specific enrichments of the heritability signal for kidney stone and CKD  
695 diseases in kidney tissue, significance level ( $q < 0.05$ ) are indicated with an asterisk.  
696 b, Dot plot show cell-type-specific enrichments of the heritability signal (y axis) for diabetes in  
697 pancreas tissue, the blue dotted line indicates significant threshold (q value of 0.05).  
698 c, Dot plot show the cell-type-specific enrichments of the heritability signal for two typical  
699 inflammatory bowel diseases across all cell types in large intestine.  
700 d, Dot plot show cell-type-specific enrichments of the heritability signal for heart traits in heart  
701 tissue.  
702



703

704 Fig. 8 | Link fetal cell type with cancer at chromatin stat level.

705 a, Rank of tumor type relevance to proliferative state cell types based on hypergeometric test.

706 x-axis is the  $-\log_{10}(p \text{ value})$ , blue dotted line is p value of 0.05.

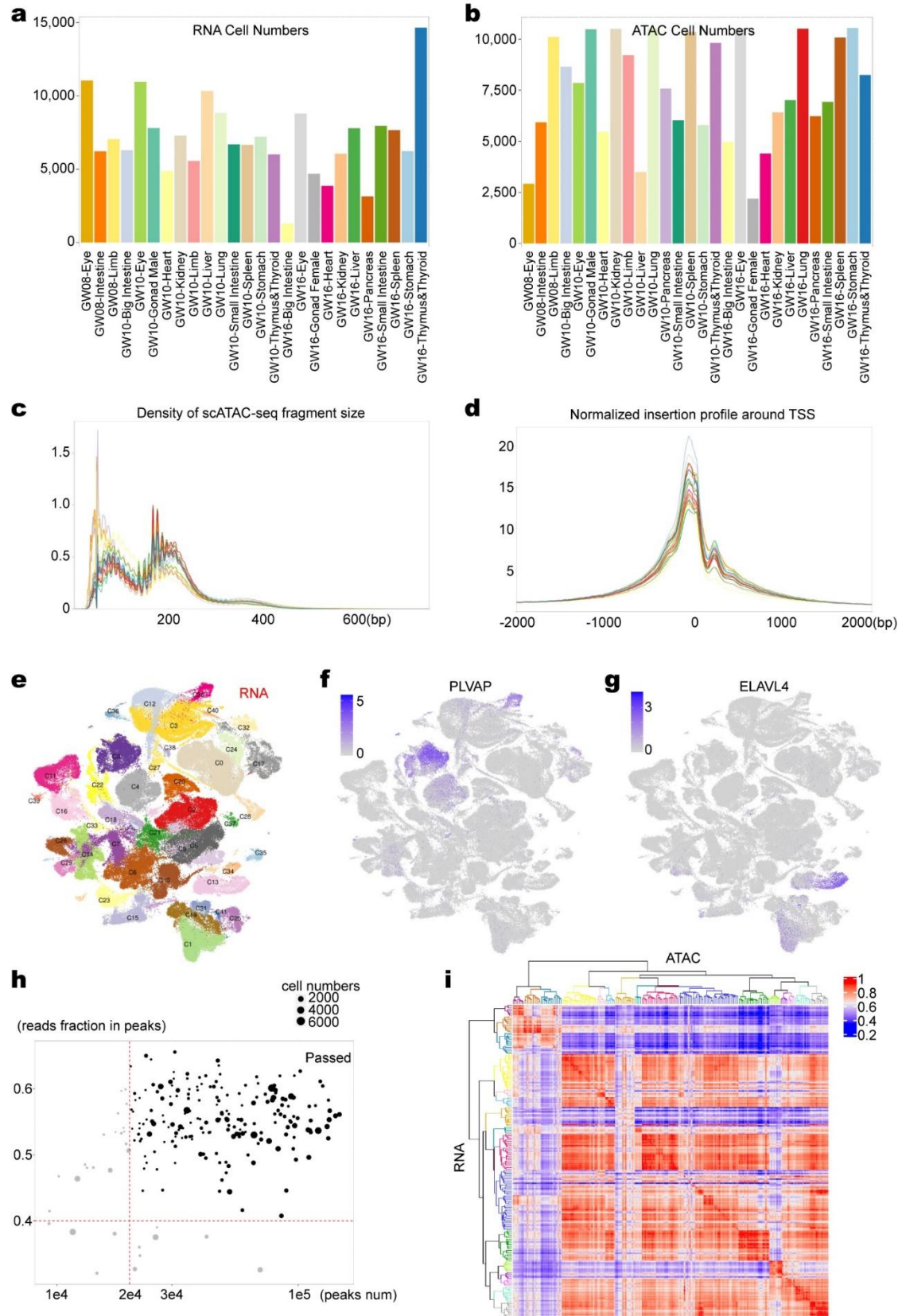
707 b, Bar plot showing the Jaccard similarity Score of top5 similar fetal cell types for Colon  
708 Adenocarcinoma (COAD) and Stomach Adenocarcinoma (STAD).

709 c, Jaccard similarities of chromatin state from 81 Colon Adenocarcinoma individuals (y axis) with  
710 the cis-elements of 20 cell types in large intestine (x axis).

711 d, Jaccard similarities of chromatin state from 41 Stomach Adenocarcinoma (STAD) individuals  
712 (y axis) with the cis-elements of 23 cell types in stomach (x axis).

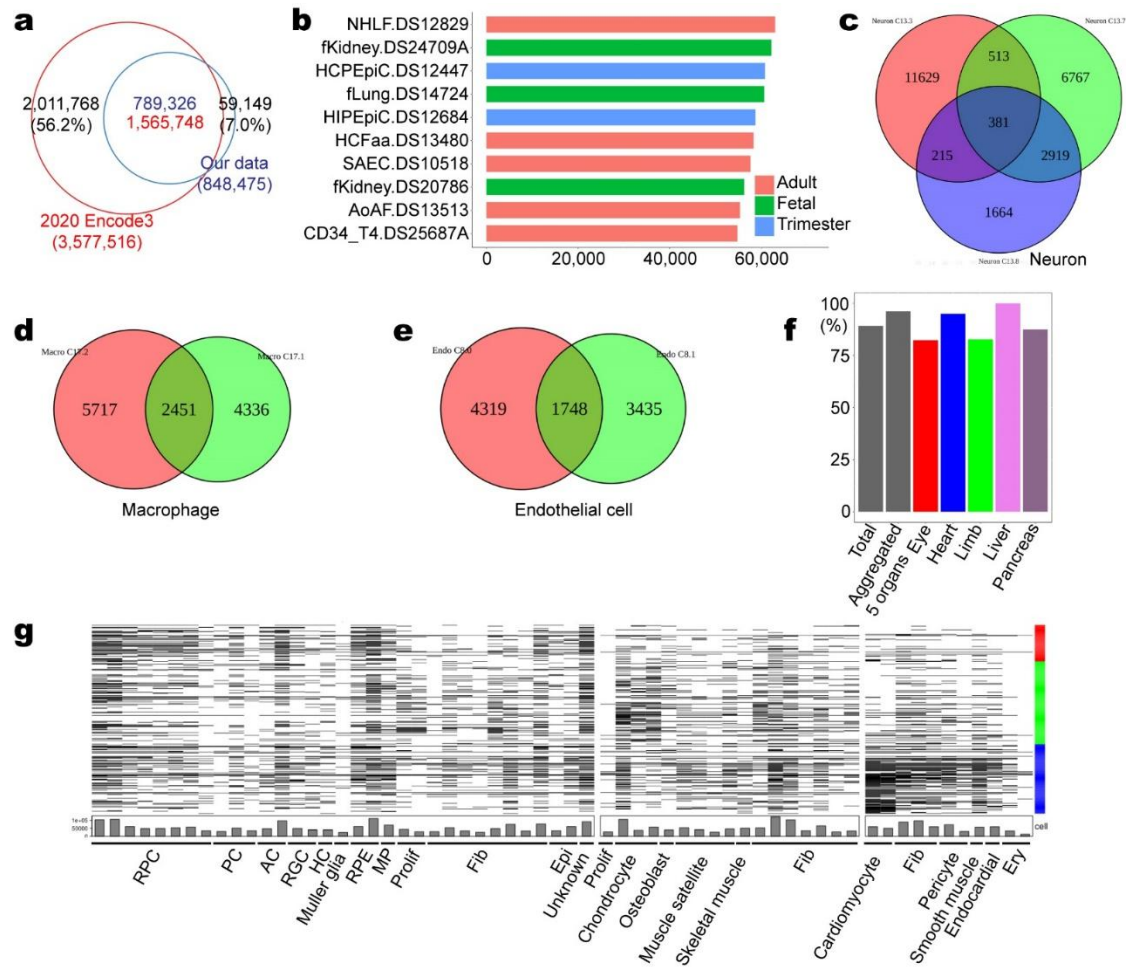
713

714



719 (c,d) Left panel: Distribution of sequenced insert sizes for each sample. Right panel: Normalized  
720 insertion profile around TSS for each sample.  
721 (e-g) (e) UMAP embedding of all 185,061 cells from the scRNA-seq data colored by 41 major  
722 clusters. (f) Normalized gene expression level of PLVAP. (g) Normalized gene expression level of  
723 ELAVL4.  
724 (h) QC of label transferring result. Bubble plot demonstrates the significant peak numbers (x axis),  
725 along with read fraction in peaks (y axis) of each cell type in scATAC-seq data. Black dots  
726 represent the cell types passing the QC filters.  
727 (i) Heatmap of spearman correlations between average gene activity score profiles (x axis) and  
728 gene expression profiles (y axis) for 225 cell types. The cell type order is the same as Fig. 1c.  
729





