

Figure #	Figure title One sentence only	Filename This should be the name the file is saved as when it is uploaded to our system. Please include the file extension. i.e.: <i>Smith_ED_Fig1.jpg</i>	Figure Legend If you are citing a reference for the first time in these legends, please include all new references in the main text Methods References section, and carry on the numbering from the main References section of the paper. If your paper does not have a Methods section, include all new references at the end of the main Reference list.
Extended Data Fig. 1	Hypothalamic and POA cell types in zebrafish and Mexican tetra	Figure E1 - Hypothalamic and POA cell types in zebrafish and Mexican tetra.eps	<b>(a)</b> UMAP of zebrafish cells coloured and labelled by annotated cell type. <b>(b)</b> UMAP of Mexican tetra surface- and cave-morphs coloured and labelled by annotated cell type. <b>(c)</b> DotPlot of the top 2 marker genes for each zebrafish cluster from (a). <b>(d)</b> DotPlot of the top 2 marker genes for each Mexican tetra cluster from (b). Examples of potentially homologous cell types and their top marker genes share a colour (blue, green, red) in (c) and (d). <b>(e)</b> UMAP of merged but not batch-corrected zebrafish and Mexican tetra single-cell datasets.
Extended Data Fig. 2	Marker genes for cell types shared between zebrafish and Mexican tetra	Figure E2. Marker genes for cell types shared between zebrafish and Mexican tetra.eps	<b>(a)</b> DotPlot of the top 5 marker genes for each integrated cluster. <b>(b)</b> Proportion of cells from each cluster by species or species-morph (height of each bar along the x-axis). Width of each bar along the y-axis indicates the proportion of that cluster in the integrated data. Red outlines indicate the Mexican tetra-specific Ciliated cluster, and the integrated Immune clusters which are over-represented in the Mexican tetra dataset. <b>(c)</b> Density plot of the number of subclusters versus the fraction of each subcluster that is either from the zebrafish or Mexican tetra dataset. Subclusters with the majority of cells from the zebrafish dataset are shown in purple, and those with the majority of cells from the Mexican tetra dataset in yellow.
Extended Data Fig. 3	Shared subclusters are highly similar due to paralogous gene expression	Figure E3. Shared subclusters are highly similar due to paralogous gene expression.eps	<b>(a)</b> Gene orthology confidence from Ensembl for all marker genes, or those marker genes which were paralogs of a marker gene in the other species. <b>(b)</b> Gene order score from Ensembl for all marker genes, or those marker genes which were paralogs of a marker gene in the other species. <b>(c)</b> The percentage of conserved, species-specific, and species-specific paralogous subcluster marker genes corrected by SCORPiOS synteny-correction. <b>(d)</b> The percentage of morph-specific marker genes for each subcluster which were paralogs of either the conserved or opposite species-specific marker gene for surface- and cave-morphs of

