

1 Gene expression patterns of the developing human face at single cell resolution reveal cell type contributions to  
2 normal facial variation and disease risk

3 Nagham Khouri-Farah<sup>1</sup>, Emma Wentworth Winchester<sup>1</sup>, Brian M. Schilder<sup>2,3</sup>, Kelsey Robinson<sup>4</sup>, Sarah W.  
4 Curtis<sup>4</sup>, Nathan G. Skene<sup>2,3</sup>, Elizabeth J. Leslie-Clarkson<sup>4</sup>, Justin Cotney<sup>5,6</sup>

## 5 Affiliations

6 1 Graduate Program in Genetics and Developmental Biology, UConn Health

7 2 Department of Brain Sciences, Faculty of Medicine, Imperial College London, London, W12 0BZ, UK

8 3 UK Dementia Research Institute at Imperial College London, London, W12 0BZ, UK

9 4 Department of Human Genetics, Emory University School of Medicine, Atlanta, GA 30322, USA

10 5 Department of Surgery, Children's Hospital of Philadelphia, Philadelphia, PA

11 6 Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA

## 12 Abstract

13 Craniofacial development gives rise to the complex structures of the face and involves the interplay of  
14 diverse cell types. Despite its importance, our understanding of human-specific craniofacial  
15 developmental mechanisms and their genetic underpinnings remains limited. Here, we present a  
16 comprehensive single-nucleus RNA sequencing (snRNA-seq) atlas of human craniofacial development  
17 from craniofacial tissues of 24 embryos that span six key time points during the embryonic period (4–8  
18 post-conception weeks). This resource resolves the transcriptional dynamics of seven major cell types  
19 and uncovers distinct major cell types, including muscle progenitors and cranial neural crest cells  
20 (CNCCs), as well as dozens of subtypes of ectoderm and mesenchyme. Comparative analyses reveal  
21 substantial conservation of major cell types, alongside human biased differences in gene expression  
22 programs. CNCCs, which play a crucial role in craniofacial morphogenesis, exhibit the lowest marker  
23 gene conservation, underscoring their evolutionary plasticity. Spatial transcriptomics further localizes  
24 cell populations, providing a detailed view of their developmental roles and anatomical context. We also  
25 link these developmental processes to genetic variation, identifying cell type-specific enrichments for  
26 common variants associated with facial morphology and rare variants linked to orofacial clefts.  
27 Intriguingly, Neanderthal-introgressed sequences are enriched near genes with biased expression in  
28 cartilage and specialized ectodermal subtypes, suggesting their contribution to modern human  
29 craniofacial features. This atlas offers unprecedented insights into the cellular and genetic mechanisms  
30 shaping the human face, highlighting conserved and distinctly human aspects of craniofacial biology.  
31 Our findings illuminate the developmental origins of craniofacial disorders, the genetic basis of facial  
32 variation, and the evolutionary legacy of ancient hominins. This work provides a foundational resource  
33 for exploring craniofacial biology, with implications for developmental genetics, evolutionary biology,  
34 and clinical research into congenital anomalies.

## 36 Main

37 Craniofacial development orchestrates the formation of the human face through the interplay of multiple  
38 cell lineages. These cell types, including mesenchyme, ectoderm, endothelium, and cranial neural crest  
39 cells (CNCCs), differentiate into a diverse array of tissues such as bone, cartilage, muscle, skin, and  
40 vasculature<sup>1-3</sup>. Together, these cells and tissues give rise to the face's essential functions like  
41 respiration, mastication, communication, and sensory perception<sup>4-7</sup>. Disruptions to craniofacial  
42 developmental processes rank amongst the most common causes of human congenital anomalies, with  
43 orofacial clefts representing a significant portion of global birth defects<sup>8-11</sup>. Thus, there exists significant  
44 need to understand the molecular, genetic, and cellular mechanisms underlying craniofacial  
45 development in humans.

46 Studies utilizing model organisms, particularly mice, have offered key insights into craniofacial  
47 development and abnormalities<sup>12-16</sup>. However, significant differences exist between mouse and human  
48 craniofacial development, including variations in timing, cellular contributions, and gene regulatory  
49 networks<sup>1,13,17-22</sup>. Furthermore, human craniofacial features exhibit evolutionary adaptations that  
50 distinguish them from other mammals and primates, underscoring the necessity for human-specific  
51 studies<sup>5,13,21,22</sup>.

52 Advances in single-cell and single-nucleus RNA sequencing (scRNA-seq and snRNA-seq respectively)  
53 technologies have enabled detailed characterization of cellular diversity and gene expression during  
54 development<sup>12,15,23-27</sup>. These tools are particularly valuable for resolving the dynamics of rare or transient  
55 cell populations, such as CNCCs, that play critical roles in craniofacial formation<sup>13</sup>. While previous  
56 efforts have developed single-cell atlases for murine craniofacial tissues, corresponding human  
57 datasets have been limited by sample availability, insufficient temporal resolution, and challenges in  
58 profiling craniofacial-specific populations<sup>12,15,23</sup>. Only a few studies have examined bulk gene expression  
59 patterns and regulation specifically during human craniofacial development<sup>13,14,28-31</sup>, and only two  
60 datasets are currently available during the embryonic period of human development<sup>13</sup>. While mapping of  
61 human genetics findings to mouse craniofacial cell types has indicated potential disease-causing  
62 subtypes<sup>32</sup>, the limited number of replicates underlying the mouse data and differences between human  
63 and mouse craniofacial development preclude confident interpretation.

64 To address these gaps in knowledge, we constructed a time-resolved gene expression atlas of human  
65 craniofacial development when the bulk of human craniofacial development occurs<sup>33,34</sup>. Using snRNA-  
66 seq on craniofacial tissues from 24 individual human embryos encompassing six key time points from 4  
67 to 8 post-conception weeks, we profiled over 42,000 nuclei and identified seven major cell types,  
68 including mesenchyme, ectoderm, endothelium, muscle progenitors, and CNCCs. Integration with  
69 human spatial transcriptomics further validated the localization of these subtypes within the developing  
70 human face. Comparative analysis with murine craniofacial datasets generated here and previously  
71 published<sup>23</sup> highlighted significant conservation of major cell types and their gene expression programs,  
72 alongside species biased markers that reflect differences in mouse and human biology.

73 Beyond the developmental biology of craniofacial formation, this study explores the genetic and  
74 evolutionary factors shaping human craniofacial features. By integrating genome-wide association  
75 studies (GWAS) with our atlas, we identified cell type-specific enrichments for genetic variants

76 associated with normal facial variation. We found that specific subtypes of ectoderm and mesenchyme,  
77 likely spatially restricted, contribute to different aspects of facial appearance and shape. We also  
78 examined rare variants associated with congenital craniofacial disorders such as orofacial clefts. We  
79 find that *de novo* protein damaging variants identified in orofacial clefting trios are enriched in genes that  
80 specify distinct cell subtypes in the face. This enrichment was heavily biased toward ectodermal  
81 subtypes that is largely obscured in previous analyses based on bulk chromatin and gene expression<sup>13,35-</sup>  
82<sup>38</sup>. We find the damaging variants coalesce in the ectodermal derived nasal placode implicating this early  
83 structure in orofacial clefts. Our analysis also uncovered evidence linking Neanderthal-introgressed  
84 sequences to genes with biased expression in specific craniofacial cell types.

85 This comprehensive atlas provides a high-resolution view of the cellular and molecular landscape of  
86 human craniofacial development, integrating gene expression, spatial mapping, and evolutionary  
87 genomics. Our work not only enhances our understanding of human craniofacial biology but also  
88 establishes a framework for future studies aimed at uncovering therapeutic targets and evolutionary  
89 insights into one of the most defining features of human anatomy. This data can be explored through an  
90 interactive web application that is accessible to most researchers:

91 [https://cotneyshiny.research.chop.edu/shiny-apps/craniofacial\\_all\\_snRNA/](https://cotneyshiny.research.chop.edu/shiny-apps/craniofacial_all_snRNA/) as well as alongside the  
92 growing number of single cell datasets hosted at the Chan-Zuckerberg CellXGene Discover resource<sup>39</sup>.

## 93 Results

### 94 Time-resolved atlas of gene-expression in the developing human face

95 To characterize the cellular landscape of human craniofacial development we performed snRNA-seq  
96 analysis of 24 individual human embryos across 6 distinct time points, encompassing major milestones  
97 of human craniofacial development from 4 to 8 post conception weeks (Fig. 1A). We profiled the entire  
98 craniofacial prominence from multiple biological replicates at each time point resulting in 86,359  
99 individual nuclei after filtering for doubles and quality of per nucleus data. While experiments in mouse  
100 offer precise control of tissue sampling for downstream processing, samples obtained from human  
101 embryos are more difficult to control what is collected. To identify potential biases or nuclei obtained  
102 from extraneous tissues we performed initial clustering of all samples to identify potential extraneous  
103 cell types. This analysis revealed a total of 13 distinct clusters (Fig. S1a). This number of clusters was  
104 substantially higher than the main cell types identified in mouse craniofacial developmental studies,  
105 suggesting that the human samples potentially contained extraneous tissues that are not part of the  
106 craniofacial complex<sup>12-14</sup>. When we examined the contributions of individual samples to these clusters,  
107 we found several clusters that were made up of nuclei derived from a small number of samples (Fig.  
108 S1b). Closer inspection of the genes strongly expressed in these clusters revealed many canonical  
109 neuronal genes, such as *TUBB3* and *MAP2* (Fig. S1c-d). We reasoned these clusters were derived from  
110 developing brain tissue not directly part of the craniofacial structures. We therefore excluded these  
111 nuclei from downstream analyses resulting in a total of 42,131 remaining nuclei with an average of 4095  
112 nuclei from each sample (Fig. S2a) and a median of 7500 counts from 2250 genes per nucleus (Fig. S2b  
113 and c). We observed that early samples had consistently higher mitochondrial reads (Fig. S2d),  
114 potentially reflecting their higher dependence on mitochondrial output or an artifact related to lower cell  
115 numbers in the processing of each sample.

116 To determine the quality of these filtered samples we sought to compare to other well  
117 characterized gene-expression profiles of craniofacial development. Our previous studies of bulk gene  
118 expression during human craniofacial development revealed a strong time related component across  
119 the samples<sup>13</sup>. When we combined gene expression profiles from all nuclei of a specific stage into  
120 pseudo-bulk gene expression profiles individual replicates were well correlated with others at the same  
121 time-point and less so with samples with greatest differences developmental time across this time  
122 course (Fig. S3). Principal component analysis of these pseudo-bulk profiles largely recapitulated our  
123 previous results with the first principal component ordering samples readily by known stage of  
124 development (Fig. 1B, Fig. S4a). Furthermore, when we performed differential expression between the  
125 pseudobulk samples we found very similar results to those obtained by bulk gene expression between  
126 the same timepoint comparisons (Fig. S4b). Specifically, the greatest number of differentially expressed  
127 genes were observed between the earliest and latest timepoints that could be compared across the two  
128 data sets (CS13 vs CS17) (Fig. S4c). Overall, these results suggest that our single-nucleus expression  
129 data closely resembles the bulk gene expression data that we have previously shown is enriched for  
130 many aspects of craniofacial biology and developmental abnormalities relative to many other tissues  
131 and cell types<sup>13</sup>.

### 132 **Identification of major cell types present in craniofacial development**

133 Having established that the single nuclei profiles at the pseudobulk level captured many of the expected  
134 aspects of craniofacial biology, we proceeded to re-cluster the filtered nuclei to first identify the major  
135 cell types present in the developing face. We identified seven major clusters and projected these high-  
136 dimensional data into two dimensions using Uniform Manifold Approximation and Projection (UMAP)<sup>40</sup>  
137 (Fig. 1C). The clusters were contributed to by samples from each of the replicates and stages in very  
138 similar proportions (Fig. 1D and E). Interestingly, this was two more distinct clusters than previously  
139 characterized in the E11.5 mouse craniofacial structures<sup>12</sup>. We reasoned that this could be due to  
140 differences in human and mouse development, but most likely related to the additional replicates and  
141 timepoints and how tissues were collected and processed. To address this, we first examined expression  
142 of the five genes examined by Li et al, *ALX1*, *EPCAM*, *HEMGN*, *CDH5*, and *FCERG1* as markers of  
143 mesenchyme, ectoderm, blood, endothelium, and immune cells respectively (Fig. 2A). *ALX1* was most  
144 strongly expressed in cluster 1, *EPCAM* in cluster 2, *HEGMN* in cluster 4, *CDH5* in cluster 5, and *FCERG1*  
145 in cluster 7, while clusters 3 and 6 did not show signal for any of these genes. While some of the  
146 timepoints were derived exclusively from female and male samples, CS12 and CS13 respectively, we did  
147 not observe any significant bias in these main cluster (Figure S2E).

148 In an attempt to characterize these unknown clusters we first identified the top 10 genes that  
149 were most strongly differentially expressed between the individual clusters (Fig. 2B, Supplemental Table  
150 1). Cluster 1 was marked by *PDGFRA*, *TWIST1*, and *PRRX1*, consistent with identifying this cluster as  
151 mesenchyme<sup>41,42</sup>. Cluster 2 was identified by *GRHL2* and *ESRP1*, genes that have been reported to be  
152 specifically active in surface ectoderm and epithelial cells<sup>43-46</sup>. Cluster 4 was marked by *SPTA1*, *ALAS2*,  
153 *RHAG*, genes involved in erythrocyte function<sup>47-51</sup>. Cluster 5 showed highly biased expression for *KDR* and  
154 *FLT1*, genes associated with the vascular system and endothelium function<sup>52,53</sup>. Cluster 7 was marked by  
155 *PTPRC*, *CD86*, and *CD136*, consistent with immune cell function<sup>54-56</sup>. These all confirmed the initial  
156 identities suggested by the markers described by Li et al in E11.5 mouse craniofacial tissue. The

157 unknown cluster 3 showed highly biased expression of *MYOG*, *MYL1*, and *MYH3*, all genes related to  
158 muscle specification and function<sup>57-59</sup>. The unknown cluster 6 showed strongly biased expression for  
159 *CDH19*, *INSC*, and *MMP17*. These genes are involved in a variety of biological processes including cell  
160 adhesion, spindle orientation during mitosis, and degradation of extracellular matrix<sup>60-62</sup>. We also noted  
161 specific expression of *FOXD3*, a developmental transcription factor which has been linked to  
162 pluripotency maintenance in stem cells and specification of neural crest in multiple species<sup>63-66</sup>.

163 We then analyzed the top 100 marker genes from each cluster for gene ontology, pathway, and  
164 disease enrichments. The genes that identified putative mesenchyme cluster 1 were enriched for a  
165 number of biological process categories related to skeletal, cartilage, and roof of mouth development  
166 (Fig. S5A, Supplemental Table 2); cellular components related to collagen processing (Fig. S5B,  
167 Supplemental Table 3); molecular functions related to gene expression, extracellular matrix, and  
168 collagen binding (Fig. S5C, Supplemental Table 4); pathways related to production of extracellular matrix  
169 (Fig. S5D, Supplemental Table 5); and diseases including cleft palate and frontonasal dysplasia (Fig. 2C,  
170 Supplemental Table 6).

171 Genes most strongly expressed in cluster 2, likely ectoderm, were enriched for biological process  
172 categories related to tight junction assembly and cell adhesion (Fig. S5A, Supplemental Table 2); cellular  
173 components related to the plasma membrane (Fig. S5B, Supplemental Table 3); molecular functions  
174 related to cadherin and laminin binding (Fig. S5C, Supplemental Table 4); pathways related to tight  
175 junction and Hippo signaling (Fig. S5D, Supplemental Table 5); and diseases including epithelioma and  
176 keratoderma (Fig. 2C, Supplemental Table 6).

177 Putative muscle progenitor cluster 3 marker genes were enriched for biological processes related  
178 to muscle cell differentiation (Fig. S5A, Supplemental Table 2); cellular components of the sarcomere  
179 (Fig. S5B, Supplemental Table 3); molecular functions related to actin filament binding (Fig. S5C,  
180 Supplemental Table 4); calcium signaling pathways (Fig. S5D, Supplemental Table 5); and diseases  
181 including myopathic abnormalities and muscular dystrophy (Fig. 2C, Supplemental Table 6).

182 The markers of red blood cell cluster 4 were enriched for biological processes related to  
183 erythrocyte homeostasis and oxygen transport (Fig. S5A, Supplemental Table 2); cellular components of  
184 the hemoglobin complex (Fig. S5B, Supplemental Table 3); molecular functions related to heme and  
185 oxygen binding (Fig. S5C, Supplemental Table 4); pathways involved in Malaria response and mineral  
186 absorption (Fig. S5D, Supplemental Table 5); and diseases including hemolytic anemia and beta  
187 thalassemia (Fig. 2C, Supplemental Table 6).

188 Genes with highest expression in putative endothelium cluster 5 were enriched for biological  
189 processes related to endothelial cell differentiation and proliferation (Fig. S5A, Supplemental Table 2);  
190 cellular components including plasma membrane rafts and caveola (Fig. S5B, Supplemental Table 3);  
191 molecular functions related to Notch and guanyl nucleotide binding (Fig. S5C; Supplemental Table 4);  
192 pathways involved in fluid shear stress and atherosclerosis (Fig. S5D; Supplemental Table 5); and  
193 diseases of the capillaries and hemangiomas (Fig. 2C, Supplemental Table 6).

194 Immune related cluster 7 marker genes were enriched for biological processes related to cytokine  
195 production and immune response (Fig. S5A, Supplemental Table 2); cellular components including



196 specific and tertiary granule membranes (Fig. S5B, Supplemental Table 3); molecular functions related  
197 to Toll-like receptor binding and immunoglobulin receptor activity (Fig. S5C, Supplemental Table 4);  
198 pathways related to the phagosome and complement and coagulation cascades (Fig. S5D,  
199 Supplemental Table 5); and diseases including many types of infections and immunodeficiencies (Fig.  
200 2C, Supplemental Table 6).

201 We then turned to the not yet concretely identified cluster 6. The marker genes we identified for  
202 this cluster were enriched for biological processes related to glial cell differentiation and myelination  
203 (Fig. S5A, Supplemental Table 2); cellular components including plasma membrane signaling receptor  
204 complexes and exocytic vesicles (Fig. S5B, Supplemental Table 3); molecular functions related to protein  
205 tyrosine kinase activator activity (Fig. S5C, Supplemental Table 4); and diseases related to central  
206 nervous system disorders and neuropathies (Fig. 2C, Supplemental Table 6). We did not detect any  
207 significant pathway enrichments for this particular cluster.

### 208 **Identification of presumptive human cranial neural crest**

209 The marker gene ontology analysis successfully confirmed the identity of six of the seven major  
210 clusters. However, cluster 6 remained difficult to identify due to the variety of enrichments identified  
211 amongst marker genes. Beyond the more nervous system-oriented enrichments listed above we also  
212 found significant biological and disease enrichments that were shared with the mesenchyme and  
213 muscle clusters. This included enrichments for extracellular matrix organization and binding, cell  
214 adhesion via plasma-membrane, skeletal muscle system development, and several types of tumors (Fig.  
215 2C, Fig. S5A-D, Supplemental Tables 2-6). When we more closely inspected biological process  
216 categories identified for cluster 6, we observed additional enrichments related to Schwann cell  
217 development and melanocyte differentiation (Supplemental Table 2). Closer inspection of full disease  
218 enrichments for this cluster revealed several types of Waardenburg Syndrome, Hirschsprung Disease,  
219 and demyelination disorders (Supplemental Table 6). The wide variety of biological functions and  
220 specific disease enrichments all suggested that this cluster might be enriched for neural crest cells.  
221 Marker genes driving these enrichments included *EDNRB*, *ERBB3*, *PAX3*, *SOX10*, *SPP1*, *TFAP2B*, and  
222 *ZEB2*, genes well known to be involved in various aspects of neural crest specification and migration<sup>67</sup>.  
223 However, while these genes are biased toward cluster 6 relative to other clusters, they are not  
224 exclusively expressed in cells found in this cluster (Fig. S6A). Amongst these *ZEB2* is more broadly  
225 expressed across all clusters except for ectoderm. Further inspection of marker genes revealed that  
226 while some of these genes are indeed strongly biased toward cluster 6, only a small percentage of cells  
227 from this cluster express each gene (Fig. S6B). Qualitatively, expression of each of these genes could be  
228 observed outside of cluster 6 and potentially in subclusters of the main clusters we have defined thus far  
229 (Fig. S6C). Given the heterogeneity of expression of each of these marker genes we reasoned that jointly  
230 analyzing expression of a module of genes might be a better indicator of cell type identity as has been  
231 shown in other single cell-based studies<sup>68</sup>. When we examined a module of genes from regulatory  
232 networks recently identified in cultured human and chimpanzee cranial neural crest cells (CNCCs)<sup>69</sup>, we  
233 found significantly higher expression in cluster 6 (Fig. 2D). Together these results strongly point to this  
234 cluster being enriched for putative CNCCs.

235 Thus far few studies have been able to identify significant populations of CNCCs from primary  
236 human tissue<sup>70-72</sup>. To better understand the gene expression programs that are active in these cells we  
237 first performed subclustering on these cells (n = 1821). We identified 11 distinct subclusters from this  
238 original population (Fig. 3A). The seven main clusters derived from the largest numbers of cells were  
239 annotated as CNCCs, while the four more punctate clusters were initially annotated as CNCC like (cnl).  
240 These clusters consisted differentially of cells derived from each of the stages profiled. Those clusters  
241 that were heavily biased toward CS12 were labeled as early (eCNCC), those that were biased toward  
242 CS13-C16 as intermediate (iCNCC), and those that were derived primarily from the CS17 and CS20  
243 timepoints as late (lCNCC) (Fig. 3B). When we examined gene expression of many CNCC markers from  
244 the literature we found variable patterns of expression. *SOX10* expression was observed in all of the  
245 CNCC and cnl clusters along with *NR2F1* and *NR2F2*, two genes identified as master regulators of CNCC  
246 fate<sup>6,73</sup> (Fig. 3C). *TFAP2A* expression was observed in all clusters but was considerably lower in the  
247 late CNCCs. Its ortholog *TFAP2B* was conspicuously absent from one late CNCC cluster (lCNCC2) and  
248 from cluster cnl2 (Fig. 3C). *ETS1* and *FOXD3* were generally expressed in most subtypes, but both  
249 expressed at very low levels in lCNCC1 (Fig. 3C). *PAX3* expression was more variable but in populations  
250 distinct from the two previously mentioned transcription factors. *SOX9* and *COL20A1* were more  
251 specifically expressed across the clusters, but again in non-overlapping patterns.

252 Overall, these genes largely confirmed that the cells we identified have neural crest character,  
253 however they did not display distinct patterns across the clusters precluding easy identification of these  
254 putative subtypes. To identify genes that readily identified each of these subtypes, we repeated the  
255 marker gene identification performed on the main types above. We identified approximately 2000 genes  
256 that were differentially expressed across these subclusters with an adjusted p-value cutoff less than  
257 0.05 and a log<sub>2</sub> fold change greater than one (Supplemental Table 7). The top five marker genes in each  
258 subtype revealed multiple transcription factors that distinguish each cluster. These included *SOX21* in  
259 early CNCCs, *MKX* in intermediate CNCCs, *HAND2* in late CNCCs, and *NKX2-1*, *ALX4*, and *FOXL1*, in  
260 clusters cnl1, cnl2, and cnl3, respectively (Fig. 3D). Identification of enriched gene ontology categories  
261 for each subtype revealed distinct functions for each. Marker genes of early CNCCs were enriched for  
262 process related to early pattern specification and axon guidance. Intermediate CNCC clusters were  
263 enriched for functions related to extracellular matrix organization and skeletal system development. Late  
264 CNCC clusters were enriched for various channel activity and sympathetic nervous system  
265 development. The cnl1 cluster was enriched for several categories shared with eCNCC1 suggesting this  
266 was an early multipotent neural crest type. The cnl4 cluster was very specifically enriched for functions  
267 related to pigment granules and melanin biosynthesis indicating these were melanocytes, a cell-type  
268 derived from neural crest (Figure S7, Supplemental Tables 8-11).

269 Thus far our analysis has lacked spatial localization, making it unclear where these cell types are  
270 derived or reside in the intact human embryo. Recently published spatial transcriptomics on two  
271 sections of a human CS13 provided an opportunity to identify such patterns of expression<sup>25</sup>. We  
272 reprocessed this data, merged all the cells from both sections, identified cell types, and confirmed their  
273 spatial locations (Figure S8). We then examined expression of marker genes that identified CNCCs  
274 versus the other major craniofacial cell types. *NR2F1* was broadly expressed across the embryo whereas  
275 *PAX3*, *TFAP2A*, and *TFAP2B* were more regionally restricted to the head and neural tube regions (Fig. 3E).

276 Genes identified as markers of CNCC subtypes showed a variety of patterns of expression. *ALX4* was  
277 generally restricted to the head region and putative frontal nasal process region. *CRABP1* was found in  
278 the anterior neural tube, eye region, and the limb. *HAND2* expression was observed in putative  
279 pharyngeal arch regions, heart, and limb. *NKX2-1* had highly restricted expression in a location that could  
280 represent a fusion zone between the lateral nasal prominence and the maxillary prominence (Fig. 3E).  
281 We then calculated module scores on these spatial data using the marker genes from each of the CNCC  
282 subtypes. We found that at this stage of development, each of these sets of marker genes were generally  
283 biased toward the neural tube region of the embryo with *cnl1* marker genes showing the most restricted  
284 pattern of expression.

## 285 **Conservation of cell types and gene expression programs in human and mouse craniofacial** 286 **development**

287 The analysis above showed compelling evidence of the identities of the major cell types in the  
288 developing human face. This included two cell types, muscle and CNCCs, not previously observed in  
289 single cell atlases of mouse craniofacial development<sup>12,15,23</sup>. We wondered whether these cell types were  
290 not present in these mouse datasets due to sampling differences in tissues and broader timepoints. To  
291 address this, we generated single-nucleus gene expression data from mouse craniofacial tissues  
292 harvested from multiple biological replicates of E10.5 to E12.5. These samples reflected the major  
293 morphological landmarks of the human tissue profiled allowing a more direct comparison of cell types.  
294 We then further combined this data with recently published single cell gene expression data from E13.5  
295 and E15.5 resulting in 79402 expression profiles after similar quality control filters applied to human data  
296 (Methods). When we clustered these data using approaches identical to the human data, we obtained  
297 the same number of main clusters with remarkably similar cluster ratios and organization in the UMAP  
298 projection space (Fig. 4A, Supplemental Table 12). When we examined gene expression of the same  
299 major markers profiled in human (Fig. 2A) we readily identified the same major mouse cell types  
300 including muscle and putative CNCCs (Fig. 4B). Roughly 70% of the tissue was of mesenchymal origin,  
301 15% was ectodermal, and the remaining 15% was distributed similarly across the remaining 5 cell types.  
302 When we projected these cell types on our recent analysis of spatial gene expression in mouse E15.5  
303 craniofacial sections, we found expected patterns of cell type localization (Fig. 4C). To determine if these  
304 cell types were specified by the same sets of genes, we compared marker gene identities obtained in the  
305 same fashion in both species. We found significant sharing of marker genes between the orthologous  
306 major cell types (Fig. 4D). The highest degree of sharing was found between mesenchyme, followed by  
307 endothelium and ectoderm. While still significant, the lowest degree of marker gene conservation was  
308 observed between CNCCs of each species (Fig. 4D). When we examined the functionally conserved  
309 marker genes for ontology enrichments, we observed distinct patterns of enrichment that confirmed the  
310 cell type assignments in each species (Fig. S9A-E, Supplemental tables 13-17). Disease enrichments  
311 related to general craniofacial abnormalities were observed in mesenchyme, while enrichment of cleft  
312 upper lip was observed in conserved markers of mesenchyme and ectoderm (Fig. S9F, Supplemental  
313 Table 18). These enrichments were driven by well-known craniofacial genes including *ALX1*, *ALX4*, *MSX1*,  
314 *RUNX2*, and *TWIST1* reinforcing their conserved role in mammalian craniofacial development  
315 (Supplemental Table 18).



## Species-specific differences in marker gene expression during human and mouse craniofacial development

While the overall craniofacial cell types and major gene expression patterns were shared between species, our analysis revealed hundreds of marker genes that were only called in a single species (Supplemental Tables 19-21). The largest fraction of species biased calls was observed for CNCCs. As expected, the shared CNCC markers were enriched for functions related to gliogenesis and nervous system development. However, the human-biased markers were biased toward genes related to ribosome biogenesis and cytoplasmic translation while mouse-biased markers were enriched for genes with functions related to oxidative phosphorylation and the electron transport chain (Figure S10A-D). When we examined the mesenchyme cluster, we found the shared markers were enriched for morphogenesis and differentiation programs for mesenchymal derived cell tissues as expected. However, human-biased markers were enriched generally for functions related to DNA replication and cell cycle while mouse-biased markers were enriched for only a few categories primarily related to MAPK signaling pathways (Figure S11 A-E). When we examined the human disease phenotypes enriched for each of these gene sets, we found general craniofacial abnormalities and isolated cleft palate among conserved genes. Mouse-biased mesenchymal markers were enriched exclusively for multiple seizure disorders. Human-biased mesenchymal markers were enriched for a number of craniofacial related phenotypes including microphthalmos and low set ears and exclusively for sloping forehead, large nose, and biparietal narrowing (Figure S11D). Given these phenotypes, genes driving these enrichments could be significant contributors to differences in skull shape, size, and function between human and mice. When we inspected these categories, we found genes with the highest levels of specificity for human mesenchyme included *ALX3*, *CTSK*, *CYP1B1*, *FOXC1*, *MAB21L1*, *MSX2*, and *TENM1* (Fig. 4E).

### Leveraging mouse craniofacial cell-type annotations to identify human craniofacial subtypes.

Having demonstrated that major cell types, including the CNCCs, could be readily identified in both species and showed significant conservation of gene expression, we reasoned we could leverage the substantial annotation resources that have been generated for mouse to identify human cell subtypes. To achieve this, we focused on the major cell types that have been extensively subclustered and characterized in previous publications<sup>12,23,74</sup>, mesenchyme and ectoderm, as well as the novel populations of CNCCs we have identified here. When we performed subclustering of mouse CNCCs, we identified 8 distinct subtypes (Fig. 4F). These had a very similar arrangement in UMAP space compared to the subclusters we identified in the human CNCCs (Fig. 3A). When we examined functional enrichment of marker genes of each of these subclusters we found similar results as in human, including a clear population of melanocytes (Fig. 4G, Supplemental Table 22). Examination of the same neural crest markers as in human CNCC subtypes revealed very similar patterns of expression (Fig. 4H). We observed that *Sox10* and *Nr2f2* were expressed across all the subtypes as well as a similar trend in variable expression of *Foxd3*, *Pax3*, *Tfap2a*, and *Tfap2b* and across subtypes. When we inspected the markers for each of these subtypes, we found many of the same genes as in human subtypes including *Alk*, *Alx4*, *Crabp1*, and *Hand2* (Fig. 4I, Supplemental Table 23). When we reprocessed mouse E11.5 spatial transcriptomic data<sup>43</sup> in a similar fashion to the human CS13 data, we found very similar patterns of expression for many of the human CNCC markers in mouse tissues (Figure S12A). Examination of module scores calculated from the mouse CNCC subtypes also revealed similar patterns across the

357 mouse embryo as observed for human (Figure S12B). To attempt to identify the orthologous CNCC  
358 subtypes across species we compared sharing of orthologous marker genes much as we did with the  
359 main cell types. When we examined a confusion matrix of comparisons of cell types we found the  
360 highest similarity amongst CNCC subclusters human iCNCC4 and mouse iCNCC2 as well as human  
361 *cnl4* and mouse *cnl*, the putative melanocyte clusters (Fig. S13A). The additional *cnl* clusters in human  
362 showed variable similarity to mouse and could reflect heterochrony, primate cell states not present in  
363 rodents, or the more genetically diverse human samples profiled.

364 Having demonstrated that even in the potentially least conserved cell type that we could readily  
365 identify shared subtypes across species we then turned to the other major cell types, mesenchyme and  
366 ectoderm. We subclustered the large mouse mesenchyme cluster and identified 19 subclusters across  
367 the mouse timeseries. Using a combination of gene ontology enrichments of marker genes and previous  
368 annotations of mouse craniofacial single cell and spatial transcriptomics we assigned functional and/or  
369 positional labels to each cluster (Fig. 4J, Supplemental Tables 22 and 24-27). For example, the well-  
370 established lateral nasal process (LNP) marker *Pax7*<sup>12,75-77</sup> and the osteoblast marker *Sp7*<sup>78,79</sup> were used  
371 to define respective clusters. In the case of osteoblasts we observed two clusters expressing similar  
372 markers but were biased in cells from different stages of development, thus we further refined these as  
373 early and late osteoblasts (Figure 4K). Two small clusters clearly represented contaminating blood  
374 derived cells or neuronal-like populations while one additional cluster could not be readily identified but  
375 had many markers associated with rapidly cycling cells (Fig. 4K). When we examined the mouse E11.5  
376 spatial transcriptomics data we had reprocessed above, we found good concordance between marker  
377 gene expression and generalized localization in the embryo (Figure S14A). In contrast to both the human  
378 and mouse CNCC analysis, projection of module scores for mouse mesenchymal clusters readily  
379 identified specific regions of the developing craniofacial structures that corresponded well to labels we  
380 had assigned them (Fig. 4L and S14B).

