# Lifelong single-cell profiling of cranial neural crest diversification

Peter Fabian[#1], Kuo-Chang Tseng[#1], Mathi Thiruppathy[#1], Claire Arata[#1], Hung-Jhen Chen[1],

Joanna Smeeton[1,2], Nellie Nelson[1], and J. Gage Crump*[1]

[1] Eli and Edythe Broad California Institute for Regenerative Medicine Center for Regenerative Medicine and Stem Cell Research, Department of Stem Cell Biology and Regenerative Medicine, University of Southern California Keck School of Medicine, Los Angeles, CA 90033, USA

[2] Department of Rehabilitation and Regenerative Medicine, Columbia University Irving Medical Center, Columbia University, New York, NY 10032, USA

[#]These authors contributed equally

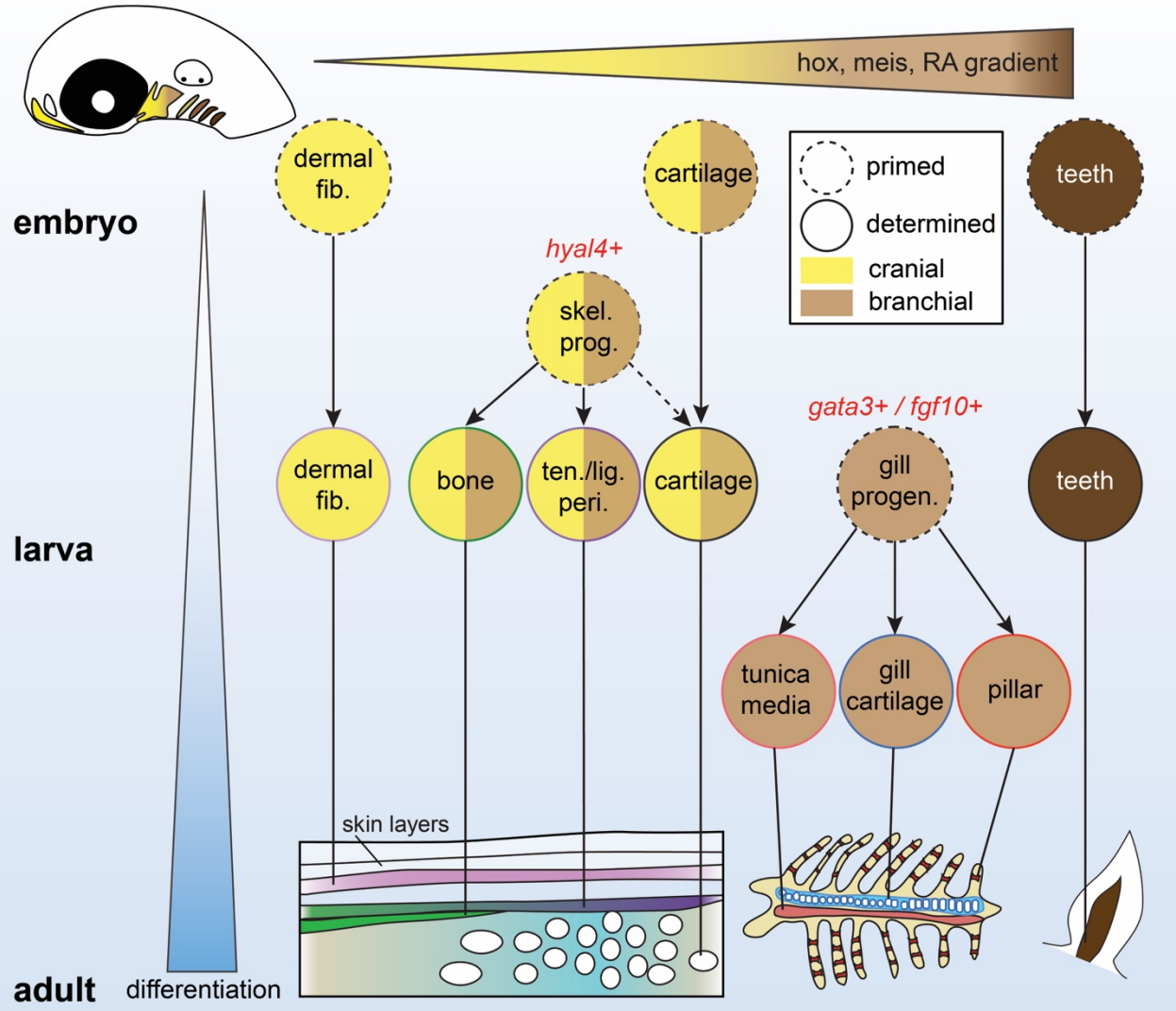*Correspondence: J. Gage Crump, gcrump@usc.edu, (323) 442-2693

**Running Title:** Single-cell profiling of neural crest

**Keywords:** Cranial neural crest; single-cell genomics; craniofacial skeleton; chromatin accessibility; zebrafish

1

## Highlights

- Single-cell transcriptome and chromatin atlas of cranial neural crest

- Progressive emergence of region-specific cell fate competency

- Chromatin accessibility mapping identifies candidate lineage regulators

- Gata3 function linked to gill-specific respiratory program

## Graphical Abstract

## Abstract

The cranial neural crest generates a huge diversity of derivatives, including the bulk of connective and skeletal tissues of the vertebrate head. How neural crest cells acquire such extraordinary lineage potential remains unresolved. By integrating single-cell transcriptome and chromatin accessibility profiles of cranial neural crest-derived cells across the zebrafish lifetime, we observe region-specific establishment of enhancer accessibility for distinct fates. Neural crest-derived cells rapidly diversify into specialized progenitors, including multipotent skeletal progenitors, stromal cells with a regenerative signature, fibroblasts with a unique metabolic signature linked to skeletal integrity, and gill-specific progenitors generating cell types for respiration. By retrogradely mapping the emergence of lineage-specific chromatin accessibility, we identify a wealth of candidate lineage-priming factors, including a Gata3 regulatory circuit for respiratory cell fates. Rather than multilineage potential being an intrinsic property of cranial neural crest, our findings support progressive and region-specific chromatin remodeling underlying acquisition of diverse neural crest lineage potential.

## Main text

Cranial neural crest-derived cells (CNCCs) are a vertebrate-specific population, often referred to as the fourth germ layer, that have extraordinary potential to form diverse cell types. In addition to pigment cells and the peripheral nervous system, CNCCs form the ectomesenchyme that populates the pharyngeal arches and gives rise to much of the skeleton and connective tissue of the jaws and face[1]. Posterior arch CNCCs contribute to a distinct set of organs, including the thymus, parathyroid, and cardiac outflow tract, and in fishes cell types important for respiration, including specialized endothelial-like pillar cells that promote gas exchange[2]. In zebrafish, teeth develop from CNCCs of the most posterior seventh arch.

59     The extent to which diverse lineage potential is an intrinsic property of CNCCs, or acquired

60     through later inductive signaling, has been investigated for over a century through labeling,

61     grafting, and extirpation experiments, yet remains unresolved[3]. Individual avian CNCCs can

62     generate multiple types of derivatives in vitro, including ectomesenchyme and neuroglial cells,

63     suggesting multilineage potential is an intrinsic property[4]. However, upon cranial transplantation,

64     trunk neural crest cells, which normally do not make mesenchymal derivatives, can contribute to

65     the facial skeleton following extended culture[5] or misexpression of key transcription factors[6]. A

66     recent study in skate also shows mesodermal contribution to the gill skeleton, a classically

67     considered CNCC-derived structure[7]. These findings point to extrinsic inductive cues for CNCC

68     fate determination. Here we take a genomics approach in zebrafish to understand when

69     enhancers linked to diverse CNCC fates first gain accessibility, thus revealing that chromatin

70     accessibility underlying multilineage potential is largely gained after CNCC migration.

71

## Single-cell atlas of CNCC derivatives across the zebrafish lifetime

73     To understand the emergence and diversification of CNCC lineages across the lifetime of a

74     vertebrate, we constructed a longitudinal single-cell atlas of gene expression and chromatin

75     accessibility of zebrafish CNCC derivatives. We permanently labeled CNCCs using *Sox10:Cre;*

76     *actab2:loxP-BFP-STOP-loxP-dsRed* (*Sox10>dsRed*) fish (Fig. 1a), in which genetic

77     recombination indelibly labels CNCCs shortly after their specification at 10 hours post-fertilization

78     (hpf)[8]. Previous single-cell analyses of CNCCs in zebrafish[9], chick[10], and mouse[11-13], and in vitro

79     human CNCC-like cells[14], had focused on CNCC establishment, migration, and early fate choices

80     between the neuroglial, pigment, and ectomesenchyme lineages. Here we investigate cellular

81     diversity and lineage progression of CNCC ectomesenchyme across embryonic (1.5 and 2 days

82     post-fertilization (dpf)), larval (3 and 5 dpf), juvenile (14 and 60 dpf), and adult (150-210 dpf)

83     stages. After fluorescence activated cell sorting (FACS) of *Sox10>dsRed*+ cells from dissected

84     heads, we performed single-cell RNA sequencing (scRNAseq) and single-nuclei assay for

85    transposase accessible chromatin sequencing (snATACseq) at each stage using the 10X

86    Genomics Chromium platform and paired-end Illumina next-generating sequencing (Fig. 1b).

87    After filtering for quality, we obtained 58,075 cells with a median of 866 genes per cell for

88    scRNAseq, and 88,177 cells with a median of 10,449 fragments per cell for snATACseq. To better

89    resolve snATACseq data, we used the SnapATAC package[15], which integrates snATACseq with

90    scRNAseq data to create "pseudo-multiome" datasets.

91

92    Analysis of CNCC cell clusters across all stages using UMAP dimensionality reduction recovered

93    most known CNCC derivatives, including Schwann cells (glia), several neuronal subtypes,

94    pigment cells, and diverse mesenchymal cell types (Fig. S1-8, Table S1). We also recovered otic

95    placode and epithelial cells, likely reflecting additional non-CNCC expression of *Sox10:Cre*[8], and

96    blood lineage cells, likely due to autofluorescence. Similar clusters were recovered using both

97    scRNAseq and SnapATAC data. We then re-clustered the CNCC ectomesenchyme sub-

98    population across stages, as this makes the most substantial and diverse cell contributions in the

99    head. To confirm ectomesenchyme identity at 1.5 dpf, we also performed scRNAseq analysis of

100    cells double-positive for the CNCC transgene *sox10:dsRed* and the ectomesenchyme transgene

101    *fli1a*:e*GFP*. Co-clustering showed high concordance between *sox10:dsRed+*; *fli1a*:e*GFP+*

102    ectomesenchyme and the *Sox10>dsRed+* ectomesenchyme subset, and between

103    *Sox10>dsRed+* ectomesenchyme scRNAseq subsets at all 7 stages (Fig. S8).