730

731

732 Supplementary Fig. 2 | Comparison of chromatin accessible sites, related to Figure 2

733 (a) The overlap between DHSs from ENCODE3 paper and our ATAC peaks. Using all DHSs.

734 (b) Top 10 tissues that contribute most to DHSs specific peaks (774,300 in Fig. 2b).

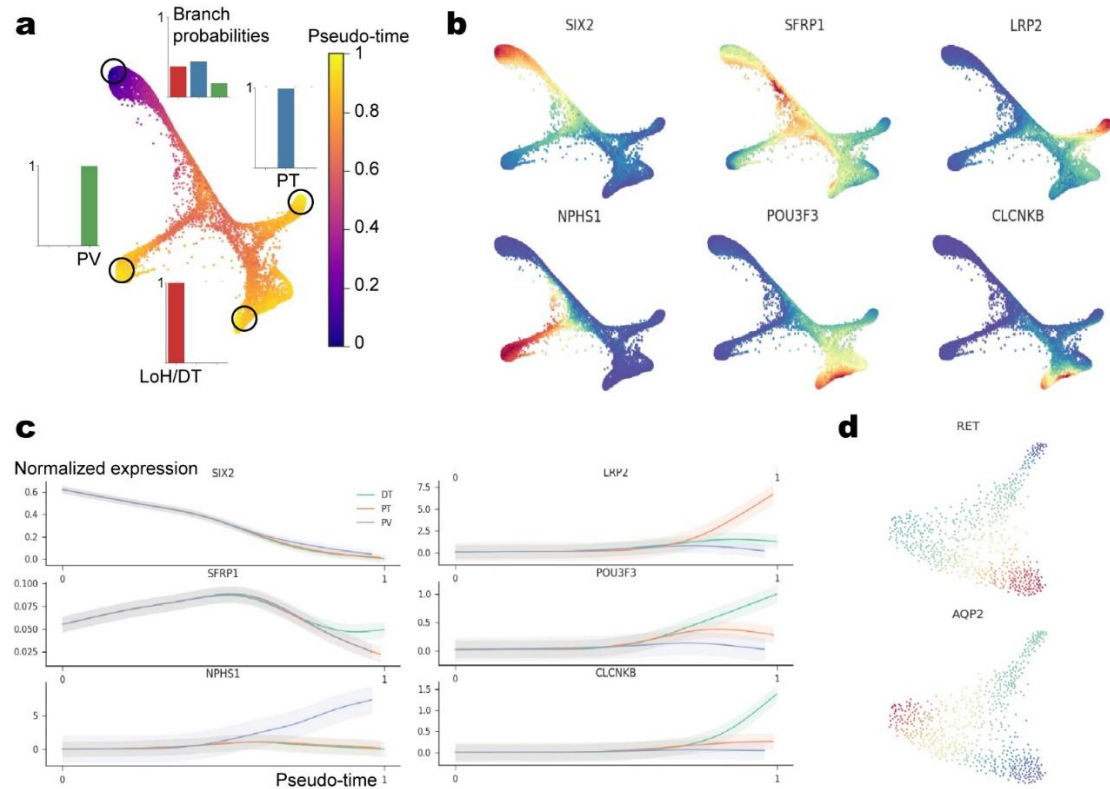
735 (c-e) Overlaps of ATAC specific peaks among sub cell types from Fig. 2c.

736 (f) Coverage of VISTA enhancers in different sets.

737 (g) Accessibility of VISTA enhancers among different cell types. Each row represents an enhancer,

738 and each column represents a cell type. The color code on right represents organ source from

739 Supplementary Fig. 2d.



740

741

742 Supplementary Fig. 3 | Trajectory analysis of kidney epithelial cells, related to Figure 3

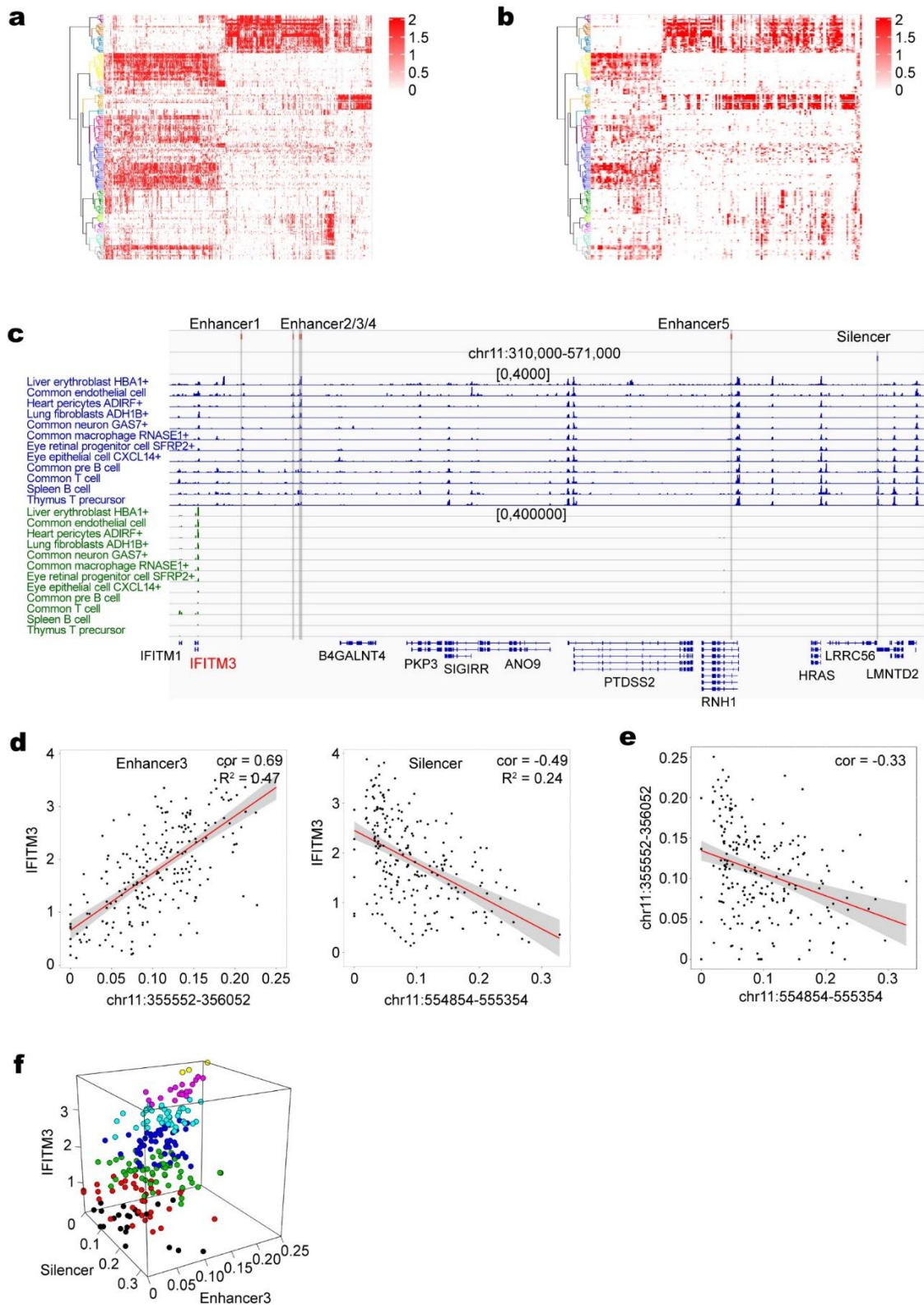
743 (a) UMAP embedding of CM derived 12,048 cells from the scATAC-seq data colored by  
744 pseudo-time. The bar chart shows the terminal state probability distributions of three selected  
745 cells.