381 We performed identical analyses for the ectodermal cluster revealing an additional 19  
382 subclusters (Fig. 4M). Applying the same analysis of marker genes from the literature, gene ontology  
383 enrichments, and expression in mouse single cell transcriptomics data we annotated each of these  
384 clusters with functional and spatial labels (Supplemental Tables 22 and 24-27). We identified highly  
385 specific ectodermal populations like periderm marked by *Gabrp*<sup>12,80,81</sup>, cells that will form structures of  
386 the inner ear marked by *Oc90* (Zhao et al 2007, Wang PNAS 1998), palate ectoderm identified by *Foxe1*<sup>82</sup>,  
387 and the putative pituitary marked with *Lhx3* and *Lhx4*<sup>83,84</sup> among others (Fig 4N). As with the  
388 mesenchyme, we found the spatially resolved expression of marker genes corresponded well to  
389 expected positions of the mouse embryo (Figure S15A). Module score calculations for each subtype  
390 resulted in refined spatial identification of subtypes that matched the labels and expected positions well  
391 (Fig. 4O and S15B).

392 While the module score analysis is indicative of the cell types and spatial locations of the labels we  
393 applied, they are calculated independently of any other cell types. We therefore sought to predict what  
394 are the dominant cell types in specific locations based on spatial transcriptomics data we had not used  
395 for any of the previous analysis. When we projected top spatial predictions for 20 of the subtypes  
396 identified across previously published E15.5 mouse head data<sup>23,85</sup> we found very good concordance for  
397 cell type labels and known anatomical features (Figure 4P). Overall, the analyses performed here

confirmed the identities of multiple cell types across the development of mouse craniofacial tissues. Moreover, the demonstration of conserved marker genes provides a framework for transferring cell type labels to subclusters identified in human data as well as putative spatial inferences from data that originally lacked that information. We explore the subtype identifications in human data below.

### Characterization of mesenchymal cell subtypes

When we subclustered the large number of mesenchymal cells, we identified 22 subtype clusters. Using the same confusion matrix-based approach from above based on orthologous gene expression in mesenchymal subtypes, we assigned cell type and/or functional labels to each of the human clusters (Fig 5A and S13B). In some cases, multiple human clusters correlated well with a single mouse cluster and were labelled as separate populations (e.g., mouse mandibular arch and human mandibular arch 1-3). The most abundant cell types were obtained from the mandibular arch and the maxillary process and were well represented from each of the timepoints. Some of the transient structures like the lateral nasal process and cells labeled as early osteoblasts were biased towards early timepoints, while later forming cell types and structures such as cartilage and palatal shelves were dominated by cells derived from CS20 samples (Fig. 5B). When we examined marker genes identifying each of these clusters we found many transcription factors including *BARX1*, *MSX1*, and *MSX2* in the maxillary process population 2 cluster (MxP2); *SHOX* in mandibular arch 1 (arch1); *PAX7* in lateral nasal process population 2 (LNP2); *SPX* in palatal fusion zone population 1 (palatal.fusion.1); *HAND1* in mandibular arch population 3 (arch3); *HOXA3*, *B3*, and *D3* in fusion mesenchyme population 1 (fusion.mes.1); *MKX* in palatal shelf population 1.1 (palatal.shelf.1.1); and *TBXT* in cartilage population 2 (cartilage.2) among many others (Fig. 5C).

Gene ontology analysis revealed many biological processes, cellular component, and molecular process categories that were relevant for these subtypes (Figure S16). For example, cartilage1 and cartilage2 were differentially enriched for hyaluronic acid and frizzled binding respectively. Cartilage 1 is primarily found in CS20 samples suggesting these are distinct stages of cartilage development. Early osteoblast markers were enriched in pathways regulating pluripotency while late osteoblast markers were enriched for PI3K-AKT signaling and parathyroid hormone response. The more regional based annotations shared many of the same functional enrichments suggesting the same underlying processes were active in these cell types. However, the maxillary process / anterior lateral nasal process derived cells (MxP.aLNP) likely from near the lambdoid junction<sup>86</sup> showed substantially higher expression of many genes related to ribosome production and cytoplasmic translation. Examination of disease enrichments across cluster marker genes revealed some tissue-specific disorders like Osteogenesis imperfecta in late osteoblasts and epiphyseal dysplasia in cartilage 1. Enrichment for genes related to isolated cleft palate were found in several clusters including MxP2, palatal.shelf2.1, palatal.shelf.2.2, and cartilage 1 (Fig 5D).

While the gene ontology analysis confirmed the labeling of some specific cell types, the more positional types remained less clear. To address this, we again turned to the CS13 human spatial transcriptomics data. When we examined some of the markers that defined the mesenchyme versus other cell types, such as *TWIST1* and *PRRX1*, we found fairly broad expression across the embryo with some bias toward the craniofacial region. Other markers like *SATB2* were much more regionally restricted and potentially specifically mark craniofacial mesenchyme versus other types (Fig. 5E). When we examined some of the

438 subtype marker genes, we found much more regionalized expression. *MSX1* was found near many  
439 surface locations with a bias toward the head. *BARX1* was rather specifically localized in the general  
440 region of the pharyngeal arches and the developing stomach. *SPX* and *CYP26C1* were both restricted to  
441 the head region of the embryo at this stage (Fig. 5F). As was observed in mouse, we found much more  
442 regionalized signals from module scores for each subtype. The mandibular arch clusters were clearly  
443 enriched in the pharyngeal arch region of the CS13 embryo and biased toward the more anterior portion  
444 of this region. The lateral nasal process clusters were enriched in distinct areas of the head with *LNP1*  
445 being more posterior and *LNP2* being more anterior. Other subtypes like *MxP2* and *palatal.shelf2.2*  
446 showed good spatial concordance with the labels that had been assigned (Fig. 5G).

## 447 **Characterization of ectodermal cell subtypes**

448 We then turned to the ectodermal cluster to identify potential subtypes. Using the same basic approach  
449 as the mesenchyme, we identified 22 distinct ectodermal clusters (Fig. 6A). Transferring of mouse labels  
450 (Figure S13C) revealed cells that would give rise to specific ectodermal-derived organs like the pituitary  
451 and thyroid, structures of the inner ear (auditory1-3), and surfaces of several structures including  
452 periderm (Fig. 6A). As with mesenchymal subclusters, many of the ectodermal subclusters annotated as  
453 early versus late had biased sample contributions (Fig. 6B). Amongst marker genes of ectodermal  
454 subtypes, transcription factors were again prominent. *LHX3*, *SIX6*, and *PITX2* were most strongly  
455 expressed in the pituitary; *GATA6* marked the palate subtype; the nasal placode (NaP) was identified by  
456 *SP8* and *FEZF1*; auditory subtypes 1-3 were marked by *SALL3*, *GRIN2A*, and *SP9* respectively; *EBF1*,  
457 *EBF2*, and *EBF3* in a single ectodermal subtype (ect.EBF); and *TBX18* marked a putative fusion zone  
458 cluster among others (Fig 6C).

459 Consistent with our findings for CNCCs and mesenchyme, gene ontology analysis revealed many  
460 biological processes, cellular component, and molecular process categories that were relevant for  
461 ectodermal subtypes (Figure S17). For example, markers for all three auditory subtypes were enriched  
462 for terms related to inner ear morphogenesis and development; genes biased for eye ectoderm were  
463 enriched for structural components of the lens; periderm marker genes were associated with the  
464 cornified envelope and skin development; markers for the thyroid cluster were enriched for thyroid  
465 hormone synthesis; and the markers of the pituitary were associated with pituitary gland development.  
466 The less specific clusters such as ectodermal surface clusters were enriched for a variety of categories  
467 suggesting they might be more regionally distinct cell states. In particular, *surface3* marker genes were  
468 biased for oxidative phosphorylation and cytoplasmic translation compared to other ectodermal  
469 subtypes (Figure S17A-E). Examination of enriched human diseases revealed many tissue- or region-  
470 specific disorders including aniridia in the NaP cluster; nonsyndromic deafness in auditory clusters 1  
471 and 2; thyroid agenesis for the thyroid cluster; anterior pituitary hypoplasia for the pituitary cluster; and  
472 congenital cataracts in the eye ectodermal cluster. Interestingly, median cleft lip and palate was only  
473 enriched in the pituitary cluster. Lastly marker genes of the ectodermal cluster expressing high levels of  
474 *EBF* genes (ect.EBF) were enriched for the largest number of disease categories suggesting this might be  
475 a particularly disease relevant cell type or state (Figure 6D).

476 Examination of overall ectodermal markers revealed relatively restricted expression to various surfaces  
477 in the human spatial transcriptomics data. One notable exception being OC90 that was strongly

478 expressed in the location of the putative inner ear (Fig. 6E). Subtype markers also showed generally  
479 restricted expression particularly for *DLX5*, *FOXE1*, and *SIX6*. *PITX2* was expressed in multiple putative  
480 fusion locations in the head but also strongly in the hindlimb (Fig. 6F). Markers of the ect.EBF subcluster,  
481 *EBF2* and *EBF3*, were biased in expression toward the head and pharyngeal arches of the CS13 human  
482 embryo. When we examined the spatial expression for both human CS13 and mouse E11.5, we found  
483 qualitatively different patterns of expression in the craniofacial regions corroborating our previous finding  
484 (Fig. S18). When we inspected module scores of each subtype, we observed exquisitely specific  
485 localization for some clusters like pituitary and auditory. Other clusters were generally enriched at  
486 surfaces of the pharyngeal arches and the putative esophagus (Fig. 6G).

### 487 **Cell-type specific enrichment of genes and variants linked to orofacial abnormalities and normal** 488 **facial variation.**

489 The analysis above demonstrated strong concordance between human and mouse cell types and  
490 subtypes, showed coherent functional and disease enrichments across these cell types, and revealed  
491 spatial enrichments consistent with functions and expected anatomical locations. The strong support of  
492 our labelling of cell types across human craniofacial development, gave us the opportunity to interrogate  
493 the cell type-specific expression profiles for enrichment of craniofacial related genetic signals. The  
494 genetic contributions of common variants to many aspects of craniofacial variation have been studied in  
495 multiple populations based on frontal and profile photographs<sup>87,88</sup>. However, the cell types and  
496 embryonic landmarks that drive these differences are currently unknown. To address this issue, we first  
497 processed the genome-wide summary statistics<sup>87,88</sup> for each craniofacial landmark measurement with at  
498 least one genome-wide significant association using the linkage disequilibrium aware approach Multi-  
499 marker Analysis of GenoMic Annotation (MAGMA<sup>89</sup>). We then calculated expression weighted cell type  
500 enrichments<sup>90</sup> (EWCE) across all the cell types identified in our study using MAGMA-Celltyping<sup>91</sup>. We  
501 observed distinct patterns of cell type enrichment related to different sections of the face. We found that  
502 profile landmark measures related to soft tissues including multiple measures of lip thickness and  
503 shape, ear size, and nose shape were enriched primarily in ectodermal subtypes. Frontal landmark  
504 measures related to aspects of these same portions of the face such as the distance of the outer edge of  
505 the eye to the nasion (ExR-N) showed similar patterns of ectodermal enrichment. Measures likely to be  
506 influenced by bone or cartilage structure such as jaw, chin, and brow protrusion as well the positioning  
507 of the eyes relative to the base of the nose (EnL-Sn) and the mouth (ExR-ChR) were enriched primarily  
508 amongst mesenchymal subtypes (Fig. 7A). Amongst mesenchymal cell types, the mandibular arch,  
509 palatal shelf, and maxillary process fusion zone subclusters had largest number of significant  
510 enrichments for facial shape. The fusion zone cluster and surprisingly the pituitary cluster had the largest  
511 number of significant enrichments amongst ectodermal subtypes. Interestingly, while CNCCs certainly  
512 give rise to many of the downstream cell types and tissues, we found relatively few shape associations  
513 for CNCC subtypes. Overall, this analysis suggests specific cell subtypes contribute differentially to  
514 individual facial differences and suggest these effects begin to manifest very early in human  
515 development.

516 We next turned to studies of the genetic underpinnings of craniofacial abnormalities. In particular, there  
517 have been dozens of genetic associations identified for risk for orofacial clefting in multiple  
518 populations<sup>92-106</sup>. However, the cell types that potentially influence risk for clefting have not been



519 identified in human development. While orofacial clefting has been examined extensively using a variety  
520 of approaches, these studies have been performed in many different populations, making cross-study  
521 comparisons challenging<sup>93,107-115</sup>. Moreover, to identify true positive signals for cell type enrichments  
522 diseases that are not expected to be related to craniofacial cell types examined in the same population  
523 are needed as negative controls. To mitigate these issues, we turned to genome wide association studies  
524 that have been systematically performed on a large cohort of Finnish ancestry<sup>116</sup>. From this resource we  
525 selected all studies annotated as a congenital abnormality by FinnGen with at least one genome-wide  
526 significant association ( $n = 45$ ) as well as two immune related diseases, Crohn's disease and systemic  
527 lupus erythematosus (SLE), that we have used as negative controls in our previous studies<sup>28,117</sup>. When we  
528 analyzed these GWAS using the same approach as for facial variation we found similar partitioning of  
529 enrichment between specific classes of cell types (Fig 7B). We found that ectodermal cell subtypes were  
530 enriched for cleft lip with cleft palate (palate.surface), ankyloglossia, and other congenital malformation  
531 of the tongue and mouth (dental, fusion.zone). Mesenchymal subtypes were enriched for cleft lip or lip  
532 and palate (MxP\_aLNP, mandibular arch 3, fusion mesenchyme subcluster 1 and 2), other congenital  
533 malformations of the ear (multiple palatal shelf subtypes), and congenital malformation of the  
534 musculoskeletal system (cartilage2) among others. CNCC subtypes were most consistently enriched for  
535 other congenital malformations of the upper alimentary tract. The immune cluster was most significantly  
536 enriched for Crohn's and SLE. Many other congenital abnormalities showed no significant enrichments  
537 for any craniofacial cell types demonstrating the specificity of our analyses. A few of these cell types,  
538 MxP\_aLNP and ect\_EBF in particular, were associated with both craniofacial disease and normal facial  
539 variation. These findings suggest that some cell types are contributors to both facial shape as well as risk  
540 for clefting. These results also point to underlying differences in how clefting phenotypes are categorized  
541 which are then in turn related to different subtypes of mesenchyme and ectoderm.

542 Our results from the marker gene ontology enrichments and common variant associations point to  
543 relevant craniofacial disease and phenotype enrichments for specific craniofacial cell types. However, it  
544 is unclear if these cell types might be generally informative for other human phenotypes. We posited that  
545 integrating continuous expression patterns instead of just binary marker gene identity may reveal  
546 additional associations. To address this, we employed a systematic examination of the entire Human  
547 Phenotype Ontology (HPO)<sup>118,119</sup> (Fig 7C). As expected, the immune cluster was systematically enriched  
548 for 90 of 253 phenotypes related to abnormality of the immune system. The red blood cell cluster was  
549 enriched for terms related to abnormality of metabolism and homeostasis (60 of 782 phenotypes). Both  
550 these cell types were enriched for phenotypes related to abnormalities of blood and blood-forming  
551 tissues (150 and 75 of 536 respectively). The endothelium cluster was enriched for abnormality of the  
552 cardiovascular system (60 of 672 phenotypes).

553 Amongst ectodermal subtypes we found the eye subcluster was strongly enriched for phenotypes  
554 related to abnormalities of the eye (75 of 717). Periderm, palate surface, and surface 2 and 3 subtypes  
555 were enriched for abnormalities of the integument. As expected, the pituitary and thyroid subtypes were  
556 associated with abnormalities of the endocrine system. Surprisingly many of the ectodermal subtypes  
557 were enriched for phenotypes related to abnormalities of the respiratory and genitourinary systems.  
558 Among the mesenchymal subclusters, many were enriched for abnormalities of the head and neck. The  
559 cartilage1 cluster showed the most diverse enrichments including phenotypes related to growth

560 abnormalities and abnormalities of the musculoskeletal system, ear, and limb. The main CNCC cluster  
561 was enriched for abnormalities of the nervous system, driven by most of the CNCC subclusters with the  
562 exception of the *cnl1,3*, and 4 subclusters. Surprisingly the specialized ectodermal subtype *ect.GDNF*  
563 was significantly associated with abnormalities of the voice. Together these results suggest that some  
564 subtypes we identified are not specific to the head and are more general states like cartilage. Moreover,  
565 this analysis revealed that while no major cell types were enriched for neoplasms, late CNCCs employ  
566 gene expression programs that likely trigger overgrowth.

### 567 **Differential enrichment of curated gene lists revealed distinct disease risk and role in skull shape** 568 **and/or function across hominid evolution.**

569 Thus far our analysis of the craniofacial cell types has leveraged annotated ontology categories and  
570 common variant associations. Other gene lists that are not part of these systematic ontology databases  
571 and potentially of use to the craniofacial field have not been interrogated. To address this were  
572 assembled multiple gene lists relevant for orofacial clefting including those compiled by CleftGeneDB<sup>120</sup>,  
573 genes co-expressed in important gene modules or prioritized for craniofacial disease in our recent  
574 work<sup>13</sup>, and genes with distinct classes of *de novo* mutations (synonymous vs protein altering) in orofacial  
575 cleft trios sequenced as part of the Gabriella Miller Kids First program<sup>121,122</sup> and CPSeq Studies<sup>123</sup>. We also  
576 curated genes at the extremes of tolerance to loss of function mutations in otherwise healthy  
577 populations that have been suggested to be enriched or depleted of disease relevant genes<sup>124</sup>. Lastly  
578 given our findings for common facial variation across humans, we wondered whether genes potentially  
579 regulated by Neanderthal derived sequences might have craniofacial cell type specific enrichments. As a  
580 control for this evolutionary analysis, we included genes near human accelerated regions, which have  
581 been reported to be enriched in neuronal related functions and expression<sup>125-128</sup>.

582 With these lists in hand, we again employed the expression weighted cell type enrichment  
583 approach. We found that the CleftGeneDB, craniofacial black co-expression modules, and our  
584 prioritized gene lists showed similar patterns of significant enrichments in mesenchymal subtypes  
585 including multiple clusters related to the maxillary process, palatal shelves, and lateral nasal process  
586 (Fig. 8A). Relatively few ectodermal and CNCC subtypes were enriched for these gene lists. The genes  
587 identified by gnomAD to have the least tolerance for loss of function mutations (LOUEF decile 1) were  
588 significantly enriched in many different subtypes identified by our analysis. In particular, *MxP.aLNP*  
589 cluster showed the most significant enrichment. This contrasted with those genes with the most  
590 tolerance for loss of function mutations (LOUEF decile 9) that showed few enrichments and were  
591 generally non-overlapping with the LOUEF decile 1 enrichments (Fig. 8A).

592 When we analyzed the genes near Neanderthal derived sequences, we found patterns of cell type  
593 enrichment distinct from the more disease-focused lists described above. The strongest enrichment was  
594 observed in *pLNP2* mesenchyme subtype. We identified the specialized *ect.EBF* and *ect.GDNF* clusters,  
595 two fusion mesenchyme subtypes, and *cartilage1*. Interestingly all three auditory types were significantly  
596 enriched in this analysis. This was contrasted by only a single cell type identified when examining HAR  
597 associated genes, consistent with their previously published association with brain cell types and  
598 neuronal function<sup>125-127</sup>. We found no consistent, significant enrichments from any of our randomly  
599 selected gene lists across the cell types in questions. We also found no enrichments for the red blood

600 cells across any of these gene lists, and only the gnomAD unconstrained genes for the immune cell types  
601 (Fig. 8A).

602 We then turned to recently identified *de novo* variants from orofacial clefting trios from the  
603 Gabriella Miller Kids First program<sup>121,122</sup> and CPSeq studies<sup>123</sup>. We found no enrichment in any cell types for  
604 genes affected by *de novo* synonymous variants. We found no enrichment in any cell types for genes  
605 affected by *de novo* synonymous variants. However, we identified multiple cell types that strongly  
606 express genes with *de novo* protein altering variants (Fig. 8A). Palate ectoderm showed the strongest  
607 enrichment from this analysis, a cell type that was not enriched for any of the community curated gene  
608 lists related to clefting nor our previous prioritized genes<sup>13</sup>. Multiple other ectodermal cell types were  
609 also identified as enriched including multiple surface ectodermal subtypes, specialized ectoderm  
610 ect.EBF and ect.GDNF, fusion zone ectoderm, and the nasal placode (NaP). Fewer mesenchymal cell  
611 type enrichments were observed but identified the MxP.aLNP and others related the lateral nasal  
612 process (pLNP2, pLNP.fusion).

613 These findings suggest that current disease associations have been biased for genes expressed in the  
614 mesenchyme and that many genes expressed in ectodermal subtypes are also substantial contributors  
615 to clefting risk. To explore this concept further we wondered whether not only the number of genes, but  
616 the total number of *de novo* variants observed in genes might reveal additional disease associations.  
617 When we applied a computational framework that examines gene lists for excess *de novo* mutational  
618 load<sup>129,130</sup> we largely confirmed the findings from the EWCE analysis. We identified 22 clusters for which a  
619 least one phenotype was significantly enriched using a Benjamin-Hochberg false discovery rate of <10%  
620 (Fig. 8B and Supplemental Table 30). These included 17 for all trios with OFCs, 19 for trios with CLP, and  
621 5 for trios with CP. We identified 14 significant enrichments across ectodermal cell types (n=22), 7  
622 enrichments from mesenchymal cell types (n=22) and a single CNCC subtype (n=11). Only 3 clusters  
623 were significantly enriched in all three categories (NaP, palate, and cartilage1), whereas there were 12  
624 shared between all OFCs and CL/P and 2 shared between all OFCs and CP. We also found 5 clusters that  
625 were only significant in the CL/P group and 1 that was only significant in the full cohort. No significant  
626 findings were observed for endothelium, muscle, red blood cell, or immune cell types in our data.

627 For the ectodermal subtypes we identified strongest enrichment *for de novo* variants identified in  
628 the whole cohort and those probands with cleft lip with cleft palate (CL/P) in the nasal placode, surface3,  
629 and palate.surface. We only identified significant enrichment *of de novo* variants from cleft palate only  
630 probands (CP) in the nasal placode and palate ectoderm. While fewer significant enrichments were  
631 observed for mesenchymal subtypes, we found cartilage1 was enriched for all analyses performed.  
632 Interestingly several subtypes were biased toward significant enrichment related to CP vs CL/P. For  
633 instance, MxP2 and palatal shelf 2.1 were enriched for the former while pLNP2 and pLNP.fusion for the  
634 latter (Fig.8B).

635 To explore the genes driving these enrichments we examined the genes with *de novo* damaging  
636 variants that were markers for the nasal placode, the most significantly enriched subtype across our  
637 analysis. As expected, these genes were all expressed in the NaP cells, but were frequently expressed in  
638 many other types of ectoderm to varying degrees (Fig. 8C). In particular, a high degree of sharing of  
639 expression was observed with periderm and multiple surface subtypes, including genes previously

640 implicated in orofacial clefting like *TP63*, *IRF6*, and *CDH1*. Among the *de novo* damaged genes those with  
641 the most biased expression in NaP were *SFRP4* and *DNAH11*. When we examined the localization of the  
642 NaP cells on the spatial transcriptomics data, we found discrete localization at the putative frontonasal  
643 and maxillary processes (Fig 8D). Finally, when we examined these genes for known disease  
644 enrichments, we found enrichment for various types of clefting and craniofacial abnormalities (Fig 8E).  
645 These were driven largely by the genes listed above related to clefting. Interestingly many of the genes we  
646 identified here are expressed in similar patterns to those known disease genes, but have not been  
647 associated with many human disease phenotypes. Amongst these *SFRP4* has the highest specificity of  
648 expression across the main cell types and ectodermal cell types (Fig 8E).

649 Compared to the shared expression and overlapping genes between the NaP and palate clusters, the  
650 genes driving enrichment in cartilage1 were more distinct. Interestingly, although both CP and CL/P were  
651 enriched to a similar degree, the makeup of genes contributing to this signal was different. For CP, the  
652 main driver of the signal was due to *COL2A1* variants, which made up half of the observed variants,  
653 where the remainder were single gene contributions (total n=10). This gene has fairly restricted  
654 expression in the head region and presumptive somites of the CS13 embryo (Figure S19). In contrast,  
655 CL/P probands collectively were enriched within cartilage1, but there were no genes that were  
656 individually overrepresented—only *KCNH5* had multiple variants (2 of total n=18), and the rest were a  
657 single variant per gene. This enrichment highlights the importance of these cells in OFC etiology, but the  
658 difference in signal drivers may provide insight into the heterogeneity of the genetic architecture between  
659 CP and CL/P.

## 660 Discussion

661 Craniofacial abnormalities are some of the most common human birth defects. Only recently have gene  
662 expression patterns active during human craniofacial development been examined<sup>13</sup>. We previously  
663 showed that genes specifically or co-expressed across craniofacial development relative to other tissues  
664 were enriched for known disease-causing genes<sup>13</sup>. However, these analyses relied on bulk gene  
665 expression data from the developing craniofacial tissues. The face is a complex structure that is derived  
666 from multiple cell lineages like ectoderm, mesenchyme, and the specialized neural crest. These major  
667 cell types undergo differentiation to become a variety of distinct cell types that make up the face  
668 including bone, cartilage, muscle, mucosa, and vasculature. Our bulk analyses showed strong bias for  
669 gene programs expressed in human and mouse mesenchyme preventing analysis of genes in ectodermal  
670 and other less abundant cell types. While other single cell atlases from human embryonic development  
671 have been described, there were few biological replicates and relatively few cells clearly derived from  
672 craniofacial regions<sup>24-27</sup>. Moreover, few craniofacial centric analyses have been previously performed on  
673 such data.

674 Our work here attempted to address these shortcomings and concentrate on cell types that are present  
675 across many of the major milestones of human craniofacial development. In this work we profiled  
676 multiple biological replicates from six distinct stages of human craniofacial development. Across these  
677 data we identified seven major cell types present in the developing human face. Most of these, including  
678 mesenchyme, ectoderm, endothelium, blood, and immune cells, have been previously identified in  
679 mouse craniofacial development<sup>12-16</sup>. However, we identified two distinct clusters not described in those

680 previous efforts or labelled as cell types not expected to exist in high levels in craniofacial tissues like glia  
681 or Schwann cells. Our thorough characterization of these clusters using curation of genes from the  
682 literature as well as extensive gene and disease ontology analyses point to these clusters being muscle  
683 progenitors and cranial neural crest. While several protocols for deriving neural crest like cells from  
684 human embryonic stem cells have been described, the primary CNCCs have remained elusive. Also,  
685 only a handful of well-known neural crest genes have been examined using immunohistochemistry in a  
686 small number of early human embryos<sup>71</sup>. Thus, it is unclear the complete repertoire of genes that are  
687 active in this cell type and how closely in vitro models reflect the primary gene expression patterns. Our  
688 analysis here not only established a large number of known marker genes as bona fide CNCC genes,  
689 including *FOXD3* and *SOX10*, but also identifies new genes that could be important for CNCC  
690 specification or function such as *INSC*, *ABCA8*, and *CTXND1*. Our identification of subclusters of the  
691 CNCC including putative melanocytes and the expression programs within them are likely to be useful to  
692 many researchers interested in these cell types. Moreover, identification of this exotic cell type and  
693 subtypes is not a fluke. Generation of data from mouse from similar tissues and stages and uniform  
694 process reveal these same populations. Upon close inspection of the gene ontology enrichments, other  
695 groups may have mistakenly labelled these cells as glia or Schwann cells simply because of biases in the  
696 ontology databases. Far more research has been performed on the human brain and related cell types  
697 than other parts of the body, likely resulting in many more brain related gene ontology annotations. While  
698 automated and machine learning based approaches are gaining traction for labelling of single cell  
699 atlases<sup>131-137</sup>, transient developmentally related cell types that are not in current databases and biases in  
700 ontology will still require close inspection and interpretation.

701 By generating comparable datasets from both mouse and human we had the unique opportunity to  
702 identify both shared and species-biased gene programs active in individual cell types. As expected, we  
703 found the main cell types identified in each species share the most significant amount of marker genes  
704 with the orthologous cell type in the other species. Among these, mesenchyme was the most  
705 functionally shared between human and mouse based on marker gene expression. Surprisingly, CNCC  
706 markers were the least shared between these species, even less than cells from the immune system that  
707 has been documented to have substantial differences across humans and mice<sup>17</sup>. This could reflect  
708 substantial functional differences in CNCC between human and mice and indicate that this cell type  
709 may be particularly labile across evolution allowing innovation of craniofacial shape as others have  
710 proposed<sup>1,18-20,70,138</sup>.

711 Although there was the largest degree of shared marker gene expression within mesenchyme, we found  
712 hundreds of differences in marker gene identity between human and mouse. While we restricted our  
713 analysis to genes with clear one-to-one orthology between these two species, some of these differences  
714 could be due to mis-annotation of orthology, substantial developmental heterochrony, or the inherit  
715 noisiness of current single nucleus gene expression data. However, by focusing on coherent gene  
716 ontologies and strongly expressed genes we identified many genes that are likely to reflect true species  
717 differences. For instance, one of the top human mesenchymal markers based on absolute and  
718 specificity of expression that was not revealed in mice was *ALX3*. Recessive mutations in human *ALX3*  
719 have been linked to frontorhiny or frontonasal dysplasia 1 (OMIM 136760)<sup>139</sup>, while the *Alx3*<sup>-/-</sup> mouse has  
720 been reported to have no phenotype<sup>140</sup>. Our analysis also identified *MSX2*, to which humans have been



721 suggested to be much more sensitive than mice to dosage of this transcription factor during craniofacial  
722 development<sup>141</sup>. Further analysis of all these subtypes and comparison with additional species could  
723 reveal novel functional differences as well as the core regulatory programs that are present in all  
724 vertebrates.

725 While we highlighted some of the species differences in major cell types that could be relevant for what  
726 human genes and diseases can be modelled in mice, our comparison framework allowed us to  
727 accurately identify subtypes of each major cell types between species. This allowed us to leverage the  
728 substantial single-cell and spatial transcriptomics resources as well cell type annotations that have  
729 been generated by many different groups<sup>23,86,142</sup>. By transferring functional and spatial labels for mouse  
730 cell subtypes to our human data we could add such information to data that were originally lacking. We  
731 confirmed these labels using a variety of gene and disease ontology analyses, but most convincingly by  
732 leveraging previously published spatial transcriptomics data for a CS13 human embryo<sup>25</sup> By using marker  
733 genes to calculate module scores across this spatial data we confirmed relevant anatomical regions  
734 from which each subtype was potentially derived. We were able to identify some exquisitely specific  
735 spatial locations for ectodermal subtypes related to the ear, eye, and pituitary. We also identified  
736 expected regionalized expression for mesenchymal subtypes putatively derived from the mandibular  
737 arch as well as important fusion zones like the lateral nasal process. Further characterization of the  
738 markers we identified in higher resolution spatial transcriptomics across multiple sections and  
739 reconstruction into a complete three-dimensional representation as has been recently described for  
740 mice will be necessary to validate these findings<sup>143</sup>.

741 One of the major goals of generating such resources is to enable better understanding of human  
742 phenotypes and disease. Not only can facial abnormalities affect our capacity for communication and  
743 feeding, but the face is also one of the most defining features of each human and is intimately tied to our  
744 sense of individuality. Thus, understanding how facial shape is encoded in our genomes is of substantial  
745 general interest. In recent years coupling of two- and three-dimensional imaging approaches with large  
746 scale genotyping has enabled the discovery of common genetic variants associated with quantitative  
747 differences in many different facial landmarks<sup>87,88,144-146</sup>. While these variants have been shown to be  
748 enriched in regulatory regions active in the developing face, the cell types that underly facial differences  
749 were unknown. Using our highly confident cell subtype annotations, we found distinct differences in  
750 enrichments for measurements across the human face. In general, the enrichments we observed were  
751 mutually exclusive, features likely driven by mesenchyme subtypes not associated with an ectodermal  
752 subtype and vice versa. As expected, mesenchyme subtypes were associated with features that are  
753 likely driven by hard structures like bone and cartilage while ectoderm subtypes were associated with  
754 some measures that are related to soft tissue shape or thickness. The most consistent associations  
755 observed were related to variation in measures of the midface. These were significantly enriched for  
756 many mesenchyme subtypes that we annotated as derived from regions that are consistent with these  
757 effects: the maxillary process, palatal shelves, and fusion zone mesenchyme. We did not observe any  
758 subtype that contributed to all aspects of the face, nor did we observe significant subtype enrichments  
759 for all measurements. These landmarks may be driven by cell types that appear later in development or  
760 be influenced by subtle gene expression differences in many cell subtypes. While the two studies we  
761 utilized were performed in populations with distinct ancestries and yielded consistent results, it is

762 possible that subtypes could influence facial variation differently in other populations. Further  
763 identification of genetic associations with more facial measures in a more diverse set of individuals and  
764 identification of cell types later in craniofacial development will be needed to address this issue.

765  
766 Craniofacial abnormalities are among the most common birth defects in humans. The most common  
767 form of these, nonsyndromic cleft lip and/or cleft palate, is thought to occur relatively early in human  
768 development between 4 and 6 weeks<sup>8,147,148</sup>. Consistent with this idea, we found that variants associated  
769 with risk for orofacial clefting are enriched in regulatory sequences active in craniofacial tissues from  
770 this developmental window<sup>28</sup>. However, the cell types in which these variants manifest their effects were  
771 unknown. Here we used uniformly generated and processed genome-wide association data for many  
772 congenital abnormalities in the Finnish population. This population has been shown to have a high  
773 incidence of clefting with interesting geographical distributions<sup>149</sup>, and we reasoned would serve as an  
774 excellent test case for subtype enrichments across relevant and unrelated diseases. Indeed, we found  
775 some subtypes of both mesenchyme and ectoderm were significantly enriched for orofacial clefting or  
776 other abnormalities of mouth. We found some overlap between phenotypes and subtypes particularly  
777 related to cardiac outflow tract abnormalities consistent with the neural crest derived nature of those  
778 structures<sup>150-154</sup>. We found expected cell type specific enrichments for immune cells in the autoimmune  
779 related diseases that we included from this cohort, SLE and Crohn's. We also did not observe  
780 enrichment for most subtypes in most abnormalities outside the craniofacial and cardiac structures.