			<p>Mexican tetra. <b>(e)</b> The odds ratio for the enrichment of paralogs in the species-specific genes for each subcluster for zebrafish and Mexican tetra. <b>(f)</b> The row-scaled <math>\Delta S/</math> for all subclusters between zebrafish and Mexican tetra. Yellow indicates the highest <math>\Delta S/</math> value between Mexican tetra and zebrafish subclusters. For all boxplots, box bounds represent the first and third quartiles and whiskers 1.5 times the interquartile range, thicker line represents the median.</p>
Extended Data Fig. 4	Paralog shifts are associated with loss of ancestral gene expression patterns	Figure E4. Paralog shifts are associated with loss of ancestral gene expression patterns.eps	<p><b>(a-b)</b> Empirical cumulative distribution function (ECDF) for expression divergence (<math>dT</math>) for paralogous gene pairs. <b>(c-d)</b> ECDF of the number of cell types that have overlapping expression patterns within ancestral cell types for paralogous genes pairs (redundancy score, orange highlight in b). <b>(e-f)</b> ECDF of the number of non-ancestral cell types expressing each individual paralogous gene. Results for c-e are grouped by the age of the duplication inferred from the last common ancestor (LCA) which had both genes - from the oldest (Opisthokonta, yellow), to the most recent common ancestor (Otophysi, red), and to those gene duplicates which are only found in either <i>Danio rerio</i> or <i>Astyanax mexicanus</i> (dark red). Results from <b>b, d, and f</b> are filtered and grouped by the originating whole genome duplication event (WGD), either vertebrate (2R) or teleost (3R).</p>
Extended Data Fig. 5	Gene regulatory networks identified by GENIE3/SCENIC	Figure E5. Gene regulatory networks identified by GENIE3.eps	<p><b>(a)</b> Comparison of the random forest weights for orthologous transcription factors in the zebrafish (y-axis) and Mexican tetra (x-axis) data for example terminal effector genes. Colours indicate whether those transcription factors are in the top 2% of transcription factors for each gene in either zebrafish (blue) and Mexican tetra (red), both (yellow), or none (black).</p>
Extended Data Fig. 6	Species-specific subcluster identities are not dependent on species-specific genes	Figure E6. Species-specific subcluster identities are not dependent on species-specific genes.eps	<p><b>(a)</b> tSNEs of cells from clusters containing a species-specific neuronal subcluster coloured by the original subcluster identity. <b>(b)</b> tSNEs of cells from clusters containing a species-specific neuronal subcluster coloured by subcluster identity derived from subclustering without species-specific genes. <b>(c)</b> Sankey diagrams illustrating the relationship between original subcluster identities and identities from subclustering without species-specific genes. Box heights and line widths are proportional to the number of cells in each subcluster and connection, respectively.</p>

			Shaded connections represent cells from species-specific subclusters.
Extended Data Fig. 7	Comparison of subcluster identities between independent and integrated analysis	Figure E7. Comparison of subcluster identities between independent and integrated analysis.eps	<p>(a) Sankey diagram of Mexican tetra surface-morph specific subclusters and their relationship to integrated subclusters, and zebrafish subclusters. Box heights and line widths are proportional to the number of cells in each subcluster and connection, respectively. (d) Sankey diagram of Mexican tetra cave-morph specific subclusters and their relationship to integrated subclusters, and zebrafish subclusters. Box heights and line widths are proportional to the number of cells in each subcluster and connection, respectively. (c) Sankey diagram of the Zebrafish species-specific subclusters (middle) and their relationship to subclusters independently identified in the zebrafish (right) or Mexican tetra datasets (left). Box heights and line widths are proportional to the number of cells in each subcluster and connection, respectively. (d) Sankey diagram of the subclusters shared by, (“Shared (147)”) or specific to, surface- and/or cave-morphs (“Cave-specific” or “Surface-specific”). The middle column depicts whether each subcluster is found in all cave-morph samples (“All Caves”), different combinations of multiple caves, or only in the datasets from specific cave-lineages (“Pachon” or “Molino”). Box heights and line widths are proportional to the number of cells in each subcluster and connection, respectively.</p>
Extended Data Fig. 8	Comparison of neuropeptides and gene regulatory networks between surface- and cave-morphs	Figure E8. Comparison of neuropeptides and gene regulatory networks between surface- and cave-morphs.eps	<p>(a) DotPlot showing expression of <i>galn</i> in the cells from the <i>galn</i> cluster (Neuronal_07), and expression of <i>oxt</i>, <i>avp</i>, and <i>ENSAMXG0000021172</i> in the Neuronal_19 cluster. Cells are grouped by species morph and cave-lineage. (b) Similarity Index between the transcription factor sets for surface- and cave-morphs of Mexican tetra for neuropeptides, neurotransmitters, synaptic genes, and ion channels. (c-f) Random forest weights for orthologous transcription factors in the Mexican tetra surface-morph (y-axis) and Mexican tetra cave-morph (x-axis) data for the neuropeptides <i>galn</i>, <i>hcrt</i>, <i>oxt</i>, and <i>avp</i>. Colours indicate whether those transcription factors are in the top 2% of transcription factors for each gene in either surface-morphs (green) and cave-morphs (yellow), both (purple), or none (black). For all boxplots, box bounds represent the first and third quartiles and whiskers 1.5 times the</p>