746 (b) Normalized gene expression level of previous known markers.

747 (c) Expression pattern of previous known marker genes in each segment along the pseudo-time  
748 path.

749 (d) UMAP embedding of UB derived 604 cells from the scATAC-seq data colored by expression  
750 of UB markers (top) and CD marker (bottom).

751

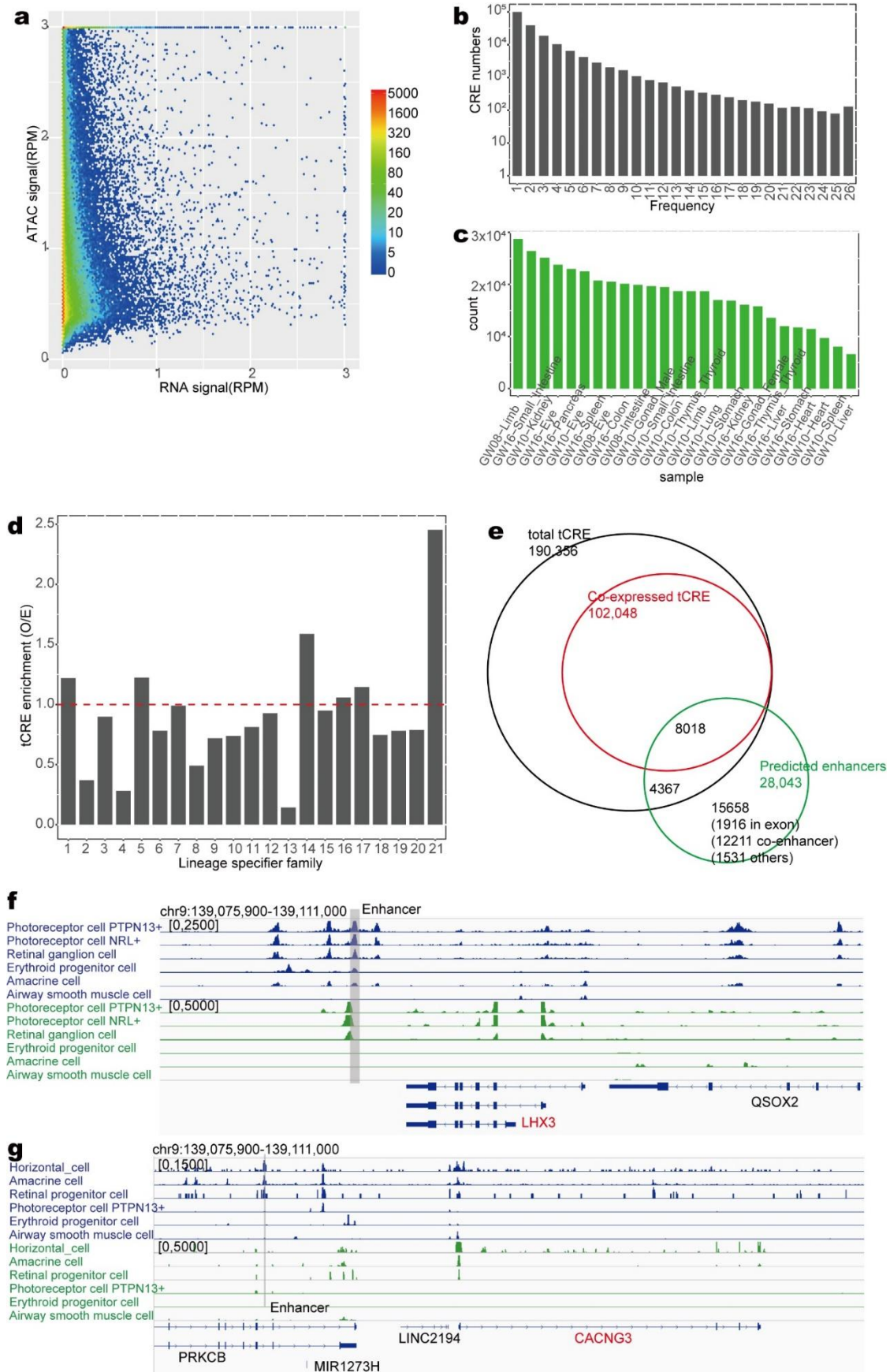


752

753 Supplementary Fig. 4 | Accessibility of positive and negative cis-regulatory elements, related to  
754 Figure 5

755 (a-b) Heatmaps showing the combinational regulation on the same gene. Same as Fig. 5b,c, but  
756 using all enhancers or silencers.

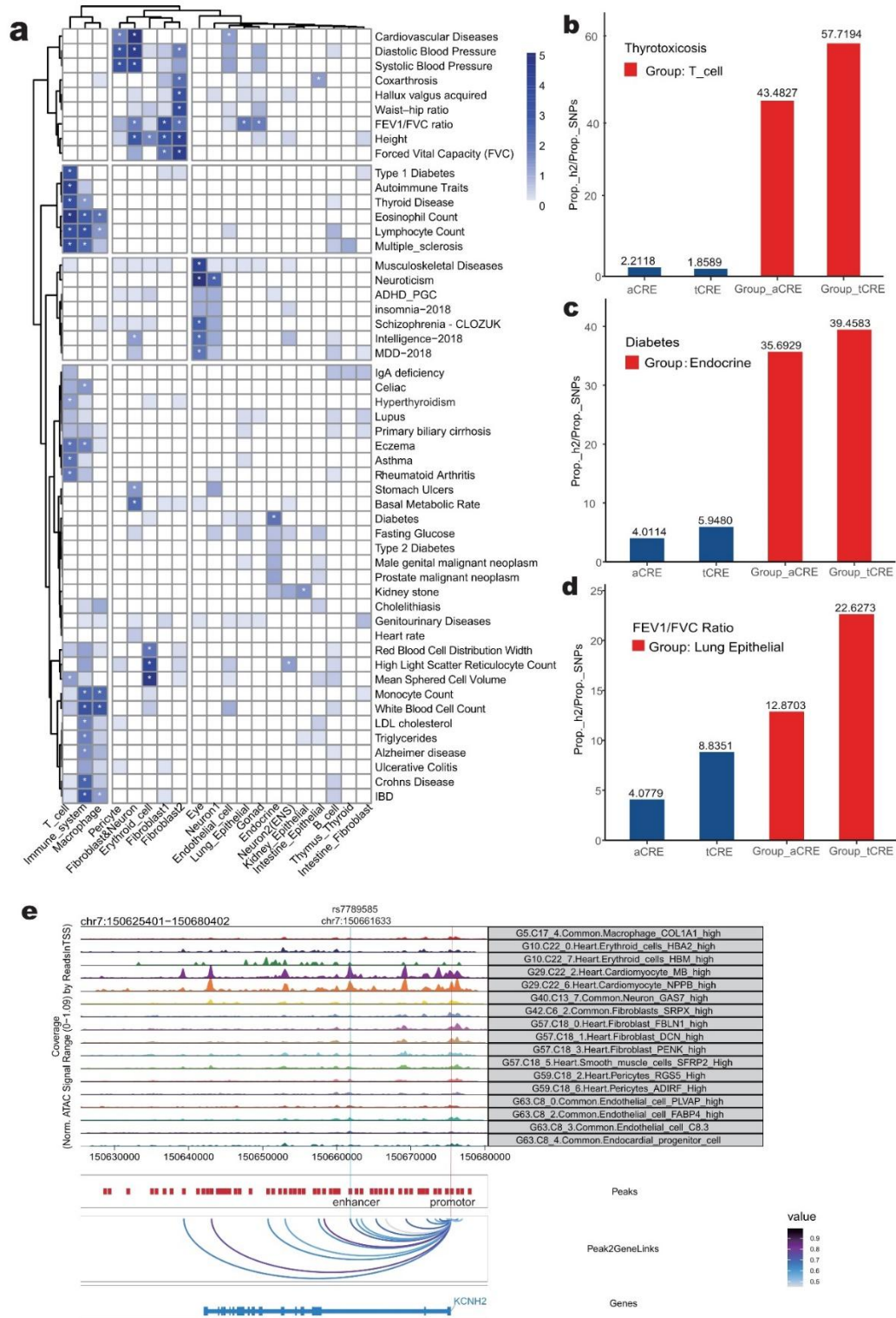
- 757 (c) Example locus around IFITM3 with annotated cis-regulatory elements on the top. Cell types  
758 are ordered according to expression level of IFITM3.
- 759 (d) Scatter plot demonstrates the peak accessibility level (x axis), along with IFITM3 expression  
760 level (y axis) of each cell type. Left is enhancer 3 and right is silencer from (c).
- 761 (e) Scatter plot demonstrates the accessibility level of the silencer (x axis), along with the  
762 accessibility level of the enhancer 3 (y axis) of each cell type.
- 763 (f) 3D scatter plot showing the relationship among the accessibility level of the enhancer 3,  
764 silencer and gene expression of IFITM3.