781 Interestingly several of the enrichments we observed for subtypes were shared across the craniofacial  
782 variation and craniofacial abnormality analyses. The MxP.aLNP and ect.EBF subtypes were examples  
783 that had several significant associations in both phenotypes. This is particularly interesting as it has been  
784 speculated that some of the same processes may be at play<sup>102,155-158</sup>. Our findings here suggest that some  
785 cell types play an outsized role in landmarks of the midface region and risk for orofacial clefting. Our  
786 analysis of marker genes for these specialized subtypes suggests these two subtypes are near one  
787 another spatially and could be located near the fusion zone termed the "lambdoid junction"<sup>86,159-161</sup>.  
788 Failure of this region to fuse in humans has been suggested to cause cleft lip that could also involve the  
789 nostril region and primary palate<sup>10,147,162,163</sup>. It is thus relatively straightforward to imagine that subtle  
790 differences in the timing of migrations and fusion of cells residing in this region could influence the shape  
791 of the midface. Interestingly, some of the major markers of the specialized ectodermal subtype are  
792 multiple members of the EBF family of transcription factors. Our previous work suggested that these  
793 genes were co-expressed more strongly in human craniofacial cell types than mouse, and found  
794 compelling evidence that *EBF3* is a *bona fide* orofacial clefting risk gene<sup>13</sup>. This EBF family of  
795 transcription factors have been linked to regulation of differentiation of multiple different tissue types  
796 and predisposition for several tumor types<sup>164-170</sup>. The timing of differentiation of cells at a fusion zone  
797 could influence the degree to which structures fuse and impact both clefting risk and facial shape.  
798 Studies leveraging the marker genes we have identified for each of these subtypes could allow more  
799 specific labelling and identification of these cells in human tissues and mouse embryos as well as  
800 experiments to test impact of facial variation.

801 Our analysis of curated gene lists that are not included in standard gene ontologies was also revealing  
802 related to both the cell type identities as well as the composition of the gene lists themselves. For  
803 instance, our previous prioritized gene list as well as the curated CleftGeneDB resource are heavily  
804 biased toward some mesenchymal subtypes. This is not surprising given the ratios of cell types we  
805 observed in the data generated here. Mesenchyme is by far the dominant major cell type, thus previous  
806 studies of gene expression and protein expression from bulk tissues were heavily biased toward this cell  
807 type. The genes identified as constrained in human populations were more broadly enriched across all  
808 the cell subtypes suggesting they play critical roles in most cell types in the body. As expected, the  
809 unconstrained genes were enriched in relatively few cell types and were not enriched in the likely  
810 craniofacial disease relevant subtypes. While both the common variant analyses for orofacial clefting  
811 and the curated craniofacial disease gene lists were biased toward mesenchyme subtypes, genes  
812 harboring rare *de novo* protein damaging variants identified in cleft probands showed much more  
813 enrichment in ectodermal subtypes. This trend was not observed for *de novo* synonymous variants  
814 suggesting this was not a population specific effect or other artifacts of sequencing. This trend was  
815 further supported when we examined the frequency of *de novo* protein altering variants, where we found  
816 significant enrichment in multiple ectoderm subtypes primarily for CL/P. While the number of CP only  
817 cases were fewer than CL/P, we found these *de novo* variants were enriched in a few mesenchymal  
818 subtypes that make sense for a spatial perspective. Overall, this points to the ectodermal subtypes, that  
819 as we discussed above make up a small proportion of craniofacial tissue, as a major contributor to  
820 clefting risk. Due to the biases of previous studies for the most abundant cell types there are likely many  
821 additional clefting risk genes that remain to be discovered. The resources we described here could help  
822 further prioritize genes that are discovered in such sequencing cohorts. For instance, the nasal placode  
823 ectodermal subtype was marked the most substantial number of genes with *de novo* damaging variants.  
824 Many known disease risk genes are expressed in this subtype thus genes with similar patterns of  
825 expression or specificity of expression could be guilty by association. In particular, our analysis  
826 highlighted the *SFRP4* gene. This gene has been linked to Pyle disease (OMIM 265900) that features bone  
827 abnormalities and fragility particularly of the long bones and GWAS of bone mineral density<sup>171-174</sup>. Similar  
828 phenotypes are observed in *Sfrp4* knockout mice<sup>175</sup>. Cell type specific dysregulation of this gene either  
829 due to somatic mosaicism or regulatory element disruption could result in bone abnormalities or other  
830 defects in a relevant part of the developing face. Further studies of this gene in a craniofacial specific  
831 context in mice as well as identification of the regulatory landscape controlling could reveal a role in  
832 clefting risk.

833 As detailed above, our analysis of craniofacial variation revealed multiple cell types that contribute to  
834 human facial shape. Beyond interindividual differences there have been reported to be substantial  
835 differences in the shape of many craniofacial features between modern humans and of closely related  
836 but extinct hominid species such as Neanderthal and Denisovans<sup>176-178</sup>. Identifying the genetic  
837 contributions to these differences and if Neanderthal derived sequences in the human genome  
838 predispose individuals to specific phenotypes or diseases has been of particular interest<sup>179-185</sup>. While  
839 Neanderthal derived variants in genes and regulatory regions active in adult bulk tissues have been  
840 linked to specific phenotypes related to brain and cranium shape, immunity, and adipose function<sup>186-189</sup>,  
841 it is unknown if any human developmental cell types might be influenced by such variants. Our analysis

842 points to Neanderthal derived regions in the European genetic background are systematically enriched  
843 near genes with biased expression in multiple cell types related to ear development, cartilage, and  
844 specialized ectodermal subtypes. *Cartilage1* and *EBF* expressing ectodermal subtype (ect.*EBF*) were  
845 also shown to be enriched for both *de novo* protein damaging variants in orofacial clefting probands and  
846 several aspects of modern human facial variation. These results could suggest that risk for orofacial  
847 clefting and facial shape could both be influenced by Neanderthal introgression events. We did not  
848 observe any such enrichments for sequences that have been shown to be accelerated on the human  
849 lineage, suggesting that findings are functionally relevant. Consistent with this idea, the Neanderthal  
850 derived analysis was the only one that demonstrated enrichment in all the ear related ectodermal  
851 subtypes. Multiple aspects of Neanderthal inner ear morphology have been shown to differ substantially  
852 from modern humans and other primates<sup>190,191</sup>. We also note that we observed enrichment of  
853 Neanderthal introgressed regions near genes with biased expression in the specialized ectodermal  
854 subcluster ect.*GDNF*. This was the lone subtype that was enriched for abnormalities of the voice.  
855 Differential DNA methylation patterns between modern humans and Neanderthals and Denisovans  
856 indicated genes related to vocal anatomy are regulated in a distinct fashion<sup>192</sup>. These findings open the  
857 distinct possibility that the degree of introgressed segments in the genomes of modern human  
858 individuals could influence ear morphology and hearing capabilities as well as vocal characteristics.

859 In summary we have provided a substantial resource for understanding the cell types and gene  
860 expression patterns that build the human and mouse face. Our analyses revealed relationships between  
861 specific cell subtypes and many aspects of human biology including facial shape and orofacial clefting  
862 risk. We also illuminated potential contributions of ancient hominids to craniofacial morphology. Future  
863 integration with cell type specific chromatin accessibility could reveal specific variants and regulatory  
864 regions that encode such phenotypic differences, risk factors, and species-specific biology. This data  
865 can be explored through an interactive web application that is accessible to most researchers:  
866 [https://cotneyshiny.research.chop.edu/shiny-apps/craniofacial\\_all\\_snRNA/](https://cotneyshiny.research.chop.edu/shiny-apps/craniofacial_all_snRNA/). The data will be deposited  
867 to other major single cell aggregation databases including the Chan-Zuckerberg CellXGene Discover  
868 resource<sup>39,193</sup>.

## 869 **Acknowledgements**

870 We would like to thank members of the UConn/JAXGM Single Cell Genomics Core for help with  
871 standardizing single-cell isolation techniques and preparing sequencing libraries. We would also like to  
872 thank members of the UConn Computational Biology Core and High-Performance Computing Facility for  
873 assistance with package installation and software/hardware support. We are grateful to Dr. Peter Tran,  
874 Dr. Sungryong Oh, and Pooja Sonawane for constructive feedback and copyediting. This work was  
875 funded by grants from the National Institutes of Health to E.J.L.C (R01-DE030342, X01-HG010835, X01-  
876 HD100701, X01-HL132363) and JC (NIDCR 1R01DE028945, NIDCR 1R03DE028588, and NIGMS  
877 5R35GM119465).

## 878 **Author contributions**

879 Conceptualization: J.C. Investigation: N.F., E.W.W. and J.C. Formal analysis: N.F., E.W.W, B.M.S. K.R.,  
880 S.W.C, J.C. Writing—original draft: J.C. Writing— review and editing: N.F., E.W.W, B.S, K.R., S.W.C,  
881 N.G.S, E.J.L.C., J.C. Funding acquisition: J.C. Supervision: J.C., E.J.L.C., S.G.K.

## 882 **Code and data availability**

883 Code for analysis and generation of figures can be found on github  
884 ([https://github.com/cotneylab/craniofacial\\_snrna](https://github.com/cotneylab/craniofacial_snrna)). An interactive website for exploring processed data  
885 is found here: [https://cotneyshiny.research.chop.edu/shiny-apps/craniofacial\\_all\\_snRNA/](https://cotneyshiny.research.chop.edu/shiny-apps/craniofacial_all_snRNA/). Raw data  
886 from mouse experiments generated in this work will be deposited in GEO. Cellranger ARC gene  
887 expression outputs for both human and mouse are available on Zenodo.

## 888 **Methods**

### 889 *Human tissue samples*

890 The use of human embryonic tissue was reviewed and approved by the Human Subjects Protection Program  
891 at UConn Health (UCHC 710-2-13-14-03) and Children’s Hospital of Philadelphia (IRB 24-022258). Human  
892 embryonic craniofacial tissues were collected via the Joint MRC/Wellcome Trust Human Developmental  
893 Biology Resource (HDBR) under-informed ethical consent with Research Tissue Bank ethical approval  
894 (18/LO/0822 and 18/NE/0290, project 200225). Donations of tissue to HDBR are made entirely voluntarily by  
895 women undergoing termination of pregnancy. Donors are asked to give explicit written consent for the fetal  
896 material to be collected, and only after they have been counseled about the termination of their pregnancy.  
897 Further documentation of all policies and ethical approvals for HDBR sample collection can be found  
898 at <https://www.hdbr.org/ethical-approvals>. Tissues were flash-frozen upon collection and stored at  $-80^{\circ}\text{C}$ .  
899 Upon thawing, the samples were quickly inspected for intactness of the general craniofacial prominences and  
900 processed for single nucleus multiomics.  
901

### 902 *Mouse embryonic tissue samples*

903 The use of mouse embryonic tissues was reviewed and approved by the UConn Health Institutional Animal  
904 Care and Use Committee (Protocol AP-2000061-0723). Eight-week-old wild-type male and female C57BL6/J  
905 mice were obtained from Jackson Laboratory. Mice were housed according to recommendations by Jackson  
906 Laboratory with 12 h light:dark cycle beginning at 7 a.m. The ambient temperature was maintained between 20  
907 and 22 °C and humidity was maintained at 40–60%. Mice were given ad libitum access to food and water.  
908 Timed matings were established by the identification of vaginal plugs the morning following the housing of a  
909 single male with multiple female mice. Embryos were harvested from pregnant mothers at mid-day either 10,  
910 11, or 12 days after identification of the vaginal plug. The staging was confirmed by counting somites and  
911 comparing overall morphology to the Theiler Staging Criteria<sup>194</sup>. All embryos from a given litter were combined  
912 for individual biological replicates, and at least three biological replicates were collected and processed for  
913 each stage. Craniofacial prominences were collected in a very similar fashion to human samples and  
914 subsequently prepared for single nucleus multiomics.  
915



## 916 *Single nucleus multiomics*

917 Primary human craniofacial tissues from CS12, CS13, CS14, CS16, CS17 and CS20, each stage represented  
918 by a minimum of 3 replicates, were obtained from HDBR. Tissue from each embryo were mechanically broken  
919 into single-cell suspensions and cells were checked for viability counted using Trypan blue staining following  
920 the 10X Genomics protocol for single-cell multiome sequencing using the ChromiumX controller. Samples  
921 were sequenced on multiple Illumina NovaSeq runs according to 10X Genomics recommendations. Raw  
922 fastqs were processed using CellRanger ARC (v2.0.2) using hg38 genome and gene annotations provided by  
923 10X Genomics.

924  
925 Primary mouse craniofacial tissues from E10.5-E12.5 from multiple (3-5 depending on stage) mixed sex  
926 C57BL/6J Mus Musculus embryos (Jackson Laboratories) were pooled. Animals were raised and sacrificed in  
927 compliance with UConn Health IACUC approval (protocol AP-200061-0723). Samples were mechanically  
928 broken into single-cell suspensions, processed for multiome using the ChromiumX controller, and sequenced  
929 in the same fashion as for human samples above. Raw fastqs were processed using CellRanger ARC (v2.0.2)  
930 using mm10 genome and gene annotations provided by 10X Genomics.

## 931 932 *Processing of snRNA and identification of major cell types.*

933 Filtered barcode matrices from each human samples generated by CellRanger ARC were individually loaded  
934 with Read10X\_h5 command in Seurat<sup>195</sup> and merged into one object. Percentage of mitochondrial reads were  
935 calculated for each cell and filtering was performed to only retain cells with less than ten percent  
936 mitochondrial derived. Further filtering was performed based on number of counts per cell ( $500 < x < 25000$ )  
937 and number of genes detected per cell ( $500 < x < 7000$ ). Filtered data were normalized with default values and  
938 cell cycle scores were calculated using Seurat. Data was scaled based on S and G2M score regression and  
939 dimensionality reduction with principal component analysis (PCA) were performed using respective  
940 commands in Seurat. The top 2000 variable features were identified and data were then further integrated with  
941 harmony R package<sup>196</sup>. Nearest neighbors based on harmony corrected embeddings were calculated with up  
942 to 30 dimensions and clusters were identified with multiple resolutions from 0.1 to 1 in Seurat. We then  
943 performed uniform manifold approximation and projection (UMAP) dimensionality reduction using harmony  
944 corrected embeddings in Seurat (dimensions = 30, minimum distance = 0.3). Resulting clusters were  
945 inspected for expression of multiple craniofacial markers from Li et al 2019 and marker genes were identified  
946 for each cluster. Cells from clusters identified with high expression of neuronal markers *TUBB3* and *MAP2*  
947 were removed and the process of normalization, harmonization, and clustering was repeated with remaining  
948 cells from all samples. Marker genes for major cell types were identified using FindAllMarkers (logfc.threshold  
949 = 0.25, min.pct = 0.1, test.use = "wilcox", min.cells.feature = 3, mi.cells.group = 3, pseudocount.use = 1, and  
950 return.thresh = 0.01). The top 100 marker genes for each cluster ( $p < 0.05$ , ranked by log2fold change versus  
951 all other clusters) were then analyzed for gene and disease ontology enrichments using compareCluster in  
952 clusterProfiler R package (v. 4.12.6). Cranial neural crest genes were compiled based on markers identified by  
953 regulatory network construction in human cultured CNCC and craniofacial tissue data<sup>69</sup>. Major cell type labels  
954 were applied to each cluster.

956 For mouse data sets, filtered barcode matrices from each mouse E10.5 to 12.5 samples generated by  
957 Cell Ranger ARC were individually loaded with Read10X\_h5 command in Seurat<sup>195</sup> and merged with E13.5 to  
958 E15.5 data from Pina et al 2023 (GSE205448). Calculation of percent mitochondrial reads and filtering were  
959 performed with similar thresholds to human data above. Subsequent harmonization, dimensionality  
960 reduction, and clustering were performed identically to those for human data above. Less significant  
961 contamination of neuronal cell types was observed in mouse data, which was identified and filtered as  
962 describe for human. Identification of marker genes and gene ontology enrichments, and CNCC modules  
963 scores were performed as above for human data. Major cell type labels were applied to each mouse cluster.  
964

### 965 *Marker gene comparisons across species*

966 Lists of all marker genes for each of the seven main subtypes for each species ( $p < 0.05$ ) were compiled and  
967 orthology based on HGNC symbol annotated by Ensembl v105 was obtained using the getLDS command in  
968 the biomaRt R package (v. 2.60.1). Only genes that had one ortholog in each species and a HGNC symbol were  
969 retained ( $n = 7504$ ). Significant overlaps between all orthologous marker gene lists were determined using the  
970 testGeneOverlap command in GeneOverlap R package<sup>197</sup> Conserved and species-specific genes were  
971 determined based on HGNC symbol and the intersection matrix obtained by getMatrix in GeneOverlap. Gene  
972 and disease ontology enrichments were calculated using clusterProfiler.  
973

974 Final Seurat objects were prepared for display in an interactive webapp using the ShinyCell R Package<sup>198</sup>.

### 975 *Subclustering of major cell types*

976 For major cell types labelled as mesenchymal, ectodermal, or CNCC further subclustering was first  
977 performed on mouse data. For each major cell type, normalization, scaling with regressed cell cycle impacts,  
978 harmonization, and subclusters were identified using same procedure as described above. Marker genes were  
979 identified for each cluster and functional enrichments were determined using clusterProfiler. Annotations for  
980 each cluster were manually assigned based on those originally described<sup>12-16</sup>. Mouse cell subtype  
981 assignments were further confirmed with ToppGene<sup>199</sup> using the scToppR package<sup>200</sup>. Mouse main and  
982 subtype annotations were further confirmed by projection on mouse E15.5 spatial transcriptomics data<sup>23</sup>  
983 (GSE245469). Links between snRNA and spatial data were determined using FindTransferAnchors and  
984 transferred using TransferData in Seurat.  
985

986 Following annotation of subtypes and for comparison with human data, an intermediate data set was created  
987 where mouse genes were reduced and converted to those to those with one to one orthology with human  
988 genes using annotations provided by Ensembl (archive dec2021) with biomaRt R package<sup>201</sup>. The intermediate  
989 dataset was used to transfer mouse subtype annotations to human subtypes by first identifying shared  
990 features across clusters using FindTransferAnchors in Seurat with log normalization and canonical correlation  
991 analysis (cca). Predicted subtype labels were transferred to human subtypes using TransferData in Seurat and  
992 further confirmed with a confusion matrix. Marker genes for human subtypes were identified as performed for  
993 major cell types and functional enrichments were characterized with compareCluster in clusterProfiler. Final  
994 Seurat main objects and subtype objects were prepared for display in an interactive webapp using the  
995 ShinyCell R Package<sup>198</sup>. Seurat objects were also converted to scanpy and anndata objects using scEasy R  
996 package (v0.0.7) for hosting at the Chan Zuckerberg CELL by GENE Discover resource.

997

## 998 *Processing of spatial transcriptomics*

999 Spatial transcriptomics data for two sections of a human CS13 embryo<sup>25</sup> were retrieved from  
1000 <https://heoa.shinyapps.io/code/>. Raw sequence count matrices were loaded using the Read10x command of  
1001 Seurat<sup>195</sup> and converted to HDF5 format. These counts were then combined with spot coordinates and section  
1002 images using CreateSeuratObject. Data from both slices were merged and variable features were determined  
1003 using Seurat. The percentage of mitochondrial reads was determined for each cell and was used to transform  
1004 all data in the merged object using SCTransform from Seurat. Data was clustered using UMAP and plotted  
1005 which revealed a strong batch effect between the two spatial objects. Data was further normalized using  
1006 Harmony (v1.2.3), projected using UMAP, and clustered with a resolution of 0.8. Marker genes for the 22  
1007 clusters were identified using FindAllMarkers in Seurat. The top 100 marker genes for each cluster ( $p < 0.05$ ,  
1008 ranked by log<sub>2</sub>fold change versus all other clusters) were then analyzed for gene and disease ontology  
1009 enrichments using compareCluster in clusterProfiler R package (v. 4.12.6). Enrichments and spatial  
1010 localization were compared to previous annotations by Xu et al 2023 and labelled accordingly. The top 100  
1011 marker genes from each of the subclusters identified in human craniofacial data were used to calculate  
1012 module scores across the merged spatial object and plotted using SpatialFeaturePlot in Seurat.

1013

1014 We chose mouse E11.5 spatial transcriptomics data as it is most morphologically similar to CS13 human  
1015 embryos. Data were retrieved all E11.5 spatial transcriptomics data from the MOSTA resource<sup>43</sup> and loaded  
1016 into Seurat as for human data above. For each section, module scores of top 100 marker genes for each main  
1017 cluster or subcluster were calculated. Gene spatial feature plots for selected genes and module scores were  
1018 then generated with Seurat.

1019

## 1020 *Facial variation and congenital abnormality GWAS enrichments*

1021 We retrieved summary statistics for facial variation<sup>87,88</sup> and all congenital abnormality GWAS summary  
1022 statistics from FinnGenn<sup>116</sup>. Raw summary stats were further processed and standardized with hg38  
1023 coordinates with MungeSumstats R package (<https://doi.org/10.1093/bioinformatics/btab665>). Variants were  
1024 mapped to genes +/- 100kb using MAGMA<sup>89</sup> Frontal facial measures and FinnGenn were processed based on  
1025 1000 genome European population while profile facial measures were processed with the 1000 genome  
1026 Middle/South American population all obtained from the MAGMA website (<https://cncr.nl/research/magma/>).  
1027 We converted the Seurat snRNA-seq expression data to a CellTypeData set with the Expression Weighted  
1028 Celltype Enrichment (EWCE) R package<sup>90</sup> and then assessed each study trait for a linear positive correlation of  
1029 cell type gene expression specificity and gene-level genetic associations using MAGMA Celltyping<sup>91</sup>. Plots  
1030 were generated using tidyheatmaps in R<sup>202</sup>.

1031

## 1032 *Gene list enrichments per cell type*

1033 The CellTypeData-formatted human craniofacial snRNA-seq objects were generated using  
1034 generate\_celltype\_data in EWCE (v1.15.0). Mean and specificity metrics for several marker genes (*SOX10*,  
1035 *TP63*, and *MSX1*) were inspected across main cell types and subtypes using plot\_ctd in EWCE. Gene lists were  
1036 compiled from multiple resources including gnomAD (v4.1), CleftGeneDB, prioritized genes and black module

1037 from craniofacial WGCNA<sup>13</sup>, and genes affected *by de novo* variation in orofacial clefting probands<sup>32,123</sup>. For  
1038 Neanderthal introgressed regions and human accelerated regions, coordinates were obtained from respective  
1039 publications<sup>188,203</sup> and assigned single nearest gene using rGREAT<sup>204</sup> with “oneClosest” association rule. Each  
1040 gene list was then tested for linear association using bootstrap enrichment test in EWCE (reps = 10,000;  
1041 geneSizeControl = TRUE). Results from all gene lists were then merged and plotted with ewce\_plot in EWCE  
1042 with correction for total number of gene lists and cell types tested using the Benjamini-Hochberg approach.  
1043

#### 1044 *Phenotype-cell type association tests*

1045 To map the relationships between cell types and phenotypes, we ran pairwise association tests between all  
1046 combinations of cell types in our snRNA-seq-derived CellTypeData and phenotypes across the Human  
1047 Phenotype Ontology (HPO)<sup>119</sup> using the run\_phenomix function from MSTExplorer (v1.0.5). In contrast to the  
1048 gene list-based approaches (e.g. EWCE) this function reframes the problem as a series of linear regressions  
1049 by leveraging continuous scores that summarize the current strength of evidence for a causal relationship  
1050 between each gene-phenotype pair (using additional data from the Gene Curation Coalition)<sup>118,205</sup>. The  
1051 continuous nature of this data allows us to more accurately capture phenotype-cell type relationships,  
1052 especially for phenotypes with large gene lists where only some genes have strong evidence of actually  
1053 causing the phenotype. The gene signature vectors for each phenotype were previously merged and shared as  
1054 a single precomputed gene (5003 unique gene symbols) x phenotype (11047 unique HPO phenotypes)  
1055 association matrix. Next, a series of linear regressions tests were performed between the gene specificity  
1056 vectors of each cell type (n=66 vectors) and the gene association vectors of each phenotype (n=11047  
1057 vectors). Finally, multiple-testing correction was applied using Benjamini-Hochberg False Discovery Rate<sup>206</sup>  
1058 (at FDR<5% significance).  
1059

1060 For the purposes of summarization and visualization, the number of significantly associated phenotypes per  
1061 cell type were then computed within each major HPO branch (Fig. 7C). Here, we define HPO branches as  
1062 groups of related phenotypes that can be labeled according to their shared ancestral term, e.g., ‘Abnormality  
1063 of the immune system’. Next, we sought to determine whether some cell types were disproportionately more  
1064 often associated with phenotypes of a particular HPO branch. To accomplish this, we performed a series of  
1065 proportion tests comparing the proportion of total phenotypes that a given cell type was significantly  
1066 associated with within a target HPO branches relative to all other HPO branches. In practice, we computed  
1067 2x2 contingency tables (number of significant phenotype association vs. number of non-significant phenotype  
1068 associations x target branch vs. non-target branches) for each cell type within each HPO branch, which were  
1069 then used as inputs to the prop\_test function within the rstatix R package (v0.7.2). This test appropriately  
1070 takes into account the different number of phenotypes across HPO branches. Only one-sided tests were  
1071 performed to test whether the target HPO branch was greater than all other (non-target) branches (set with the  
1072 alternative = "greater" parameter). All proportion tests were then corrected for multiple testing at FDR<5%.  
1073

#### 1074 *Orofacial Clefting de novo variant analysis*

1075 We used the R package ‘DenovolyzeR’ (version 0.2.0) to test enrichment of *de novo* variants (DNs) in a  
1076 dataset of OFC case-parent trios. Enrichment is calculated by comparing the expected number of  
1077 variants, as determined by mutation models described by Samocha, et al<sup>130</sup>, to the observed number of

1078 variants in a given gene or group of genes using the ‘DenovolyzeByClass’ and ‘includeGenes’ functions.  
1079 Using our dataset of 2031 DNs in 1171 genes identified in 1676 trios with OFCs, we first compared this  
1080 list of genes to those with calculated mutational rates in the R package ‘DenovolyzeR’ (version 0.2.0)  
1081 using the ‘viewProbabilityTable()’ function. There were 12 trios in which DNs were identified, but no  
1082 mutational rates for the affected genes were present; thus, we ultimately tested 1662 trios with OFCs,  
1083 broken down by subtype including 1180 cleft lip with or without cleft palate (CL/P; 226 cleft lip (CL), 954  
1084 cleft lip and palate (CLP)), and 482 cleft palate (CP) trios. We then tested enrichment of all OFC trios and  
1085 by subtype within the top 20% of genes by log<sub>2</sub>FC derived from single nucleus RNA sequencing of human  
1086 craniofacial tissue at CS20.

1087



## Figure Legends

Figure 1. Generation of single nucleus gene expression atlas of human craniofacial development.

A). Anatomical regions of the developing craniofacial region from 4 to 8 weeks post conception. Individual Carnegie Stages (CS) and replicates at stage are indicated below images. B). Pseudo-bulk gene expression of tissues from each stage displayed in principal component (PC) space based on the first two PCs. Progression of developmental time is indicated along PC1 dimension. C) UMAP projection and cluster identification of all human craniofacial cells after filtering of neurons. D). Number of cells obtained at each CS stage for each cluster identified in C. E). Distribution of samples from each sample across the UMAP projection.

Figure 2. Identification of main cell types in the developing human face.

A). Gene expression feature plots for indicated genes across UMAP projection. B) Average and percent expression for the top 10 marker genes for each main cluster. C). Disease ontology enrichments of categories curated by DisGeNet for each indicated cluster. D). Identification of CNCC cluster based on module score of curated neural crest genes and labelling of all remaining clusters. Violin plots and individual values for all cells of a given cluster type based on CNCC module score calculated by Seurat.

Figure 3. Identification of CNCC subtypes in the developing human face.

A). UMAP projection of subclustered CNCC main cell type. B). Contribution of cells from each CS timepoint to each CNCC subcluster. C) Violin plots of published neural crest marker genes across each subcluster. D). Average and percent expression for the top 5 marker genes for each of the CNCC subclusters. E). Gene expression spatial feature plot for indicated CNCC marker genes in two sections from a CS13 human embryo. F). Gene expression spatial feature plot for indicated CNCC subtype marker genes in same sections as E. G). Spatial feature plot for module scores calculated from top 100 marker genes from indicated CNCC subtype.

Figure 4. Single-nucleus gene expression in the developing mouse face.

A). UMAP projection of all cells profiled by this study and combined with published studies. Major cell types are indicated. B) Heatmap of expression for indicated marker genes across each cluster. C). Spatial prediction of major cell types across E15.5 craniofacial section from Pina et al 2023<sup>23</sup>. D). Heatmap of sharing of marker genes between each major cell type in human and mouse. (P-values calculated by GeneOverlap R package). E). Network plot of human specific mesenchymal markers related to selected ontology categories. Shading of individual gene nodes based on fold change in expression of cells in the mesenchymal main cell type versus all other cell types. F). UMAP projection of subclustered CNCC cells from mouse. G). Gene ontology enrichments for indicated categories across each CNCC subcluster. H). Violin gene expression plots across CNCC subcluster for neural crest gene orthologous to human genes plotted in Fig. 3C. I). Average and percent expression for the top 5 marker genes for each of the mouse CNCC subclusters. J). UMAP projection of subclustered mesenchymal cells from mouse. K). Average and percent expression for the top 5 marker genes for each of the mouse mesenchymal subclusters. L). Spatial feature plot for module scores of top 100 marker genes of the PalatalShelf2 subcluster on a section of a mouse E11.5 embryo<sup>43</sup>. M). UMAP projection of subclustered ectodermal cells from mouse. N). Average and percent expression for the top 5 marker genes for each of the mouse ectodermal subclusters. O). Spatial feature plot for module scores of top 100 marker genes of the palate surface subcluster on a section of a mouse E11.5 embryo<sup>43</sup>. P). Spatial predictions of selected craniofacial subtypes on E15.5 craniofacial section from Pina et al 2023<sup>23</sup>.

1128

1129 Figure 5. Identification of mesenchymal subtypes in human craniofacial development.

1130 A). UMAP projection of subclustered mesenchymal main cell type. Subtype labels based on transfer of  
1131 mouse mesenchymal subtypes to human. B). Contribution of cells from each CS timepoint to each  
1132 mesenchymal subcluster. C) Average and percent expression for the top 5 marker genes for each of the  
1133 mesenchymal subclusters. D). Disease ontology enrichments for each of the indicated mesenchymal  
1134 subcluster. E.) Gene expression spatial feature plot for indicated mesenchymal marker genes in two  
1135 sections from a CS13 human embryo. F). Gene expression spatial feature plot for indicated  
1136 mesenchymal subtype marker genes in same sections as E. G). Spatial feature plot for modules scores  
1137 calculated from top 100 marker genes from indicated mesenchymal subtype.

1138 Figure 6. Identification of ectodermal subtypes in human craniofacial development.

1139 A). UMAP projection of subclustered ectodermal main cell type. Subtype labels based on transfer of  
1140 mouse ectodermal subtypes to human. B). Contribution of cells from each CS timepoint to each  
1141 ectodermal subcluster. C) Average and percent expression for the top 5 marker genes for each of the  
1142 ectodermal subclusters. D). Disease ontology enrichments for each of the indicated ectodermal  
1143 subcluster. E.) Gene expression spatial feature plot for indicated ectodermal marker genes in two  
1144 sections from a CS13 human embryo. F). Gene expression spatial feature plot for indicated ectodermal  
1145 subtype marker genes in same sections as E. G). Spatial feature plot for modules scores calculated from  
1146 top 100 marker genes from indicated ectodermal subtype.

1147 Figure 7. Enrichment of common variation associate with facial shape differences and congenital abnormality risk.

1148 A). Clustered heatmap of significance values calculated by MAGMA Celltyping<sup>91</sup> for each facial variation trait and cell subtype.  
1149 Profile landmark diagram adapted from Bonfante et al 2021<sup>88</sup>. Frontal landmark diagram adapted from Xiong et al 2019<sup>87</sup>.  
1150 Colors along top of heatmap indicate main cell type classification. Shaded gray indicators along left of heatmap indicate  
1151 study origin. Colors along left of heatmap indicate general region of the face each landmark is located. Hyphenated trait  
1152 measures are obtained from Xiong et al 2019 and combinatorial code is indicated in coded legend (e.g., EnL-ALL indicates  
1153 landmark segment 5 to 7). Descriptive named traits obtained from Bonfante et al 2021<sup>88</sup>. Levels of significance indicated by  
1154 asterisks or period according to figure. B). Clustered heatmap of significance values calculated by MAGMA Celltyping<sup>91</sup> for  
1155 each congenital abnormality or disease and cell subtype. Colors along top of heatmap indicate main cell type classification.  
1156 Levels of significance indicated by asterisks or period according to figure. C) Barplot showing the number of enriched human  
1157 phenotypes (max-normalized from 0-1 within each branch) for main cell types and subtypes as calculated by  
1158 MSTExplorer::run\_phenomix. Significance of the proportion tests, testing for disproportionate numbers of phenotype  
1159 enrichments for a given cell type within a given HPO branch, is denoted with asterisks (FDR<0.001=\*\*\*, FDR<0.01=\*\*,  
1160 FDR<0.05=\*) as well as black outlines around the bars.

1161 Figure 8. Genes associated with orofacial clefting, constraint in human populations, and Neanderthal introgression show distinct cell  
1162 subtype enrichments.