104

105    At the adult stage, we recovered 17 distinct clusters using scRNAseq that corresponded to 16

106    clusters using SnapATAC; these were largely associated with the jaw skeleton or gills (Fig. 1c-e).

107    Skeletal derivatives include bone, cartilage, teeth, and a population with properties of periosteum,

108    tendon, and ligament. Gills are composed of primary filaments containing cartilage rods and

109    primary veins surrounded by a tunica media, and numerous secondary filaments housing

110    endothelial-like "pillar" cells that promote gas exchange. Unexpectedly, we recovered a

111    specialized type of gill cartilage distinct from that in the rest of the head, as well as pillar and

112    tunica media cells and putative gill progenitors. We also recovered smooth muscle, perivascular,

113    and stromal cells (see Table S1 for cluster marker genes and Fig. S9-10 for in situ validation).

114

115    In addition to skeletal and gill populations, we recovered a distinct type of fibroblast enriched for

116    the cell adhesion molecule *chl1a* and *wnt5a*. Strikingly, these fibroblasts are also enriched for

117    genes encoding enzymes for all steps of phenylalanine and tyrosine breakdown (Fig. 1f, Fig. S11).

118    In situ hybridization for two of these genes (*hpdb* and *pah*) reveals that these fibroblasts are in

119    the dermis between the skin epidermis and *runx2b+/sp7+* osteoblast lineage cells (Fig. 1g,h).

120    Humans with mutations in *HGD*, which encodes an intermediate enzyme in the Phe/Tyr catabolic

121    pathway, develop Alkaptonuria, or black bone disease, due to accumulation and pathological

122    aggregation of homogentisic acid[16]. As the abundant melanocytes in the zebrafish skin use high

123    levels of Tyr to synthesize melanin, one possibility is that these specialized dermal fibroblasts

124    function to protect the skeleton by removing damaging Phe/Tyr metabolites.

125

126    **Progressive emergence of CNCC derivatives and region-specific progenitors**

127    To understand lineage decisions of CNCC mesenchyme across time, we first used the STITCH

128    algorithm[17] to connect individual stages into developmental trajectories for scRNAseq and

129    snATACseq datasets (Fig. 2a,b). As early as 3 dpf (particularly apparent for snATACseq), we

130    observe divergence of CNCCs into skeletogenic versus gill lineages. A *hyal4+* perichondrium

131    population precedes branches for tendon/ligament, periosteum, and osteoblasts (Fig. S9), and an

132    *fgf10b+* gill progenitor population appears at 5 dpf and precedes branches for gill cartilage, pillar,

133    and tunica media cells (Fig. S10). We also observe a distinct trajectory to dermal fibroblasts by 3

134    dpf (Fig. S11), as well as to *cxcl12a+* stromal cells (Fig. S9) and teeth. We do not observe CNCC

135    contributions to cardiomyocytes (Fig. S8), in contrast to reports for amniotes[18]. By creating an

136    index for ectomesenchyme-enriched gene expression at 1.5 dpf, a stage preceding the onset of

137    differentiation, we found no evidence for retention of ectomesenchyme identity at later stages, as

138    shown by aggregated ectomesenchyme gene expression and the early ectomesenchyme marker

139    *nr2f5*[19] (Fig. S12). Although formation of CNCC ectomesenchyme involves a reacquisition of the

140    pluripotency network[14], we also did not observe expression of pluripotency genes *pou5f3* (*oct4*),

141    *sox2*, *nanog*, and *klf4* at any stage of post-migratory ectomesenchyme, with the exception of

142    *lin28aa* that displays broad expression at 1.5 dpf and is rapidly extinguished by 2 dpf (Fig. S12).

143    Rather than maintenance of a multipotent ectomesenchyme population, our data point to

144    progressive emergence of specialized *hyal4*+ perichondrium, *cxcl12a*+ stromal, and *fgf10a/b*+ gill

145    populations at 3 dpf and beyond (Fig. S12).

146

147    To further dissect region-specific lineages, we used Monocle3[20] on scRNAseq datasets to

148    construct pseudotime trajectories of anterior arch (i.e. skeletogenic) versus posterior arch (i.e. gill,

149    *hoxb3a+/gata3*+) CNCC mesenchyme at 5-14 dpf (Fig. S13, dermal fibroblasts and teeth were

150    removed). For skeletogenic clusters, cell distribution from 5 to 14 dpf suggested two distinct

151    lineages: one involving chemokine-expressing stromal cells (*cxcl12a+/ccl25b*+) and a second

152    emanating from *hyal4*+ cells (Fig. 2c-e, Fig. S13). In situ hybridization at 14 dpf revealed broad

153    mesenchymal expression of *cxcl12a*, and expression of *hyal4* in perichondrium in a largely

154    complimentary pattern to *postnb* and *col10a1a* expression in periosteum (Fig. S9). Branches from

155    *hyal4*+ perichondrium led to periosteum, tendon and ligament cells, chondrocytes, and

156    osteoblasts, consistent with studies showing perichondrium to be the precursor of the periosteum

157    in endochondral bones[21,22].

158

159    For gill clusters, cell distribution from 5 to 14 dpf revealed two primary trajectories (Fig. 2f-h, Fig.

160    S13). In the first, *cxcl12a+/ccl25b*+ stromal cells give rise to mesenchyme associated with retinoic

161    acid metabolism (*aldh1a2+/rdh10a*+), with in situ hybridization revealing these cell types restricted

162    to the base of secondary filaments (Fig. S10). In the second, *fgf10a*+ cells are connected to

163    *fgf10b+* cells, which then diverge into gill cartilage, pillar, tunica media, and perivascular

164    populations. To test whether *fgf10b+* cells are progenitors for specialized gill subtypes, we used

165    CRISPR/Cas9 to insert a photoconvertible nuclear EOS protein into the endogenous *fgf10b* locus.

166    We found *fgf10b*:nEOS to be robustly expressed in the forming gills, with expression becoming

167    progressively restricted to the tips of gill filaments over time, similar to endogenous *fgf10b*

168    expression (Fig. S10, S14). We then used UV light to convert *fgf10b*:nEOS fluorescence from

169    green to red in a small number of filaments at 7 dpf and observed contribution to gill chondrocytes

170    and pillar cells 3 days later, with new *fgf10b*:nEOS cells (i.e. green only) being generated at the

171    tips of growing filaments (Fig. 2i). Similar results were seen in adult gill filaments (Fig. S14). These

172    data support *fgf10b+* cells being progenitors for gill-specific cell types from larval through adult

173    stages.

174

175    To understand how CNCC mesenchyme changes from embryogenesis to adulthood, we next

176    interrogated patterns of gene usage and chromatin accessibility (Fig. 2j, Fig. S15-16, Table S2).

177    Gene ontogeny (GO) analysis of ectomesenchyme at 1.5 and 2 dpf revealed terms linked to cell

178    division and metabolism, consistent with early expansion of this population. We also find

179    enrichment of transcription factors for early ectomesenchyme (*dlx2a*, *twist1a*, *nr2f6b*) and arch

180    patterning (*pou3f3b*, *hand2*), as well as transcription factor binding motifs for several types of

181    nuclear receptors, in accordance with known roles of Nr2f members in ectomesenchyme

182    development[19]. The *hyal4+* population contains skeletal-associated terms (collagen fibril

183    organization, skeletal system development, regulation of ossification, cartilage development),

184    consistent with being a common progenitor for cartilage, tendon, ligament, and bone in

185    pseudotime analysis. The *hyal4+* population is enriched for transcription factors implicated in

186    perichondrium biology (*mafa*, *foxp2*, *foxp4*)[23,24] and cartilage formation (*barx1*, *sox6*, *emx2*)[25-27],

187    and motifs for Bmp signaling (SMAD) and transcription factors (NFAT, RUNX) known to regulate

188    cartilage and bone[28]. For gill *fgf10a/b+* progenitors, we recover terms for general growth (e.g.

8

189 translation, cellular biosynthetic process), response to Fgf signaling, and respiratory system

190 development, consistent with lineage tracing showing *fgf10b:nEOS*-labeled cells giving rise to gill

191 respiratory cell types through adult stages. We also observe enrichment of *gata2a*, *gata3*, and

192 GATA motif accessibility, suggesting important roles of Gata factors in gill-specific lineages.

193

194 In contrast to *hyal4+* and *fgf10a/b+* populations that display hallmarks of progenitors, the *cxcl12a+*

195 stromal population is associated with terms for regeneration, response to injury and wounding,

196 negative regulation of the Wnt signaling pathway, and, particularly at adult stages, response to

197 stress and modulation of the immune response. This population is enriched for *osr1*, early

198 response genes of the Fos/Jun family, C/EBP family members implicated in response to

199 inflammation[29], and *egr1* that has recently been linked to injury-induced regenerative responses

200 across the animal kingdom[30]. Recovery of motifs for STAT and C/EBP also point to immune

201 system interactions. As *Cxcl12+* stromal cells in murine bone marrow have been shown to only

202 contribute to osteoblasts during bone regeneration[31], it will be interesting to test whether the

203 *cxcl12a+* stromal population in animals such as zebrafish that lack bone marrow also plays a role

204 in skeletal regeneration[32].

205

206 **Highly resolved embryonic spatial expression domains from integrated datasets**

207 We next sought to understand the developmental origins of distinct cell types and lineage

208 programs in CNCC ectomesenchyme. To do so, we first examined the ability of integrated

209 transcriptomic and chromatin accessibility datasets to predict the expression patterns of potential

210 ectomesenchyme patterning genes at 1.5 dpf, a stage before overt cell type differentiation.

211 Compared to scRNAseq (Fig. 3a) or snATACseq alone (Fig. S17), SnapATAC pseudo-multiome

212 analysis (Fig. 3b) was better able to separate CNCCs along the major positional axes, including

213 the dorsal-ventral axis and the anterior-posterior axis (frontonasal, mandibular (arch 1), hyoid

214 (arch 2), branchial (arch 3-6), and tooth-bearing (arch 7)).

215

216    Comparison of the predicted SnapATAC expression of known region-specific genes - *pou3f3b*

217    (dorsal arches 1 and 2), *dlx5a* (intermediate arches), *hand2* (ventral arches), *meis2b* (arch 7),

218    and *pitx1* (oral mandibular)[25,33,34] - revealed tight correlation to reported expression, including

219    zebrafish-specific overlap of *dlx5a* and *hand2* in the mandibular arch (Fig. 3d). We also identified

220    a previously unappreciated oral-aboral axis of the mandibular arch in zebrafish, marked by *pitx1*

221    and *nr5a2* respectively, which we validated by in situ hybridization for *nr5a2* (Fig. 3e). Re-

222    examination of genes identified from previous bulk RNA sequencing of zebrafish arches further

223    revealed strong correlation of SnapATAC domains with reported expression for 23 of 27 genes

224    (Fig. S18), with SnapATAC suggesting frontonasal and tooth-domain expression for two genes

225    previously annotated as false positives[25]. We also observed correlation of the transcription factor

226    binding motifs enriched in cluster-specific accessible chromatin with the activities of transcription

227    factors of the same family, including POU3F3, MEIS2, HAND, DLX5, PITX1, and NR5A2 (Fig.