765

766 Supplementary Fig. 5 | Assessing properties of transcribed cis elements, related to Figure 6.

- 767 (a) Smooth scatter plot demonstrates the peak transcription level (x axis), along with ATAC signal  
768 intensity (y axis) of transcribed cis elements in GW10 limb. Only cis elements with open  
769 chromatin state are shown.
- 770 (b) Frequency distribution of transcribed cis elements in all samples.
- 771 (c) Counts of identified transcribed cis elements in each sample.
- 772 (d) Enrichment for transcribed cis elements in different peak groups from Fig 2A.
- 773 (e) The overlap between transcribed cis elements, co-expressed cis elements and putative  
774 enhancers from peak-to-gene links.
- 775 (f) Example locus of transcription-dependent enhancer of LHX3 with annotated cis-regulatory  
776 elements on the top. Cell types are ordered according to expression level of LHX3.
- 777 (g) Example locus of transcription-independent enhancer of CACNG3 with annotated  
778 cis-regulatory elements on the top. Cell types are ordered according to expression level of  
779 CACNG3.
- 780



781

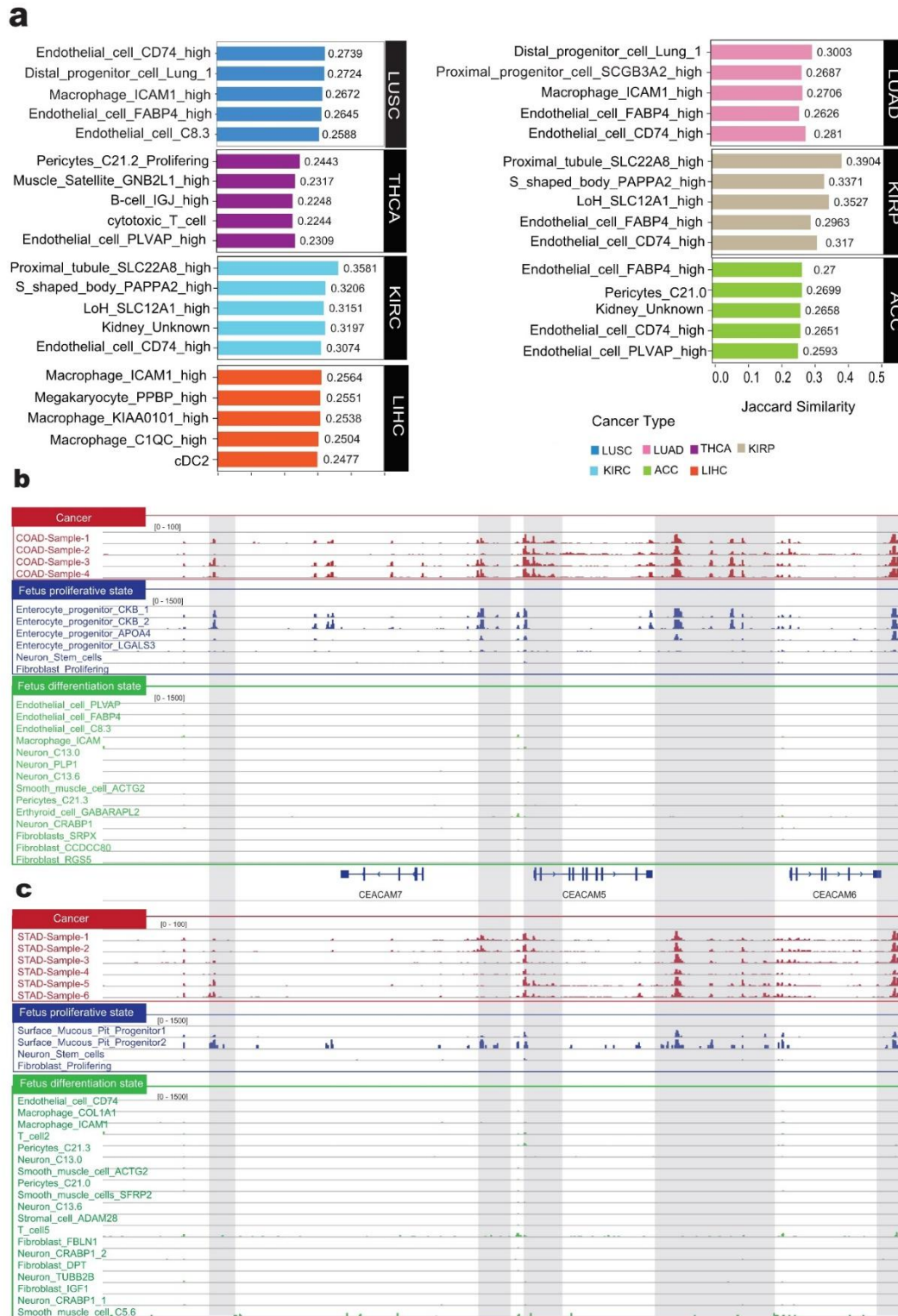
782 Supplementary Fig. 6 | S-LDSC results from 52 traits show heritability enrichment in 20 LSFs.

783 (a) Heatmap displaying the  $-\log_{10}(q)$  value of the enrichment for 20 peak groups across 52 traits  
784 analyzed (Except LSF 21). 20 LSFs were classified and colored by broader cell-type category, that  
785 met the across 20 LSFs, significance level ( $q < 0.05$ ) are indicated with an asterisk.

786 (b-d) Bart plots displaying Enrichment of heritability in various CRE-types.

787 X-axis, from left to right are 'aCRE' represent total cis-elements detected (840K); 'tCRE' represent  
788 total transcribed cis-elements (190K); 'Group\_aCRE' represent specific peak group of total  
789 cis-elements; Group\_tCRE' represent specific peak group of transcribed cis-elements. Y-axis,  
790 Heritability enrichment  $\Pr(h^2)/\Pr(\text{SNPs})$ , estimated by LDSC. Red bar shows heritability  
791 enrichment of assigned group peak, the blue bar shows bulk level.  
792 (e) Genome browser tracks for scATAC-seq (top; scale, RPM) and indicated one AF-associated  
793 risk variant. Co-accessibility track shows linkages between the AF variant-containing CRE and  
794 promoters.  
795





796

797 Supplementary Fig. 7 |

798 (a) Bar plot showing the Jaccard similarity Score of top5 similar fetal cell types for 7 cancer types  
799 (See Supplement table).

800 (b) Regulatory landscape around the CEACAM family genes (CEACAM5, CEACAM6,

801 CEACAM7), indicating GENCODE gene annotations, ATAC seq tracks for each cell type of

802 Colon (blue and green), and top 6 from COAD sample (Red). Fetal cell types in colon have been  
803 classified two parts (See Method), those cell types which are labeled by green color belong to  
804 differentiation state cell types, other blue cell types are proliferative state cell types.

805 (c) The Same as FigureS7B, highlight the chromatin profile between COAD and Surface Mucous  
806 Pit Progenitor cells.

807

## 808 Supplementary information

809

810 Legends for Supplementary Files

811

812 **File S1** | Metadata of cells in scRNA-seq data. Includes sample metadata, per-cell QC  
813 stats, cluster id and cell type annotation.

814 **File S2** | Metadata of cells in scATAC-seq data. Includes sample metadata, cell type  
815 annotation and various Cell\_ID information for each software.

816 **File S3** | Gene count matrix of cells in scRNA-seq data in RDS format. Includes  
817 expression UMI values for each gene in each cell.

818 **File S4** | Peak count matrix of cells in scATAC-seq data in RDS format. Includes  
819 insertion counts within each peak in each cell, while the maximum value was set to 4.

820 **File S5** | Normalized peak by cell type matrix in RDS format. Includes normalized  
821 peak accessible values (reads per million reads/100) for each cell type.