1163 A). Bar plot of standard deviations from the mean of bootstrapping tests performed by EWCE method<sup>90</sup> for each indicated  
1164 gene list and cell subtype. Asterisks indicate significant subtype enrichments corrected for number of gene lists and cell  
1165 subtypes performed for entire figure. B). Bubbleplot of  $-\log_{10}$  transformed significance and fold enrichment values for each  
1166 cell subtype from denovolzeR<sup>129</sup> analysis of protein *damaging de novo* variation in orofacial cleft trios from the Gabriella Miller  
1167 Kids First program<sup>121,122</sup>. Colored circles indicate variants identified in whole cohort (Any), in cleft lip with cleft palate probands  
1168 (CL/P), or cleft palate only probands (CP). Cell subtypes are clustered by main cell type. C) Average and percent expression  
1169 across all ectodermal subtypes of genes identified in nasal placode subtype *with de novo* protein damaging mutations for B.  
1170 D). Spatial feature plot of modules scores calculated from top 100 marker genes from nasal placode ectodermal subtype on

CS13 human embryo. E). Heatmap of fold differences in expression of each indicated gene in nasal placode subtype versus all other ectodermal cells. Presence or absence of box indicates membership in indicated disease ontology category indicated as significantly enriched in nasal placode marker genes.

## References

1. Minoux, M. & Rijli, F.M. Molecular mechanisms of cranial neural crest cell migration and patterning in craniofacial development. *Development* **137**, 2605-2621 (2010).
2. Cordero, D.R. *et al.* Cranial neural crest cells on the move: Their roles in craniofacial development. *American Journal of Medical Genetics Part A* **155**, 270-279 (2011).
3. Tang, W. & Bronner, M.E. Neural crest lineage analysis: from past to future trajectory. *Development* **147**(2020).
4. Guo, J. *et al.* Variation and signatures of selection on the human face. *J Hum Evol* **75**, 143-52 (2014).
5. Mitteroecker, P., Gunz, P., Bernhard, M., Schaefer, K. & Bookstein, F.L. Comparison of cranial ontogenetic trajectories among great apes and humans. *J Hum Evol* **46**, 679-97 (2004).
6. Naqvi, S. *et al.* Decoding the Human Face: Progress and Challenges in Understanding the Genetics of Craniofacial Morphology. *Annu Rev Genomics Hum Genet* **23**, 383-412 (2022).
7. Smith, D.W. Recognizable patterns of human malformation: genetic, embryologic, and clinical aspects. *Major Probl Clin Pediatr* **7**, 1-368 (1970).
8. Leslie, E.J. & Marazita, M.L. Genetics of cleft lip and cleft palate. *American journal of medical genetics. Part C, Seminars in medical genetics* **163C**, 246-258 (2013).
9. Mc Goldrick, N. *et al.* A multi-program analysis of cleft lip with cleft palate prevalence and mortality using data from 22 International Clearinghouse for Birth Defects Surveillance and Research programs, 1974–2014. *Birth Defects Research* **115**, 980-997 (2023).
10. Mossey, P.A., Little, J., Munger, R.G., Dixon, M.J. & Shaw, W.C. Cleft lip and palate. *Lancet* **374**, 1773-85 (2009).
11. Mai, C.T. *et al.* National population - based estimates for major birth defects, 2010-2014. *Birth defects research* **111**, 1420-1435 (2019).
12. Li, M. *et al.* Integrative functional genomic analysis of human brain development and neuropsychiatric risks. *Science (New York, NY)* **362**, eaat7615 (2018).
13. Yankee, T.N. *et al.* Integrative analysis of transcriptome dynamics during human craniofacial development identifies candidate disease genes. *Nat Commun* **14**, 4623 (2023).
14. Rajderkar, S.S. *et al.* Dynamic enhancer landscapes in human craniofacial development. *Nat Commun* **15**, 2030 (2024).
15. Sun, J. *et al.* Single-cell RNA-Seq reveals transcriptional regulatory networks directing the development of mouse maxillary prominence. *J Genet Genomics* **50**, 676-687 (2023).
16. Han, X. *et al.* Runx2-Twist1 interaction coordinates cranial neural crest guidance of soft palate myogenesis. *Elife* **10**(2021).
17. Mestas, J. & Hughes, C.C. Of mice and not men: differences between mouse and human immunology. *J Immunol* **172**, 2731-8 (2004).
18. Martik, M.L. & Bronner, M.E. Riding the crest to get a head: neural crest evolution in vertebrates. *Nat Rev Neurosci* **22**, 616-626 (2021).
19. York, J.R., Yuan, T. & McCauley, D.W. Evolutionary and Developmental Associations of Neural Crest and Placodes in the Vertebrate Head: Insights From Jawless Vertebrates. *Frontiers in Physiology* **11**(2020).
20. Donoghue, P.C.J., Graham, A. & Kelsh, R.N. The origin and evolution of the neural crest. *BioEssays* **30**, 530-541 (2008).

- 1218 21. Selleri, L. & Rijli, F.M. Shaping faces: genetic and epigenetic control of craniofacial  
1219 morphogenesis. *Nature Reviews Genetics* **24**, 610-626 (2023).
- 1220 22. Martínez-Abadías, N. *et al.* The Developmental Basis of Quantitative Craniofacial Variation in  
1221 Humans and Mice. *Evolutionary Biology* **39**, 554-567 (2012).
- 1222 23. Pina, J.O. *et al.* Multimodal spatiotemporal transcriptomic resolution of embryonic palate  
1223 osteogenesis. *Nat Commun* **14**, 5687 (2023).
- 1224 24. Cao, J. *et al.* A human cell atlas of fetal gene expression. *Science* **370**(2020).
- 1225 25. Xu, Y. *et al.* A single-cell transcriptome atlas profiles early organogenesis in human embryos.  
1226 *Nature Cell Biology* **25**, 604-615 (2023).
- 1227 26. Wang, C. *et al.* Single-cell RNA sequencing analysis of human embryos from the late Carnegie to  
1228 fetal development. *Cell & Bioscience* **14**, 118 (2024).
- 1229 27. Zeng, B. *et al.* The single-cell and spatial transcriptional landscape of human gastrulation and  
1230 early brain development. *Cell Stem Cell* **30**, 851-866.e7 (2023).
- 1231 28. Wilderman, A., VanOudenhove, J., Kron, J., Noonan, J.P. & Cotney, J. High-Resolution Epigenomic  
1232 Atlas of Human Embryonic Craniofacial Development. *Cell Rep* **23**, 1581-1597 (2018).
- 1233 29. Vieille-Grosjean, I., Hunt, P., Gulisano, M., Boncinelli, E. & Thorogood, P. Branchial HOX Gene  
1234 Expression and Human Craniofacial Development. *Developmental Biology* **183**, 49-60 (1997).
- 1235 30. Cai, J. *et al.* Gene expression in pharyngeal arch 1 during human embryonic development.  
1236 *Hum.Mol.Genet.* **14**, 903-912 (2005).
- 1237 31. Samuels, B.D. *et al.* FaceBase 3: analytical tools and FAIR resources for craniofacial and dental  
1238 research. *Development (Cambridge)* **147**(2020).
- 1239 32. Bishop, M.R. *et al.* Genome-wide Enrichment of De Novo Coding Mutations in Orofacial Cleft  
1240 Trios. *American journal of human genetics* **107**, 124-136 (2020).
- 1241 33. Schoenwolf, G.C., Bleyl, S.B., Brauer, P.R. & Francis-West, P. *Larsen's Human Embryology*,  
1242 (Elsevier, Philadelphia, PA, 2021).
- 1243 34. Jirasek, J.E. *An Atlas of Human Prenatal Developmental Mechanics: Anatomy and Staging*, 312  
1244 (CRC Press, London, 2004).
- 1245 35. Brunskill, E.W. *et al.* A gene expression atlas of early craniofacial development. *Developmental*  
1246 *biology* **391**, 133-146 (2014).
- 1247 36. Hooper, J.E., Jones, K.L., Smith, F.J., Williams, T. & Li, H. An Alternative Splicing Program for  
1248 Mouse Craniofacial Development. *Front Physiol* **11**, 1099 (2020).
- 1249 37. Feng, W. *et al.* Spatial and temporal analysis of gene expression during growth and fusion of the  
1250 mouse facial prominences. *PLoS One* **4**, e8066 (2009).
- 1251 38. Hooper, J.E. *et al.* Systems biology of facial development: contributions of ectoderm and  
1252 mesenchyme. *Developmental biology* **426**, 97-114 (2017).
- 1253 39. Program, C.S.-C.B. *et al.* CZ CELL×GENE Discover: A single-cell data platform for scalable  
1254 exploration, analysis and modeling of aggregated data. *bioRxiv*, 2023.10.30.563174 (2023).
- 1255 40. Leland, M., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for  
1256 Dimension Reduction. *arXiv* **1802.03426**(2020).
- 1257 41. Martin, J.F., Bradley, A. & Olson, E.N. The paired-like homeo box gene M<sub>Hox</sub> is required for early  
1258 events of skeletogenesis in multiple lineages. *Genes Dev* **9**, 1237-49 (1995).
- 1259 42. Bartoletti, G., Dong, C., Umar, M. & He, F. Pdgfra regulates multipotent cell differentiation  
1260 towards chondrocytes via inhibiting Wnt9a/beta-catenin pathway during chondrocranial cartilage  
1261 development. *Developmental Biology* **466**, 36-46 (2020).
- 1262 43. Chen, A. *et al.* Spatiotemporal transcriptomic atlas of mouse organogenesis using DNA nanoball-  
1263 patterned arrays. *Cell* **185**, 1777-1792.e21 (2022).
- 1264 44. Nikolopoulou, E. *et al.* Spinal neural tube closure depends on regulation of surface ectoderm  
1265 identity and biomechanics by Grhl2. *Nature Communications* **10**, 2487 (2019).



- 1266 45. Bebee, T.W. *et al.* The splicing regulators Esrp1 and Esrp2 direct an epithelial splicing program  
1267 essential for mammalian development. *Elife* **4**(2015).
- 1268 46. Revil, T. & Jerome-Majewska, L.A. During Embryogenesis, Esrp1 Expression Is Restricted to a  
1269 Subset of Epithelial Cells and Is Associated With Splicing of a Number of Developmentally  
1270 Important Genes. *Developmental Dynamics* **242**, 281-290 (2013).
- 1271 47. Knowles, W.J., Bologna, M.L., Chasis, J.A., Marchesi, S.L. & Marchesi, V.T. Common structural  
1272 polymorphisms in human erythrocyte spectrin. *J Clin Invest* **73**, 973-9 (1984).
- 1273 48. Bennett, V. & Stenbuck, P.J. The membrane attachment protein for spectrin is associated with  
1274 band 3 in human erythrocyte membranes. *Nature* **280**, 468-73 (1979).
- 1275 49. Gardner, L.C., Smith, S.J. & Cox, T.M. Biosynthesis of delta-aminolevulinic acid and the regulation  
1276 of heme formation by immature erythroid cells in man. *J Biol Chem* **266**, 22010-8 (1991).
- 1277 50. Huang, C.H. The human Rh50 glycoprotein gene. Structural organization and associated splicing  
1278 defect resulting in Rh(null) disease. *J Biol Chem* **273**, 2207-13 (1998).
- 1279 51. Vallese, F. *et al.* Architecture of the human erythrocyte ankyrin-1 complex. *Nat Struct Mol Biol* **29**,  
1280 706-718 (2022).
- 1281 52. Millauer, B. *et al.* High affinity VEGF binding and developmental expression suggest Flk-1 as a  
1282 major regulator of vasculogenesis and angiogenesis. *Cell* **72**, 835-46 (1993).
- 1283 53. Kendall, R.L., Wang, G. & Thomas, K.A. Identification of a natural soluble form of the vascular  
1284 endothelial growth factor receptor, FLT-1, and its heterodimerization with KDR. *Biochem Biophys  
1285 Res Commun* **226**, 324-8 (1996).
- 1286 54. Justement, L.B., Campbell, K.S., Chien, N.C. & Cambier, J.C. Regulation of B cell antigen receptor  
1287 signal transduction and phosphorylation by CD45. *Science* **252**, 1839-42 (1991).
- 1288 55. Engel, P. *et al.* The B7-2 (B70) costimulatory molecule expressed by monocytes and activated B  
1289 lymphocytes is the CD86 differentiation antigen. *Blood* **84**, 1402-7 (1994).
- 1290 56. Wang, M.H. *et al.* Identification of the ron gene product as the receptor for the human  
1291 macrophage stimulating protein. *Science* **266**, 117-9 (1994).
- 1292 57. Weintraub, H. *et al.* The myoD gene family: nodal point during specification of the muscle cell  
1293 lineage. *Science* **251**, 761-6 (1991).
- 1294 58. Seidel, U. & Arnold, H.H. Identification of the functional promoter regions in the human gene  
1295 encoding the myosin alkali light chains MLC1 and MLC3 of fast skeletal muscle. *J Biol Chem* **264**,  
1296 16109-17 (1989).
- 1297 59. Karsch-Mizrachi, I., Travis, M., Blau, H. & Leinwand, L.A. Expression and DNA sequence analysis  
1298 of a human embryonic skeletal muscle myosin heavy chain gene. *Nucleic Acids Res* **17**, 6167-79  
1299 (1989).
- 1300 60. Takahashi, M. & Osumi, N. Identification of a novel type II classical cadherin: Rat cadherin19 is  
1301 expressed in the cranial ganglia and Schwann cell precursors during development.  
1302 *Developmental Dynamics* **232**, 200-208 (2005).
- 1303 61. Izaki, T., Kamakura, S., Kohjima, M. & Sumimoto, H. Two forms of human Inscuteable-related  
1304 protein that links Par3 to the Pins homologues LGN and AGS3. *Biochemical and Biophysical  
1305 Research Communications* **341**, 1001-1006 (2006).
- 1306 62. Puente, X.S., Pendás, A.M., Llano, E., Velasco, G. & López-Otín, C. Molecular cloning of a novel  
1307 membrane-type matrix metalloproteinase from a human breast carcinoma. *Cancer Res* **56**, 944-9  
1308 (1996).
- 1309 63. Sasai, N., Mizuseki, K. & Sasai, Y. Requirement of FoxD3-class signaling for neural crest  
1310 determination in *Xenopus*. *Development* **128**, 2525-2536 (2001).
- 1311 64. Kos, R., Reedy, M.V., Johnson, R.L. & Erickson, C.A. The winged-helix transcription factor FoxD3 is  
1312 important for establishing the neural crest lineage and repressing melanogenesis in avian  
1313 embryos. *Development* **128**, 1467-1479 (2001).



- 1314 65. Lukoseviciute, M. *et al.* From Pioneer to Repressor: Bimodal foxd3 Activity Dynamically Remodels  
1315 Neural Crest Regulatory Landscape *In Vivo*. *Developmental Cell* **47**, 608-  
1316 628.e6 (2018).
- 1317 66. Simões-Costa, M.S., McKeown, S.J., Tan-Cabugao, J., Sauka-Spengler, T. & Bronner, M.E.  
1318 Dynamic and Differential Regulation of Stem Cell Factor FoxD3 in the Neural Crest Is Encrypted in  
1319 the Genome. *PLOS Genetics* **8**, e1003142 (2012).
- 1320 67. Simões-Costa, M. & Bronner, M.E. Establishing neural crest identity: a gene regulatory recipe.  
1321 *Development (Cambridge)* **142**, 242-257 (2015).
- 1322 68. Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by single-cell  
1323 RNA-seq. *Science* **352**, 189-96 (2016).
- 1324 69. Feng, Z. *et al.* hReg-CNCC reconstructs a regulatory network in human cranial neural crest cells  
1325 and annotates variants in a developmental context. *Communications Biology* **4**, 442 (2021).
- 1326 70. Thomas, S. *et al.* Human neural crest cells display molecular and phenotypic hallmarks of stem  
1327 cells. *Human Molecular Genetics* **17**, 3411-3425 (2008).
- 1328 71. Betters, E., Liu, Y., Kjaeldgaard, A., Sundström, E. & García-Castro, M.I. Analysis of early human  
1329 neural crest development. *Dev Biol* **344**, 578-92 (2010).
- 1330 72. O’Rahilly, R. & Müller, F. The development of the neural crest in the human. *Journal of Anatomy*  
1331 **211**, 335-351 (2007).
- 1332 73. Rada-Iglesias, A. *et al.* Epigenomic annotation of enhancers predicts transcriptional regulators of  
1333 human neural crest; PMC3751405. *Cell Stem Cell* **11**, 633-648 (2012).
- 1334 74. Sun, K.Y. *et al.* A deep catalogue of protein-coding variation in 983,578 individuals. *Nature* (2024).
- 1335 75. Aoto, K., Nishimura, T., Eto, K. & Motoyama, J. Mouse GLI3 Regulates Fgf8 Expression and  
1336 Apoptosis in the Developing Neural Tube, Face, and Limb Bud. *Developmental Biology* **251**, 320-  
1337 332 (2002).
- 1338 76. Baker, J.L., Wood, B., Karpinski, B.A., LaMantia, A.S. & Maynard, T.M. Testicular receptor 2, Nr2c1,  
1339 is associated with stem cells in the developing olfactory epithelium and other cranial sensory and  
1340 skeletal structures. *Gene Expr Patterns* **20**, 71-9 (2016).
- 1341 77. Mansouri, A., Hallonet, M. & Gruss, P. Pax genes and their roles in cell differentiation and  
1342 development. *Curr Opin Cell Biol* **8**, 851-7 (1996).
- 1343 78. Nakashima, K. *et al.* The Novel Zinc Finger-Containing Transcription Factor Osterix Is Required for  
1344 Osteoblast Differentiation and Bone Formation. *Cell* **108**, 17-29 (2002).
- 1345 79. Hojo, H., Ohba, S., He, X., Lai, L.P. & McMahon, A.P. Sp7/Osterix Is Restricted to Bone-Forming  
1346 Vertebrates where It Acts as a Dlx Co-factor in Osteoblast Specification. *Dev Cell* **37**, 238-53  
1347 (2016).
- 1348 80. Sur, A. *et al.* Single-cell analysis of shared signatures and transcriptional diversity during zebrafish  
1349 development. *Dev Cell* **58**, 3028-3047.e12 (2023).
- 1350 81. Van Otterloo, E. *et al.* AP-2 $\alpha$  and AP-2 $\beta$  cooperatively function in the craniofacial surface  
1351 ectoderm to regulate chromatin and gene expression dynamics during facial development. *Elife*  
1352 **11**(2022).
- 1353 82. De Felice, M. *et al.* A mouse model for hereditary thyroid dysgenesis and cleft palate. *Nat Genet*  
1354 **19**, 395-8 (1998).
- 1355 83. Sheng, H.Z. *et al.* Specification of pituitary cell lineages by the LIM homeobox gene Lhx3. *Science*  
1356 **272**, 1004-7 (1996).
- 1357 84. Raetzman, L.T., Ward, R. & Camper, S.A. Lhx4 and Prop1 are required for cell survival and  
1358 expansion of the pituitary primordia. *Development* **129**, 4229-39 (2002).
- 1359 85. Pina, J.O. *et al.* Spatial Multiomics Reveal the Role of Wnt Modulator, Dkk2, in Palatogenesis.  
1360 *bioRxiv* (2024).
- 1361 86. Li, H., Jones, K.L., Hooper, J.E. & Williams, T. The molecular anatomy of mammalian upper lip and  
1362 primary palate fusion at single cell resolution. *Development (Cambridge)* **146**(2019).

- 1363 87. Xiong, Z. *et al.* Novel genetic loci affecting facial shape variation in humans. *eLife* **8**, e49898  
1364 (2019).
- 1365 88. Bonfante, B. *et al.* A GWAS in Latin Americans identifies novel face shape loci, implicating VPS13B  
1366 and a Denisovan introgressed region in facial variation. *Science advances* **7**, eabc6160 (2021).
- 1367 89. de Leeuw, C.A., Mooij, J.M., Heskes, T. & Posthuma, D. MAGMA: generalized gene-set analysis of  
1368 GWAS data. *PLoS Comput Biol* **11**, e1004219 (2015).
- 1369 90. Skene, N.G. & Grant, S.G. Identification of Vulnerable Cell Types in Major Brain Disorders Using  
1370 Single Cell Transcriptomes and Expression Weighted Cell Type Enrichment. *Front Neurosci* **10**, 16  
1371 (2016).
- 1372 91. Skene, N.G. *et al.* Genetic identification of brain cell types underlying schizophrenia. *Nature*  
1373 *Genetics* **50**, 825-833 (2018).
- 1374 92. Ludwig, K.U. *et al.* Imputation of Orofacial Clefting Data Identifies Novel Risk Loci and Sheds Light  
1375 on the Genetic Background of Cleft Lip ± Cleft Palate and Cleft Palate Only. *Human Molecular*  
1376 *Genetics* **26**, 829-842 (2017).
- 1377 93. Leslie, E.J. *et al.* Association studies of low-frequency coding variants in nonsyndromic cleft lip  
1378 with or without cleft palate. *American journal of medical genetics Part A* **100**, 493-8 (2017).
- 1379 94. Butali, A. *et al.* Genomic analyses in African populations identify novel risk loci for cleft palate.  
1380 *Hum Mol Genet* **28**, 1038-1051 (2019).
- 1381 95. Huang, L. *et al.* Genetic factors define CPO and CLO subtypes of nonsyndromic orofacial cleft.  
1382 *PLoS Genet* **15**, e1008357 (2019).
- 1383 96. Leslie, E.J. *et al.* A multi-ethnic genome-wide association study identifies novel loci for non-  
1384 syndromic cleft lip with or without cleft palate on 2p24.2, 17q23 and 19q13. *Human molecular*  
1385 *genetics* **25**, 2862-2872 (2016).
- 1386 97. He, M. *et al.* Genome-wide Analyses Identify a Novel Risk Locus for Nonsyndromic Cleft Palate. *J*  
1387 *Dent Res* **99**, 1461-1468 (2020).
- 1388 98. Curtis, S.W. *et al.* The PAX1 locus at 20p11 is a potential genetic modifier for bilateral cleft lip.  
1389 *HGG Adv* **2**(2021).
- 1390 99. Yu, Y. *et al.* Genome-wide analyses of non-syndromic cleft lip with palate identify 14 novel loci  
1391 and genetic heterogeneity. *Nat Commun* **8**, 14364 (2017).
- 1392 100. Wu, T. *et al.* Evidence of gene-environment interaction for two genes on chromosome 4 and  
1393 environmental tobacco smoke in controlling the risk of nonsyndromic cleft palate. *PLoS One* **9**,  
1394 e88088 (2014).
- 1395 101. Sun, Y. *et al.* Genome-wide association study identifies a new susceptibility locus for cleft lip with  
1396 or without a cleft palate. *Nat Commun* **6**, 6414 (2015).
- 1397 102. Howe, L.J. *et al.* Investigating the shared genetics of non-syndromic cleft lip/palate and facial  
1398 morphology. *PLoS Genet* **14**, e1007501 (2018).
- 1399 103. Ray, D. *et al.* Pleiotropy method reveals genetic overlap between orofacial clefts at multiple novel  
1400 loci from GWAS of multi-ethnic trios. *PLoS Genet* **17**, e1009584 (2021).
- 1401 104. Curtis, S.W. *et al.* FAT4 identified as a potential modifier of orofacial cleft laterality. *Genet*  
1402 *Epidemiol* **45**, 721-735 (2021).
- 1403 105. Haaland Ø, A. *et al.* A genome-wide scan of cleft lip triads identifies parent-of-origin interaction  
1404 effects between ANK3 and maternal smoking, and between ARHGEF10 and alcohol consumption.  
1405 *F1000Res* **8**, 960 (2019).
- 1406 106. Carlson, J.C. *et al.* Genome-wide interaction studies identify sex-specific risk alleles for  
1407 nonsyndromic orofacial clefts. *Genet Epidemiol* **42**, 664-672 (2018).
- 1408 107. Grosen, D. *et al.* Risk of oral clefts in twins. *Epidemiology (Cambridge, Mass.)* **22**, 313-319 (2011).
- 1409 108. Christensen, K. & Andersen, P.F. Isolated Cleft Palate in Danish Multiple Births, 1970-1990. *The*  
1410 *Cleft Palate Craniofacial Journal* **30**, 469-474 (1993).

- 1411 109. Christensen, K. The 20th century Danish facial cleft population--epidemiological and genetic-  
1412 epidemiological studies. *Cleft Palate Craniofac J* **36**, 96-104 (1999).
- 1413 110. Diaz Perez, K.K. *et al.* Rare variants found in clinical gene panels illuminate the genetic and allelic  
1414 architecture of orofacial clefting. *Genetics in Medicine* **25**(2023).
- 1415 111. Auslander, A. *et al.* The International Family Study of Nonsyndromic Orofacial Clefts: Design and  
1416 Methods. *The Cleft Palate Craniofacial Journal* **59**, S37-S47 (2021).
- 1417 112. Marazita, M.L. *et al.* Meta-Analysis of 13 Genome Scans Reveals Multiple Cleft Lip/Palate Genes  
1418 with Novel Loci on 9q21 and 2q32-35. *The American Journal of Human Genetics* **75**, 161-173  
1419 (2004).
- 1420 113. Ludwig, K.U. *et al.* Genome-wide meta-analyses of nonsyndromic cleft lip with or without cleft  
1421 palate identify six new risk loci. *Nature genetics* **44**, 968-971 (2012).
- 1422 114. Ludwig, K.U. *et al.* Meta-analysis Reveals Genome-Wide Significance at 15q13 for Nonsyndromic  
1423 Clefting of Both the Lip and the Palate, and Functional Analyses Implicate GREM1 As a Plausible  
1424 Causative Gene. *PLoS Genetics* **12**, e1005914 (2016).
- 1425 115. Yu, Y. *et al.* Genome-wide meta-analyses identify five new risk loci for nonsyndromic orofacial  
1426 clefts in the Chinese Han population. *Molecular Genetics & Genomic Medicine* **11**, e2226 (2023).
- 1427 116. Kurki, M.I. *et al.* FinnGen provides genetic insights from a well-phenotyped isolated population.  
1428 *Nature* **613**, 508-518 (2023).
- 1429 117. VanOudenhove, J., Yankee, T., Wilderman, A. & Cotney, J. Epigenomic and Transcriptomic  
1430 Dynamics During Human Heart Organogenesis. *Circulation research* **127**, e184-e209 (2020).
- 1431 118. Murphy, K.B. *et al.* Identification of cell type-specific gene targets underlying thousands of rare  
1432 diseases and subtraits. *medRxiv*, 2023.02.13.23285820 (2023).
- 1433 119. Gargano, M.A. *et al.* The Human Phenotype Ontology in 2024: phenotypes around the world.  
1434 *Nucleic Acids Research* **52**, D1333-D1346 (2023).
- 1435 120. Xu, H. *et al.* CleftGeneDB: a resource for annotating genes associated with cleft lip and cleft  
1436 palate. *Science bulletin* **66**, 2340-2342 (2021).
- 1437 121. Mukhopadhyay, N. *et al.* Whole genome sequencing of orofacial cleft trios from the Gabriella  
1438 Miller Kids First Pediatric Research Consortium identifies a new locus on chromosome 21.  
1439 *Human Genetics* **139**, 215-226 (2020).
- 1440 122. Curtis, S.W. *et al.* Rare genetic variants in SEC24D modify orofacial cleft phenotypes. *medRxiv*  
1441 (2023).
- 1442 123. Robinson, K. *et al.* Trio-based GWAS identifies novel associations and subtype-specific risk  
1443 factors for cleft palate. *medRxiv*, 2023.03.01.23286642 (2023).
- 1444 124. Karczewski, K.J. *et al.* The mutational constraint spectrum quantified from variation in 141,456  
1445 humans. *Nature* **581**, 434-443 (2020).
- 1446 125. Doan, R.N. *et al.* Mutations in Human Accelerated Regions Disrupt Cognition and Social Behavior.  
1447 *Cell* **167**, 1-27 (2016).
- 1448 126. Levchenko, A., Kanapin, A., Samsonova, A. & Gainetdinov, R.R. Human Accelerated Regions and  
1449 Other Human-Specific Sequence Variations in the Context of Evolution and Their Relevance for  
1450 Brain Development. *Genome Biol Evol* **10**, 166-188 (2018).
- 1451 127. Driessens, S.L.W. *et al.* Genes associated with cognitive ability and HAR show overlapping  
1452 expression patterns in human cortical neuron types. *Nature Communications* **14**, 4188 (2023).
- 1453 128. Prabhakar, S. *et al.* Human-Specific Gain of Function in a Developmental Enhancer. *Science (New*  
1454 *York, NY)* **321**, 1346-1350 (2008).
- 1455 129. Ware, J.S., Samocha, K.E., Homsy, J. & Daly, M.J. Interpreting de novo Variation in Human Disease  
1456 Using denovolyzeR. *Curr Protoc Hum Genet* **87**, 7.25.1-7.25.15 (2015).
- 1457 130. Samocha, K.E. *et al.* A framework for the interpretation of de novo mutation in human disease.  
1458 *Nat Genet* **46**, 944-50 (2014).

- 1459 131. Luecken, M.D. *et al.* Benchmarking atlas-level data integration in single-cell genomics. *Nature*  
1460 *Methods* **19**, 41-50 (2022).
- 1461 132. Luecken, M.D. & Theis, F.J. Current best practices in single - cell RNA - seq analysis: a tutorial.  
1462 *Molecular Systems Biology* **15**, e8746 (2019).
- 1463 133. Heumos, L. *et al.* Best practices for single-cell analysis across modalities. *Nature Reviews*  
1464 *Genetics* **24**, 550-572 (2023).
- 1465 134. Amezquita, R.A. *et al.* Orchestrating single-cell analysis with Bioconductor. *Nature Methods* **17**,  
1466 137-145 (2020).
- 1467 135. Abdelaal, T. *et al.* A comparison of automatic cell identification methods for single-cell RNA  
1468 sequencing data. *Genome Biology* **20**, 194 (2019).
- 1469 136. Domínguez Conde, C. *et al.* Cross-tissue immune cell analysis reveals tissue-specific features in  
1470 humans. *Science* **376**, eabl5197 (2022).
- 1471 137. Fischer, F. *et al.* scTab: Scaling cross-tissue single-cell annotation models. *Nature*  
1472 *Communications* **15**, 6611 (2024).
- 1473 138. Prescott, Sara L. *et al.* Enhancer Divergence and cis-Regulatory Evolution in the Human and  
1474 Chimpanzee Neural Crest. *Cell* **163**, 68-83 (2015).
- 1475 139. Twigg, S.R. *et al.* Frontorhiny, a distinctive presentation of frontonasal dysplasia caused by  
1476 recessive mutations in the ALX3 homeobox gene. *Am J Hum Genet* **84**, 698-705 (2009).
- 1477 140. Beverdam, A., Brouwer, A., Reijnen, M., Korving, J. & Meijlink, F. Severe nasal clefting and  
1478 abnormal embryonic apoptosis in Alx3/Alx4 double mutant mice. *Development* **128**, 3975-86  
1479 (2001).
- 1480 141. Wilkie, A.O. *et al.* Functional haploinsufficiency of the human homeobox gene MSX2 causes  
1481 defects in skull ossification. *Nat Genet* **24**, 387-90 (2000).
- 1482 142. Ye, Q., Bhojwani, A. & Hu, J.K. Understanding the development of oral epithelial organs through  
1483 single cell transcriptomic analysis. *Development* **149**(2022).
- 1484 143. Qiu, X. *et al.* Spatiotemporal modeling of molecular holograms. *Cell*.
- 1485 144. Weinberg, S.M., Cornell, R. & Leslie, E.J. Craniofacial genetics: Where have we been and where  
1486 are we going? *PLoS Genet* **14**, e1007438 (2018).
- 1487 145. Claes, P. *et al.* Genome-wide mapping of global-to-local genetic effects on human facial shape.  
1488 *Nat Genet* **50**, 414-423 (2018).
- 1489 146. White, J.D. *et al.* Insights into the genetic architecture of the human face. *Nat Genet* **53**, 45-53  
1490 (2021).
- 1491 147. Dixon, M.J., Marazita, M.L., Beaty, T.H. & Murray, J.C. Cleft lip and palate: understanding genetic  
1492 and environmental influences. *Nature reviews Genetics* **12**, 167-178 (2011).
- 1493 148. Marazita, M.L. The evolution of human genetic studies of cleft lip and cleft palate. *Annual Review*  
1494 *of Genomics and Human Genetics* **13**, 263-283 (2012).
- 1495 149. Rahimov, F. *et al.* High incidence and geographic distribution of cleft palate in Finland are  
1496 associated with the IRF6 gene. *Nat Commun* **15**, 9568 (2024).
- 1497 150. Kirby, M.L. & Waldo, K.L. Role of neural crest in congenital heart disease. *Circulation* **82**, 332-340  
1498 (1990).
- 1499 151. Farrell, M., Waldo, K., Li, Y.-X. & Kirby, M.L. A Novel Role for Cardiac Neural Crest in Heart  
1500 Development. *Trends in Cardiovascular Medicine* **9**, 214-220 (1999).
- 1501 152. Bronner, M.E. Formation and migration of neural crest cells in the vertebrate embryo.  
1502 *Histochemistry and Cell Biology* **138**, 179-186 (2012).
- 1503 153. Etchevers, H.C., Dupin, E. & Le Douarin, N.M. The diverse neural crest: from embryology to human  
1504 pathology. *Development* **146**(2019).
- 1505 154. Yamagishi, H. Cardiac Neural Crest. *Cold Spring Harbor Perspectives in Biology* **13**(2021).
- 1506 155. Rahimov, F. *et al.* Disruption of an AP-2alpha binding site in an IRF6 enhancer is associated with  
1507 cleft lip. *Nature genetics* **40**, 1341-1347 (2008).