228    3c,d). This approach shows the power of integrated scRNAseq and snATACseq data to predict

229    the spatial expression domains of the vast majority of CNCC ectomesenchyme genes at

230    pharyngeal arch stages.

231

232    **Chromatin accessibility predicts cell type competency in early arches**

233    We next sought to understand how the establishment of cell fate competency is linked to the

234    earlier activity of arch patterning genes. To do so, we first computed unique patterns of chromatin

235    accessibility ("peaks") for each cell cluster at 14 dpf (Fig. 4a, Table S3). Modules of the top

236    enriched peaks for each cell type were then mapped onto UMAP projections of SnapATAC data

237    at 1.5, 2, 3, and 5 dpf (Fig. S19). To understand when cluster-specific peaks become established,

238    as well as cluster relatedness, we developed the bioinformatics pipeline "Constellations". First,

239    we calculated whether projections of cluster-specific peak modules are skewed toward particular

240    regions of UMAP space at each earlier time-point, suggesting establishment of cluster-specific

10

241  chromatin accessibility (a proxy for cell type competency).  We then computed the relatedness of

242  peak module projections in two dimensions for each mapped cluster at each stage (Fig. 4b).

243  Analysis of cell competency trajectories shows that cell types can be grouped into five main

244  classes: skeletogenic cells (including *hyal4+* perichondral and *postnb+* periosteal cells), stromal

245  cells, dermal fibroblasts, gill cell types, and cartilage. Constellations analysis also reveals a

246  temporal order of cell type competency establishment, with unique chromatin accessibility for

247  cartilage and dermal fibroblast lineages emerging at 1.5 dpf; bone and perichondrium at 2 dpf;

248  and periosteum, tendon and ligament, and gill progenitors and pillar cells at 3 dpf (Fig. 4c). This

249  analysis suggests that chromatin accessibility prefiguring diverse CNCC cell types is

250  progressively established rather than being inherited from earlier multipotent CNCCs.

251

252  **Constellations analysis reveals candidate transcription factors for lineage priming**

253  To discover potential transcription factors for establishing cell type competency, we analyzed the

254  Constellations dataset for transcription factors whose expression and predicted binding motifs

255  were co-enriched in particular clusters. We identified 287 transcription factor expression/motif

256  pairs showing enrichment (Fig. S20, Table S4). The FOXC1 motif and *foxc1b* gene body activity

257  were highly enriched in the cartilage trajectory, and LEF1/*lef1* in the dermal fibroblast trajectory

258  (Fig. 5a). Projection of FOX motifs and merged Fox gene activity (*foxc1a*, *foxc1b*, *foxf1*, *foxf2a*,

259  *foxf2b*) and LEF1/*lef1* onto SnapATAC UMAPs at 1.5 dpf reveals close correlation to mapping of

260  the 14 dpf peak modules for cartilage and dermal fibroblasts at this stage (Fig. 5b,c), as well as

261  the known fate map of cartilage precursors in the arches[35] (Fig. 5d,e). This confirms genetic

262  evidence for roles of Foxc1 and Foxf1/2 in cartilage formation in zebrafish and mouse[36,37], and

263  more specifically Foxc1 in establishing accessibility of cartilage enhancers in the developing

264  face[28]. It also raises the possibility that Wnt signaling, mediated in part by Lef1, may play a role

265  in early dermal fibroblast specification, consistent with enrichment of *wnt5a* in this population (Fig.

266  S11).

11

267

268  We also find GATA3/*gata3* to be highly enriched in gill populations, with SnapATAC UMAP

269  projections of GATA3 motif and *gata3* gene body activity at 5 dpf correlating with 14 dpf gill

270  progenitor peaks (Fig. 5f). The enrichment of ETS2/*ets2*, which plays a role in endothelial

271  development[38], in the gill pillar trajectory is consistent with ETS factors driving a mesenchyme-to-

272  endothelia transition during formation of these vascular cells. Skeletogenic trajectories are

273  uniquely marked by IRF8/*irf8*. Whereas loss of bone in mouse *Irf8*[-/-] mutants has been attributed

274  to increased osteoclastogenesis[39], our analysis suggests that Irf8 may also have an early function

275  in priming the skeletal lineage. Enrichment of CEBPA/*cebpa* in stromal trajectories may reflect

276  the immunomodulatory role of this mesenchymal population[29]. These findings show the power of

277  Constellations analysis to reveal potential factors for establishing regional chromatin accessibility

278  important for later cell type differentiation.

279

280  **Gill-specific lineages distinguished by early Gata3 activity**

281  Given the selective enrichment of GATA3 motifs and *gata3* activity in gill lineages, we further

282  investigated the presence of a Gata3 regulatory circuit directing CNCCs to gill fates. Whereas

283  previous work had shown that *gata3* is expressed in and required for initial gill bud formation in

284  zebrafish, larval lethality had precluded analysis of gill subtype differentiation[40]. We find *gata3*

285  expression to be maintained in gill populations through adult stages in scRNAseq data, which we

286  validated by in situ hybridization at 14 dpf and 2 years of age (Fig. S21). We then identified a non-

287  coding region ~143kb downstream of the *gata3* gene, itself containing a predicted GATA3 binding

288  site, that was selectively accessible in posterior arch CNCCs by 3 dpf, gill progenitors and pillar

289  cells by 5 dpf, and gill cartilage cells by 14 dpf (Fig. 6a, Fig. S22). This *gata3-P1* element was

290  sufficient to drive highly restricted GFP expression in posterior arch CNCCs starting at 1.5 dpf,

291  which continued in gill progenitors, pillar cells, and chondrocytes through 60 dpf (Fig. 6c-e, Fig.

292  S21).

293

294    Gill cartilage has a markedly distinct expression and chromatin accessibility profile from hyaline

295    cartilage of the jaw, as shown by selective expression of *ucmaa* in gill cartilage versus *ucmab* in

296    hyaline cartilage (Fig. S23). We identified a non-coding region ~5kb upstream of the *ucmaa* gene

297    that was selectively accessible in gill cartilage starting at 14 dpf and contained a predicted GATA3-

298    binding site (Fig. 6b, Fig. S22). This *ucmaa-P1* element drives highly restricted GFP expression

299    in gill chondrocytes at 11 and 23 dpf, in contrast to a previously described *ucmab* enhancer[28]

300    driving GFP expression in hyaline but not gill cartilage (Fig. 6f, Fig. S23). Although functional

301    assays are needed to confirm Gata3 dependence, our findings are consistent with GATA factors

302    establishing a positive autoregulatory circuit in posterior arch CNCCs that maintains *gata3*

303    expression and promotes the later differentiation of gill-specific cell types (Fig. 6g).

304

305    **Conclusion**

306    Integration of transcriptome and chromatin accessibility data of the CNCC lineage has allowed us

307    to connect patterning along major development axes to the emergence of the wide diversity of

308    CNCC-derived cell types. Rather than lineage-specific chromatin accessibility being an intrinsic

309    property of CNCCs, our Constellations analysis points to the progressive remodeling of chromatin

310    underlying diverse cell type differentiation. Roles for inductive signaling in establishing enhancer

311    accessibility in post-migratory arch CNCCs would help explain reports of mesodermal cells

312    contributing to classically considered CNCC structures such as the skate gill skeleton[7]. Further,

313    retrograde mapping of cell type-specific chromatin accessibility, combined with our highly

314    resolved atlas of pharyngeal arch gene expression, reveals candidate transcription factors priming

315    distinct CNCC lineages. Consistent with recent reports of organ-specific fibroblast

316    heterogeneity[41], we also uncover a CNCC-derived dermal fibroblast population characterized by

317    Phe/Tyr metabolism genes, which may be induced by early Wnt/Lef1 activity. Expression of some

318    of the same Phe/Tyr catabolic genes is observed in a subset of axolotl limb fibroblasts[42], with the

319    blackening of bone and cartilage in Alkaptonuria patients defective in Phe/Tyr breakdown

320    suggesting general roles for these specialized dermal fibroblasts in protecting the skeleton. In the

321    gill region, we identify a *fgf10*-expressing progenitor population characterized by sustained Gata3

322    activity, with later emergence of Ets2 activity in pillar cells providing a potential mechanism for the

323    mesenchyme-to-endothelia transition of these specialized vascular cells. The presence of a

324    similar Fgf10-expressing mesenchyme population in the mammalian lung[43] raises the possibility

325    that an ancestral CNCC-derived gill respiratory program may have been co-opted by the

326    mesoderm during later lung evolution. Single-cell profiling of transcriptome and chromatin

327    accessibility across time thus provides a blueprint for understanding the diversification and post-

328    embryonic production of the huge variety of CNCC-derived cell types throughout the head.

329

## 330    Materials and methods

### 331    Zebrafish lines

332    The Institutional Animal Care and Use Committee of the University of Southern California

333    approved all animal experiments (Protocol 20771). Published lines include *Tg(Mmu.Sox10-*

334    *Mmu.Fos:Cre)$^{zf384}$* [8]; *Tg(actab2:loxP-BFP-STOP-loxP-dsRed)$^{sd27}$* [44]; and *Tg(ucmab_p1:GFP)$^{el806}$,*

335    *Tg(fli1a:eGFP)$^{y1}$*, and *Tg(sox10:DsRedExpress)$^{el10}$* [28]. Five transgenic lines were generated as

336    part of this study: *Tg(fgf10b:nEOS)$^{el865}$*, *Tg(gata3_p1:GFP)$^{el857}$*, *Tg(gata3_p1:GFP)$^{e,858}$*,

337    *Tg(ucmaa_p1:GFP)$^{el851}$* and *Tg(ucmaa_p1:GFP)$^{el854}$*. The *fgf10b:nEOS* knock-in line was made

338    using CRISPR/Cas9-based integration [45]. Three gRNAs targeting sequences upstream of the

339    *fgf10b*      translational      start      site      (5'-CATGATAACCCTTCCTAGAT-3',      5'-

340    GAGCTCTTTGATAGCGGGCT-3', and 5'-GTTGAGCAGCATGTCCCATG-'3) were co-injected at

341    100 ng/uL into wild-type embryos with Cas9 RNA (100 ng/uL), an mbait-NLS-EOS plasmid (20

342    ng/uL) [46], and the published gRNA targeting the mbait sequence [45] to linearize the plasmid. A

343    germline founder was identified based on nEOS fluorescence in the progeny of injected animals.