822 **File S6** | Binary peak by cell type matrix in RDS format. Includes binary values for  
823 each cell type, where 1 denotes accessible and 0 denotes inaccessible.

824 **File S7** | Seurat object of 185,061 high-quality cells in scRNA-seq data. Includes  
825 count matrix, low-dimension embedding and cell informations from the global  
826 perspective.

827 **File S8** | Seurat object of average profiles of 335 cell types in scRNA-seq data.

828 **File S9** | tCRE transcription intensity matrix of each cell type in scRNA-seq data in  
829 RDS format. Includes RPM value for each peak in each cell type.

830 **File S10** | 3D animated scatter plot representing relationship between gene expression  
831 level and enhancer/silencer activity in gif format. Pattern 1 is related to Fig. 5g, while  
832 pattern 2 is related to Supplementary Fig. 4f.

833

834 The Supplementary Table S1 can be downloaded from

835 <https://figshare.com/ndownloader/files/30790600>

836 The Supplementary Table S2 can be downloaded from

837 <https://figshare.com/ndownloader/files/30790603>

838 The Supplementary Table S3 can be downloaded from

839 <https://figshare.com/ndownloader/files/30790606>

840 The Supplementary Table S4 can be downloaded from

841 <https://figshare.com/ndownloader/files/30790609>

842 The Supplementary Table S5 can be downloaded from

843 <https://figshare.com/ndownloader/files/30790612>

844 The Supplementary Table S6 can be downloaded from

845 <https://figshare.com/ndownloader/files/30790615>

846 The Supplementary Table S7 can be downloaded from

847 <https://figshare.com/ndownloader/files/30790588>

848 The Supplementary Table S8 can be downloaded from

849 <https://figshare.com/ndownloader/files/30790591>

850 The Supplementary Table S9 can be downloaded from

851 <https://figshare.com/ndownloader/files/30790594>

852 The Supplementary Table S10 can be downloaded from

853 <https://figshare.com/ndownloader/files/30790597>

854

855 These processed files are also uploaded to Open Archive for Miscellaneous Data  
856 (OMIX) database: <http://ngdc.cncb.ac.cn/omix/preview/MCawh0yL>.

857

858

## 859 **METHOD DETAILS**

### 860 **Tissue acquisition and processing**

861 The study of human embryos was approved by the Reproductive Study Ethics  
862 Committee in Peking Union Medical College Hospital, Beijing, China. All tissue  
863 samples used for this study were obtained with written informed consent from all  
864 participants. Samples from surgically removed aborted fetal tissues were collected  
865 into Leibovitz's L-15 (11415064, Gibco) plus with 10% fetal bovine serum (FBS)  
866 right after resection and immediately transported on ice from hospital to the  
867 laboratory in less than 1 h.

868 We collected 4 individual ranging from: 6 PCW (post conception weeks), 10 PCW to  
869 16 PCW and a total of 28 samples (15 organs or tissues): spleen, pancreas, liver,  
870 thymus, thyroid, lung, stomach, small intestine, big intestine, kidney, male gonad,  
871 female gonad, fore-limb, heart, and eye were including (Supplementary Table S1).  
872 Each organ was dissected and washed with DPBS twice, then collected in 1.5 mL EP  
873 tubes.

### 874 **Single cell preparation and Nuclei Isolation**

875 Tissues were minced into pieces (~1 mm) on ice using scissors, and digested into  
876 single-cell suspensions with 1 mg/ml type II collagenase (17101015, GIBCO) and 1  
877 mg/ml type IV collagenase (17104019, GIBCO) for 30min at 37 °C with intermittent  
878 shaking. The dissociated cells were separated and remaining undigested tissue were

879 digested again with fresh digestion buffer. Digested suspension was passed through  
880 70um strainer (Biologix).

881 Dissociated cells were centrifuged at 300 g for 5 min at 4 °C, then re-suspended in 1  
882 mL of cold DPBS with 0.1% BSA. After passing through a 40um cell strainer  
883 (Biologix), cells were washed twice, centrifuged at 300 g for 5 min at 4 °C,  
884 re-suspended in cold DPBS with 0.1% BSA at a density of  $1 \times 10^5$  cells/ml, and stored  
885 on ice before scRNA-Seq and nuclei isolation.

886 To isolate nuclei, the half of the cell pellets were re-suspended in 100 uL chilled lysis  
887 buffer (10 mM Tris-HCl, pH 7.4, 10 mM NaCl, 3 mM MgCl<sub>2</sub>, 0.1% NP40, 0.1%  
888 Tween-20, and 0.01% digitonin - from <sup>1</sup> supplemented with 1% BSA), and pipette  
889 mix 10X. After incubation for 5 min on ice, add 1 ml chilled Wash Buffer ((10 mM  
890 Tris-HCl, pH 7.4, 10 mM NaCl, 3 mM MgCl<sub>2</sub>, 0.1% Tween-20, 1% BSA) to the  
891 lysed cells. Pipette mix 5x, then centrifuged at 300 g for 5 min at 4 °C. Based on  
892 number of cells used for isolation and assuming ~50% nuclei loss during cell lysis,  
893 resuspend in chilled Diluted Nuclei Buffer (PN-2000153, 10x Genomics). If cell  
894 debris and large clumps are observed, pass through a cell strainer. For low volume,  
895 use a 40 μm Flowmi Cell Strainer (H13680-0040, Bel-Art) to minimize volume loss.

#### 896 **sc-RNA-seq Libraries Construction and sequencing**

897 Single cell RNA-seq was performed using the Single Cell 5' RNA Reagent Kits (10x  
898 Genomics, Pleasanton, California) according to the manufacturer's instruction. The  
899 aimed target cell recovery for each library was ~9,000 cell per sample. In brief,

900 cellular suspensions were loaded on the sample chip in the Chromium Controller  
901 instrument (10X Genomics) to generate single-cell Gel Bead-In-Emulsions (GEMs).  
902 GEM-reverse transcription (RT) was performed in a Veriti 96-well thermal cycler  
903 (BioRad, 1851197). After RT, GEMs were harvested and the cDNAs were amplified,  
904 and cleaned up with SPRIselect Reagent Kit (Beckman Coulter, Pasadena, CA).  
905 Indexed sequencing libraries were constructed using Chromium Single-Cell 3' Library  
906 Kit or Single Cell 5' Library kit based for enzymatic fragmentation, end-repair,  
907 A-tailing, adaptor ligation, ligation cleanup, sample index PCR, and PCR cleanup.  
908 Libraries were quantified using Bioanalyzer (Agilent) and QuBit (Thermofisher)  
909 analysis and then sequenced in NovaSeq 6000 (Illumina, San Diego, CA) with a  
910 150-bp paired-end read length, targeting a depth of 50,000–100,000 reads per cell.

#### 911 **sc-ATAC-seq Libraries preparation and sequencing**

912 The scATAC library was prepared using the 10x Genomics platform with the  
913 Chromium Single Cell ATAC Library & Gel Bead Kit (10x Genomics, Pleasanton,  
914 California) as instructed by the manufacturer. A total of 15,000 nuclei per sample  
915 were used as input for single-cell ATAC-seq following the manufacturer's  
916 instructions. Briefly, after tagmentation, the cells were loaded on a Chromium  
917 Controller Single-Cell instrument to generate single-cell Gel Bead-In-Emulsions  
918 (GEMs) followed by linear PCR as described in the 10X scATAC-seq protocol using  
919 a Veriti 96-well thermal cycler (BioRad, 1851197). After breaking the GEMs, the  
920 barcoded tagmented DNA was purified with SPRIselect Reagent Kit (Beckman

921 Coulter, Pasadena, CA) and further amplified to enable sample indexing and  
922 enrichment of scATAC-seq libraries. The final libraries were quantified using  
923 Bioanalyzer (Agilent) and QuBit (Thermofisher) analysis and then sequenced in  
924 Nextseq 550AR or NovaSeq 6000 (Illumina, San Diego, CA) with a 50-bp paired-end  
925 read length, or MGISEq-2000FCL (MGI Tech Co., Ltd., China) with 100-bp  
926 paired-end read length targeting a depth of 30,000–50,000 reads per cell.