			interquartile range, thicker line represents the median.
Extended Data Fig. 9	Transcriptional signatures of neuroinflammation resistance in cave-morphs	Figure E9. Transcriptional signatures of neuroinflammation resistance in cave-morphs.eps	(a) tSNE reduction of immune clusters (Tcells, Bcells, Microglia, Macrophages, Mast cells, Thrombocytes, Neutrophils, and Erythorcytes) from surface- and cave-morph Mexican tetra coloured and labelled by species-morph. (b) tSNE reduction of immune cell types from surface- and cave-morph Mexican tetra coloured by cluster. (c) Marker genes for surface- and cave-morph versions of each immune cell type. Red outlines indicate differential expression of neuroinflammation associated genes in cave-morph immune cells. Gene expression is quantified by both the percentage of cells which express each gene (dot size) and the average expression in those cells (colour scale). (d) tSNE reduction showing expression of <i>ccr9a</i> in Mexican tetra immune cells. (e) Proportion of cells within each immune subcluster which come from Choy surface-morphs, or Molino, Tinaja, or Pachon cave-morphs.
Extended Data Fig. 10	A permanent stress-response in a cave-morph specific neuronal subcluster	Figure E10. A permanent stress-response in a cave-morph specific neuronal subcluster.eps	(a) tSNE reduction of Neuronal_03 cluster from Mexican tetra coloured and labelled by subcluster. (b) tSNE reduction of Neuronal_03 cluster from Mexican tetra coloured by species-morph. (c) DotPlot of the top 5 marker genes for each subcluster of the Neuronal_03 cell type (x-axis), and their expression across all subclusters (y-axis). Gene expression is quantified by both the percentage of cells which express each gene (dot size) and the average expression in those cells (colour scale). (d) Dendrogram of the Neuronal_03 subclusters based on the Variable Features of the Neuronal_03 cluster, and the proportion barplot of cells from each species-morph per subcluster. (e) GO analysis of genes differentially expressed between Neuronal_03_1 and Neuronal_03_4. (f) tSNE reduction of Neuronal_03 cluster from Mexican tetra coloured by <i>hspb1</i> expression. Neuronal_03_4 subcluster is highlighted by a dotted line. (g) Sankey diagram of the relationships between the Mexican tetra subclusters (left-hand side), integrated subclusters (middle), and zebrafish subclusters (right-hand side). Box heights and line widths are proportional to the number of cells in each subcluster and connection, respectively.



Item	Present?	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. The extension must be .pdf	A brief, numerical description of file contents. i.e.: <i>Supplementary Figures 1-4, Supplementary Discussion, and Supplementary Tables 1-4.</i>
Supplementary Information	Yes	Supplementary_Information_Shafer.pdf	Combined Supplementary Methods/Results, Supplementary References, Supplementary Figures 1-9
Reporting Summary	No		
Peer Review Information	Yes	<i>Shafer_PFfile</i>	

2  
3

Type	Number If there are multiple files of the same type this should be the numerical indicator . i.e. "1" for Video 1, "2" for Video 2, etc.	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. i.e.: <i>Smith_Supplementary_Video_1.mov</i>	Legend or Descriptive Caption Describe the contents of the file
Supplementary Data		Supplementary_data.zip	<p>Supplementary data for Cavefish single-cell sequencing publication =====</p> <p>This archive contains the supplementary data for the paper "Gene family evolution underlies cell type diversification in the hypothalamus of teleosts", which includes all of the raw and partially processed data produced by the analyses presented.</p> <p>The archive contains:</p> <ol style="list-style-type: none"> <li>1) The raw count data for both the zebrafish (<i>Danio rerio</i>) and Mexican tetra (<i>Astyanax mexicanus</i>) single-cell experiments, as compressed .csv files.</li> <li>2) The Seurat object meta data for the zebrafish (<i>Danio rerio</i>), Mexican tetra (<i>Astyanax mexicanus</i>), and integrated Seurat objects, containing sample, species, and cell type cluster labels for each cell.</li> </ol>