344    For enhancer transgenic lines, we synthesized peaks for *gata3* (chr4:24918100-24918770) and

14

345    *ucmaa* (chr4:7836670-783720) using iDT gBlocks and cloned these into a modified pDest2AB2

346    construct containing E1b minimal promoter, GFP, and an eye-CFP selectable marker [28] using In-

347    Fusion cloning (Takara Bio). We injected plasmids and Tol2 transposase RNA (30 ng/uL each)

348    into one-cell stage zebrafish embryos, raised these animals, and screened for founders based on

349    eye CFP expression in the progeny. Two independent germline founders were identified for each

350    that showed similarly specific activity in the gills.

351

352    **In situ hybridization and immunohistochemistry**

353    All samples were prepared by fixation in 4% paraformaldehyde and embedded in paraffin, with

354    decalcification for one week in 20% EDTA if over 14 dpf. All in situ patterns were confirmed in at

355    least 3 independent animals. RNAscope probes were synthesized by Advanced Cell Diagnostics

356    in channels 1 through 4. Channel 1 probes: *ifitm5, ucmaa, col10a1a*. Channel 2 probes: *postnb,*

357    *myh11a, cxcl12a, sp7, gata3*. Channel 3 probes: *pah, lum, fgf10b, sox9a.* Channel 4 probe: *hyal4,*

358    *acta2, ncam3*. Paraformaldehyde-fixed paraffin-embedded sections were deparaffinized, and the

359    RNAscope Fluorescent Multiplex V2 Assay was performed according to manufacturer's protocols

360    with the ACD HybEZ Hybridization oven. Colorimetric in situ hybridization was performed as

361    described [32]. The *hpdb* riboprobe was generated by cloning a purchased gBlock fragment

362    designed from transcript hpdb-201 using nucleotides 679-1395 (tggatga...gactccc) into pCR-

363    BluntII-TOPO (Life Technologies). The *nr5a2* riboprobe was generated by PCR amplification of

364    cDNA with primers 5'-ATGGGGAACAGGGGCATATG-3' and 5'-AGGGGTCGGGATACTCTGAT-3',

365    the *ucmaa* riboprobe with primers 5'-TGGTACCAGCTCAAGACACT-3' and 5'-

366    ATAGTACTGGCGGTGGTGAG-3', the *ucmab* riboprobe with primers 5'-

367    ATGTCCTGGACTCAACCTGC-3' and 5'-GTTATCTCCCAGCGTGTCCA-3', and the

368    *thbs4a* riboprobe with primers 5'-CCCATGTTTCTTCGGTGTGA-3' and 5'-

369    GGTTTGGTACCAGCCTACAG-3'. Amplified products were cloned into pCR-BluntII-TOPO.

370    pCR-BluntII-TOPO plasmids were linearized by restriction digest (enzyme dependent on direction

371    of blunt insertion), and RNA probe was synthesized using either T7 or Sp6 polymerase (Roche)

372    depending on direction of blunt insertion. Immunohistochemistry for dsRed was performed with a

373    7 minute -20°C 100% acetone target retrieval and blocking in 2% normal goat serum (Jackson

374    ImmunoResearch, cat. no. 005-000-121). Primary antibodies include rabbit anti-mCherry (1:100,

375    Rockland Immunochemicals, cat. no. RL600-401-P16) and rabbit anti-mCherry (1:100, Novus

376    Biologicals, cat. no. NBP2-25157) used at the same time. The secondary antibody was goat anti-

377    rabbit Alexa Fluor 546.

378

379    **Imaging**

380    Confocal images of whole-mount or section fluorescent in situ hybridizations and live images of

381    transgenic fish were captured on a Zeiss LSM800 microscope using ZEN software. Colorimetric

382    in situs were imaged on a Zeiss AxioScan Z.1 For *fgf10:nEOS* experiments, we used the ROI

383    function on the confocal microscope to specifically convert nEOS-expressing cells in the gill

384    filaments of live animals using targeted UV irradiation, prior to the emergence of gill filament

385    cartilage. At the specified days post-conversion, we euthanized the animal, fixed it in 4% PFA for

386    1 hour, and dissected the gill arches. We stained the gills with DRAQ5 nuclear dye (Abcam) for

387    30 min and imaged at 40X to locate converted cells. For all transgenic imaging experiments,

388    expression patterns were confirmed in at least 5 independent animals.

389

390    **Single-cell analysis and statistics**

391    *scRNAseq and snATACseq library preparation and alignment*

392    Dissected heads from converted *Sox10:Cre; actab2:loxP-BFP-STOP-loxP-dsRed* fish were

393    incubated in fresh Ringer's solution 5-10 min, followed by mechanical and enzymatic dissociation

394    by pipetting every 5 minutes in protease solution (0.25% trypsin (Life Technologies, 15090-046),

395 1 mM EDTA, and 400 mg/ml Collagenase D (Sigma, 11088882001) in PBS) and incubated at

396 28.5°C for 20-30 minutes or until full dissociation. Reaction was stopped by 6X stop solution (6

397 mM CaCl2and 30% fetal bovine serum (FBS) in PBS). Cells were pelleted (2000 rpm, 5 min, 4

398 °C) and resuspended in suspension media (1% FBS, 0.8 mM CaCl2, 50 U/ml penicillin, and 0.05

399 mg/ml streptomycin (Sigma-Aldrich, St. Louis, MO) in phenol red-free Leibovitz's L15 medium

400 (Life Technologies)) twice. Final volumes of 500 µl resuspended cells were placed on ice and

401 sorted by fluorescence activated cell sorting (FACS) to isolate live cells that excluded the

402 cytoplasmic stain Zombie green (BioLegend, 423111). For scRNAseq library construction,

403 barcoded single-cell cDNA libraries were synthesized using 10X Genomics Chromium Single Cell

404 3' Library and Gel Bead Kit v.2 per manufacturer's instructions. Libraries were sequenced on

405 Illumina NextSeq or HiSeq machine at a depth of at least 1,000,000 reads per cell for each library.

406 Read2 was extended to 126 cycles for higher coverage. Cellranger v3.0.0 (10X Genomics) was

407 used for alignment against GRCz11 (built with GRCz11.fa and GRCz11.98.gtf) and gene-by-cell

408 count matrix generation with default parameters.

409 For snATACseq library construction, we used the same cell dissociation and sorting protocol as

410 for scRNAseq, with isolation of live cells that excluded the cytoplasmic stain Zombie green

411 (BioLegend, 423111) and collected live cells in 0.04% BSA/PBS. Nuclei isolation was performed

412 per manufacturer's instructions (10X Genomic, protocol CG000169). Cells were incubated with

413 lysis buffer on ice for 90 s, followed by integrity check of nuclei under fluorescence microscope

414 with DAPI before library synthesis. Barcoded single-nuclei ATAC libraries were synthesized using

415 10X Genomics Chromium Single Cell ATAC Reagent Kit v1.1 per manufacturer's instructions.

416 Libraries were sequenced on Illumina NextSeq or HiSeq machine at a depth of at least 75,000

417 reads per nucleus for each library. Both read1 and read2 were extended to 65 cycles. Cellranger

418 ATAC v1.2.0 (10X Genomics) was used for alignment against genome (built with GRCz11.fa,

419 JASPAR2020, and GRCz11.98.gtf), peak calling, and peak-by-cell count matrix generation with

420 default parameters.

421 We included biological replicates at several stages to test the reproducibility of library preparation

422 and increase depth of data. For scRNAseq, we performed two replicates at 5 and 14 dpf, and

423  three replicates at 3 and 150 dpf. For snATACseq, we performed two replicates at 2, 3, and 14

424  dpf.

425

426  *SnapATAC for peak refinement and gene activity matrix imputation*

427  To refine the peak profile for better representation of diverse cell types across libraries, we

428  performed a second round of peak calling using package Snaptools (v1.2.7) and SnapATAC

429  (v1.0.0) [15]. We first removed low-quality cell and cell doublets by setting cutoffs based on

430  percentage of reads in peaks (> 30 for 60 dpf, > 45 for 210 dpf, and > 50 for the rest) and fragment

431  number within peaks (5,000 – 30,000 for 5 dpf, 1,000 – 11,000 for 14 dpf, and 1,000 – 20,000 for

432  the rest). Potential cell debris or low-quality cells were removed by setting hard fragments-in-peak

433  number cutoffs. Using the SnapATAC package, we then generated "pseudo-multiome" data at

434  each stage. To recover every aligned fragment, we binned the genome into 5 kb sections and

435  constructed the bin-by-cell matrices (bmats) for each library by Snaptools from the positional-

436  sorted bam files generated by Cellranger ATAC v1.2.0. The cells were filtered, dimensionally

437  reduced by diffusion map, and clustered with inputs of the first 34 dimensions followed the

438  SnapATAC vignette. The specific peaks were called for each cluster by the wrapped MACS2

439  function in SnapATAC with parameter gsize = 1.5e9, shift = 100, ext = 200, and qval = 5e-2. The

440  finalized and refined peak profile was derived by collapsing and merging all 175 individual peak

441  files to 445,307 peaks. To impute the gene activity with the corresponding scRNAseq data, the

442  bmats of each time point were used to calculate gene-activity-by-cell matrices (gmats) by

443  SnapATAC. The gmats were used to find anchors within the scRNAseq data at the same time

444  point by Seurat. We then transferred the expression data from scRNAseq through the anchors to

445  derive the imputed gene-activity-by-cell matrices for each time point.

446  *Data processing of scRNAseq and snATACseq*

447  The count matrices of both scRNAseq and snATACseq data were analyzed by R package Seurat

448  (v3.2.3) and Signac (v1.0.0). The count matrices of each sample were aggregated where

449  replicates were available. For scRNAseq data, the matrices are normalized (NormalizeData) and

450    scaled for the top 2,000 variable genes (FindVariableFeatures and ScaleData). The scaled

451    matrices were dimensionally reduced to 50 principal components (60 components for 150 dpf),

452    and then subjected to neighbor finding (FindNeighbors, k = 20) and clustering (FindClusters,

453    resolution = 0.8). The data were visualized through UMAP with 50 principal components as input.

454    For snATACseq data, the matrices are dimensionally reduced to 30 latent semantic indices (LSIs)

455    through RunTFIDF and RunSVD functions. The neighbor finding, clustering, and visualization are

456    performed as for scRNAseq data (algorithm = 3 for FindClusters) with input of the second to

457    thirtieth LSIs. To calculate the motif accessibility, the enrichment of motifs in JASPAR2020 [47] was

458    calculated by chromVAR [48] through function RunChromVar. To test the enriched genes and gene

459    activities in both scRNAseq and snATACseq data as shown in Table S1, two-sided likelihood-

460    ratio test is performed through FindAllMarkers function (min.pct = 0.25) with cutoff of adjusted p

461    value smaller than 0.001.