### 927 **scRNA-seq Data processing**

928 FASTQ files generated from sequencing were used as inputs to the 10X Genomics  
929 Cellranger (3.1.0) RNA pipeline using default arguments. Briefly, de-multiplexed  
930 reads were mapped to the hg19 genome by STAR. Filtered feature-barcode matrix  
931 containing feature, barcode list and matrix was generated and as input to Seurat  
932 (version 3.2.3). Cells with low complexity (fewer than 400 expressed genes) were  
933 excluded; cells with mitochondrial read fraction outside 10 percent were also cleared  
934 out. The Seurat (version 3.2.3) workflow were run separately on each sample, most of  
935 these parameters have default setting, and the resulting files were used for further  
936 processing. Doublet was estimated for each 10x sample by applying the  
937 ‘doubletFinder\_v3’ function in the DoubletFinder package (version 2.0.2), which is  
938 implemented to interface with Seurat. This function predicts doublets according to  
939 each real cell’s proximity in gene expression space to artificial doublets created by  
940 averaging the transcriptional profile of randomly chosen cell pairs.

### 941 **scRNA-seq clustering and cell type annotation**



942 Each dataset was integrated together using the ‘merge’ function in the Seurat  
943 package. High quality cells from all samples were merged and normalized  
944 (normalization.method = "LogNormalize", scale.factor = 10000). Highly variable  
945 genes (HVGs) had significantly variance were retained (selection.method = "vst",  
946 nfeatures = 10000). Notably, we regressed out the difference between the G2M and S  
947 phase scores (vars.to.regress = S.Score - G2M.Score) to mitigate the effects of cell  
948 cycle heterogeneity in scRNA-seq data. Next, batch effects were removed by  
949 harmony on 75 principal components computed from the HVGs only. Correction was  
950 performed between the samples of each time point, this method was carried out on the  
951 whole atlas dataset, and Harmony embeddings calculated from this batch-corrected  
952 principal component analysis were used for all further analysis steps. We used  
953 shared-nearest-neighbours (SNN) and Louvain method to cluster cells and identified  
954 42 distinct major clusters (dims = 1:75 and resolution = 0.3). To identify finer  
955 substructure from these major clusters, each cluster underwent a second round of  
956 clustering using the same methods as above with resolution range from 0.2 to 0.6,  
957 respectively. We further remove 3192 cells from 21 sub clusters with doublet ratio  
958 previous calculated higher than 55%. Finally, we identified a total of 331 sub clusters.  
959 Differential expression analysis for each cluster was performed by using the  
960 “FindAllMarkers” function with default Wilcoxon rank-sum test. Cell types were  
961 assigned to each sub cluster based on the enrichment of cell type of Human Cell

962 Landscape (HCL) and the expression of known marker genes. Details of cell type  
963 annotation information are listed in Supplementary Tables 2.

#### 964 **scATAC-seq Data processing**

965 After sequencing, FASTQ files were processed with 10X Genomics Cellranger-atac  
966 (1.2.0) pipeline with default parameters. Briefly, the reads were aligned to hg19 using  
967 BWA to generate fragment files. Only fragments with MAPQ > 30 on both reads  
968 were retained. Each unique fragment is associated with a single cell barcode. After  
969 filtering low quality barcodes and removing PCR duplicates, a total of ~3.1 billion  
970 read pairs were retained from scATAC-seq. These reads constitute 269,920 valid cells.  
971 The output HTML files containing metrics and library information are organized into  
972 a table (Table S1). The output fragment files were loaded into ArchR to generate  
973 cell-bin matrix. Briefly, we exclude low-quality cell barcode based on loose quality  
974 control parameters: 200 unique fragments per cell and a transcription start site (TSS)  
975 enrichment score of 4. Then, we used computational framework bap (bead based  
976 ATAC processing) to combine cells which have similar fragments but with different  
977 barcodes. New fragment files generated by bap2 were loaded into ArchR again. We  
978 picked the top 12,000 cells with the highest TSS enrichment score to remove the  
979 effects of cell numbers per organ and adopted a strict quality control parameter: 1000  
980 unique fragments per cell. Finally, we filtered the doublets with addDoubletScores  
981 function in ArchR and attained final cell-bin matrix for further analysis. Finally,

982 230,732 high-quality cells with balanced sample sources are used for downstream  
983 analysis.

#### 984 **Cell type identity assignment of scATAC-seq data**

985 To annotate cell types for scATAC-seq data, we transferred cell type labels from  
986 scRNA-seq to scATAC-seq data within paired assays. First of all, we arranged 331  
987 cell types by transcriptomic similarity and pre-divided them into 65 groups by using  
988 the R package dendextend. Then, we performed two rounds of label transferring using  
989 ArchR, which utilize Seurat's canonical correlation analysis (CCA) based integration  
990 infrastructure. For the first round, we transfer 65 cell type group labels with  
991 unconstrained integration mode. For the second round, we transferred 331 cell type  
992 labels with constrained integration mode. Briefly, dimensionality reduction of whole  
993 scATAC-seq dataset was performed by using Latent Semantic Indexing (LSI). Cells  
994 were clustered by Louvain algorithm with  $r=7$  (seurat's FindClusters) and visualized  
995 by UMAP. Through first-turn label transferring, we identified which cell type group  
996 labels from the scRNA-seq data are most abundant in each of scATAC-seq clusters.  
997 We constructed a "groupList" which contains 65 pair of lists of cell IDs across  
998 scRNA-seq and scATAC-seq dataset. Then we pass this list to the 'groupList'  
999 parameter of the 'addGeneIntegrationMatrix()' function in ArchR and performed  
1000 second-turn label transferring constrained in each group and sample. We achieved a  
1001 median prediction score of 0.58-1.0 across 28 samples. 283 cell types were  
1002 successfully transferred. The cell types with cell number higher than 50 were

1003 performed peak calling by using macs2. Totally, 848,475 non-overlapped 501bp  
1004 fixed-width master peaks was generated. Any peak that directly overlaps with most  
1005 significant peak was removed. After filtering cell types with less than 50 cells or  
1006 20,000 peaks, we got 225 cell types with paired pseudo-bulk profiles of gene  
1007 expression and chromatin accessibility.

### 1008 **Genome browser visualization of two assays**

1009 Firstly, we used samtools to merge sample bam files together. Secondly, we used  
1010 filterbarcodes command in the Python package sinto (v0.1,  
1011 <https://github.com/timoast/sinto>) to get bam file for each cell type. Finally, we  
1012 generated bigWig files using bamCoverage program in Deeptools2 with parameter  
1013 “-noralizeUsingRPKM” and visualized them in IGV (version 2.8.13) (Fig 1D).

### 1014 **Generate DNA accessibility patterns using binary peak-by-cell type matrix**

1015 We constructed a binary matrix  $M_{p2ct}$  consisting of the presence or absence calls of the  
1016 master peak list ( $n = 848,475$ ) across 225 cell types.  $M_{p2ct}$  ( $225 \times 848,475$ ) was  
1017 clustered by rows and columns separately. Firstly, we selected top 200,000 most  
1018 variable peaks across cell types as features. Secondly, we calculated distance between  
1019 each cell type using (1-pearson correlation). Thirdly, we did hierarchical clustering  
1020 using calculated distance using ward.D2 algorithm (Fig 1C). For column clustering,  
1021 we unitized 2-norm of each column of  $M_{p2ct}$  to 1 and got a normalized matrix  $M_{nor}$ .  
1022 Then we took cell types as features and applied K-means to 848,475 columns of  $M_{nor}$   
1023 in Hartigan-Wong algorithm. We tested different K according to an arithmetical

1024 sequence, and selected satisfactory one ( $K = 21$ ) based on internal structure of  $M_{p2ct}$   
1025 heatmap organized in clustering results. Lastly, we manually adjusted peak group  
1026 orders to visualize the binary matrix in a fashion of neatly arranged blocks on the  
1027 diagonal. Note that the same procedures were also applied to identify sub patterns of  
1028 cell types of kidney epithelial.

### 1029 **Overlap of the ATAC peaks with consensus human DHSs**

1030 To assess the overlap between our ATAC peaks and DHSs from large-scale bulk  
1031 DNase-seq, we obtained index of consensus human DHSs from ENCODE Project and  
1032 computed intersection as well as subtraction between two datasets. The comparison  
1033 were made in two cases: whole dataset level (Fig S2A); among corresponding primary  
1034 tissues (Fig 2B). To explore differences between datasets in case two, we also  
1035 calculated tissues/cell types contributions to datasets specific peaks (Fig 2C and S2B).  
1036 Note that one peak may be calculated repeatedly, but only a limited overlap exists  
1037 between sub cell types (Fig S2C). Lastly, two-tailed Student's t test was conducted  
1038 between contributions from common cell types and contributions from organ specific  
1039 cell types (Fig 2D).