- 1508 156. Indencleef, K. *et al.* Six NSCL/P Loci Show Associations With Normal-Range Craniofacial  
1509 Variation. *Frontiers in Genetics* **9**(2018).
- 1510 157. Indencleef, K. *et al.* The Intersection of the Genetic Architectures of Orofacial Clefts and Normal  
1511 Facial Variation. *Frontiers in Genetics* **12**(2021).
- 1512 158. Weinberg, S.M. What's Shape Got to Do With It? Examining the Relationship Between Facial  
1513 Shape and Orofacial Clefting. *Frontiers in Genetics* **13**(2022).
- 1514 159. Tamarin, A. & Boyde, A. Facial and visceral arch development in the mouse embryo: a study by  
1515 scanning electron microscopy. *J Anat* **124**, 563-80 (1977).
- 1516 160. Depew, M.J. & Compagnucci, C. Tweaking the hinge and caps: testing a model of the organization  
1517 of jaws. *J Exp Zool B Mol Dev Evol* **310**, 315-35 (2008).
- 1518 161. Losa, M. *et al.* Face morphogenesis is promoted by Pbx-dependent EMT via regulation of Snail1  
1519 during frontonasal prominence fusion. *Development* **145**(2018).
- 1520 162. Jiang, R., Bush, J.O. & Lidral, A.C. Development of the upper lip: morphogenetic and molecular  
1521 mechanisms. *Dev Dyn* **235**, 1152-66 (2006).
- 1522 163. Wang, K.H. *et al.* Evaluation and integration of disparate classification systems for clefts of the lip.  
1523 *Front Physiol* **5**, 163 (2014).
- 1524 164. Lin, H. & Grosschedl, R. Failure of B-cell differentiation in mice lacking the transcription factor  
1525 EBF. *Nature* **376**, 263-7 (1995).
- 1526 165. Jin, S. *et al.* Ebf factors and MyoD cooperate to regulate muscle relaxation via Atp2a1. *Nature*  
1527 *communications* **5**, 3793 (2014).
- 1528 166. Garcia-Dominguez, M., Poquet, C., Garel, S. & Charnay, P. Ebf gene function is required for  
1529 coupling neuronal differentiation and cell cycle exit. *Development* **130**, 6013-25 (2003).
- 1530 167. Parsons, D.W. *et al.* An integrated genomic analysis of human glioblastoma multiforme. *Science*  
1531 **321**, 1807-12 (2008).
- 1532 168. Jones, D.L. & Wagers, A.J. No place like home: anatomy and function of the stem cell niche. *Nat*  
1533 *Rev Mol Cell Biol* **9**, 11-21 (2008).
- 1534 169. Zardo, G. *et al.* Integrated genomic and epigenomic analyses pinpoint biallelic gene inactivation in  
1535 tumors. *Nature Genetics* **32**, 453-458 (2002).
- 1536 170. Liao, D. Emerging roles of the EBF family of transcription factors in tumor suppression. *Mol*  
1537 *Cancer Res* **7**, 1893-901 (2009).
- 1538 171. Cho, Y.S. *et al.* A large-scale genome-wide association study of Asian populations uncovers  
1539 genetic factors influencing eight quantitative traits. *Nat Genet* **41**, 527-34 (2009).
- 1540 172. Kiel, D.P. *et al.* Genome-wide association with bone mass and geometry in the Framingham Heart  
1541 Study. *BMC Med Genet* **8 Suppl 1**, S14 (2007).
- 1542 173. Koller, D.L. *et al.* Genome-wide association study of bone mineral density in premenopausal  
1543 European-American women and replication in African-American women. *J Clin Endocrinol Metab*  
1544 **95**, 1802-9 (2010).
- 1545 174. Stykarsdottir, U. *et al.* Multiple genetic loci for bone mineral density and fractures. *N Engl J Med*  
1546 **358**, 2355-65 (2008).
- 1547 175. Kiper, P.O.S. *et al.* Cortical-Bone Fragility--Insights from sFRP4 Deficiency in Pyle's Disease. *N*  
1548 *Engl J Med* **374**, 2553-2562 (2016).
- 1549 176. Bilsborough, A. Cranial Morphology of Neanderthal Man. *Nature* **237**, 351-352 (1972).
- 1550 177. MORANT, G.M. STUDIES OF PALAEO-LITHIC MAN. *Annals of Eugenics* **2**, 318-381 (1927).
- 1551 178. Falk, D. Comparative anatomy of the larynx in man and the chimpanzee: Implications for language  
1552 in Neanderthal. *American Journal of Physical Anthropology* **43**, 123-132 (1975).
- 1553 179. Prüfer, K. *et al.* The complete genome sequence of a Neanderthal from the Altai Mountains.  
1554 *Nature* **505**, 43-9 (2014).
- 1555 180. Noonan, J.P. Neanderthal genomics and the evolution of modern humans. *Genome research* **20**,  
1556 547-553 (2010).



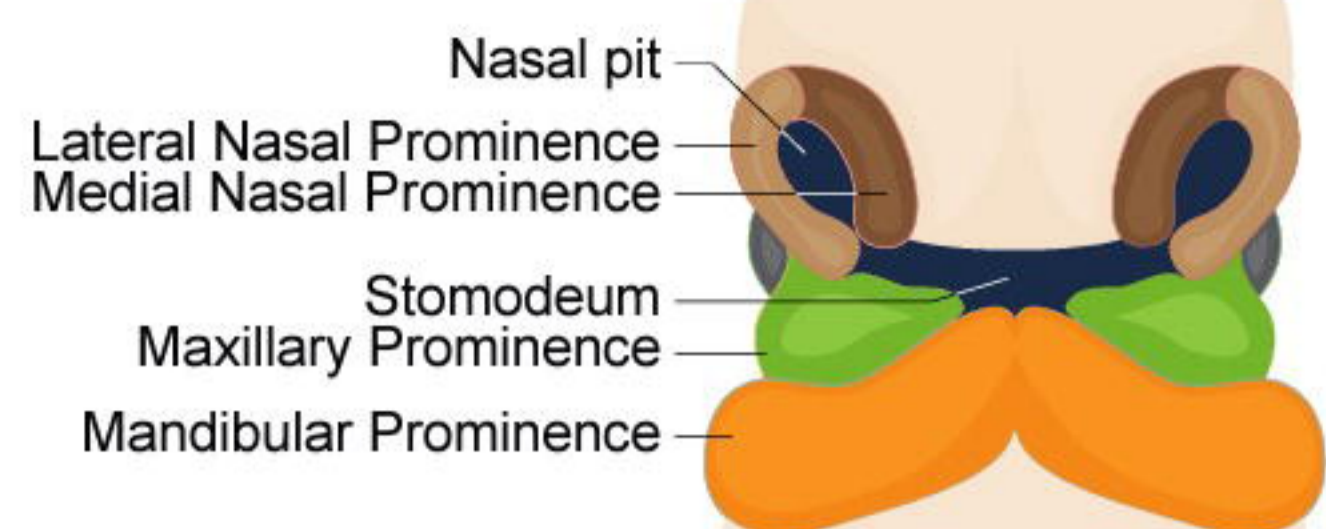
- 1557 181. Noonan, J.P. *et al.* Sequencing and analysis of Neanderthal genomic DNA. *Science* **314**, 1113-8  
1558 (2006).
- 1559 182. Sankararaman, S., Patterson, N., Li, H., Pääbo, S. & Reich, D. The date of interbreeding between  
1560 Neandertals and modern humans. *PLoS genetics* **8**, e1002947 (2012).
- 1561 183. Moorjani, P. *et al.* A genetic method for dating ancient genomes provides a direct estimate of  
1562 human generation interval in the last 45,000 years. *Proceedings of the National Academy of*  
1563 *Sciences* **113**, 5652-5657 (2016).
- 1564 184. Green, R.E. *et al.* A draft sequence of the Neandertal genome. *Science (New York, NY)* **328**, 710-  
1565 722 (2010).
- 1566 185. Vernot, B. *et al.* Excavating Neandertal and Denisovan DNA from the genomes of Melanesian  
1567 individuals. *Science* **352**, 235-239 (2016).
- 1568 186. van Galen, P. *et al.* Reduced Lymphoid Lineage Priming Promotes Human Hematopoietic Stem  
1569 Cell Expansion. *Cell Stem Cell* **14**, 94-106 (2014).
- 1570 187. Telis, N., Aguilar, R. & Harris, K. Selection against archaic hominin genetic variation in regulatory  
1571 regions. *Nature Ecology & Evolution* **4**, 1558-1566 (2020).
- 1572 188. McArthur, E., Rinker, D.C. & Capra, J.A. Quantifying the contribution of Neanderthal introgression  
1573 to the heritability of complex traits. *Nature Communications* **12**, 4481 (2021).
- 1574 189. Gregory, M.D. *et al.* Neanderthal-Derived Genetic Variation Shapes Modern Human Cranium and  
1575 Brain. *Scientific Reports* **7**, 6308 (2017).
- 1576 190. Stoessel, A. *et al.* Morphology and function of Neandertal and modern human ear ossicles. *Proc*  
1577 *Natl Acad Sci U S A* **113**, 11489-11494 (2016).
- 1578 191. Spoor, F., Hublin, J.J., Braun, M. & Zonneveld, F. The bony labyrinth of Neanderthals. *J Hum Evol*  
1579 **44**, 141-65 (2003).
- 1580 192. Gokhman, D. *et al.* Differential DNA methylation of vocal and facial anatomy genes in modern  
1581 humans. *Nat Commun* **11**, 1189 (2020).
- 1582 193. Megill, C. *et al.* cellxgene: a performant, scalable exploration platform for high dimensional  
1583 sparse matrices. *bioRxiv*, 2021.04.05.438318 (2021).
- 1584 194. Theiler, K. *The House Mouse : Atlas of Embryonic Development*, (Springer Science+Business  
1585 Media, 1989).
- 1586 195. Hao, Y. *et al.* Dictionary learning for integrative, multimodal and scalable single-cell analysis.  
1587 *Nature Biotechnology* **42**, 293-304 (2024).
- 1588 196. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony.  
1589 *Nature Methods* **16**, 1289-1296 (2019).
- 1590 197. Shen, L.S.I. GeneOverlap: Test and visualize gene overlaps. (2024).
- 1591 198. Ouyang, J.F., Kamaraj, U.S., Cao, E.Y. & Rackham, O.J.L. ShinyCell: simple and sharable  
1592 visualization of single-cell gene expression data. *Bioinformatics* **37**, 3374-3376 (2021).
- 1593 199. Chen, J., Bardes, E.E., Aronow, B.J. & Jegga, A.G. ToppGene Suite for gene list enrichment analysis  
1594 and candidate gene prioritization. *Nucleic Acids Research* **37**, W305-W311 (2009).
- 1595 200. Granger, B. & Berto, S. scToppR: a coding-friendly R interface to ToppGene. *Bioinformatics*  
1596 **40**(2024).
- 1597 201. Durinck, S. *et al.* BioMart and Bioconductor: a powerful link between biological databases and  
1598 microarray data analysis. *Bioinformatics* **21**, 3439-3440 (2005).
- 1599 202. Engler, J. tidyheatmaps: Heatmaps from Tidy Data. (2024).
- 1600 203. Keough, K.C. *et al.* Three-dimensional genome rewiring in loci with human accelerated regions.  
1601 *Science* **380**, eabm1696 (2023).
- 1602 204. Gu, Z. & Hübschmann, D. rGREAT: an R/bioconductor package for functional enrichment on  
1603 genomic regions. *Bioinformatics* **39**(2023).
- 1604 205. DiStefano, M.T. *et al.* The Gene Curation Coalition: A global effort to harmonize gene-disease  
1605 evidence resources. *Genet Med* **24**, 1732-1742 (2022).

- 1606 206. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate - a Practical and Powerful  
1607 Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B-Methodological* **57**,  
1608 289-300 (1995).