462

463    *STITCH network construction and force directed layout*

464    To identify the overall cell trajectories in our scRNAseq and snATACseq data, we used the

465    STITCH algorithm [17] to construct cell neighbor networks. As the dimensional reduction space of

466    snATACseq data are LSIs, we modified the *stitch_get_knn* and *stitch_get_link* function of the

467    STITCH package to make it compatible to LSI. For *stitch_get_knn* function of snATACseq data,

468    we used the LSI matrix to find the k nearest neighbor of each cell for each time point. For

469    *stitch_get_link* function of snATACseq data, we projected the LSI space of time point *t* to time

470    point *t-1* by solving the right orthogonal matrix of the singular vector decomposition (SVD) for *t-1*.

471    SVD ($M = U\textstyle\sum V^T$) is the initial step of latent semantic analysis where *M* is the peak-by-cell matrix

472    and *U* will be later transformed to LSI. For *t* and *t-1*, we have $M_t = U_t \textstyle\sum_t V^T_t$ and $M_{t-1} = U_{t-1} \textstyle\sum_{t-1} V^T_{t-1}$. To

473    project the space from *t* to *t-1*, we derived a projected $U_{t-1}$ as $U^p_{t-1}$ through solving the equation

474    $M_{t-1} = U^p_{t-1} \textstyle\sum_t V^T_t$. Both $U_t$ and $U^p_{t-1}$ were further combined, normalized, and subjected to the default

475    neighbor finding as *stitch_get_link*. To visualize the STITCH networks of both scRNAseq and

476    snATACseq data, we used the force directed layout by ForceAtlas2 in Gephi (v0.9.2) to derive

477    the 2-dimensional layouts.

478

479    *Pseudotime analysis*

480    We used the R package monocle3 (v0.2.3.0) to predict the pseudotemporal relationships within

481    skeletogenic or gill populations. We first merged 5 and 14 dpf scRNAseq data, including an

482    additional scRNAseq library of *sox10>dsRed+* cells sorted from the dissected ceratohyal

483    endochondral bone at 14 dpf (to further enrich skeletogenic populations), and performed

484    clustering and dimensionality reduction. After removing dermal fibroblast (*pah+*) and teeth

485    (*spock3+*) populations, we placed *hoxb3a+/gata3+* cells into a "gill" cluster and all other cells into

486    a "jaws" cluster. Cell paths were predicted by the *learn_graph* function of monocle3. We set the

487    origin of the cell paths based on the enriched distribution of 5 dpf cells.

488

489    *Gene ontology, motif family, and TF analysis of CNCC mesenchyme*

490    Analysis was performed on ectomesenchyme, perichondrium, gill progenitor, and stromal

491    populations at each stage based on markers from the scRNAseq data (Fig. S1-7). The enriched

492    genes of each cluster are tested by running a two-sided Wilcoxon rank sum test against all other

493    clusters using an adjusted p value ≤ 0.001. These enriched gene sets are subjected to gene

494    ontology analysis for terms of biological processes (BP) by R package ViSEAGO (v1.2.0) [49]. The

495    heatmap is generated by GOterms_heatmap function using values of log10(adjusted p.value). To

496    generate the heatmap of motif families for each cluster, we first averaged and aggregated the

497    motif accessibilities for each cell according to the motif family by TRANSFAC [50]. The means of

498    each motif family are used for the heatmap. To generate the heatmap of TFs for each cluster, we

499    subsetted the TFs from the enriched gene sets, and used the mean of each TF for every cluster

500    for the heatmap.

501

502    *Constellations analysis and calculation of cluster skewedness and correlation*

20

503   The tissue module scores of the snATACseq data were calculated based on the enriched peak

504   sets and their module scores for each cluster identified at 14 dpf by R packages Seurat and

505   Signac. The enriched peak sets were calculated by the FindAllMarkers function using two-sided

506   likelihood ratio test with fragment numbers in peak region as latent variables. We used the peaks

507   with adjusted p values smaller than 0.001 as the enriched peaks for a cluster. As there are 23

508   clusters (tissues) identified at 14 dpf, we ended up with 23 peak sets, which we applied to

509   calculate the tissue module scores to earlier time points (1.5, 2, 3, and 5 dpf) using

510   AddChromatinModule function. To determine whether a tissue score at a time point distributes in

511   a statistically significant, and hence biologically interesting, way, we calculated the skewedness

512   of the distribution of a tissue score by the R package parameter (v0.12.0). We considered a tissue

513   score to be distributed in a meaningful way if it was strongly right skewed by a hard cutoff of

514   skewedness greater than 1. For 1.5 dpf, the cutoff of skewedness was lowered to 0.4 to

515   accommodate overall lower skewedness at that time point, but with additional filter of max module

516   score > 15 to avoid tissue module scores with extremely low values.

517   To profile the relationship of all tissue scores, we constructed a distance matrix of all 23 tissue

518   module scores across all the time points (1.5, 2, 3, 5, and 14 dpf). For the distance $D$ between

519   score of tissue $A$ at time point $t_1$ and score of tissue $B$ at time point $t_2$, the distance $D$ can be

520   described as $D = D(\text{tissue}) + (a \times D(\text{time point}))$. $D(\text{tissue})$ stands for the distance between the

521   score of tissue $A$ and $B$ by averaging the Euclidian distance between score $A$ and $B$ at time point

522   $t_1$ and the Euclidian distance between score $A$ and $B$ at time point $t_2$. $D(\text{time point})$ stands for the

523   distance between time point $t_1$ and $t_2$ derived by the distance between the dendrogram of all the

524   tissue scores at $t_1$ and the dendrogram at $t_2$. Since $D(\text{time point})$ is relatively smaller than

525   $D(\text{tissue})$, we multiply $D(\text{time point})$ by 8 to make the distance between time points comparable

526   to the distance between tissue scores. The distance matrix was dimensionally reduced and

527   visualized by UMAP.

528   To detect the potential factors that contribute to the patterning of tissue-specific peaks, we

529   performed linear regression of each tissue module score against all motif accessibilities and the

530   gene activities of transcription factors. We used ZFIN and JASPAR2020, converted by homology

21

531    data from MGI, to build up a list of transcription factors in zebrafish. We then curated and paired

532    every motif in JASPAR2020 with its potential binding transcription factors. The coefficients of

533    regression results were used as indications of whether a motif or transcription factor is positively

534    correlated with a tissue module score with upper cutoff of adjusted p value 0.05. We transformed

535    the coefficients of all the negative related motifs and transcription factors to 0 to filter out irrelevant

536    motifs and transcription factors. To visualize the correlation of each pair of motif and transcription

537    factor, we plotted the coefficient magnitudes of motifs by dot sizes and transcription factor gene

538    body activities by a red color scale on Constellation maps.

539

## Acknowledgments

544

## Competing interests

546    No competing interest declared.

547

## Funding

551

## Data availability

553    The data that support the findings of this study are available from the corresponding author upon

554    reasonable request. Processed and raw sequencing data have been deposited at GEO as

555    GSE178969.

## Author Contributions

556

557 P.F., K.-C.T., M.T., C.A., H.-J.C., J.S., N.N., and J.G.C. performed the experiments. J.G.C.

558 oversaw the project and wrote the manuscript.

559

## References

560

561 1 Platt, J. B. Ectodermic origin of the cartilages of the head. *Anat. Anz.* **8**, 506-509 (1893).

562 2 Mongera, A. *et al.* Genetic lineage labeling in zebrafish uncovers novel neural crest
563 contributions to the head, including gill pillar cells. *Development* **140**, 916-925,
564 doi:10.1242/dev.091066 (2013).

565 3 Dupin, E., Calloni, G. W., Coelho-Aguiar, J. M. & Le Douarin, N. M. The issue of the
566 multipotency of the neural crest cells. *Dev Biol* **444 Suppl 1**, S47-S59,
567 doi:10.1016/j.ydbio.2018.03.024 (2018).

568 4 Baroffio, A., Dupin, E. & Le Douarin, N. M. Common precursors for neural and mesectodermal
569 derivatives in the cephalic neural crest. *Development* **112**, 301-305 (1991).

570 5 McGonnell, I. M. & Graham, A. Trunk neural crest has skeletogenic potential. *Curr Biol* **12**,
571 767-771, doi:10.1016/s0960-9822(02)00818-7 (2002).

572 6 Simoes-Costa, M. & Bronner, M. E. Reprogramming of avian neural crest axial identity and
573 cell fate. *Science* **352**, 1570-1573, doi:10.1126/science.aaf2729 (2016).

574 7 Sleight, V. A. & Gillis, J. A. Embryonic origin and serial homology of gill arches and paired fins
575 in the skate, Leucoraja erinacea. *Elife* **9**, doi:10.7554/eLife.60635 (2020).

576 8 Kague, E. *et al.* Skeletogenic fate of zebrafish cranial and trunk neural crest. *PLoS One* **7**,
577 e47394, doi:10.1371/journal.pone.0047394 (2012).

578 9 Mitchell, J. M. *et al.* The alx3 gene shapes the zebrafish neurocranium by regulating
579 frontonasal neural crest cell differentiation timing. *Development*, doi:10.1242/dev.197483
580 (2021).

581 10 Williams, R. M. *et al.* Reconstruction of the Global Neural Crest Gene Regulatory Network In
582 Vivo. *Dev Cell* **51**, 255-276 e257, doi:10.1016/j.devcel.2019.10.003 (2019).

583 11 Soldatov, R. *et al.* Spatiotemporal structure of cell fate decisions in murine neural crest.
584 *Science* **364**, doi:10.1126/science.aas9536 (2019).

585 12 Xu, J. *et al.* Hedgehog signaling patterns the oral-aboral axis of the mandibular arch. *Elife* **8**,
586 doi:10.7554/eLife.40315 (2019).

587 13 Yuan, Y. *et al.* Spatiotemporal cellular movement and fate decisions during first pharyngeal
588 arch morphogenesis. *Sci Adv* **6**, doi:10.1126/sciadv.abb0119 (2020).

589 14 Zalc, A. *et al.* Reactivation of the pluripotency program precedes formation of the cranial
590 neural crest. *Science* **371**, doi:10.1126/science.abb4776 (2021).

591 15 Fang, R. *et al.* Comprehensive analysis of single cell ATAC-seq data with SnapATAC. *Nat
592 Commun* **12**, 1337, doi:10.1038/s41467-021-21583-9 (2021).

593 16 Fernandez-Canon, J. M. *et al.* The molecular basis of alkaptonuria. *Nat Genet* **14**, 19-24,
594 doi:10.1038/ng0996-19 (1996).

595 17 Wagner, D. E. *et al.* Single-cell mapping of gene expression landscapes and lineage in the
596 zebrafish embryo. *Science* **360**, 981-987, doi:10.1126/science.aar4362 (2018).

597 18 Tang, W., Martik, M. L., Li, Y. & Bronner, M. E. Cardiac neural crest contributes to
598 cardiomyocytes in amniotes and heart regeneration in zebrafish. *Elife* **8**,
599 doi:10.7554/eLife.47929 (2019).