		<p>3) CSVs for all marker gene lists used in the publication.</p> <p>4) CSVs for all pseudobulk expression data for all cell type labels.</p> <p>5) The raw data used for calculating the SI for each cluster and subcluster identity in the integrated data</p> <p>6) Results from SCENIC/GENIE3 analysis, including the Linklists and tfModules outputs from SCENIC.</p> <p>7) Results of the weir fst analysis between cave and surface populations, for both INDELS and SNPs</p> <p>8) Ensembl biomart export files for determine paralogy relationships between genes within and across species</p> <p>9) Results of trinarization of gene expression across all identities, an the uniquely expressed genes per identity. These are provided as R object files (.rds)</p> <p>Supplemental_data/3-marker_gene_lists  =====</p> <p>This folder contains marker gene lists for clusters and subclusters (".sub"), for the zebrafish ("Drerio"), Mexican tetra ("Amexicanus"), or integrated ("Integrated") datasets.</p> <p>Supplemental_data/4-pseudobulk_expression  =====</p> <p>This folder contains psuedobulk expression profiles for clusters ("Clusters") and subclusters ("Subclusters"), for the zebrafish ("Drerio"), combined Mexican tetra, and the surface and cave morphs of Mexican tetra, and integrated datasets.</p>
--	--	--

4

5

6 **Gene family evolution underlies cell type diversification in the hypothalamus of teleosts**

7 Maxwell E.R. Shafer<sup>1,2,\*†</sup>, Ahilya N. Sawh<sup>1,2,\*</sup>, Alexander F. Schier<sup>1,2,3,4,5,6 \*†</sup>

8

9 <sup>1</sup> Biozentrum, University of Basel, Switzerland

10 <sup>2</sup> Department of Molecular and Cellular Biology, Harvard University, Cambridge, Massachusetts, USA

11 <sup>3</sup> Allen Discovery Center for Cell Lineage Tracing, Seattle, Washington, USA

12 <sup>4</sup> Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

13 <sup>5</sup> Harvard Stem Cell Institute, Harvard University, Cambridge, Massachusetts, USA

14 <sup>6</sup> Center for Brain Science, Harvard University, Cambridge, Massachusetts, USA

15

16 \* Current address: Biozentrum, University of Basel, Switzerland

17 † Corresponding authors: MERS: max.shafer@gmail.com; AFS: alex.schier@unibas.ch

18

19 **ABSTRACT:**

20 Hundreds of cell types form the vertebrate brain, but it is largely unknown how similar cellular  
21 repertoires are between or within species or how cell type diversity evolves. To examine cell type diversity  
22 across and within species, we performed single-cell RNA sequencing of ~130,000 hypothalamic cells  
23 from zebrafish (*Danio rerio*) and surface- and cave-morphs of Mexican tetra (*Astyanax mexicanus*). We  
24 found that over 75% of cell types were shared between zebrafish and Mexican tetra, which diverged from  
25 a common ancestor over 150 million years ago. Shared cell types displayed shifts in paralog expression

26 that were generated by sub-functionalization after genome duplication. Expression of terminal effector  
27 genes, such as neuropeptides, was more conserved than the expression of their associated  
28 transcriptional regulators. Species-specific cell types were enriched for the expression of species-specific  
29 genes, and characterized by the neo-functionalization of expression patterns of members of recently  
30 expanded or contracted gene families. Comparisons between surface- and cave-morphs revealed  
31 differences in immune repertoires and transcriptional changes in neuropeptidergic cell types associated  
32 with genomic differences. The single-cell atlases presented here are a powerful resource to explore  
33 hypothalamic cell types, and reveal how gene family evolution and shifts in paralog expression contribute  
34 to cellular diversity.

35

## 36 **INTRODUCTION:**

37 The homology of neuronal cell types was first revealed by Ramón y Cajal, who observed that  
38 morphologically similar neurons were present in the brains of many species<sup>1</sup>. Since then, the comparison  
39 of cell types has largely relied on morphological criteria and, more recently, data from select marker  
40 genes<sup>2</sup>. These studies have led to the definition of major neuronal classes and subclasses<sup>2,3</sup> but it is still  
41 unknown how molecularly similar or different brain cell types are between species. Moreover, it is unclear  
42 how cell types diversify during evolution or adaptation to extreme environments. Biological novelty may  
43 arise as a result of gene expansion<sup>4</sup>, but it is unknown how the evolution of gene families influences the  
44 diversification of cell types in the brain.