1609



A



5 weeks



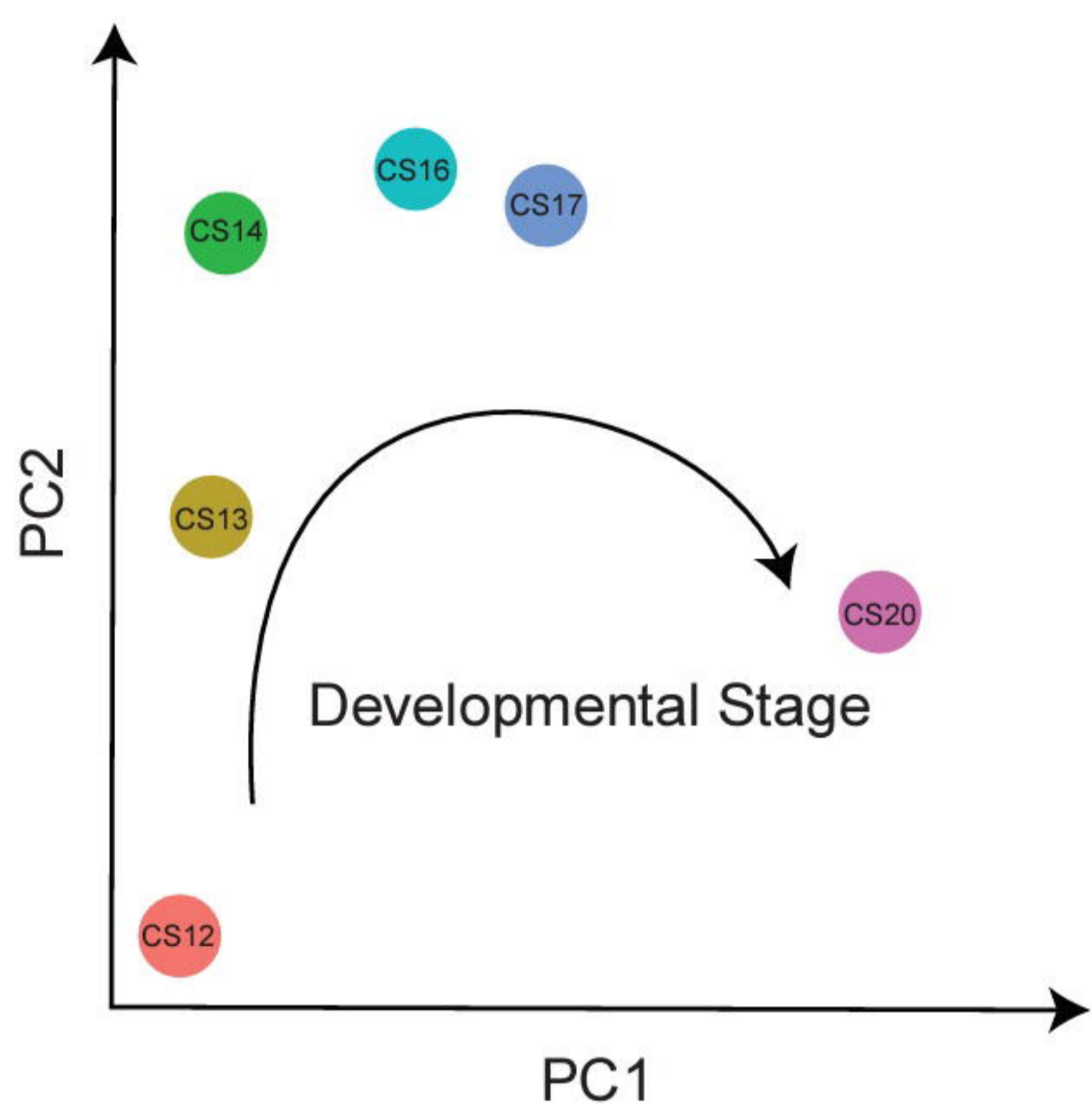
6 weeks



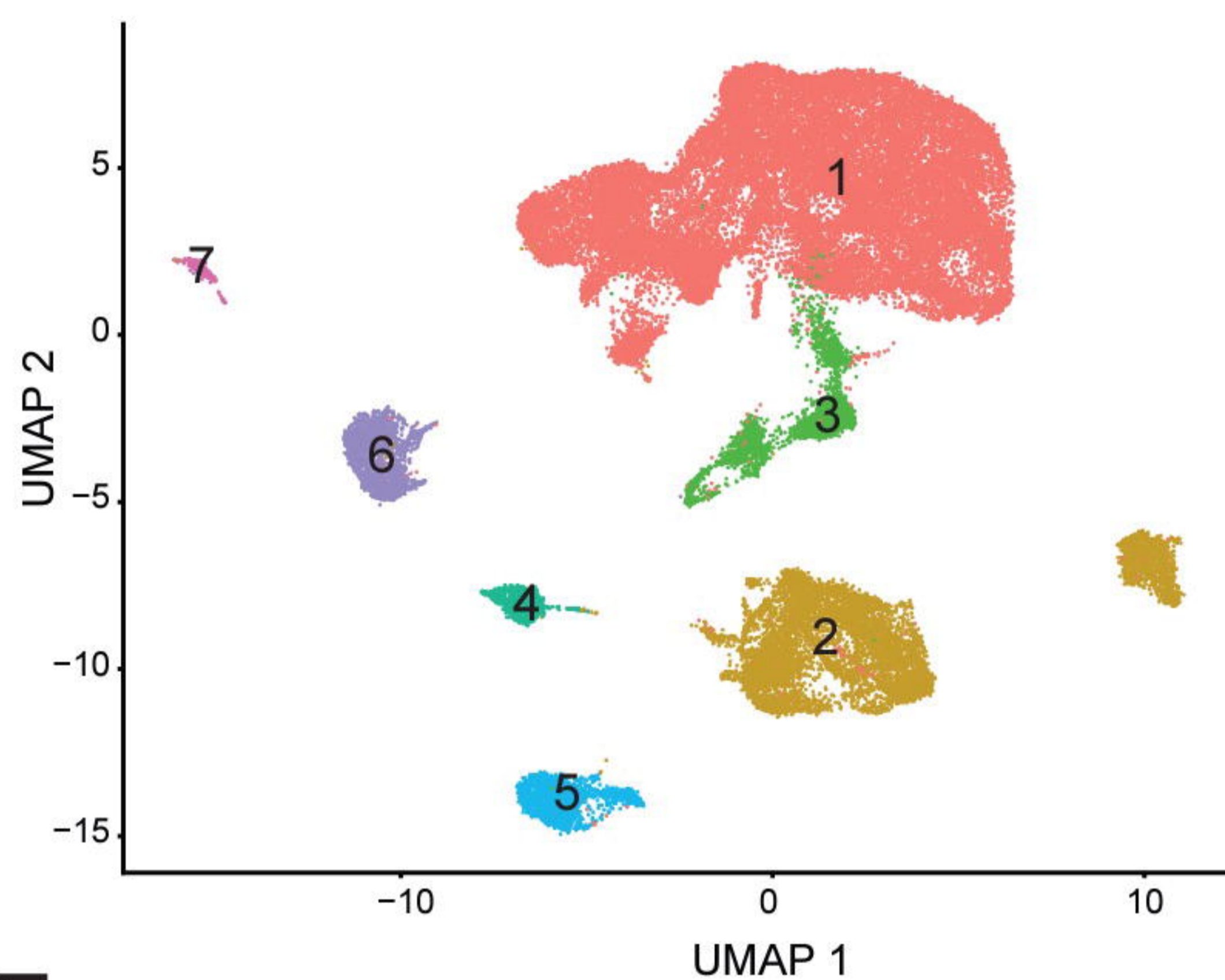
8 weeks

CS12  
n=4CS13  
n=6CS14  
n=3CS16  
n=3CS17  
n=3CS20  
n=5

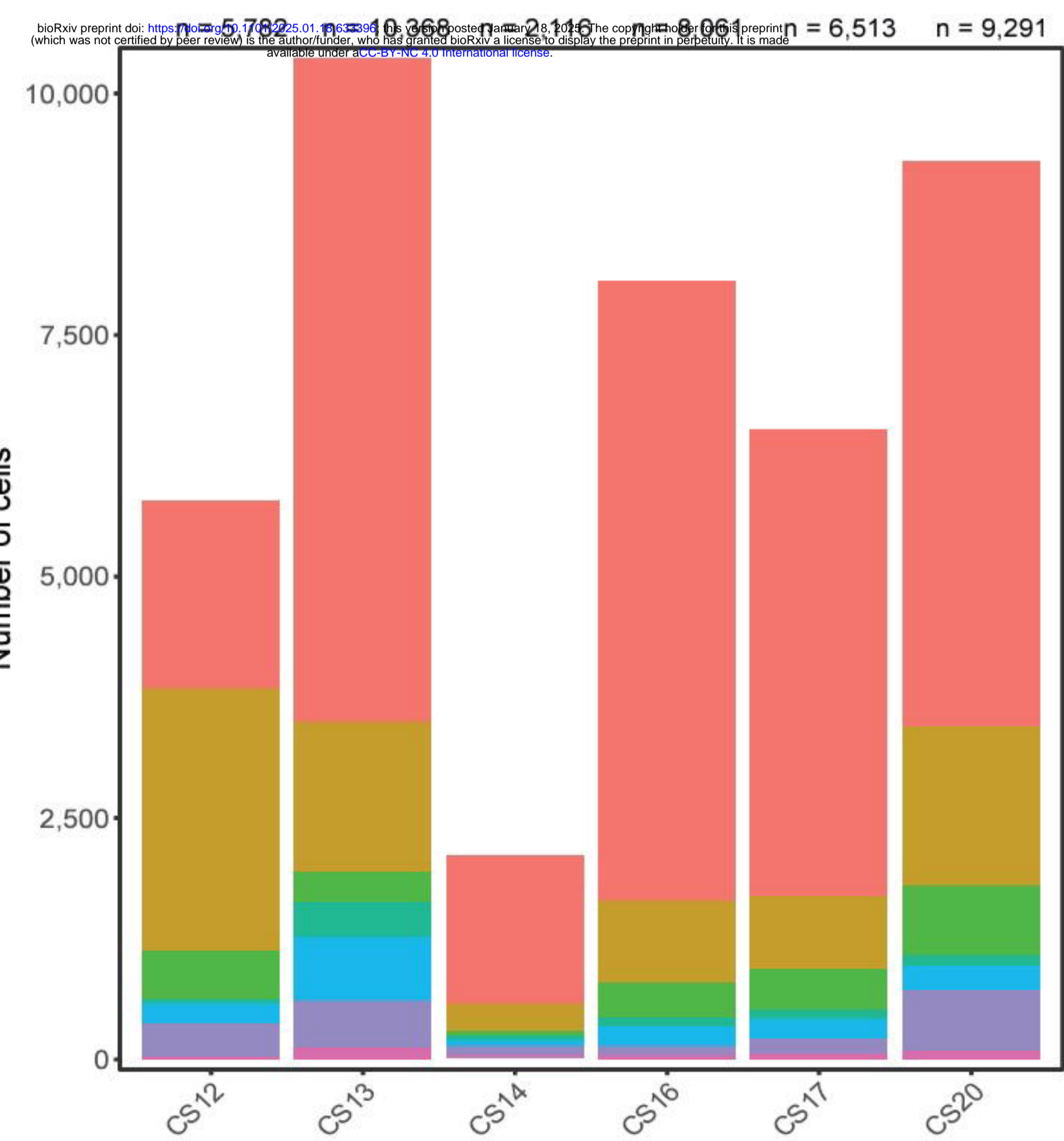
B



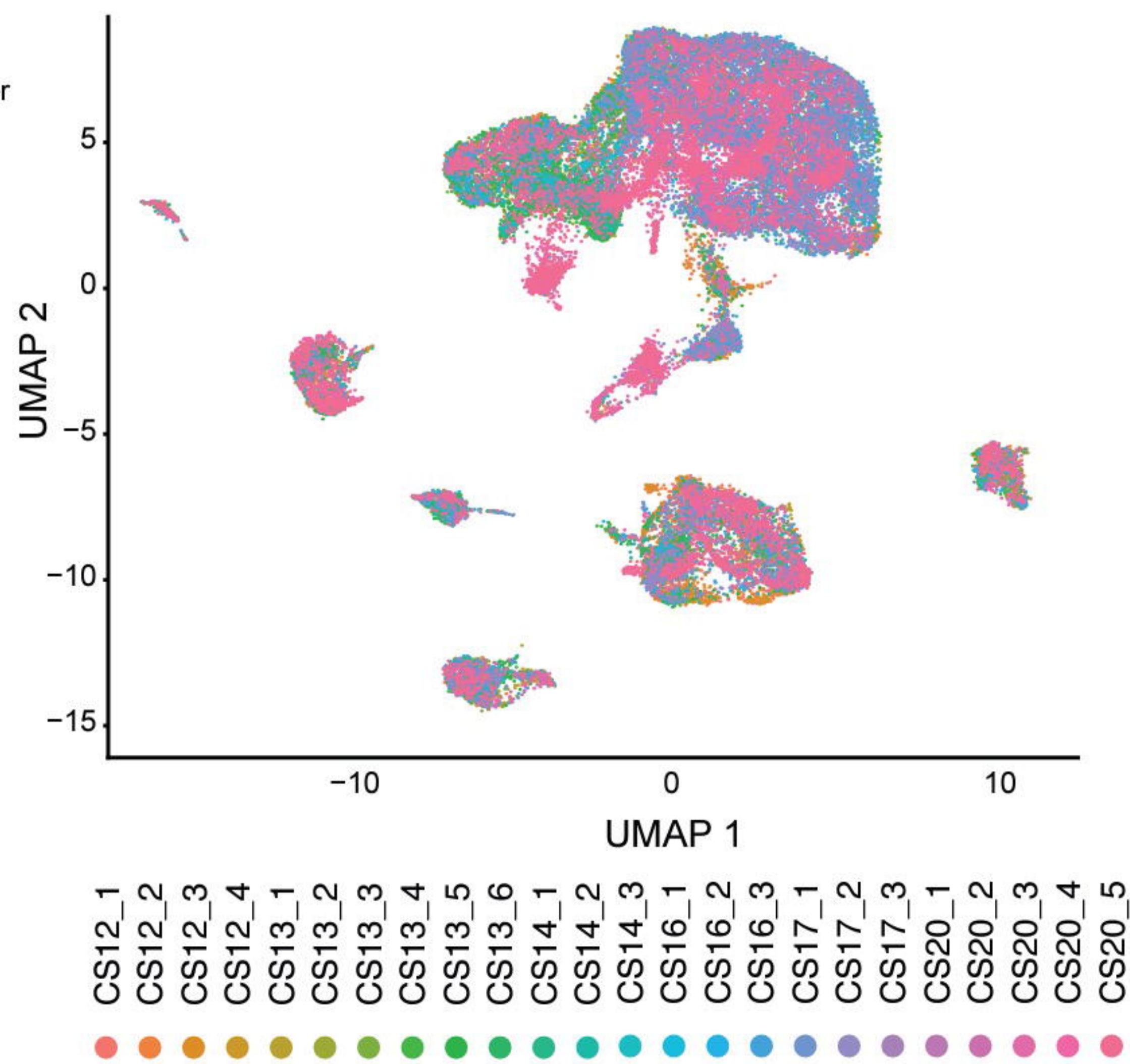
C



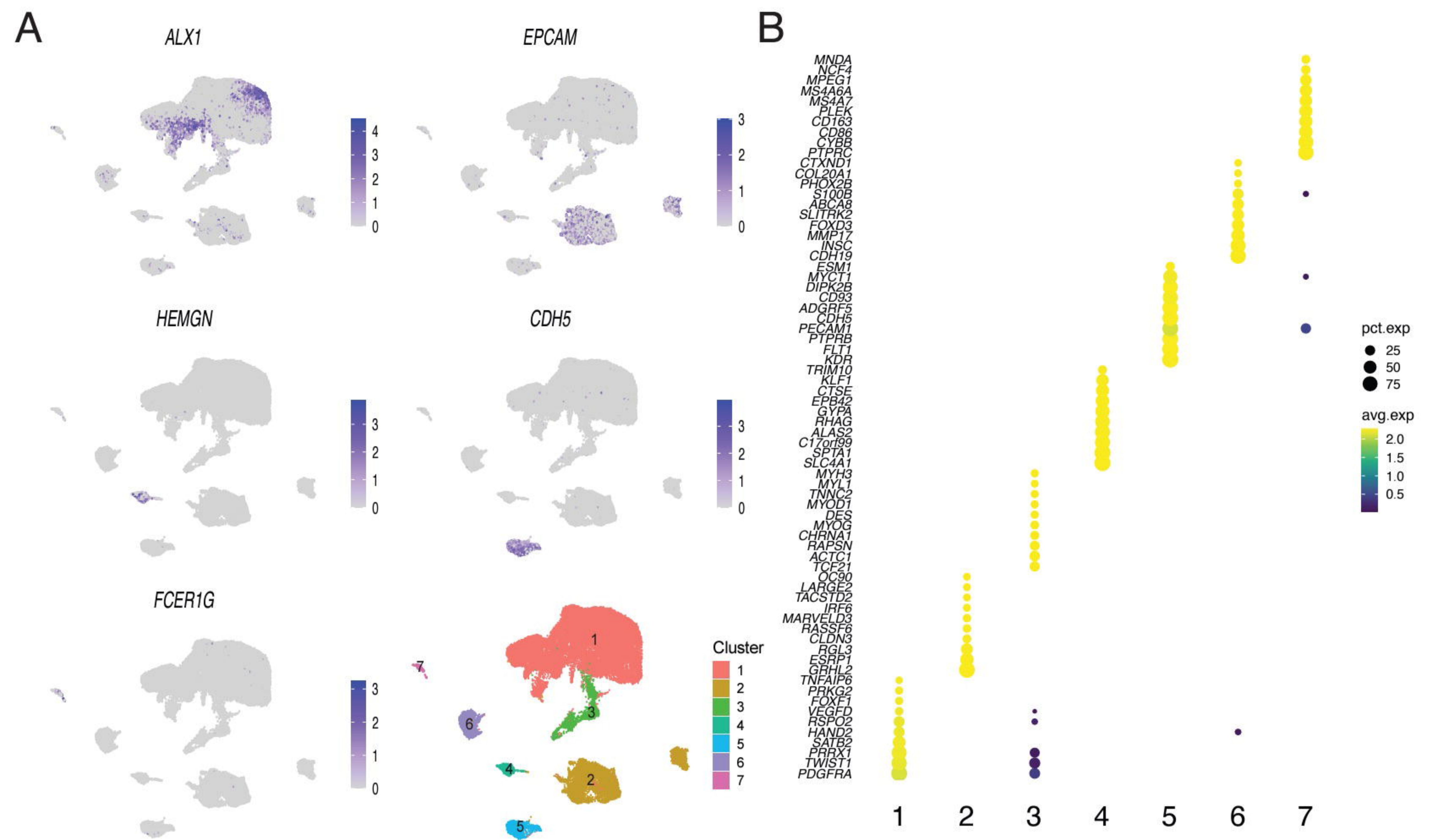
D



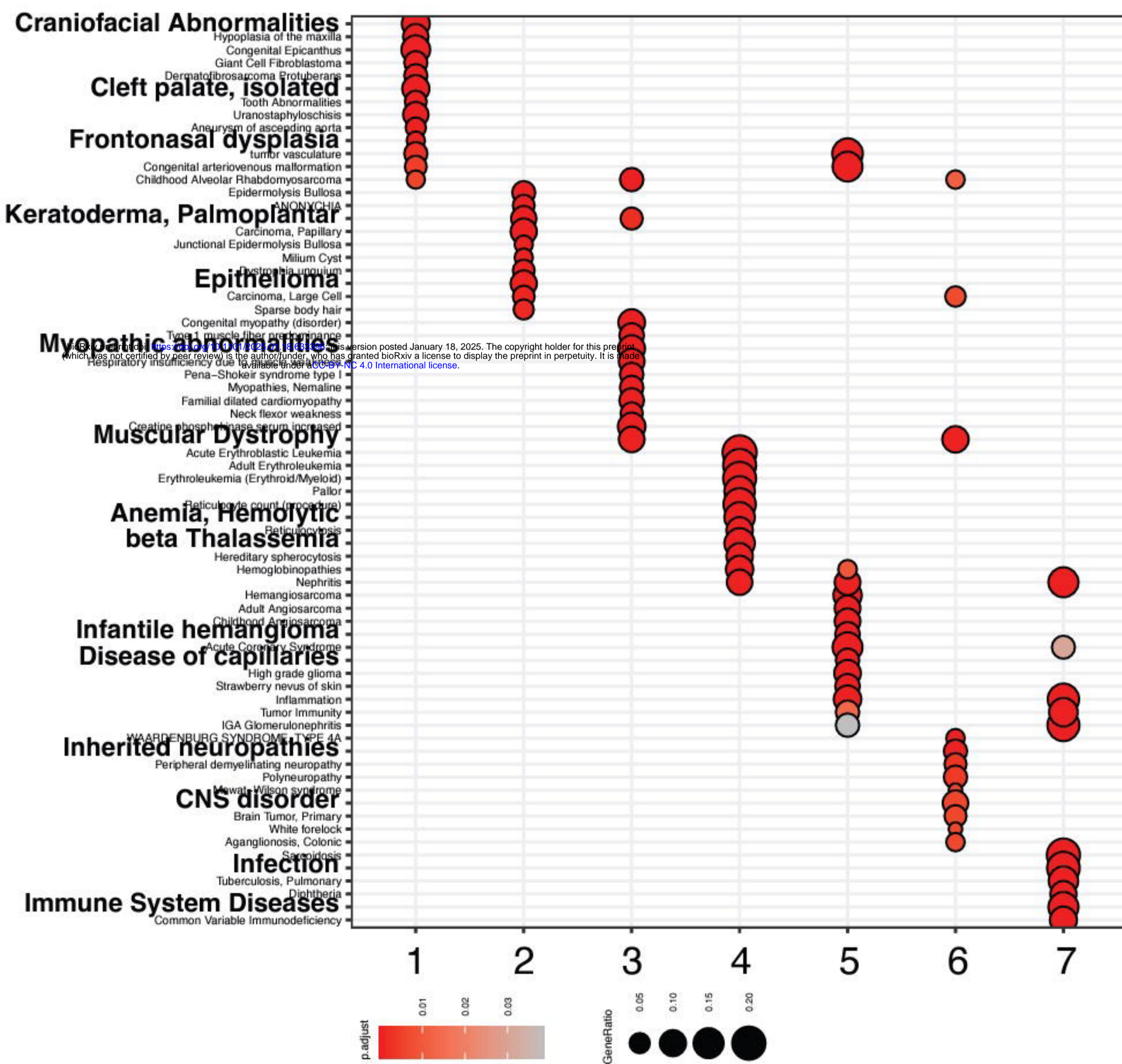
E



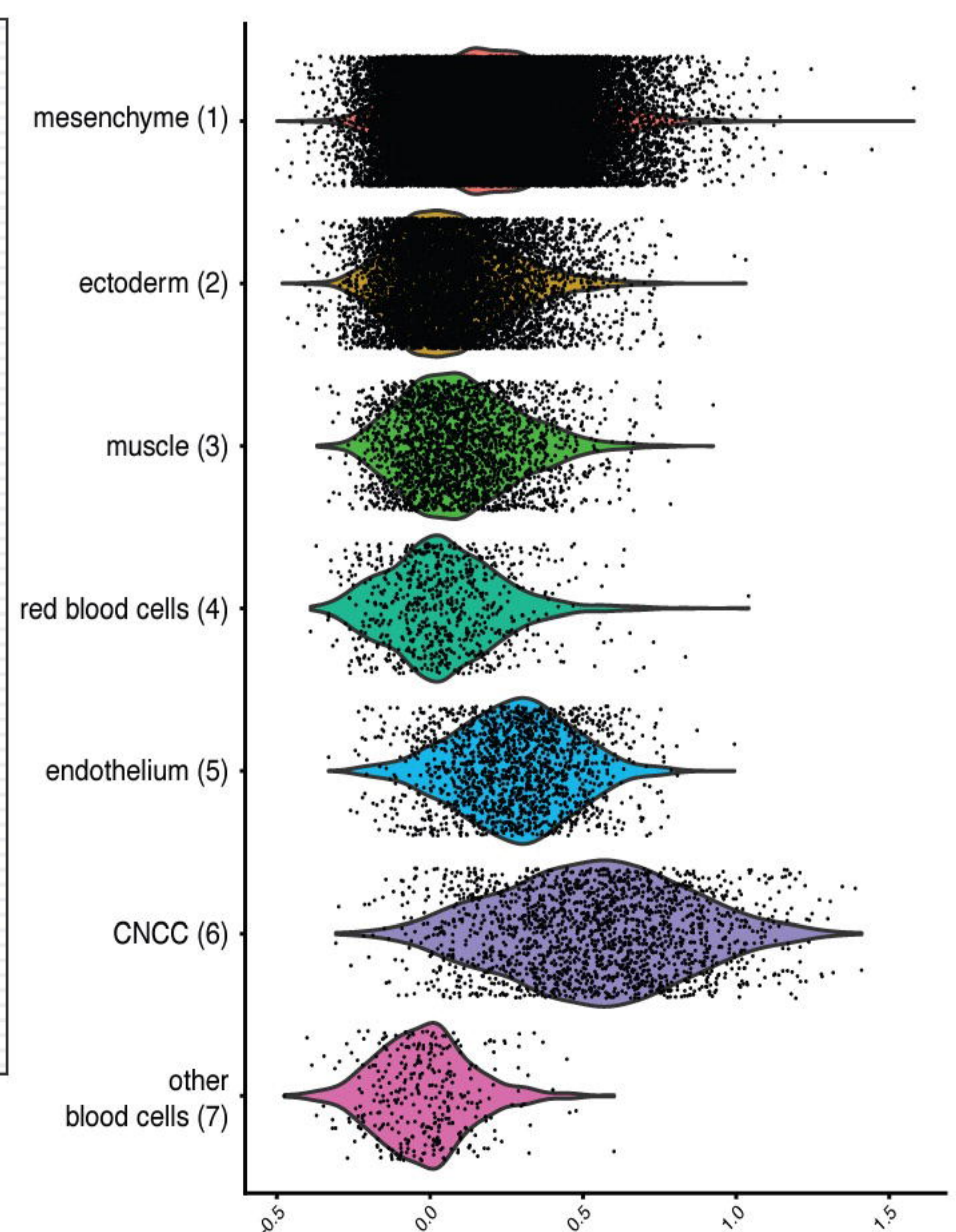




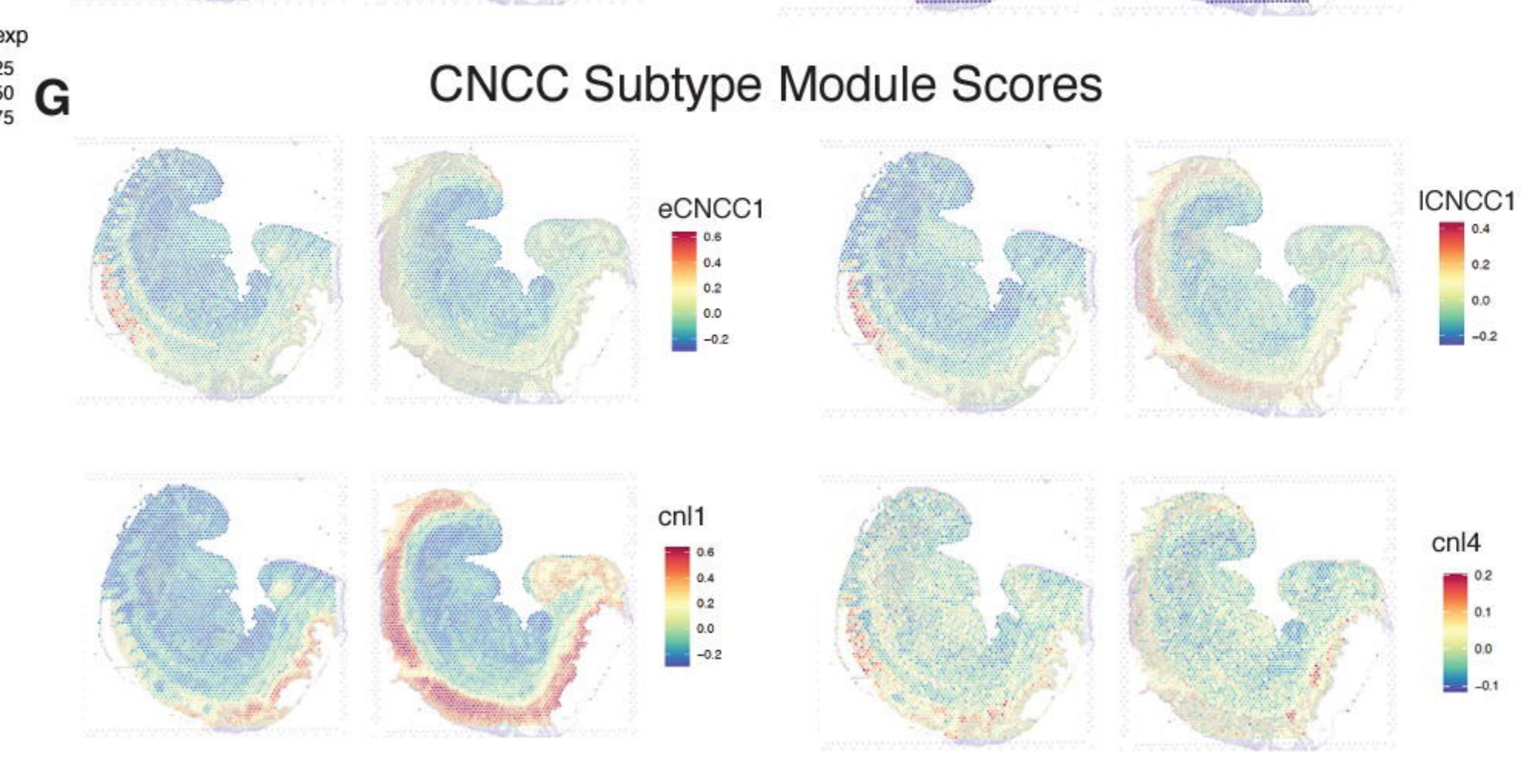
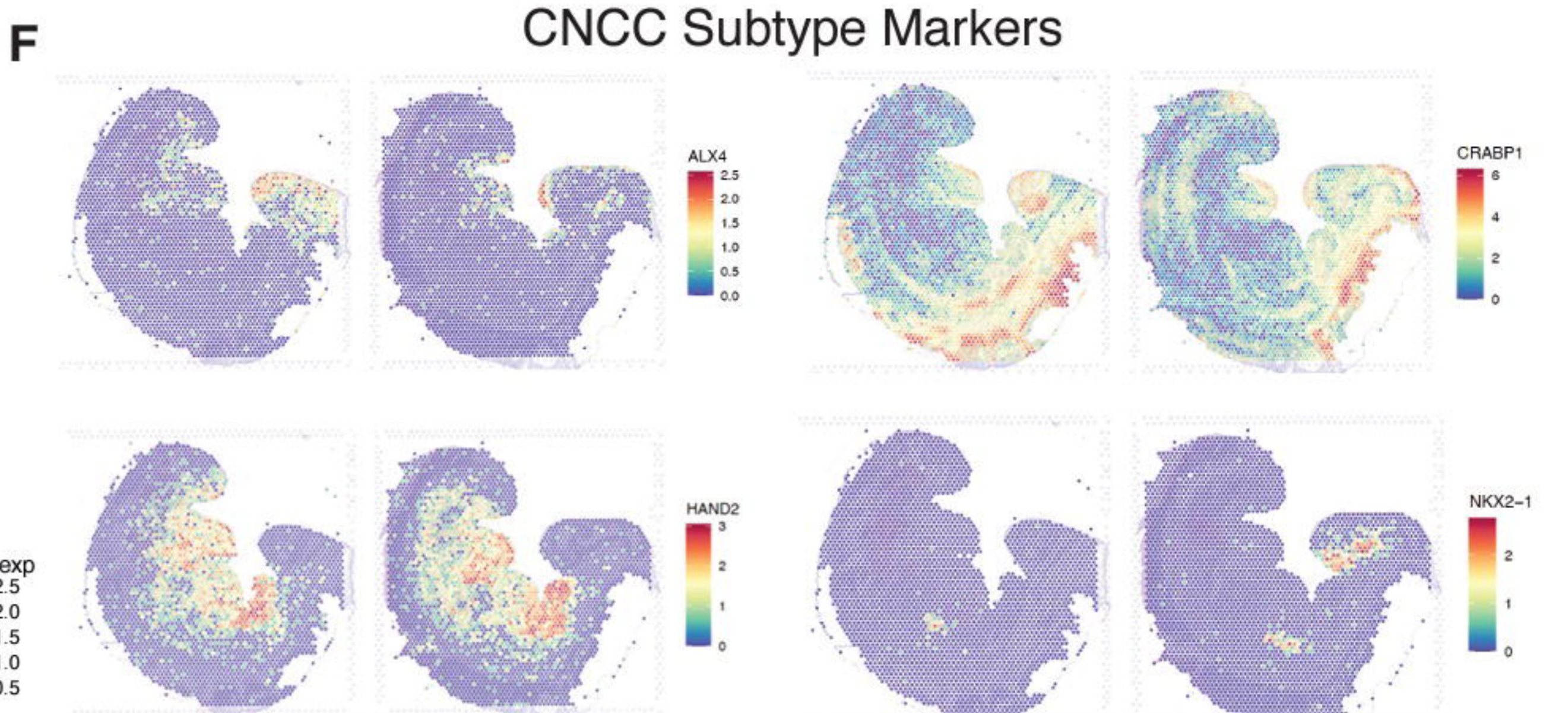
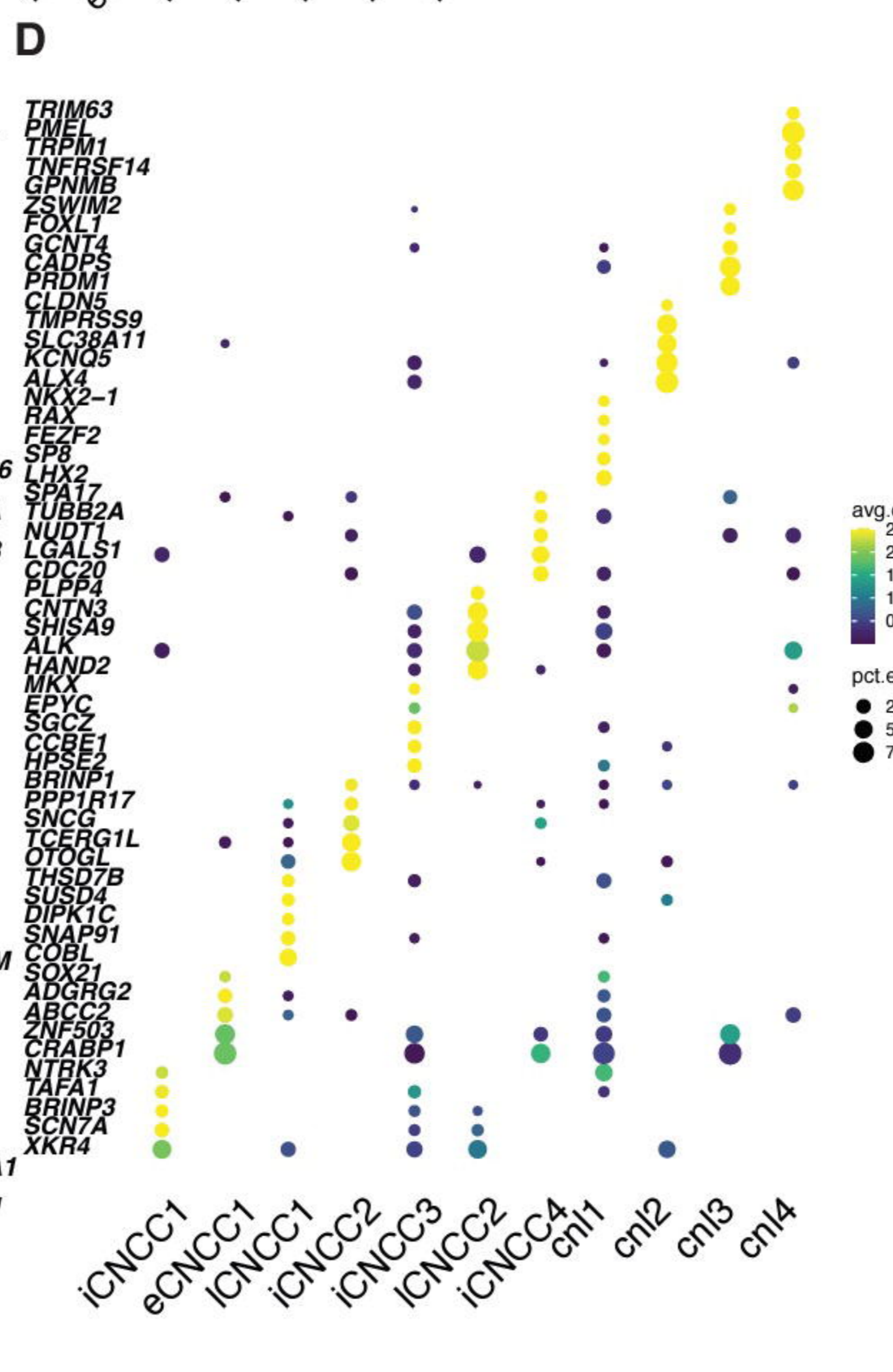
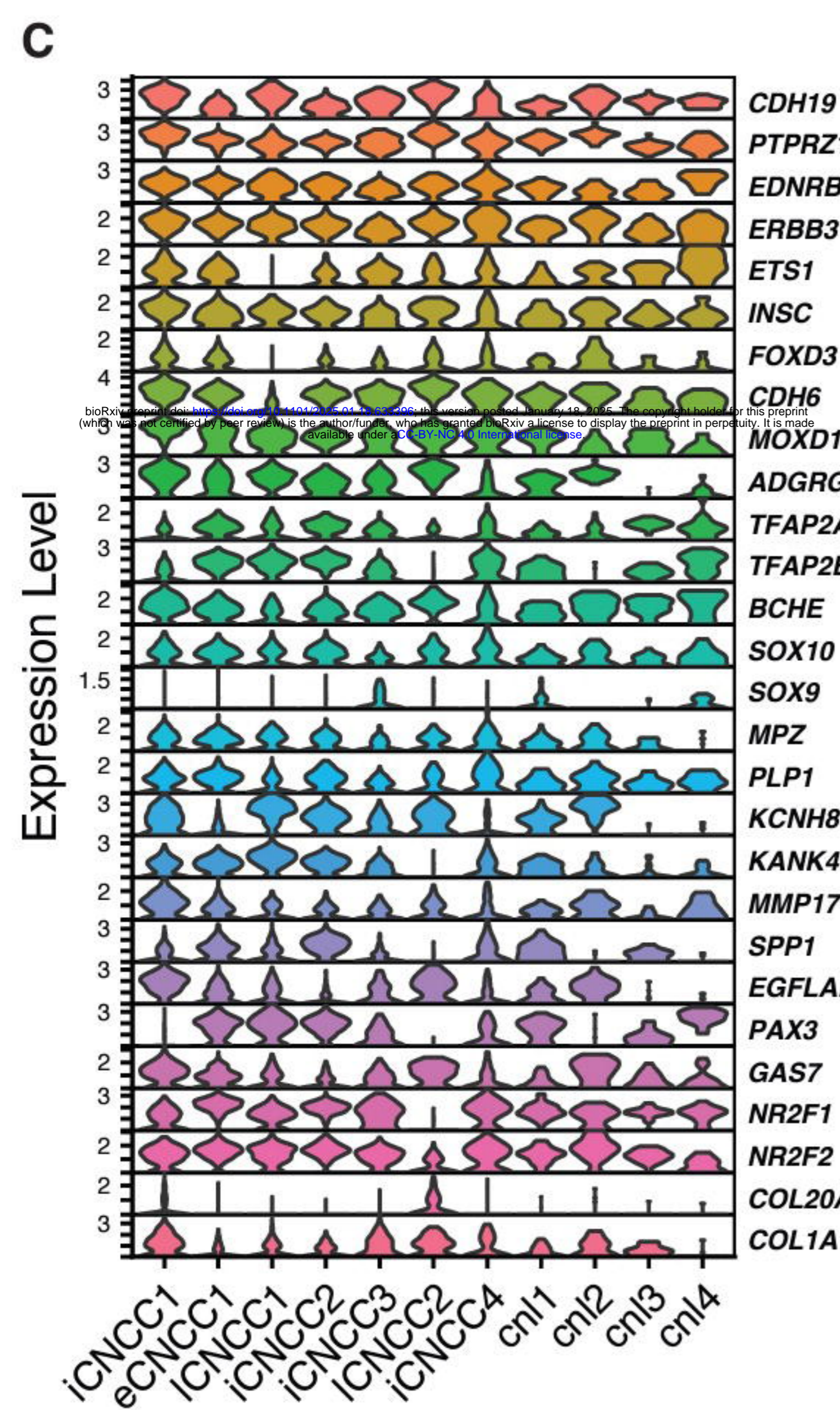
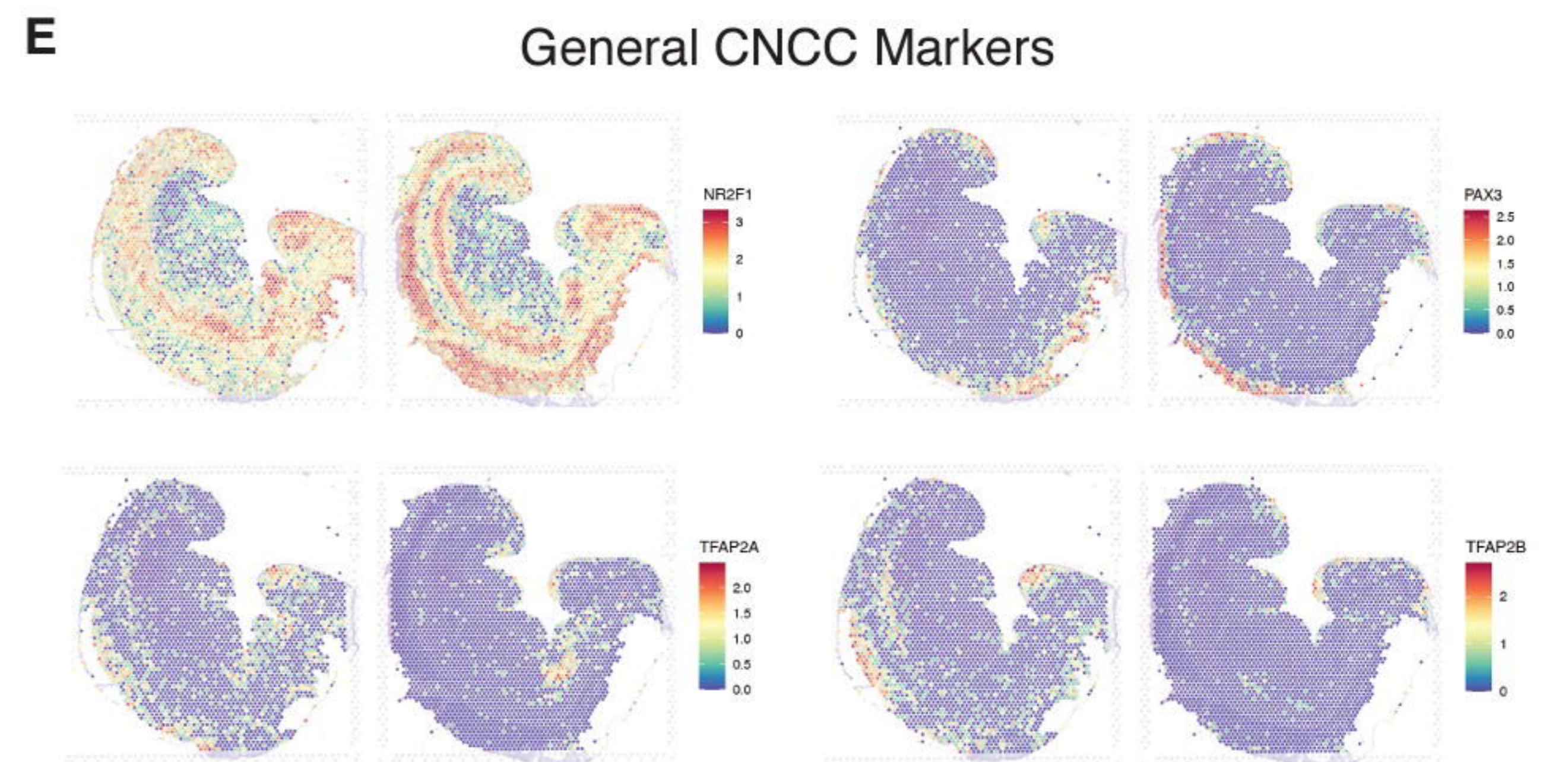
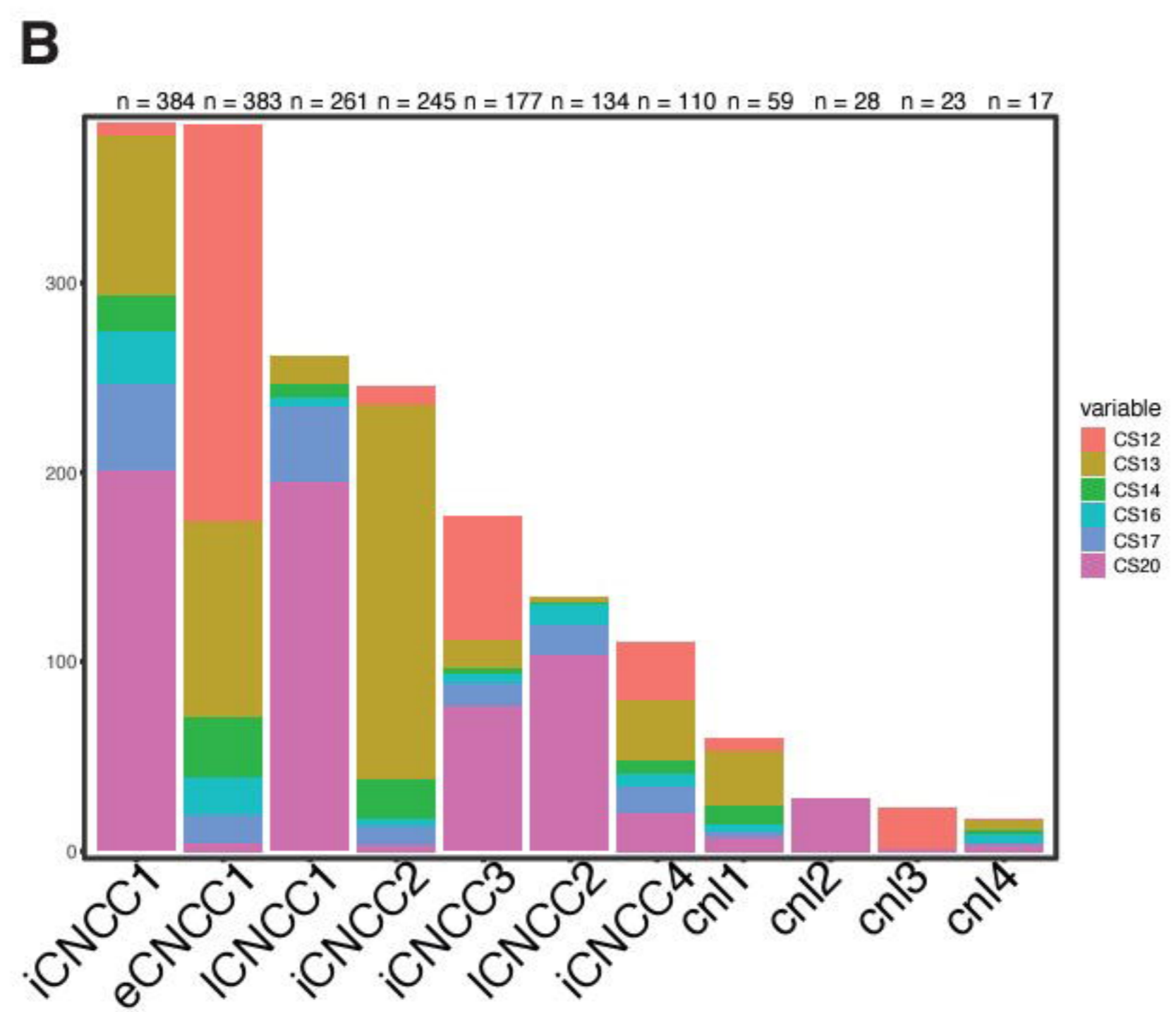
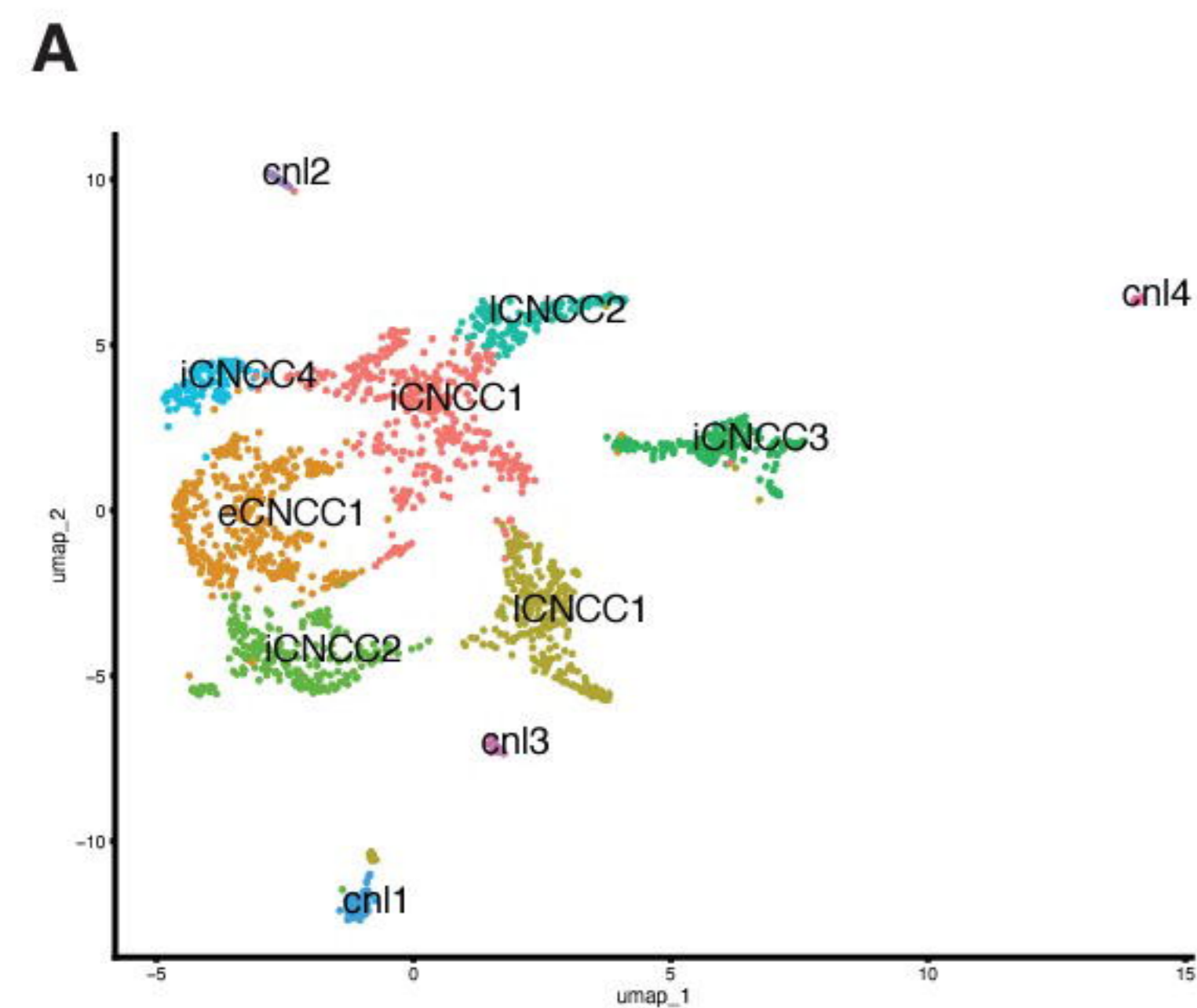
**C** DisGeNet Cluster Enrichments



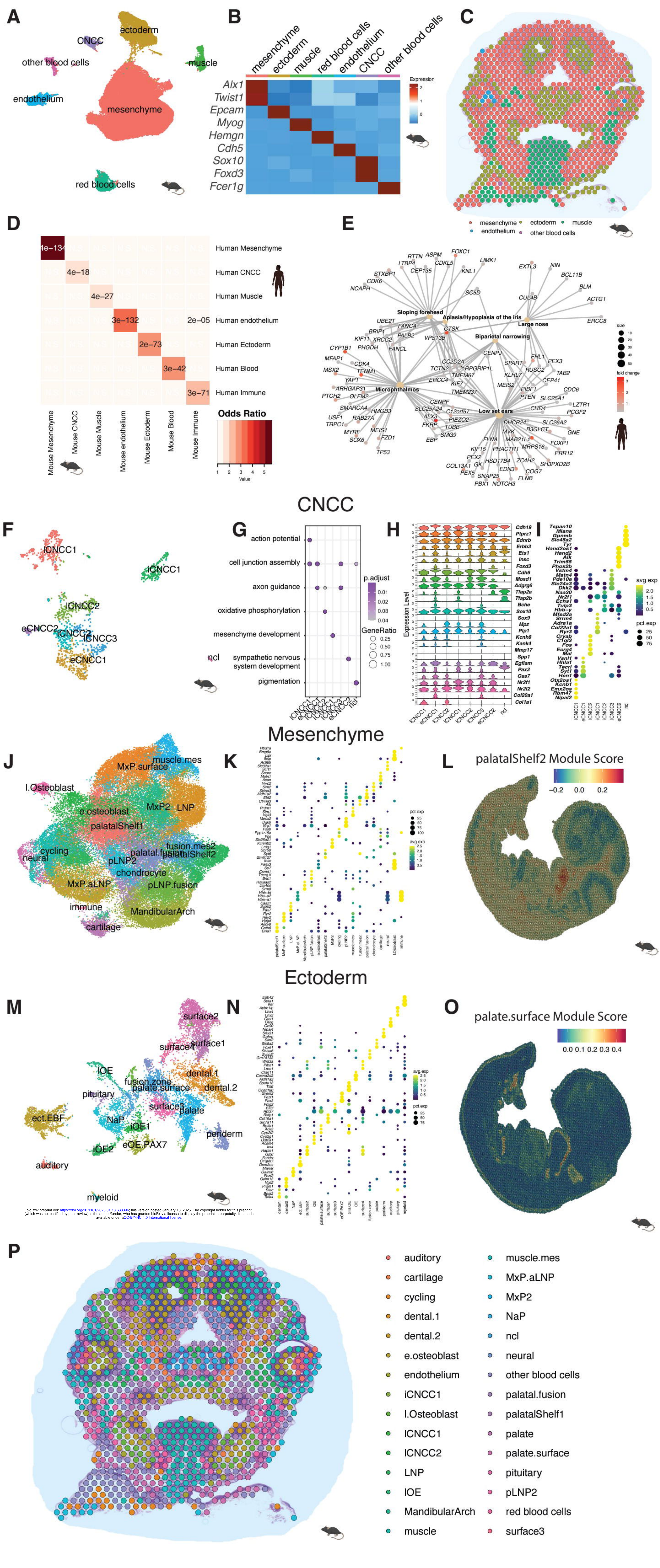
**D** Cranial Neural Crest Genes Module Score