600 19 Barske, L. *et al.* Essential Role of Nr2f Nuclear Receptors in Patterning the Vertebrate Upper
601 Jaw. *Dev Cell* **44**, 337-347 e335, doi:10.1016/j.devcel.2017.12.022 (2018).

20  Cao, J. *et al.* The single-cell transcriptional landscape of mammalian organogenesis. *Nature* **566**, 496-502, doi:10.1038/s41586-019-0969-x (2019).

21  Bandyopadhyay, A., Kubilus, J. K., Crochiere, M. L., Linsenmayer, T. F. & Tabin, C. J. Identification of unique molecular subdomains in the perichondrium and periosteum and their role in regulating gene expression in the underlying chondrocytes. *Dev Biol* **321**, 162-174, doi:10.1016/j.ydbio.2008.06.012 (2008).

22  Colnot, C., Lu, C., Hu, D. & Helms, J. A. Distinguishing the contributions of the perichondrium, cartilage, and vascular endothelium to skeletal development. *Dev Biol* **269**, 55-69, doi:10.1016/j.ydbio.2004.01.011 (2004).

23  Huang, W., Lu, N., Eberspaecher, H. & De Crombrugghe, B. A new long form of c-Maf cooperates with Sox9 to activate the type II collagen gene. *J Biol Chem* **277**, 50668-50675, doi:10.1074/jbc.M206544200 (2002).

24  Zhao, H. *et al.* Foxp1/2/4 regulate endochondral ossification as a suppresser complex. *Dev Biol* **398**, 242-254, doi:10.1016/j.ydbio.2014.12.007 (2015).

25  Askary, A. *et al.* Genome-wide analysis of facial skeletal regionalization in zebrafish. *Development* **144**, 2994-3005, doi:10.1242/dev.151712 (2017).

26  Lefebvre, V., Li, P. & de Crombrugghe, B. A new long form of Sox5 (L-Sox5), Sox6 and Sox9 are coexpressed in chondrogenesis and cooperatively activate the type II collagen gene. *EMBO J* **17**, 5718-5733, doi:10.1093/emboj/17.19.5718 (1998).

27  Nichols, J. T., Pan, L., Moens, C. B. & Kimmel, C. B. barx1 represses joints and promotes cartilage in the craniofacial skeleton. *Development* **140**, 2765-2775, doi:10.1242/dev.090639 (2013).

28  Xu, P. *et al.* Foxc1 establishes enhancer accessibility for craniofacial cartilage differentiation. *Elife* **10**, doi:10.7554/eLife.63595 (2021).

29  Poli, V. The role of C/EBP isoforms in the control of inflammatory and native immunity functions. *J Biol Chem* **273**, 29279-29282, doi:10.1074/jbc.273.45.29279 (1998).

30  Gehrke, A. R. *et al.* Acoel genome reveals the regulatory landscape of whole-body regeneration. *Science* **363**, doi:10.1126/science.aau6173 (2019).

31  Matsushita, Y. *et al.* A Wnt-mediated transformation of the bone marrow stromal cell identity orchestrates skeletal regeneration. *Nat Commun* **11**, 332, doi:10.1038/s41467-019-14029-w (2020).

32  Paul, S. *et al.* Ihha induces hybrid cartilage-bone cells during zebrafish jawbone regeneration. *Development* **143**, 2066-2076, doi:10.1242/dev.131292 (2016).

33  Barske, L. *et al.* Evolution of vertebrate gill covers via shifts in an ancient Pou3f3 enhancer. *Proc Natl Acad Sci U S A* **117**, 24876-24884, doi:10.1073/pnas.2011531117 (2020).

34  Bessa, J. *et al.* meis1 regulates cyclin D1 and c-myc expression, and controls the proliferation of the multipotent cells in the early developing zebrafish eye. *Development* **135**, 799-803, doi:10.1242/dev.011932 (2008).

35  Crump, J. G., Swartz, M. E., Eberhart, J. K. & Kimmel, C. B. Moz-dependent Hox expression controls segment-specific fate maps of skeletal precursors in the face. *Development* **133**, 2661-2669, doi:10.1242/dev.02435 (2006).

36  Hong, H. K., Lass, J. H. & Chakravarti, A. Pleiotropic skeletal and ocular phenotypes of the mouse mutation congenital hydrocephalus (ch/Mf1) arise from a winged helix/forkhead transcriptionfactor gene. *Human molecular genetics* **8**, 625-637 (1999).

37  Xu, P. *et al.* Fox proteins are modular competency factors for facial cartilage and tooth specification. *Development* **145**, doi:10.1242/dev.165498 (2018).

38  Wei, G. *et al.* Ets1 and Ets2 are required for endothelial cell survival during embryonic angiogenesis. *Blood* **114**, 1123-1130, doi:10.1182/blood-2009-03-211391 (2009).

39  Zhao, B. *et al.* Interferon regulatory factor-8 regulates bone metabolism by suppressing osteoclastogenesis. *Nat Med* **15**, 1066-1071, doi:10.1038/nm.2007 (2009).

40  Sheehan-Rooney, K., Swartz, M. E., Zhao, F., Liu, D. & Eberhart, J. K. Ahsa1 and Hsp90 activity confers more severe craniofacial phenotypes in a zebrafish model of

hypoparathyroidism, sensorineural deafness and renal dysplasia (HDR). *Dis Model Mech* **6**, 1285-1291, doi:10.1242/dmm.011965 (2013).

41 Muhl, L. *et al.* Single-cell analysis uncovers fibroblast heterogeneity and criteria for fibroblast and mural cell identification and discrimination. *Nat Commun* **11**, 3953, doi:10.1038/s41467-020-17740-1 (2020).

42 Gerber, T. *et al.* Single-cell analysis uncovers convergence of cell identities during axolotl limb regeneration. *Science* **362**, doi:10.1126/science.aaq0681 (2018).

43 Malpel, S., Mendelsohn, C. & Cardoso, W. V. Regulation of retinoic acid signaling during lung morphogenesis. *Development* **127**, 3057-3067 (2000).

44 Kobayashi, I. *et al.* Jam1a-Jam2a interactions regulate haematopoietic stem cell fate through Notch signalling. *Nature* **512**, 319-323, doi:10.1038/nature13623 (2014).

45 Kimura, Y., Hisano, Y., Kawahara, A. & Higashijima, S. Efficient generation of knock-in transgenic zebrafish carrying reporter/driver genes by CRISPR/Cas9-mediated genome engineering. *Sci Rep* **4**, 6545, doi:10.1038/srep06545 (2014).

46 Thomas, E. D. & Raible, D. W. Distinct progenitor populations mediate regeneration in the zebrafish lateral line. *Elife* **8**, doi:10.7554/eLife.43736 (2019).

47 Fornes, O. *et al.* JASPAR 2020: update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res* **48**, D87-D92, doi:10.1093/nar/gkz1001 (2020).

48 Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat Methods* **14**, 975-978, doi:10.1038/nmeth.4401 (2017).

49 Brionne, A., Juanchich, A. & Hennequet-Antier, C. ViSEAGO: a Bioconductor package for clustering biological functions using Gene Ontology and semantic similarity. *BioData Min* **12**, 16, doi:10.1186/s13040-019-0204-1 (2019).

50 Wingender, E., Dietze, P., Karas, H. & Knuppel, R. TRANSFAC: a database on transcription factors and their DNA binding sites. *Nucleic Acids Res* **24**, 238-241, doi:10.1093/nar/24.1.238 (1996).
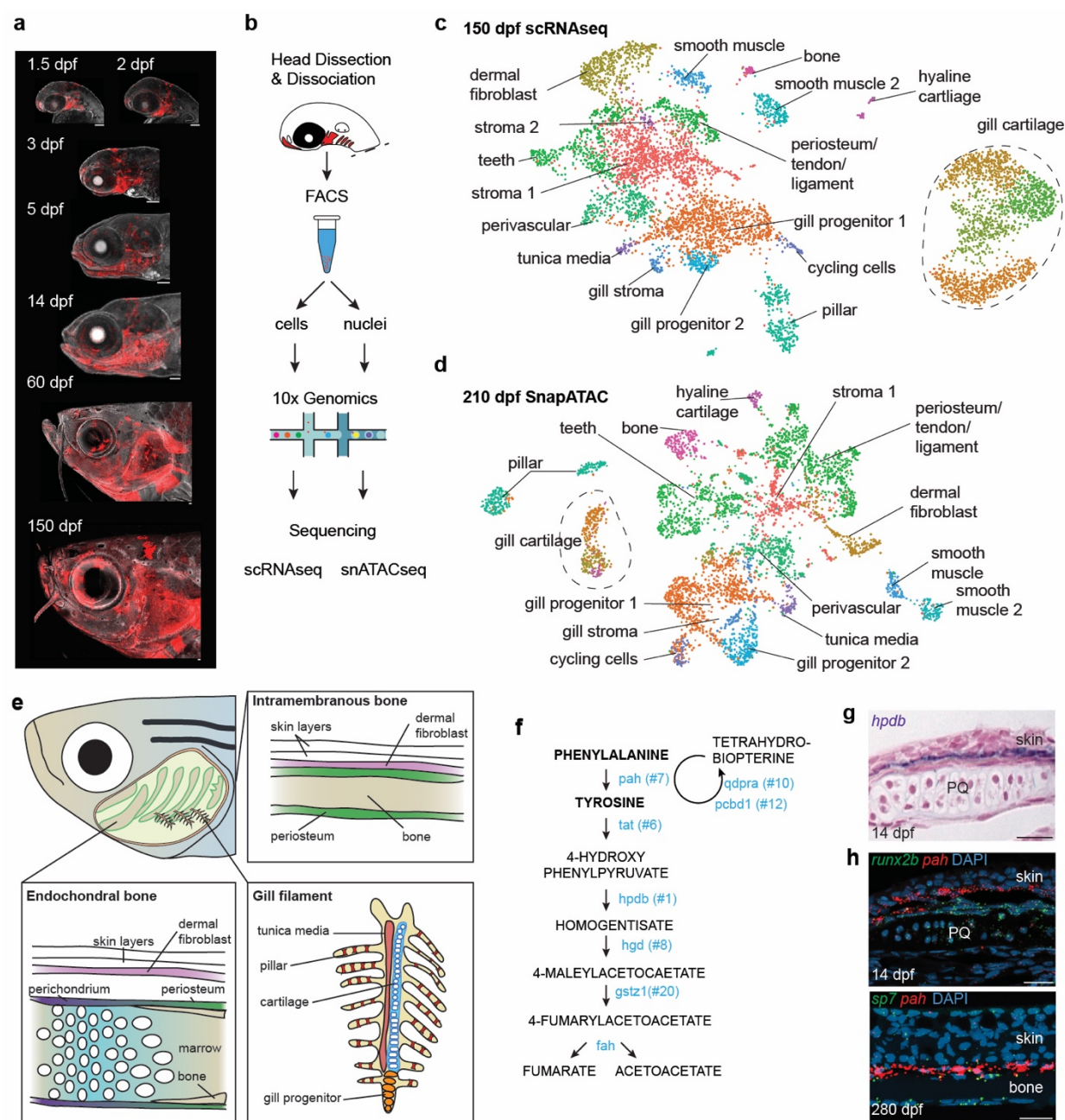
**Figure 1. Single-cell transcriptomes and chromatin accessibility of CNCCs across the zebrafish lifetime. a,** *Sox10:Cre; actab2:loxP-BFP-STOP-loxP-dsRed* labeling of CNCCs (red) in zebrafish heads across 7 stages. **b,** Scheme of cell or nuclei dissociation, fluorescence activated cell sorting (FACS), and processing on the 10X Genomics platform for sequencing. **c,d,** UMAPs of scRNAseq and SnapATAC datasets at adult stages. **e,** Diagram of an adult zebrafish head cut out to show major cell types of the interior skeletal elements and gill system. **f,** Pathway