45 Single-cell sequencing has recently emerged as a powerful tool to study and map the cell types of  
46 individual species, and has allowed the identification of hundreds of transcriptionally unique cell types in  
47 vertebrate tissues, including the brain<sup>5</sup>. Recently, cross-species comparisons using single-cell RNA-seq  
48 have identified shared and species-specific cell types, as well as mechanisms for neuronal evolution<sup>5,6</sup>.  
49 These studies have identified conserved cell types during vertebrate development<sup>7</sup> and mammalian  
50 neurogenesis<sup>8,9</sup>, as well as primate-specific adaptations<sup>10,11</sup>. Extension of these approaches to more  
51 diverse phylogenies is necessary for understanding the molecular and evolutionary basis of cell type  
52 conservation and diversification across the tree of life.

53 A powerful model for comparative studies of biological diversification are the teleosts. This group  
54 of nearly 30,000 described ray-finned fish species represents the largest clade within vertebrates and has  
55 undergone a taxon-specific whole genome duplication (WGD)<sup>12,13</sup>. It has been hypothesised that the vast  
56 diversification in morphology, physiology, and behaviour observed across teleost species was driven by  
57 gene family expansions associated with the teleost-specific WGD<sup>4,12,14</sup>. Most duplicated genes lose their  
58 functions through deleterious mutations (non-functionalization), but genes that are retained may undergo  
59 either sub-functionalization (partitioning of functions or gene expression patterns), or neo-functionalization  
60 (gain of novel functions or gene expression patterns). Little is known about the fate of these duplicated  
61 genes in teleosts, their roles in the vertebrate brain, or their links to cellular diversification.

62 In this study we analyze the conservation and diversification of teleost brain cell types using the  
63 zebrafish (*Danio rerio*) and the Mexican tetra (*Astyanax mexicanus*) as model systems. Zebrafish is the  
64 leading fish model system in developmental and neurobiology, whereas Mexican tetra is a powerful  
65 system for comparative studies. Mexican tetra has two morphs, an eyed surface-morph, and an eye-less  
66 and pigment-less cave-morph<sup>15-17</sup>. Comparisons between species, and between species-morphs  
67 represent two informative evolutionary distances: between distantly related species (150-200 million  
68 years, zebrafish and Mexican tetra), and within a species with large phenotypic differences (250-500  
69 thousand years, species-morphs of Mexican tetra), that have been linked to changes in the development  
70 and gene expression patterns of the nervous system<sup>15,18</sup>.

71 To characterise cell type diversity at a high resolution in both zebrafish and Mexican tetra, we  
72 focus on the hypothalamus. The hypothalamus is a highly conserved forebrain region that is responsible  
73 for the generation and secretion of hormones and neuropeptides involved in diverse behaviours. Within  
74 the hypothalamus these functions are partitioned into specific neuropeptidergic cell populations regulating  
75 sleep/wake (*hcrt+* and *galn+* neurons), food intake (*agrp*, *npv*, *pomc*), aggression and sexual behaviours  
76 (*oxt*, *avp*, *npv*), and physiological homeostasis<sup>19-21</sup>. It is thought that hormone-secreting brain centres are  
77 ancient, and were present in the last common ancestor of all metazoans<sup>22</sup>. However, the level of  
78 homology in the cellular populations of the hypothalamus has not been comprehensively compared  
79 between species.

80 We used single-cell transcriptomics followed by high resolution clustering and cross-species  
81 integration to systematically identify the molecular similarities in the cellular repertoire of the teleost  
82 hypothalamus. First, we observe high conservation of cell types between species over 150 million years  
83 of evolution. Second, our results suggest that shared cell types have undergone shifts in paralog gene  
84 expression, and divergence in gene regulatory networks. Third, we link cellular novelty with genetic

85 novelty and the species-specific expression of paralogous genes. Fourth, we identify transcriptional and  
86 genomic differences between surface- and cave-morphs of Mexican tetra that are candidates to be  
87 associated with behavioural phenotypes of cave-adaptation.