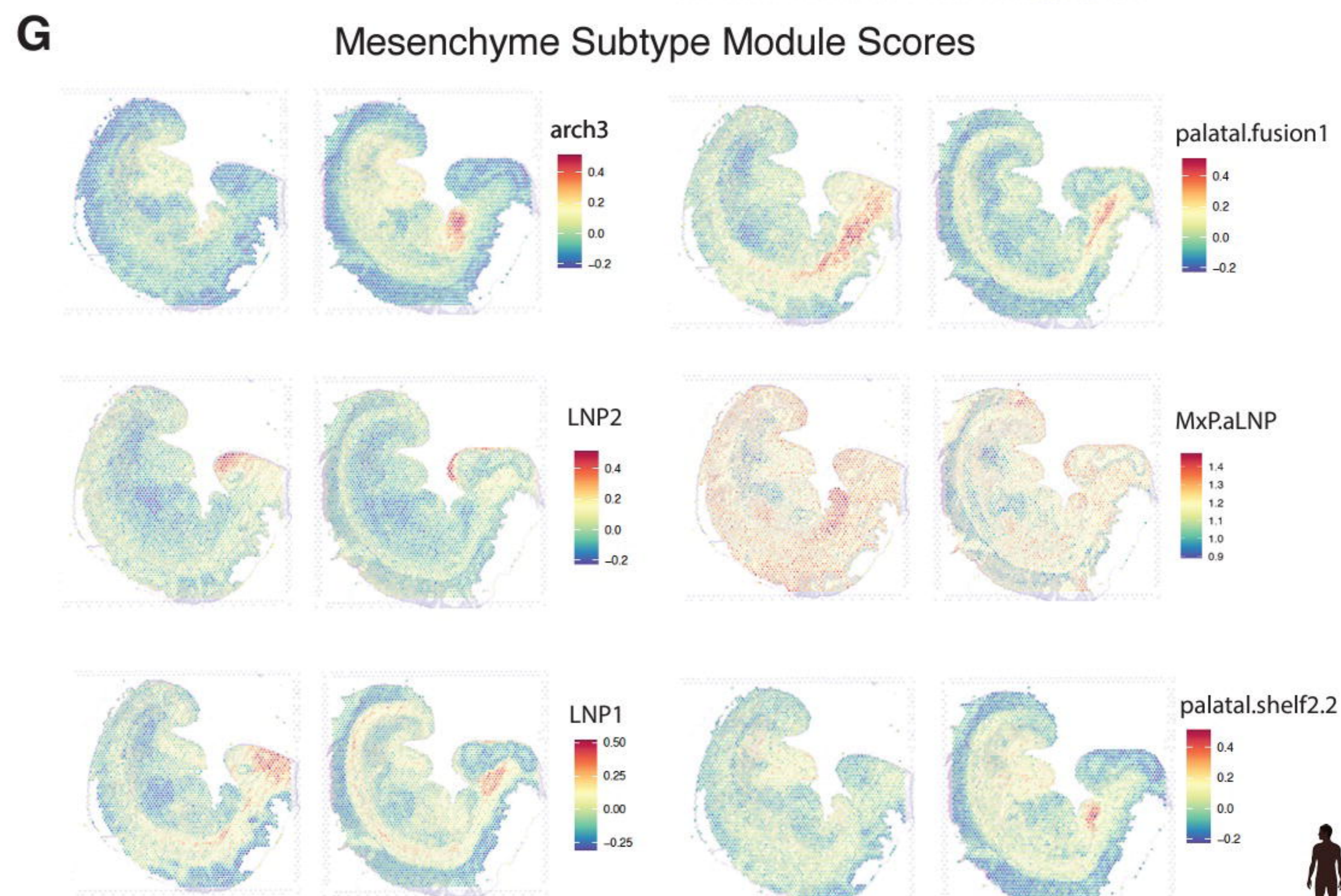
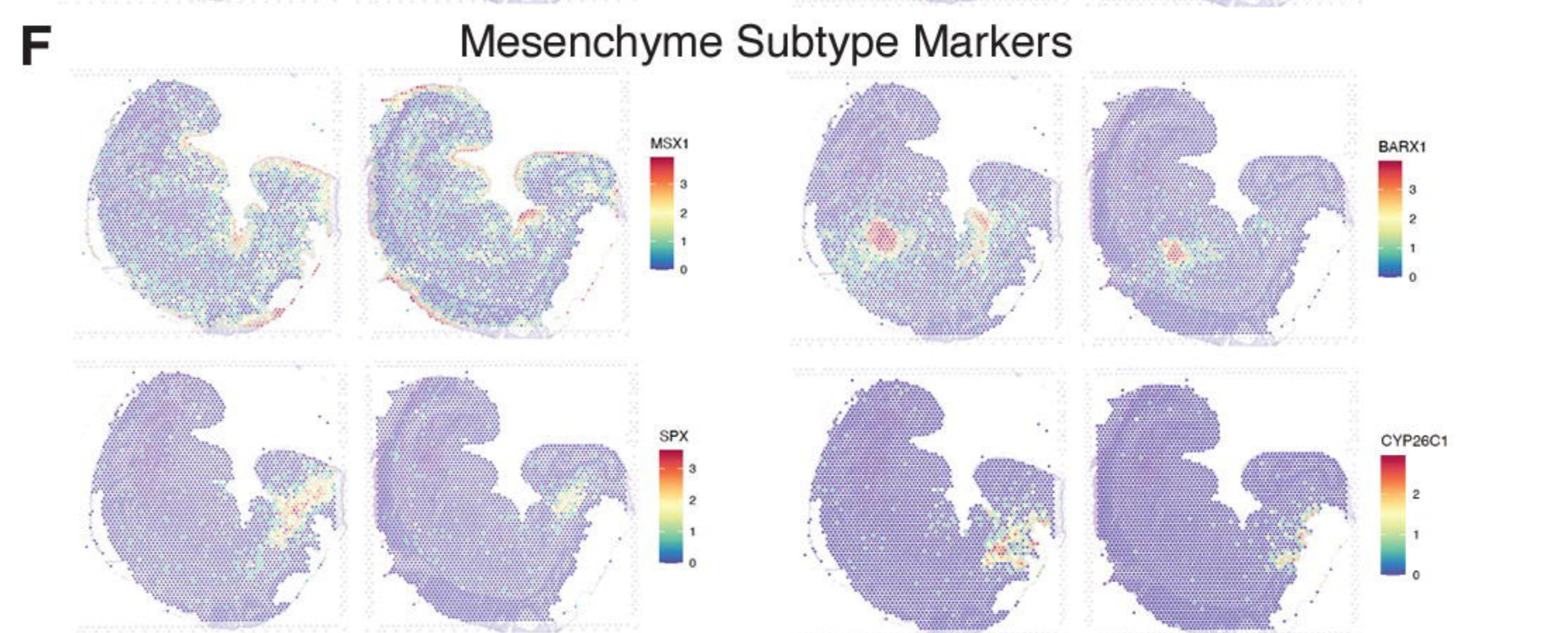
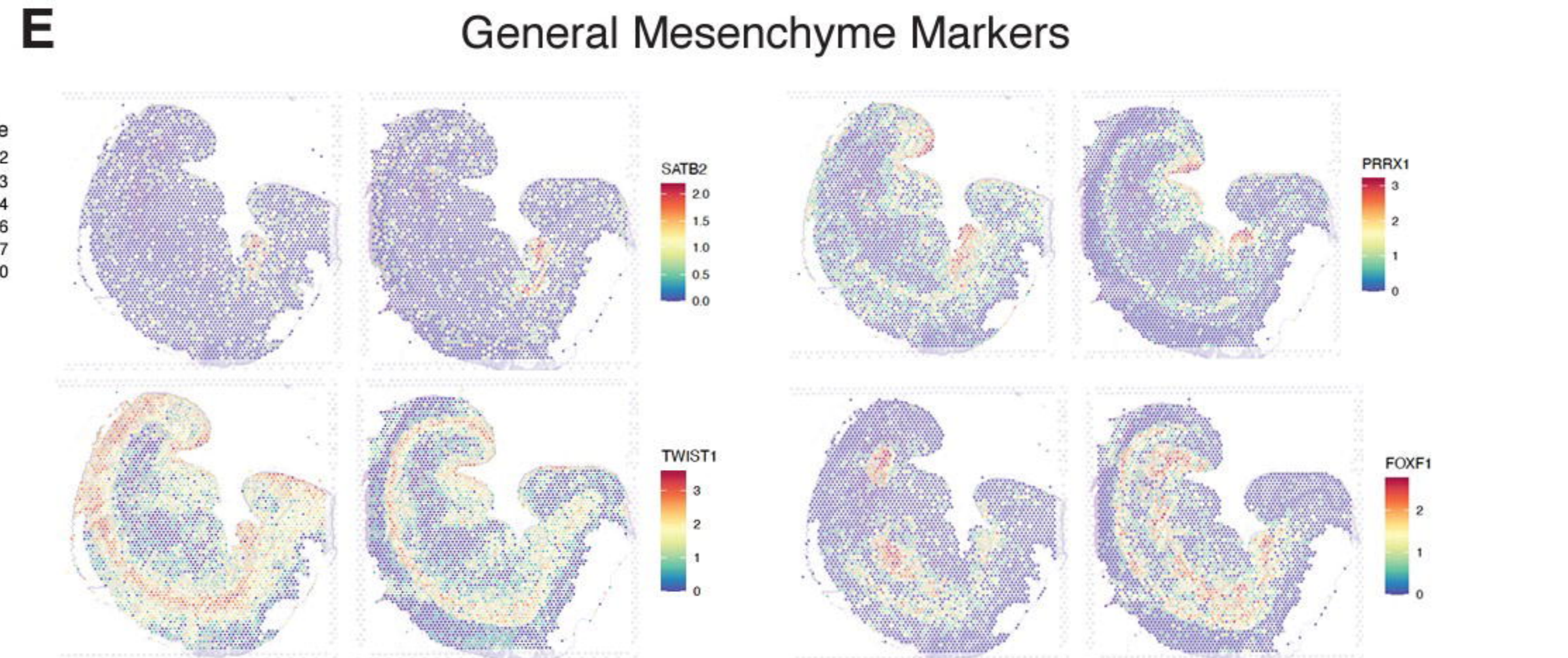
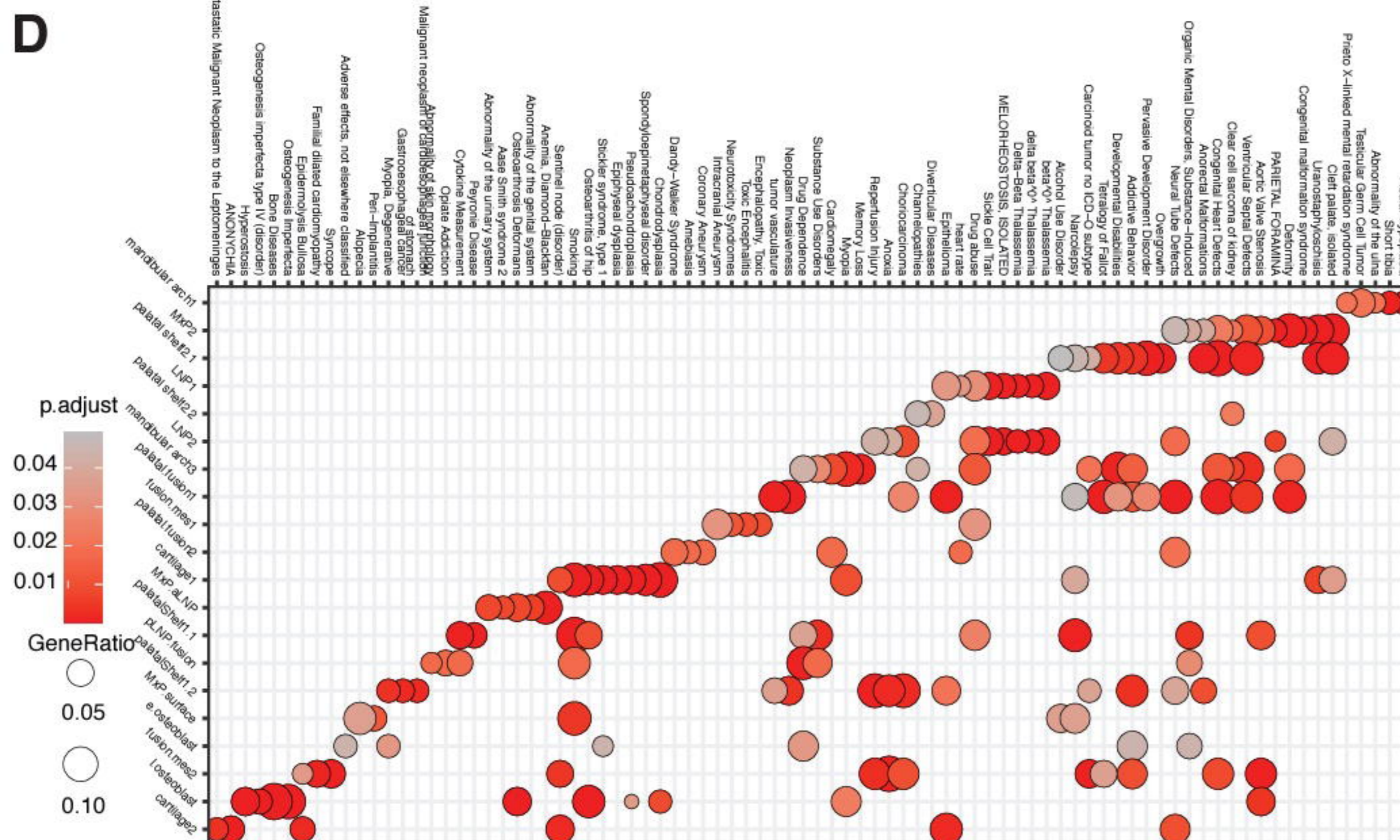
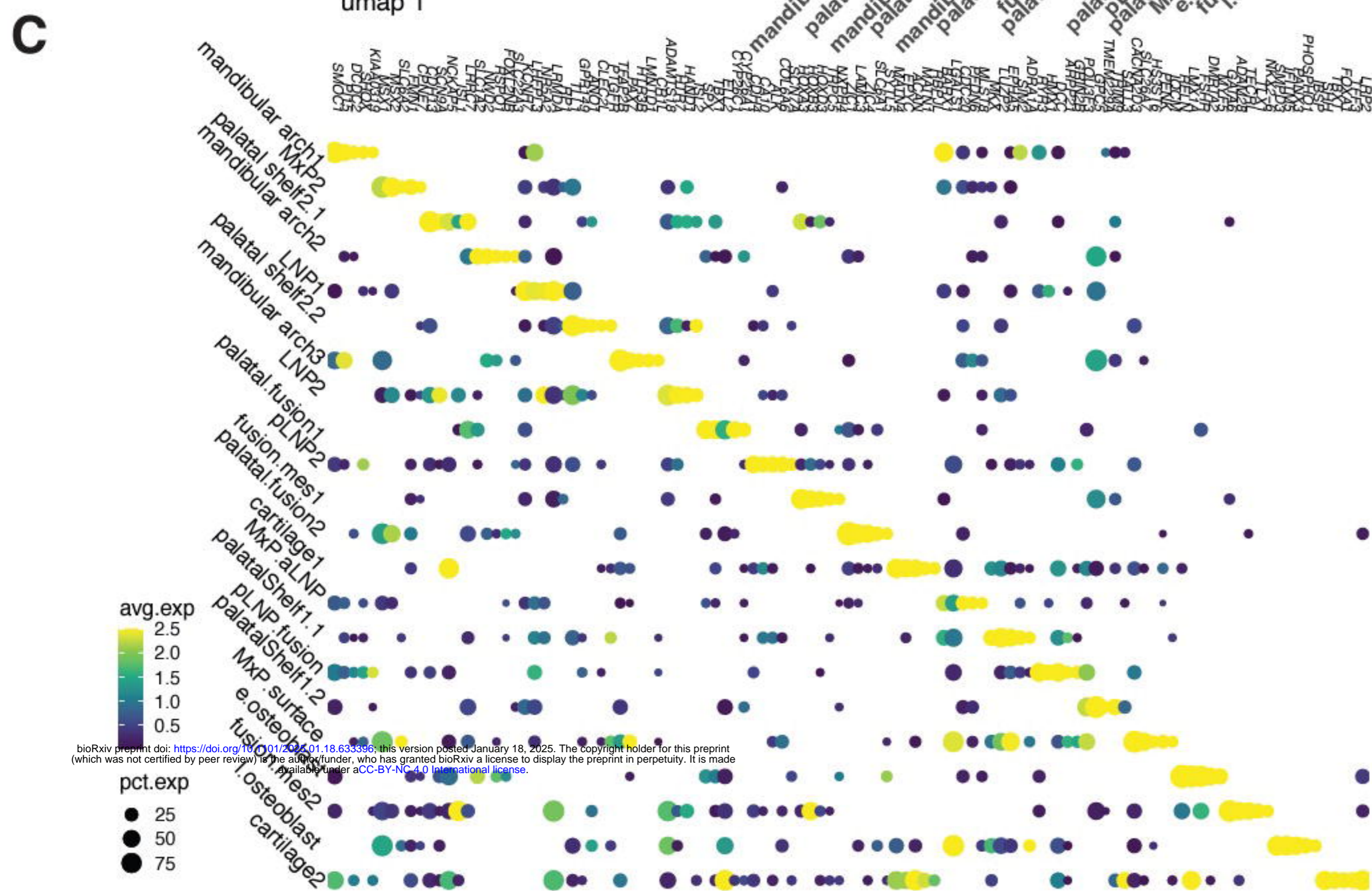
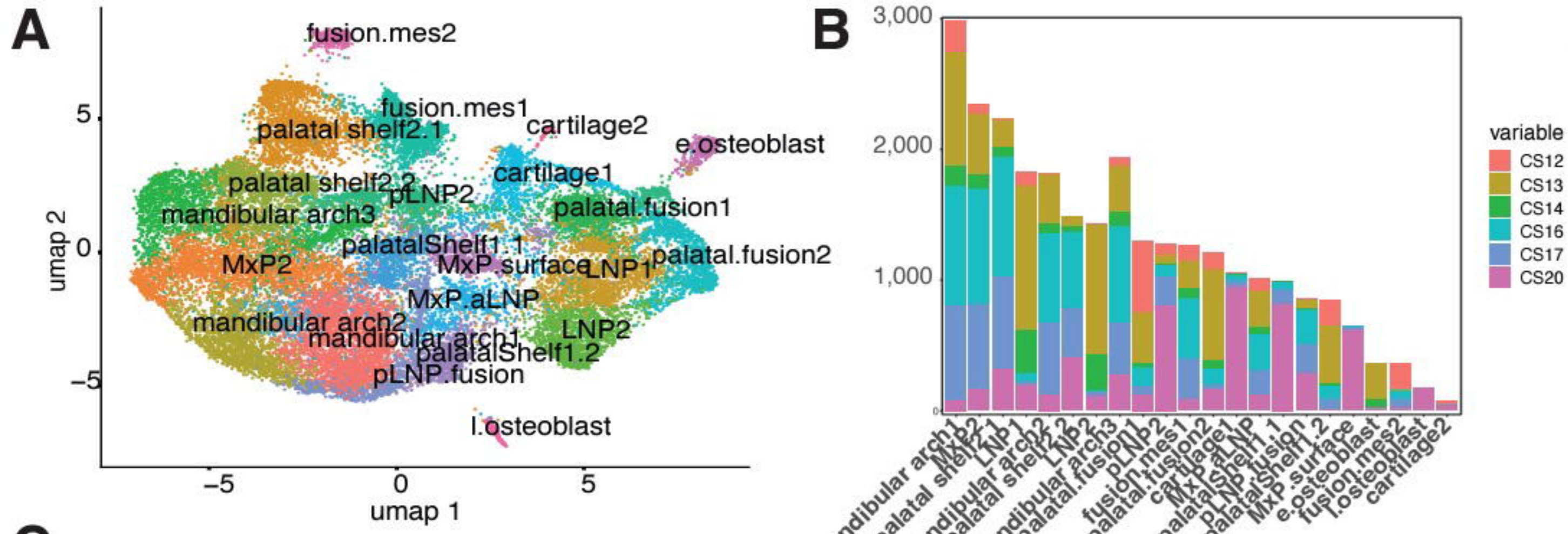








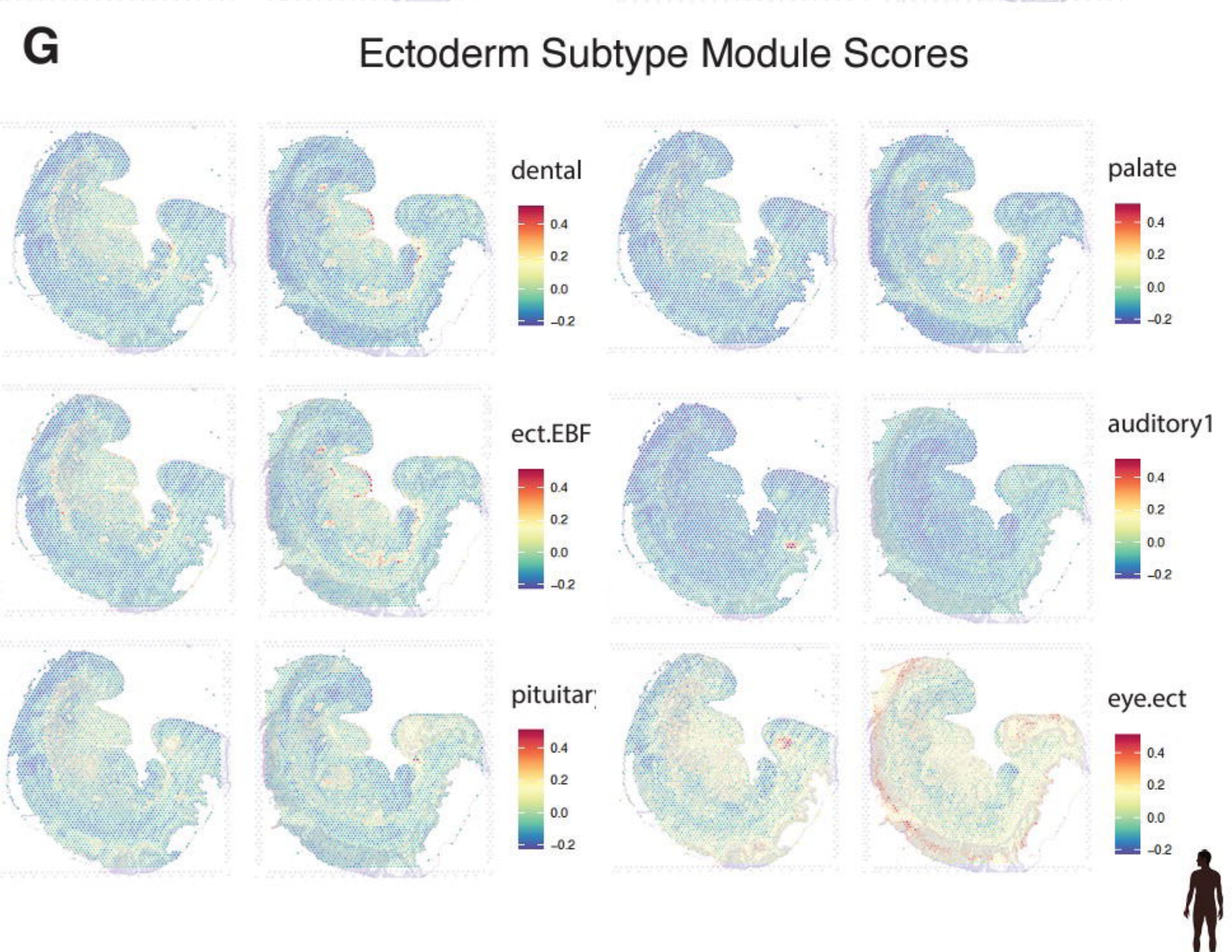
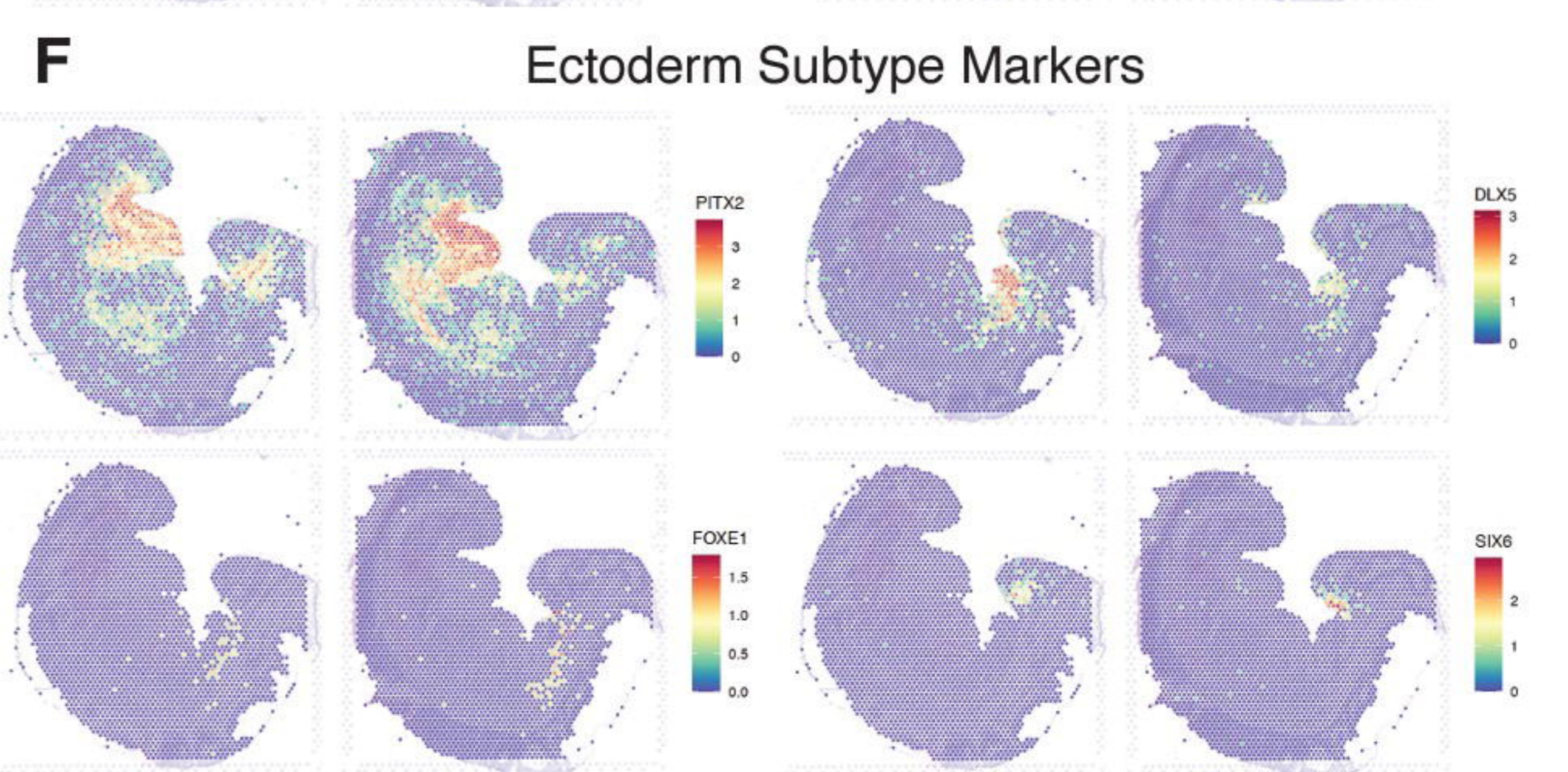
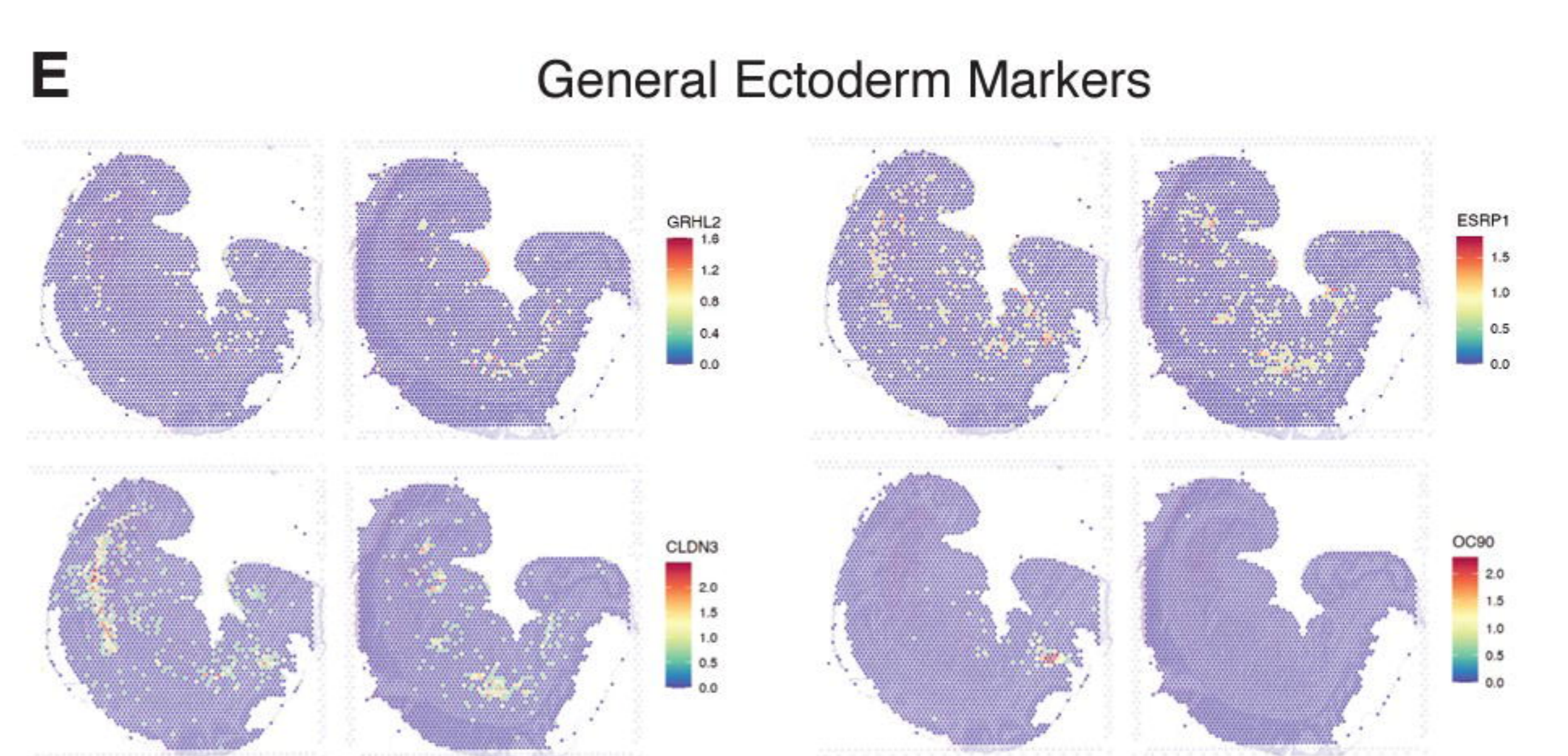
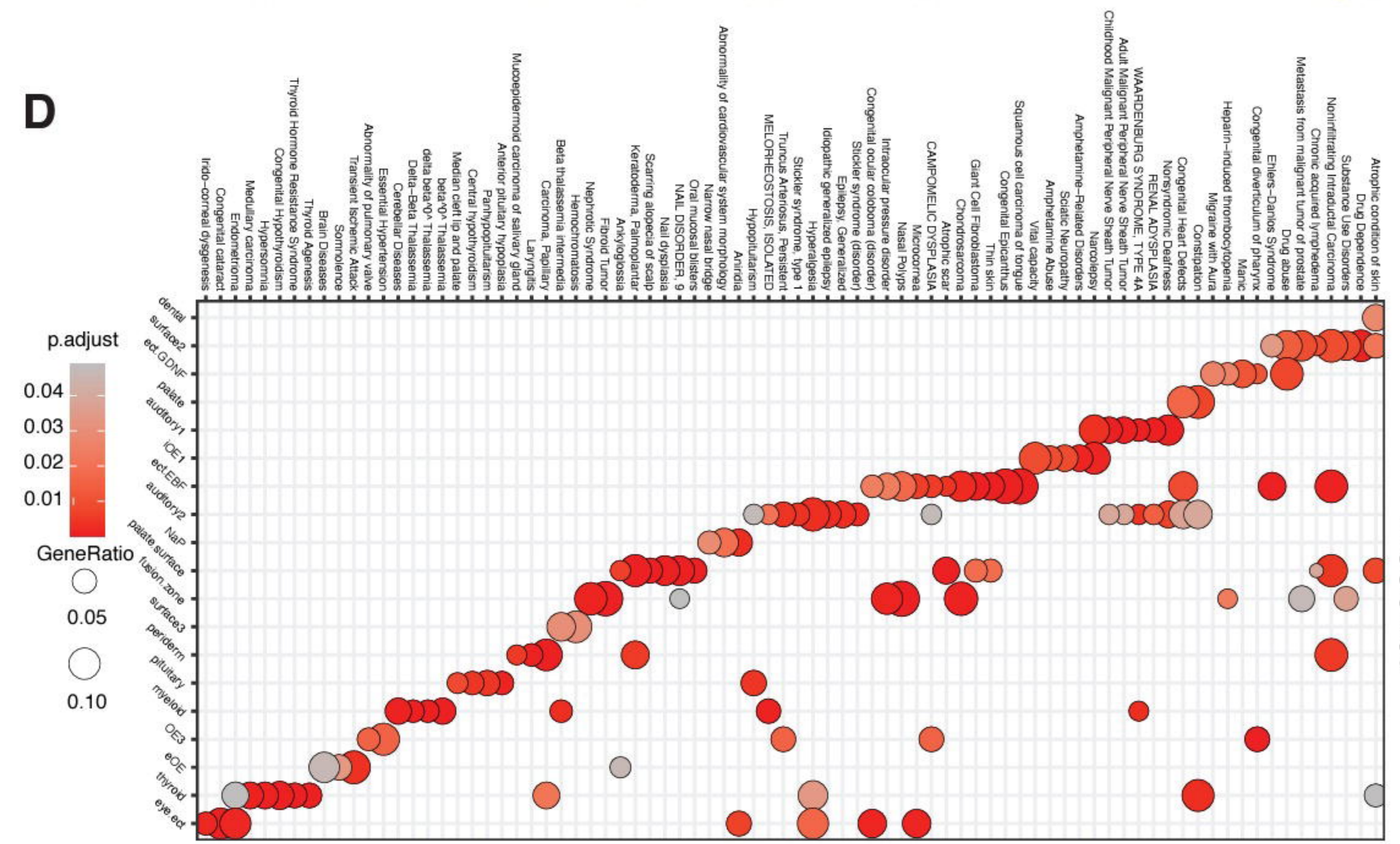
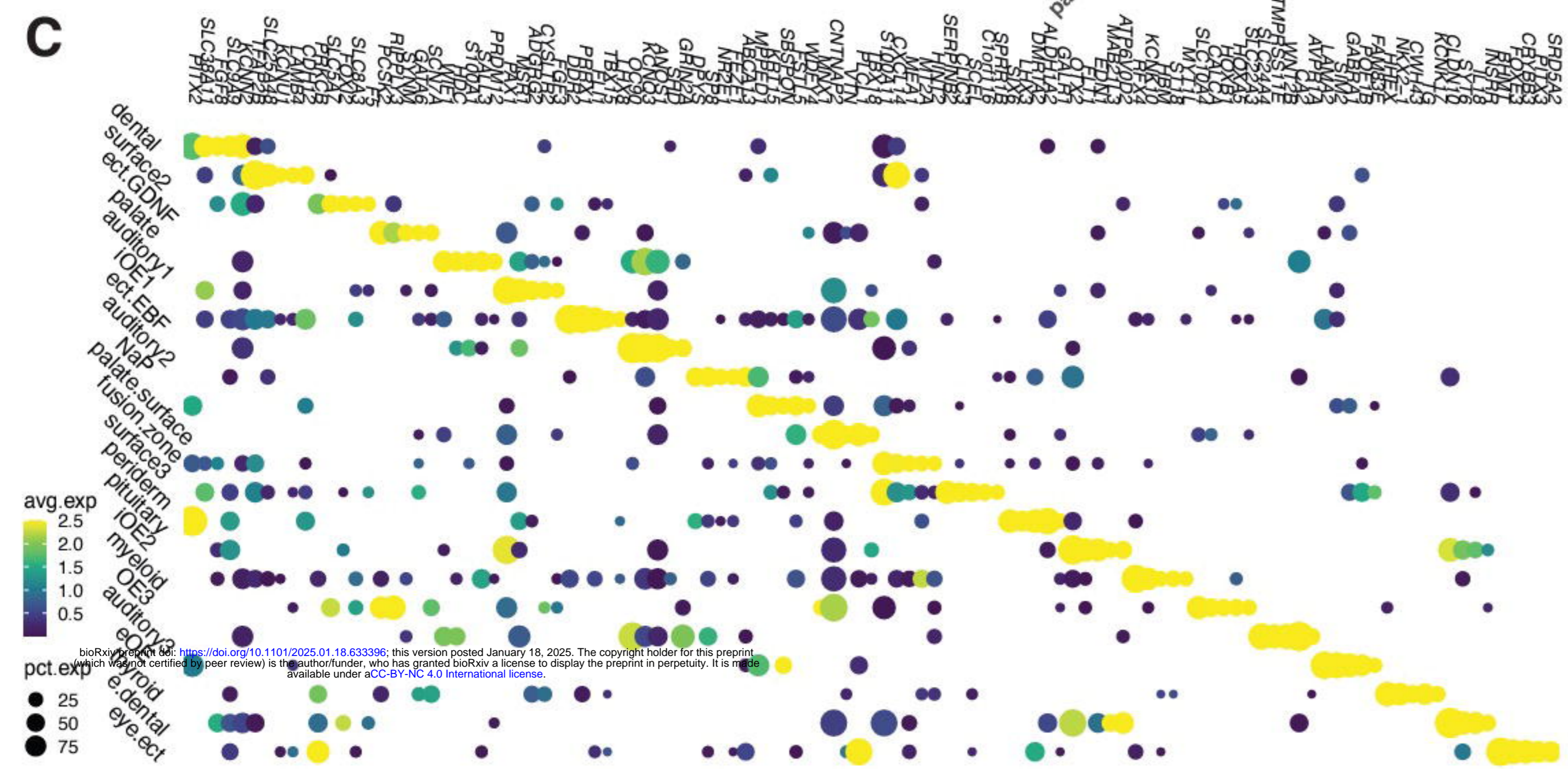
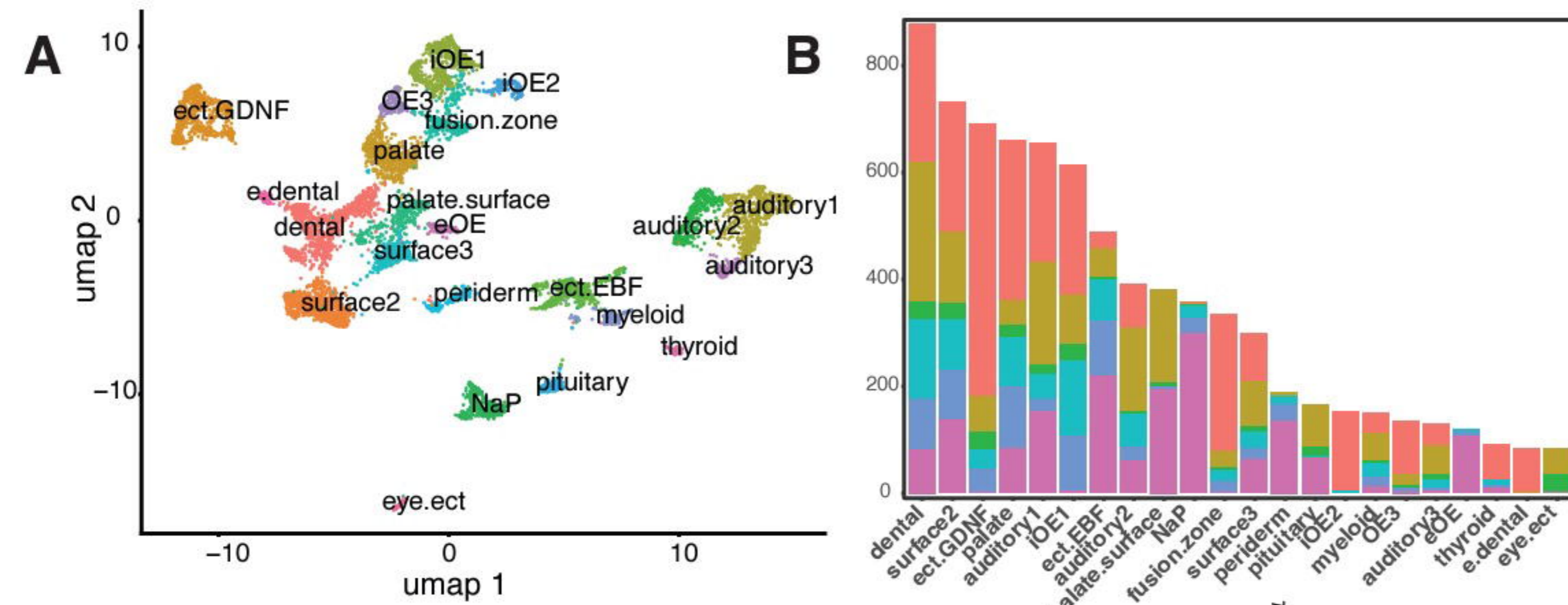