689    for Phe and Tyr breakdown. 7 of 8 genes encoding catabolic enzymes are in the top 20 selectively

690    enriched genes for the dermal fibroblast cluster (numbers show rank). **g**, Section colorimetric in

691    situ hybridization for *hpdb* RNA shows expression in the dermis between the skin and

692    palatoquadrate (PQ) endochondral bone. **h**, Section RNAscope in situ hybridizations show *pah*

693    expression in dermal fibroblasts between the *runx2b+* periosteum of PQ and skin at 14 dpf, and

694    between *sp7+* osteoblasts of intramembranous bone and skin at 280 dpf. DAPI labels nuclei in

695    blue. Scale bars, 100 μm (a), 20 μm (g,h).

696



697

698    **Figure 2. Progressive emergence of region-specific lineage programs. a,b**, STITCH plots

699 connect individual scRNAseq and snATACseq datasets across the zebrafish lifetime. Cell type

700 annotations show divergence of CNCC ectomesenchyme into skeletogenic, gill, dermal fibroblast,

701 and stromal branches. **c-e**, Pseudotime analysis using Monocle3 of the skeletogenic subset

702 (combined 5 and 14 dpf) shows a *cxcl12a+* stromal branch, and a *hyal4+* branch connected to

703 cartilage, bone, periosteum, and tendon and ligament. **f-h**, Pseudotime analysis of gill subsets

704 (combined 5 and 14 dpf) shows a *cxcl12a+* stromal branch, and a *fgf10a/b+* branch connected to

705 gill pillar, tunica media, and perivascular cells, as well as a distinct type of gill cartilage. **i**, Following

706 UV-mediated photoconversion of *fgf10b:nEOS* from green to magenta in a subset of filaments,

707 re-imaging 3 days later revealed contribution of converted cells to pillar cells (yellow arrow, white

708 reflects mixture of converted magenta and new unconverted green *fgf10b:nEOS*) and gill

709 chondrocytes (white arrow) (*n* = 3). Draq5 labels nuclei in grey. **j**, Gene Ontology (GO) terms of

710 biological processes for the respective clusters at the stages indicated. Heatmap reflects the

711 negative log of the adjusted p value. Ectomesenchyme is defined as the average of all cells at
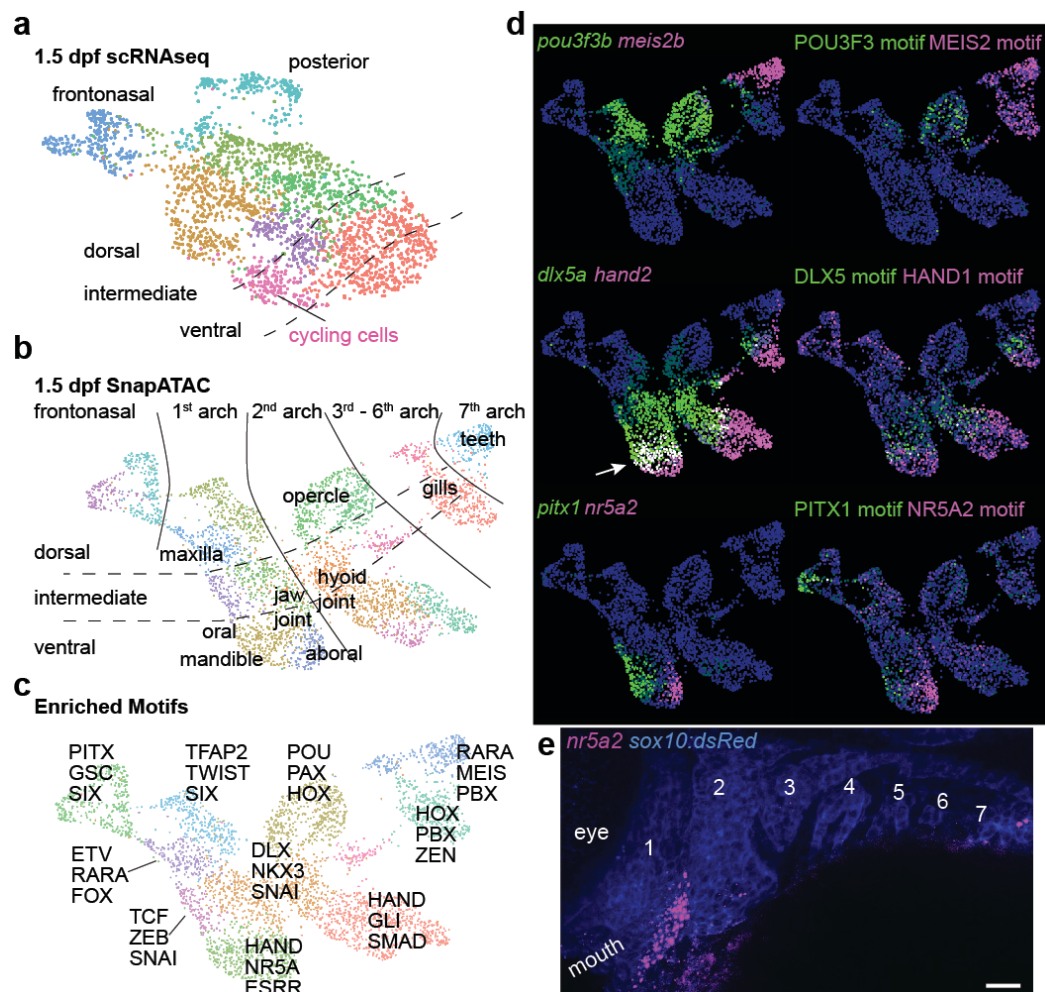
712 1.5 and 2 dpf stages. Scale bar, 100 μm.

713

**Figure 3. Highly resolved embryonic spatial expression domains from integrated datasets.**

**a,b**, UMAPs at 1.5 dpf generated by scRNAseq versus integration of scRNAseq and snATACseq datasets using SnapATAC. SnapATAC outperforms scRNAseq in resolving dorsoventral (vertical), anterior-posterior (horizontal), and major arch landmarks including a previously unappreciated oral-aboral axis in zebrafish. **c**, Select enriched transcription factor binding motifs for each cluster. **d**, Gene body activities for transcription factors and their corresponding DNA-binding motifs reveal tight correspondence with published expression patterns, including zebrafish-specific overlap of *dlx5a* and *hand2* in the mandibular arch (arrow). **e**, Fluorescent in situ hybridization shows restricted expression of *nr5a2* in the aboral domain of the mandibular arch as predicted by SnapATAC. *sox10:dsRed* labels CNCCs of the arches (numbered, anti-dsRed antibody stain). Scale bar, 20 μm.
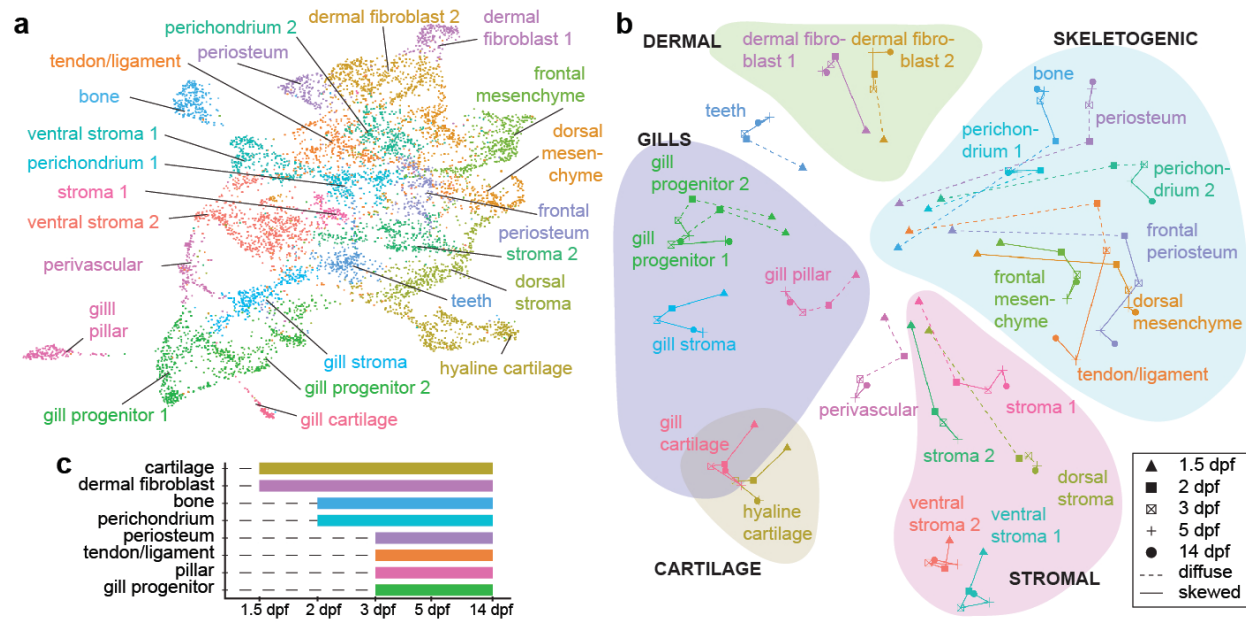
29