88

## 89 **RESULTS**

### 90 **Iterative clustering identifies shared and divergent cell types in the hypothalamus and preoptic** 91 **area of *D. rerio* and *A. mexicanus***

92 To characterise similarities and differences of brain cell types between and within species, we  
93 performed scRNA-seq on ~130,000 cells from the hypothalamus and preoptic area (POA) of *D. rerio*  
94 (zebrafish), and from surface and 3 different cave species-morphs of *A. mexicanus* (Mexican tetra)  
95 (**Figure 1a-c, Extended Data Figure 1-2, Supplementary Information, and Supplementary Figures**  
96 **1**). To resolve cell populations at high resolution, we performed iterative subclustering resulting in 194  
97 subclusters with distinct gene expression patterns (**Figures 1d, see Supplementary Information,**  
98 **Supplementary Figures 2-4** and Methods for details on our clustering approach and the comparison  
99 between clusters and subclusters). Subclustering resolved rare cell populations such as the *hcrf*<sup>+</sup>  
100 subcluster which also expressed the neuropeptide *npvf* and the transcription factor *lhx9*  
101 (Neuronal\_01\_10) (Supplementary **Figure 2e**)<sup>19,26</sup>. The majority of subclusters were shared between  
102 species (151 out of 194 subclusters composed of > 10% of cells from both species), whereas 43  
103 subclusters were specific to either zebrafish or Mexican tetra (species-specific, > 90% of cells from either  
104 species in our analysis) (**Figure 1d**). Analysis of the similarity between subclusters from each of the cave-  
105 morphs reflected the known phylogenetic relationship between surface-morphs, and Pachon, Tinaja, and  
106 Molino cave-morphs (Figure 1e). Thus, subclustering identified both similar and divergent cell types  
107 between zebrafish and Mexican tetra.

108 In the following sections we analyze this dataset to (1) identify the similarities and differences in  
109 gene expression for cell types shared between zebrafish and Mexican tetra (**Figure 2-3**), (2) compare  
110 gene regulatory networks across species (**Figure 4**), (3) define gene expression signatures associated  
111 with species-specific cell types (**Figure 5**), and (4) examine cell type similarities and differences between  
112 the surface- and cave-morphs of Mexican tetra (**Figure 6**). We then discuss our findings in relation to the  
113 molecular and evolutionary basis of cell type conservation and diversification.

114

### 115 **Transcriptional Similarity Index measures gene expression similarities and differences between** 116 **cell types**

117 To quantify the similarities and differences in the subclusters shared between species we  
118 developed a transcriptional Similarity Index (*SI*) that compares the presence or absence of cell type  
119 specific marker genes (Extended Data Figure 2a, Supplementary Information). We calculated the *SI*  
120 between the marker gene sets for zebrafish and Mexican tetra for each of the 151 shared subclusters  
121 (Figure 2a-b). *SI* was consistently the highest between shared subclusters, and between subclusters from  
122 the same or related clusters (**Figure 2a-b**).

123 Analysis of the patterns of *SI* between subclusters revealed that progenitor subclusters had  
124 significantly higher levels (average *SI* of 0.327) than differentiated neuronal cells (average *SI* of 0.224)  
125 (**Figure 2c**). The lower divergence for progenitor cells between species could be due to their function in  
126 generating many different cell types, with pleiotropic effects expected from changing gene expression  
127 patterns in progenitor cell types. Indeed, the *SI* for the larger parental clusters (marker genes shared by  
128 multiple subclusters) was significantly higher than the similarity for subclusters (marker genes specific to  
129 subclusters) also suggesting pleiotropic effects (**Figure 2d**). For example, the cluster Neuronal\_01 had a  
130 *SI* of 0.305, whereas its subclusters, Neuronal\_01\_0 - Neuronal\_01\_11, had *SIs* between 0.286 and  
131 0.146 (mean of 0.206). To test whether the amount of transcriptional similarity scales with the  
132 evolutionary time between species, we also calculated *SI* for the same subclusters between cave and  
133 surface morphs of Mexican tetra (Figure 2e). Higher *SI* was observed for surface versus cave morph cells

134 than between zebrafish and Mexican tetra, indicating that *SI* reflects divergence time. These results  
135 highlight that the transcriptomes of progenitors (versus differentiated cells), and of cell clusters (versus  
136 subclusters), have changed the least during evolution.