bioRxiv preprint doi: <https://doi.org/10.1101/2025.01.18.633356>; this version posted January 18, 2025. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

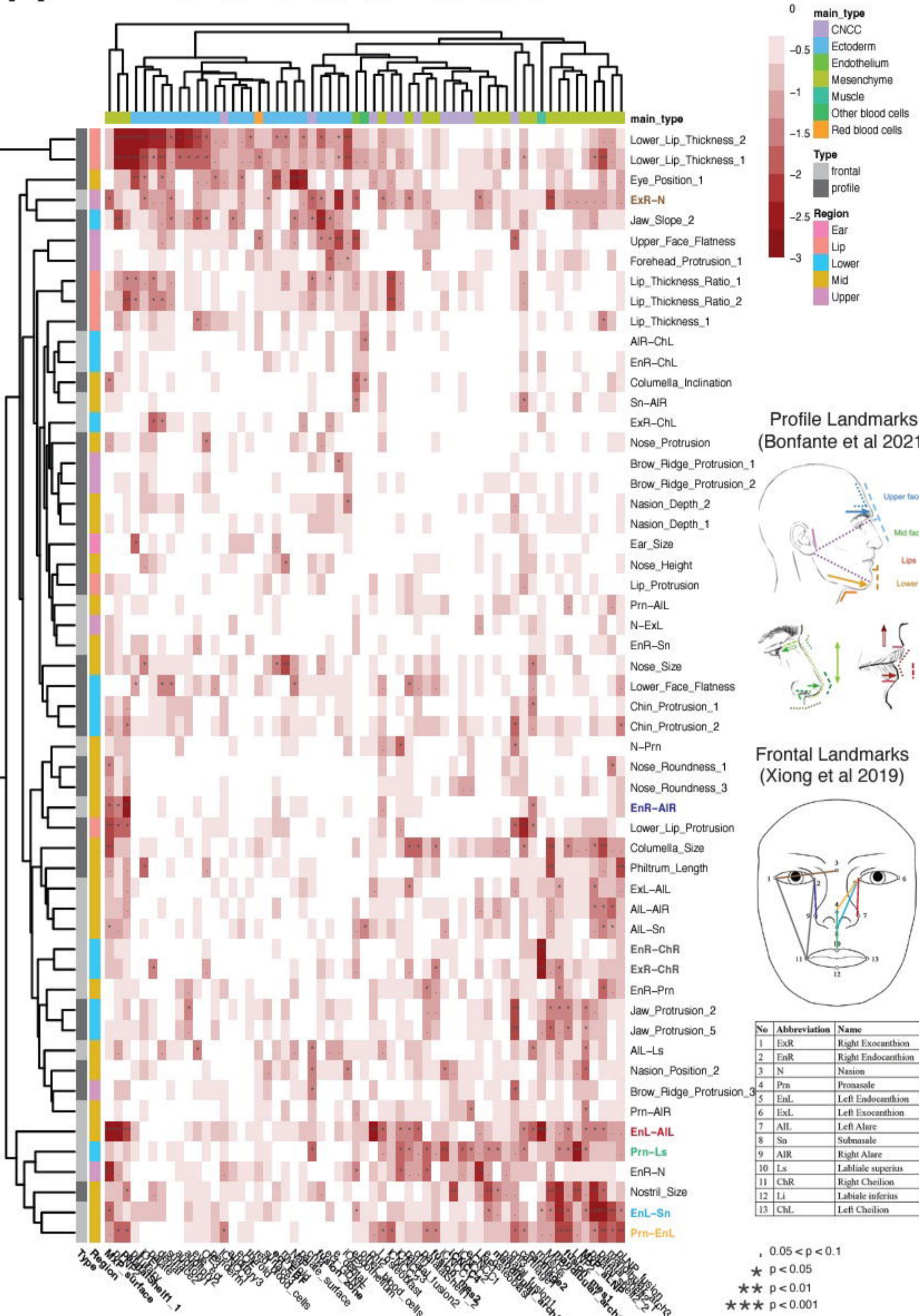




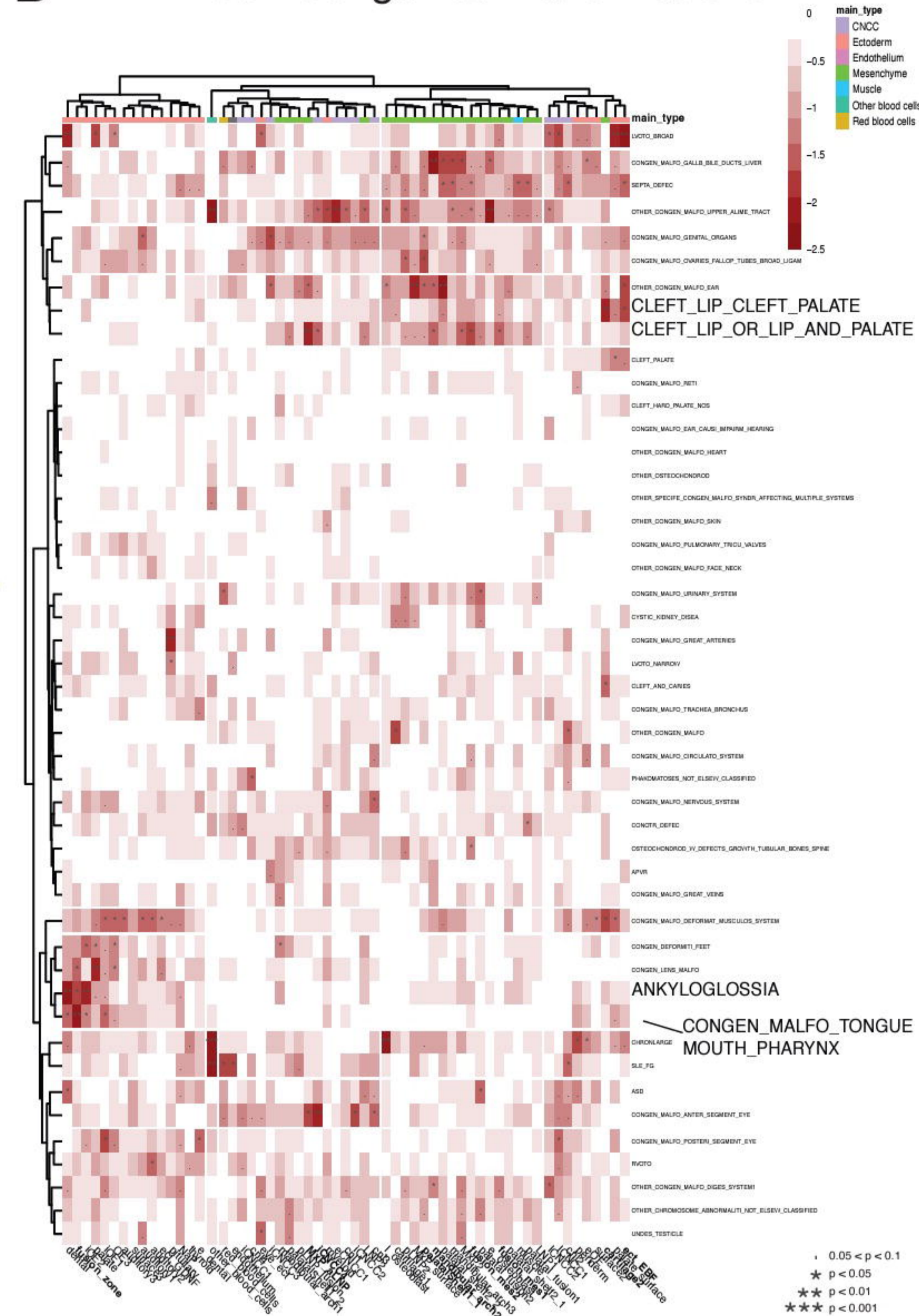




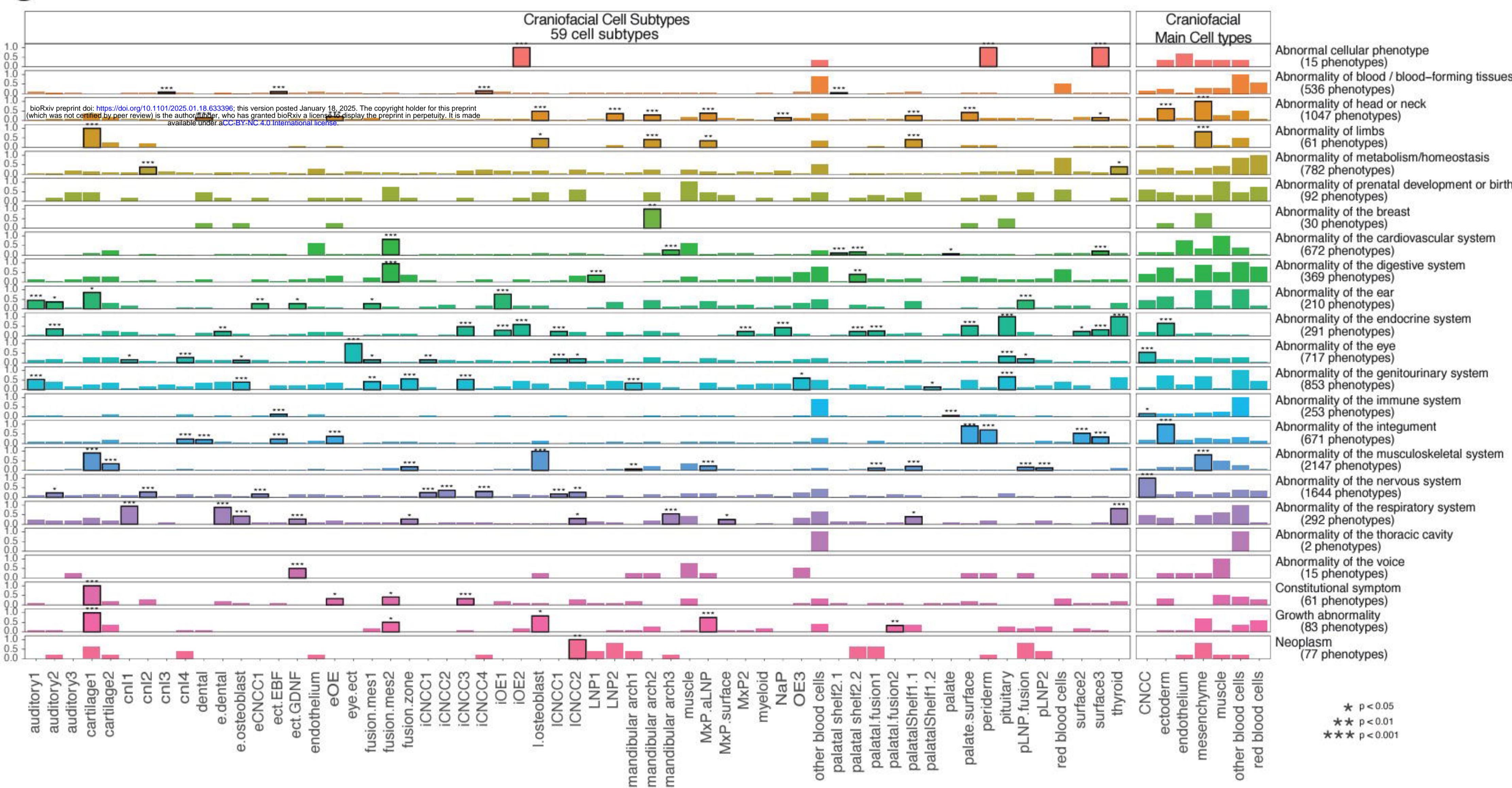
# A Craniofacial Variation



# B FinnGen Congenital Malformations

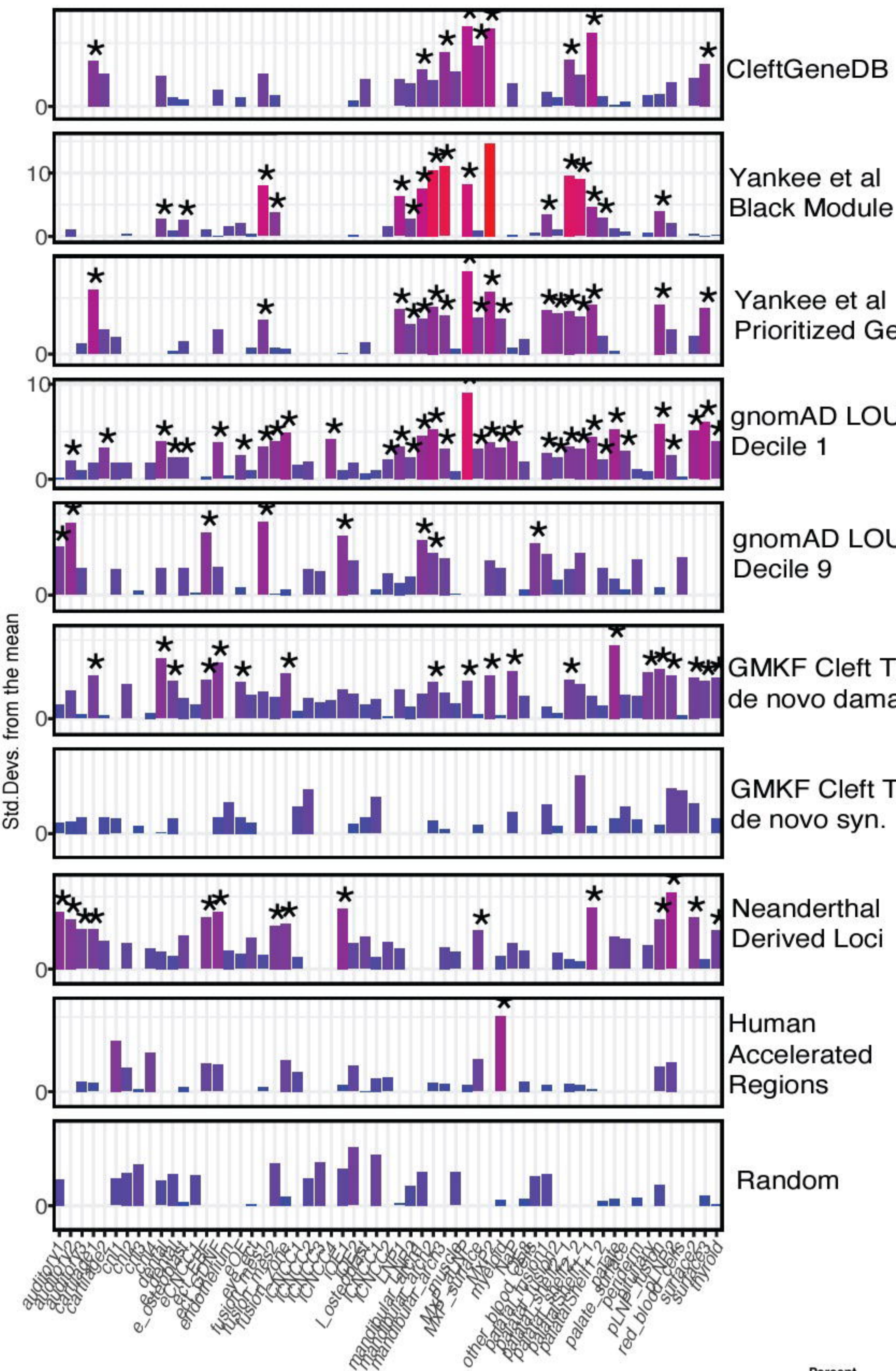


# C HPO Enrichments

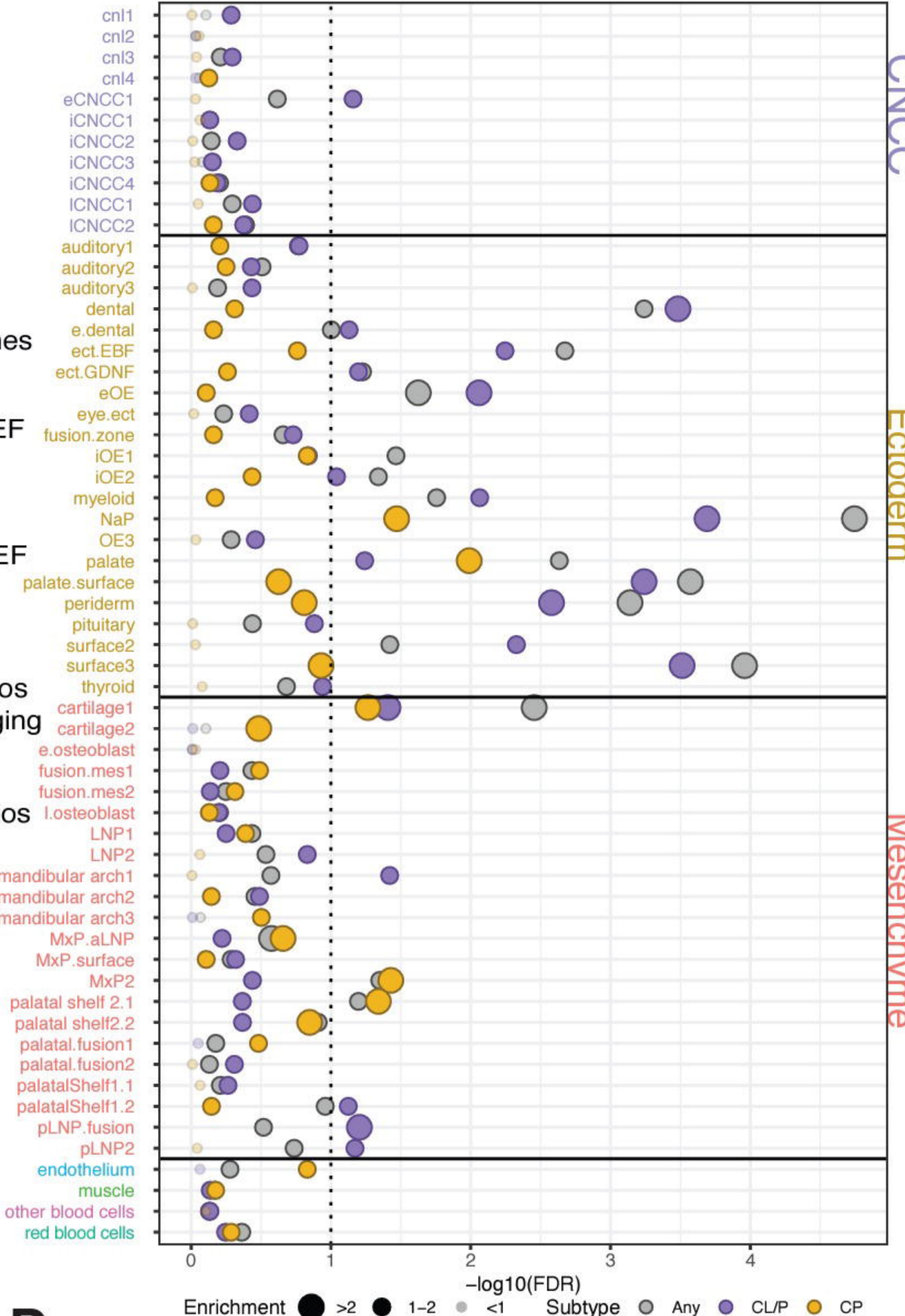




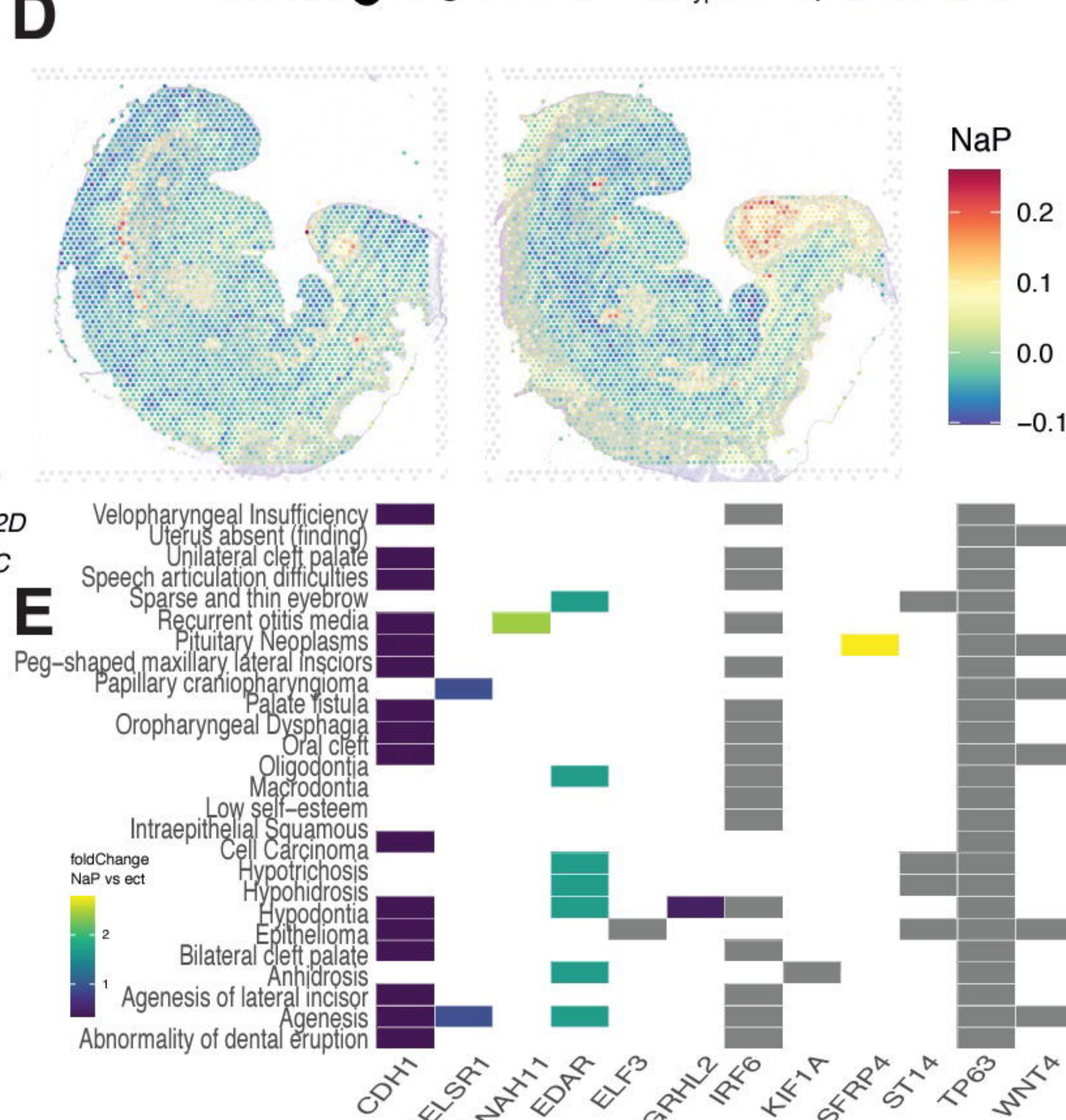
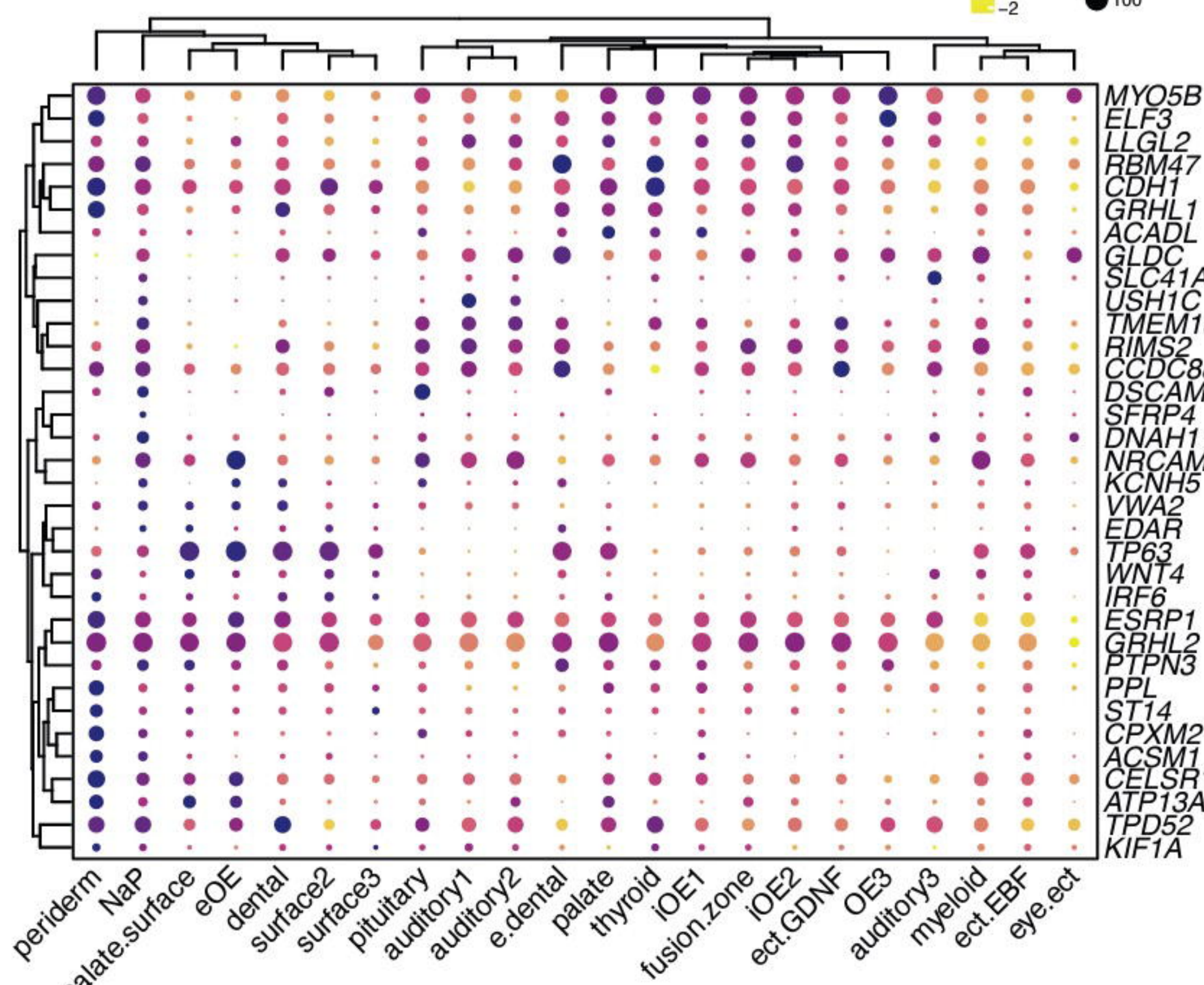
### A Expression Weighted Cell Subtype Enrichment



### B GMKF Cleft Trios de novo protein altering variants



### C Nasal Pacode de novo Protein Altered Genes



bioRxiv preprint doi: <https://doi.org/10.1101/2025.01.18.633396>; this version posted January 18, 2025. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.