### 1040 **Enrichment analyses for enhancers from the VISTA enhancer database**

1041 VISTA validated elements were downloaded from <https://enhancer.lbl.gov> on 27  
1042 September 2020. To attain the expression pattern of each enhancer, we used advanced  
1043 search on the website and downloaded the enhancers from corresponding organs (eye,  
1044 heart, limb, liver and pancreas) in turn. Firstly, a global comparison was made

1045 regardless of organ source (Fig S2D). Secondly, we characterized accessibility pattern  
1046 of enhancers across different cell types using binary matrix (Fig S2E). Finally, we  
1047 shuffled organ peaks 3 times as background for each test, and calculated observed to  
1048 expected (median value of overlaped peaks in random situation) ratio as enrichment to  
1049 eliminate quantity effects (Fig 2E). We repeated the above operation and got  
1050 enrichment in cell type level (Fig 2F).

### 1051 **Transcription factor motif enrichment and expression analysis**

1052 The findMotifsGenome.pl in HOMER was used to calculate TF motif enrichments in  
1053 different peak groups (Fig 2A and 3A) with parameter “-size 400”. Only the top 10  
1054 motifs of each peak groups were selected to perform visualization and annotate peak  
1055 groups. Gene expression levels of TFs were normalized across cell types by Z-score  
1056 and visualized using `DotPlot()` function in Seurat. Note that a gap exists between TF  
1057 names from HOMER and official gene symbols. We filled the gap by taking two  
1058 strategies: convert lower-case characters to upper-case to see if matching any official  
1059 gene symbol; manually search the TF names on GeneCards database to see if  
1060 matching any aliases of a gene. An organized csv file was available on the website.

### 1061 **Finding Instance of Specific Motifs**

1062 To recover the locations of each motif found in the motif discovery process, we ran  
1063 the findMotifsGenome.pl again with parameter: -find SIX2.motif. The recovered  
1064 peaks were defined as TF target peaks.

### 1065 **Linking regulatory elements to cognate genes**

1066 By ArchR, we leveraged the gene expression data and created a correlation-based  
1067 map between chromatin accessibility peaks and their cognate genes directly. Briefly,  
1068 an approach introduced by Cicero is adopted to create low-overlapping aggregates of  
1069 single-cell profiles. Aggregates with greater than 80% overlap with any other  
1070 aggregate are filtered in order to reduce bias. Then we leveraged scATAC-seq data  
1071 and integrated scRNA-seq data to look for correlations between peak accessibility and  
1072 gene expression. These putative gene regulatory interactions were predicted using the  
1073 "getPeak2GeneLinks" function with default parameters in ArchR. We searched a  
1074 region of  $\pm 250$ kb for each gene and filtered peaks which were proximal to TSS ( $\pm$   
1075 1kb). Links with absolute value of correlation larger than 0.45 or less than -0.40 were  
1076 used for downstream analysis. Positive links are defined as enhancer-gene links, and  
1077 negative links are defined as silencer-gene links.

1078 We repeated these procedure in whole organism level as well as within each organ.  
1079 To retain reliable linkages against random noise, we filtered links that only shows in  
1080 one condition and merged the leftovers into 155,620 positive peak-to-gene links  
1081 (associated with 108,699 peaks and 12,783 genes) and 34,287 negative peak-to-gene  
1082 links (associated with 23,392 peaks and 7,628 genes).

### 1083 **Association with ReSE-identified silencers**

1084 To validate our data, we did overlap between correlation-based silencers and  
1085 ReSE-identified silencers by using intersectBed. Then, we applied the same  
1086 correlation-based methods linking ReSE-identified silencers to cognate genes. Only

1087 the negative correlated links were taken into consideration. We set the region as  
1088  $\pm 500\text{kb}$  for each silencer, and assigned the gene with the smallest correlation to this  
1089 silencer. 2113 of 5472 ReSE-identified silencers were assigned with a target gene.

1090

### 1091 **Classification of silencers**

1092 To determine the class of each silencer, we focus on the distribution of gene  
1093 expression with different peak accessibility and see if there is a sharp decline once the  
1094 peak accessibility reach a critical value.

1095 For each silencer, we simply take a list of value of  $1/10, 2/10, \dots, 8/10 * \text{Max}$ , where  
1096 Max denotes the max value of the peak accessibility. For each value  $i$ , we separate  
1097 cell types into two group, one with peak accessibility more than  $i$  (group  $i1$ ), and one  
1098 with peak accessibility no more than  $i$  (group  $i2$ ). If either of the groups has less than  
1099 6 cell types, we skip the value  $i$ . Then we calculate mean value and variance for each  
1100 group ( $E_{i1}, \text{Var}_{i1}$  for group  $i1$ ;  $E_{i2}, \text{Var}_{i2}$  for group  $i2$ ). A silencer is classified as  
1101 strong silencer only if  $E_{i2}/E_{i1} > 3$  and  $\text{Var}_{i2}/\text{Var}_{i1} > 3$  for any of the value  $i$ . We  
1102 tested the classifier on both correlation-based silencers and ReSE-identified silencers,  
1103 and got the same result with the independent man-made result.

### 1104 **Trajectory inference with Palantir**

1105 The Palantir workflow consists of three core steps to align cells along differentiation  
1106 trajectories. Palantir also includes visualization tools to help explore trajectories and  
1107 capture the stochasticity in cell fate determination.



1108 Dimensionality reduction with force-directed layouts (FDL). Firstly, we exported  
1109 cell-peak matrix and cell-gene matrix from ArchR and transferred it into mtx format.  
1110 Secondly, the matrices were loaded into Palantir via ``scanpy.read_10x_mtx()``  
1111 function. To settle the high sparsity of scATAC-seq data, we searched 50 nearest  
1112 neighbors for each cell via ``scanpy.pp.neighbors`` function, and aggregated single-cell  
1113 profiles using following formula:

$$1114 \quad agg(A_i) = A_i + \left[ \left( \sum_{k=1}^{49} A_{n_k} - 0.49 \right) / 5 \right]$$

1115 , where  $A_i$  denotes count number of cell  $i$  on a peak and  $n_k (k = 1 \dots 49)$  denotes  
1116 neighborhood cells. Thirdly, the aggregated ATAC profiles were used for FDL  
1117 visualization via ``harmony.plot.force_directed_layout()`` function.

1118 Integration with scRNA-seq data. To integrate transcriptome into the Palantir  
1119 framework, we took the diffusion maps of the scATAC-seq data from  
1120 ``palantir.utils.run_diffusion_maps()`` function. Using the same diffusion maps, we can  
1121 visualize gene expression levels on the same FDL plot. Then we plotted marker genes  
1122 on FDL to attain cell type locations.

1123 Grouping cells into different trajectories. We first specifying an approximate early cell  
1124 and terminal cells based on marker genes. Next we ran Palantir core function on  
1125 scATAC-seq data by ``palantir.core.run_palantir()``. Palantir generates the following  
1126 results: pseudo time ordering of each cell; terminal state probabilities of each cell; a  
1127 quantitative measure of the differentiation potential of each cell. We partitioned cells

1128 into trunk and branches according to terminal state probabilities. Cells with balanced  
1129 probabilities are defined as trunk and are used for start of lineage differentiation.  
1130 Pseudo time from 0 to 1 is used to order cells.

### 1131 **Characterization of TF related enhancer elements and genes along** 1132 **differentiation**

1133 We characterized chromatin accessibility of TF related enhancer elements and  
1134 expression level of TF related genes by using Locally Weighted Linear Regression  
1135 (Loess). Briefly, we extracted profiles of chromatin accessibility/gene expression  
1136 from cell-peak/cell-gene matrix and ordered cells according to pseudo time. We  
1137 truncated the top 5% and bottom 5% among all cells and applied Min-Max  
1138 normalization to each profile to make cross-data comparison. Finally, each profile of  
1139 chromatin accessibility/gene expression combined with pseudo time was fitted with  
1140 Loess model by `'geom_smooth()'` function.

### 1141 **Generating paired DNA accessibility patterns and gene expression patterns**

1142 To visualize DNA accessibility patterns and gene expression patterns, we firstly  
1143 calculated average gene expression levels/DNA accessibility for each cell type. For  
1144 scRNA-seq data, we used `'Seurat::AverageExpression()'` function to average gene  
1145 expression by cell types. For scATAC-seq data, the read count of each cell in the  
1146 cell-peak matrix was normalized to 10,000. All cells with the same cell type label  
1147 were pooled together to get the average DNA accessibility. Then we took  
1148 enhancers/silencers related to the same gene as an unit, and used average value to

1149 represent accessibility of enhancers/silencers to the gene. Next we drew heatmap for  
1150 gene expression patterns of 108 genes identified above, and clustered genes using R  
1151 package ComplexHeatmap with parameter: `cluster_columns = T`. Enhancers and  
1152 silencers are in the same order as their linked genes, and were visualized with  
1153 heatmaps. Gene expression levels was normalized across cell types in Z-score and  
1154 limited from -2 to 2 for the visualization. DNA accessibility was normalized across  
1155 cell types in Z-score and limited from 0 to 2 for the visualization.