**Figure 4. Cell type competency mapping through retrograde chromatin accessibility analysis. a**, UMAP of 14 dpf SnapATAC data shows cell clusters for which top accessible peak modules were calculated for Constellations analysis. **b,** Constellations analysis involves mapping of cluster-specific chromatin accessibility from 14 dpf back to earlier stages and then plotting relatedness of mapped accessibility in two dimensions. Diffuse refers to a stage where cluster-specific chromatin accessibility does not map to a discrete portion of UMAP space, and skewed where it does (see Methods for details). Groups of related cell types are color-coded. **c**, Graphical representation from the Constellations analysis of when chromatin accessibility of major cell types first shows a skewed distribution in UMAP space, suggestive of establishment of competency. **a**, UMAP of 14 dpf SnapATAC data shows cell clusters for which top accessible peak modules were calculated for Constellations analysis. **b**, Constellations analysis involves mapping of cluster-specific chromatin accessibility from 14 dpf back to earlier stages and then plotting relatedness of mapped accessibility in two dimensions. Diffuse refers to a stage where cluster-specific chromatin accessibility does not map to a discrete portion of the UMAP space, and skewed to where it does (see Methods for details). Groups of related cell types are color-coded. **c**, Graphical representation from the Constellations analysis of when chromatin accessibility of major cell types first shows a skewed distribution in UMAP space, suggestive of establishment of competency.
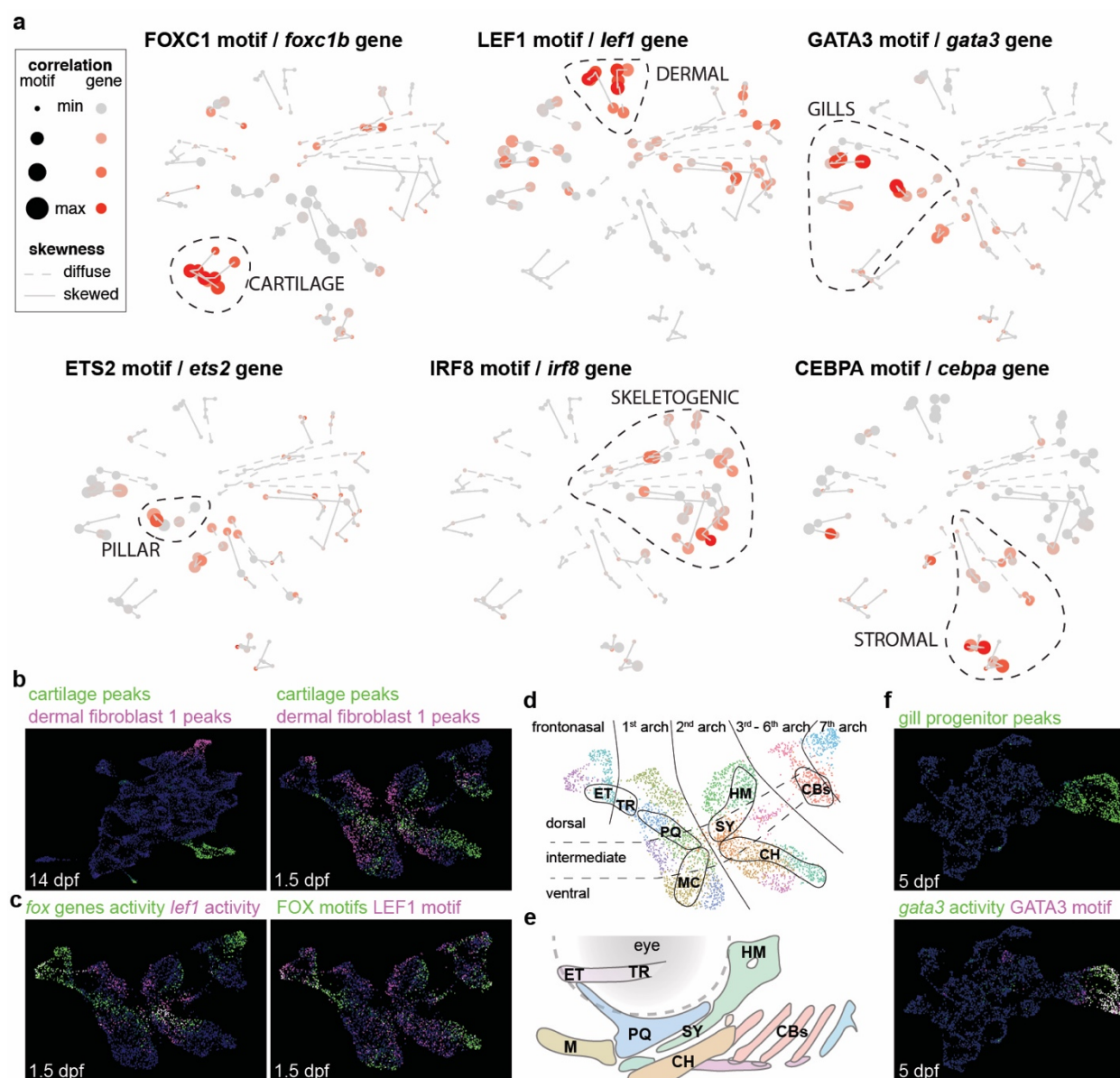
742

**Figure 5. Constellations analysis reveals candidate transcription factors for lineage priming. a**, Mapping onto the Constellations plot of transcription factors with correlated gene body activity and DNA-binding motif enrichment in specific clusters. Sizes of circles denote correlation of peak module mapping to motifs, and red color to gene body activities. **b**, Top peak modules for hyaline cartilage and dermal fibroblast 1 clusters at 14 dpf mapped onto 14 dpf and 1.5 dpf SnapATAC UMAPs. **c**, Summed gene body activities of *foxc1a*, *foxc1b*, *foxf1*, *foxf2a*, *foxf2b* and summed FOXC1, FOXF1, and FOXF2 motifs at 1.5 dpf correlate with cartilage peak mapping,

750 and *lef1* gene body activity and LEF1 motif with dermal fibroblast peak mapping. **d**,**e**, Retrograde

751 cartilage accessibility mapping at 1.5 dpf allows predictions of the arch origins of the individual

752 cartilaginous elements of the week-old skeleton: ceratobranchials (CBs), ceratohyal (CH),

753 ethmoid (ET), hyomandibula (HM), Meckel's (M), palatoquadrate (PQ), symplectic (SY), and

754 trabecula (TR). **f,** Mapping of the top peak module for 14 dpf gill progenitor clusters onto 5 dpf

755 SnapATAC UMAP shows correlation with *gata3* gene body activity and GATA3 motif.
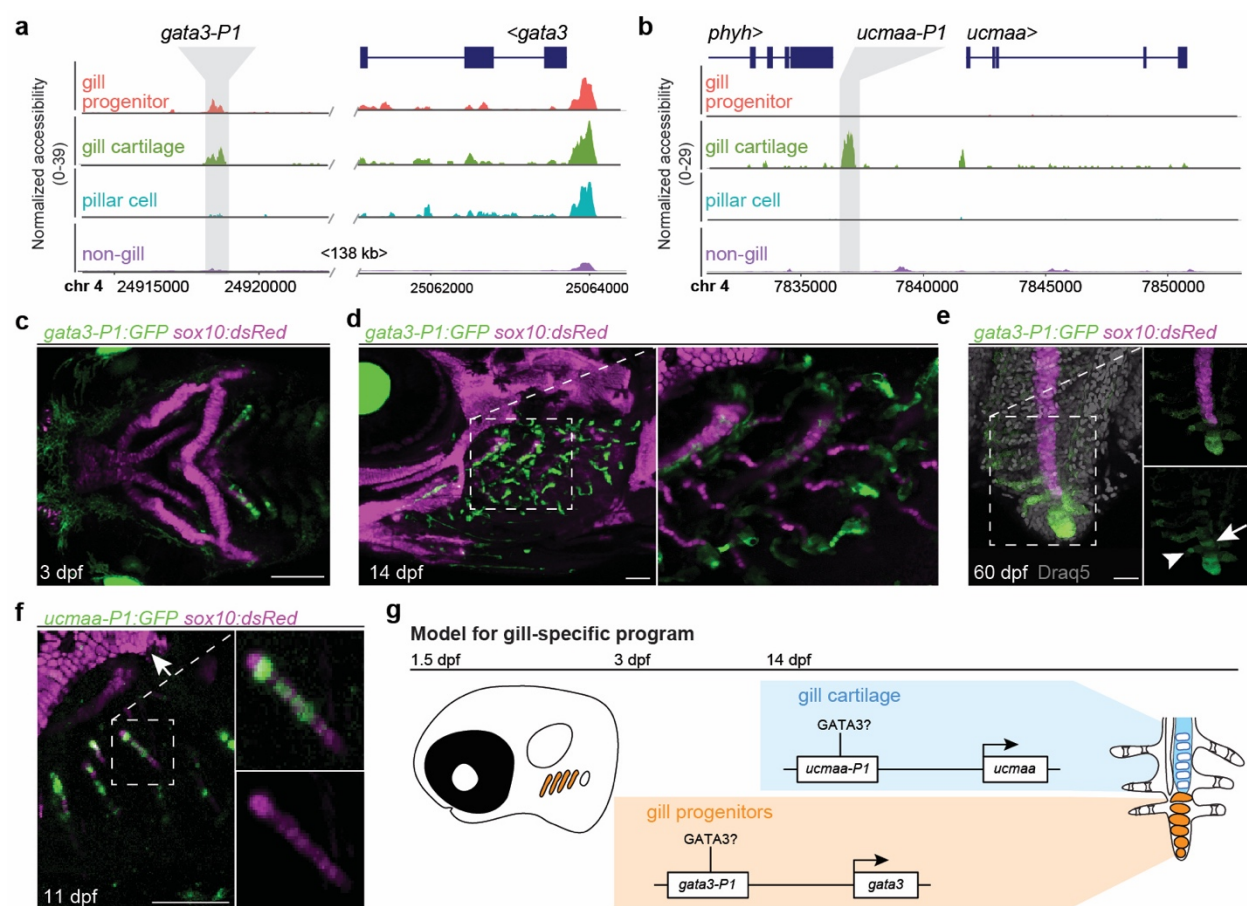
756



757

758 **Figure 6. Gata3 activity distinguishes the gill-specific lineage. a,b**, Genome tracks show

759 chromatin accessibility for cell clusters from snATACseq data at 14 dpf. Grey shading shows gill-

760 specific accessible regions near the *gata3* and *ucmaa* genes. Chromosome positions refer to the

761 GRCz11 genome assembly. **c-f**, *gata3-P1:GFP* drives expression in the posterior gill-forming

762     arches at 3 dpf, the gill filament system at 14 dpf, and gill progenitors at the tips of primary

763     filaments at 60 dpf, as well as some pillar cells (arrowhead) and gill chondrocytes (arrow) near

764     the growing tips. *sox10:dsRed* labels cartilage and Draq5 nuclei. **f**, *ucmaa-P1:GFP* drives highly

765     restricted expression in *sox10:dsRed+* gill chondrocytes (boxed region shown in merged and

766     single channels to the right) but not hyaline cartilage (top left). **g**, Model shows initiation of *gata3*

767     expression in the posterior gill-forming arches through the *gata3-P1* enhancer, maintenance of

768     *gata3* in gill progenitors at the tips of growing filaments, and expression of *ucmaa* in gill cartilage

769     through the *ucmaa-P1* enhancer. Both *gata3-P1* and *ucmaa-P1* contain predicted Gata3 binding

770     sites. Scale bars, 100 μm (c,d,f), 20 μm (e).