### 1156 **Colocalization of scATAC-seq signal and 5' scRNA-seq signal**

1157 To distinguish transcription at CRE from mRNAs, we firstly filtered scRNA-seq reads  
1158 proximal to TSS ( $\pm 1\text{kb}$ ) or overlapped with any exon. We have 8 samples with  
1159 paired-end sequencing and 18 samples with single-end sequencing on read 2 (median  
1160 fragment size: 350bp). To uncover transcription start sites precisely, we focused on  
1161 read 1 for paired-end sequencing, and shifted upstream 200bp for single-end  
1162 sequencing. Only the very beginning 50bp of each read are used for downstream  
1163 analysis. We calculated scATAC-seq signal and 5' scRNA-seq signal per distal  
1164 ATAC peak and prepared an intermediate file via `computeMatrix` in `deeptools`  
1165 (version 3.3.0). Finally, we visualized all the results in paired heatmaps via  
1166 `plotHeatmap` (Fig 5A).

### 1167 **Identifying transcribed cis-regulatory elements**

1168 To identify transcribed cis-regulatory elements, we started from sample levels and  
1169 chose representative characteristics. For each sample, transcribed cis-regulatory

1170 elements are defined as open and significant transcribed. We used ATAC data to call  
1171 peaks via macs2. Then we calculated local RNA signal enrichment by using the ratio  
1172 between core read count and average background read counts ( (upstream 500bp +  
1173 downstream 500bp)/2 ). Only the peaks with more than 5 read count and more than  
1174 1.5 local RNA signal enrichment are considered as significant transcribed. Finally, we  
1175 merged transcribed cis elements from each sample into a master list of 190,356 peaks.  
1176 For each cell type, any of the 190,356 peaks with open state and read count larger  
1177 than 3 are considered as transcribed cis elements.

#### 1178 **Identifying Cell type Specific aCREs**

1179 In each organ, we calculate specificity score for every cell type based on the cells  
1180 versus 84K aCREs matrix by ‘Specificity scores’ preprint protocol V1.01<sup>2,3</sup> which  
1181 provided by Silvia et al. Then rank these aCREs based on the specificity score, the  
1182 top 10,000 most specific CREs per cell type is used in downstream analysis.

#### 1183 **Enrichment analysis of Heritability**

1184 Partitioned heritability was measured using LD Score Regression v1.0.0<sup>4,5</sup> to identify  
1185 enrichment of GWAS summary statistics among lineage specifier families (LSF). To  
1186 do so, first all necessary data set needed to run S-LDSC including baseline scores,  
1187 PLINK files, frequency files, weights, and SNPs, were downloaded from the Broad  
1188 Institute. All files were ‘1000G\_Phase3’ versions (See TableS6). Additionally,  
1189 Roadmap Epigenetic Project LDSC files were used as additions to the baseline model  
1190 as was done in a previous application of LDSC on ATAC seq data. We obtained

1191 GWAS summary statistics data from the UK Biobank project as processed by the  
1192 Neale lab (<http://www.nealelab.is/uk-biobank/>). Summary statistics for 52 GWAS  
1193 were obtained from have been processed into LDSC-format using the  
1194 ‘munge\_sumstats.py’ script.

1195 Firstly, annotation file was created which marked all HapMap3 SNPs that fell within  
1196 top 10K CREs for each cell type, which were ranked by cell type-specificity scores.

1197 Then LD-scores were calculated for these SNPs within 1 cM windows using the 1000

1198 Genomes data with the ‘ldsc.py’ script. These LD-scores were included

1199 simultaneously with the baseline distributed annotation file from 1000 Genome

1200 project phase 3 with population code EUR and another baseline model from Roadmap

1201 Epigenetic Project LDSC files. Subsequently, the heritability explained by these

1202 annotated regions of the genome was assessed from these genome-wide association

1203 studies: The enrichment was calculated as the heritability explained for each

1204 phenotype within a given annotation divided by the proportion of SNPs in the genome

1205 and Benjamini–Hochberg FDR correction (Benjamini and Hochberg, 1995) was used

1206 to correct for multiple comparisons. Partitioned heritability calculations for all traits

1207 were combined and analyzed in R. The creation of plots was carried out using custom

1208 R scripts. The level of significance was set for LDSC results as the Bonferroni

1209 corrected P-value when take into account all summary statistics and cell populations

1210 tested.

1211 Heritability enrichment analysis workflow in 20 LSFs were similar. Each LSF has  
1212 two types, one is classified within all accessible peaks (84K, aCRE), another input set  
1213 of peaks are derived from transcribed peaks (19K, tCRE). Firstly, we collected 80  
1214 traits to do downstream analysis, only traits with an estimated heritability were carried  
1215 forward for analysis. (q value >0.2).

1216 For some significant traits, we compared the heritability enrichment level among four  
1217 conditions (all tCRE, all aCRE, significant group's tCRE, significant group's aCRE).

1218 We calculate  $\frac{\text{Pr}(h_g^2)}{\text{Pr}(\text{SNPs})}$  to measure four LSFs' genetic associations and  
1219 heritability.

### 1220 **Jaccard Similarity Analysis**

1221 Based on CytoTRACE inference and unique gene expression, the 225 fetal cell types  
1222 in our study can be grouped into two general categories with respect to cell  
1223 proliferation. Most differentiated cells, such as cardiac muscle cells in humans, are no  
1224 longer capable of cell division. These cells are produced during embryonic  
1225 development, differentiate, and are then retained throughout the life of the organism.  
1226 In contrast, 54 cell types have been annotated as proliferative state, sustain  
1227 proliferation.

1228 We obtained each patient's raw atac counts in those cancer type specific peak sets  
1229 (<https://gdc.cancer.gov/about-data/publications/ATACseq-AWG>), filtered out those  
1230 low detected peaks (reads counts < 20) and generate each patient accessible peak set as  
1231 bed format file, then convert to hg19 chromosome using LiftOver. Meanwhile, we use

1232 bedtools to generate evenly-sized 1000-bp bins across genome, score the chromatin  
1233 accessibility similarity between patients and cell types by calculated Jaccard similarity  
1234 coefficients using peak signal overlap those windows.

1235 The process to evaluate the cancer type and cell types chromatin state similarity are  
1236 basically same. We obtain each cancer type specific strongest peak sets and produce a  
1237 binary bin matrix for cancer and cell types in correspond organ, Jaccard index was  
1238 computed, and these results were summarized using heatmap.

1239 Then rank these cell types based on the cell type's Jaccard similarity coefficient, we  
1240 can evaluate chromatin accessible similarity among 9 malignant cancer types and  
1241 proliferate cell types by calculate proliferate cell type proportion in top 10% cell types.  
1242 Specifically, we calculated the hyper-geometric p-value testing the overlap within  
1243 each cancer's top 10% similar cell types compared to the proliferate cell type set using  
1244 “phyper” in R. (See TableS7)

1245

1246

1247

## 1248 **Reference**

- 1249 1. Corces, M.R. *et al.* An improved ATAC-seq protocol reduces background  
1250 and enables interrogation of frozen tissues. *Nat Methods* **14**, 959-962  
1251 (2017).
- 1252 2. Domcke, S. *et al.* A human cell atlas of fetal chromatin accessibility.  
1253 *Science* **370**(2020).
- 1254 3. Domcke, S. & Hill, A. Specificity scores. *Bio-protocol*  
1255 [bio-protocol.org/prep644](https://www.bio-protocol.org/prep644).(2020).

- 1256 4. Bulik-Sullivan, B.K. *et al.* LD Score regression distinguishes confounding  
1257 from polygenicity in genome-wide association studies. *Nat Genet* **47**,  
1258 291-5 (2015).
- 1259 5. Finucane, H.K. *et al.* Heritability enrichment of specifically expressed genes  
1260 identifies disease-relevant tissues and cell types. *Nat Genet* **50**, 621-629  
1261 (2018).
- 1262 6. Gulati, G.S. *et al.* Single-cell transcriptional diversity is a hallmark of  
1263 developmental potential. *Science* **367**, 405-411 (2020).
- 1264