137

### 138 **Expression of paralogous and functionally similar genes contributes to transcriptional divergence** 139 **in shared cell types**

140 To determine what factors contribute to the transcriptional divergence of shared cell types, we  
141 compared the identity of marker genes that were specific to one or the other species. Examination of the  
142 zebrafish-specific marker genes revealed that many were paralogs of Mexican tetra-specific marker  
143 genes, and vice versa. For example, erythrocytes were marked by *hbaa2* in Mexican tetra, but its paralog  
144 *ba1* in zebrafish. Neuronal\_00 cells were marked by *zic2b/zic5* in Mexican tetra, but their paralogs  
145 *zic1/zic3* in zebrafish (Extended Data Figure 1). The best markers for GABAergic and Glutamatergic cell  
146 types were also paralogs of each other in the two species. Zebrafish cells expressed *slc17a6b*, and *gad2*,  
147 whereas Mexican tetra cells expressed the paralogs *slc17a6a* and *gad1b*, for Glutamatergic and  
148 GABAergic cells respectively (Supplementary Figure 4). These results suggest that instead of expressing  
149 orthologous genes, shared cell types often express paralogous genes.

150 To determine how frequent such shifts in paralog expression are, we calculated for each  
151 subcluster the percentage of zebrafish-specific marker genes that were paralogs of either a shared  
152 marker gene or a Mexican tetra-specific marker gene, and vice versa (Figure 2a). Up to 14% of the  
153 species-specific marker genes for each subcluster were paralogous to another gene expressed in the  
154 same subcluster in either species (Figure 2f). These levels represented a 4-15 fold or higher enrichment  
155 for paralogous genes (odds ratio, Fisher's test) (**Figure 2g**), and were not due to mis-identification of  
156 orthology/paralogy relationships between genes (Extended Data Figure 3a-c). The vast majority of the  
157 observed changes in paralog expression between zebrafish and Mexican tetra were also conserved  
158 between surface- and cave-morphs of Mexican tetra (**Figure 2h**). Differences between surface- and cave-  
159 morph marker genes were also enriched for paralogous genes (Extended Data Figure 3d-e) Progenitor  
160 cells, which had the least amount of transcriptional divergence, also had the highest enrichment for  
161 paralogous genes (Figure 2i). Additionally, the enrichment for paralog expression was positively  
162 correlated with the *SI* across all subclusters (**Figure 2i**). Non-paralogous differentially expressed genes  
163 between species were also enriched for similar gene ontology terms suggesting conservation of function  
164 across species (Supplementary Information and Supplementary Figure 6). To account for shifts in paralog  
165 expression and to produce a more accurate estimation of the similarity of subclusters between species  
166 we calculated a corrected-*SI*, which considers paralogs as functionally equivalent. We also calculated the  
167 difference between corrected-*SI* and *SI* ( $\Delta SI$ ) (Extended Data Figure 3f). The mean of the corrected-*SI* of  
168 subclusters was twice that of the mean of the *SI* for the same subclusters (**Figure 2j**). The  $\Delta SI$  was  
169 highest between shared subclusters (Extended Data Figure 3f). These results indicate that gene  
170 expression differences between shared cell types are largely due to shifts in the expression of functionally  
171 similar paralogs both between species and between species-morphs.

172

### 173 **Enrichment of paralogous genes in shared subclusters is due to gene duplication followed by** 174 **differential retention of expression patterns**

175 Possible explanations for the observed enrichment for paralogous gene expression in shared  
176 subclusters include the differential retention or gain of expression patterns in the two species following  
177 gene duplication. In this scenario, the expression patterns of newly duplicated genes are initially identical  
178 or highly similar, but over time undergo differential sub- or neo-functionalization in gene expression in  
179 different species. Consider the hypothetical Gene X, which is duplicated in the ancestor of Species 1 and  
180 2. After time, expression of Gene Xa is retained in the cell types of interest in Species 1, whereas  
181 expression of Gene Xb may be retained in Species 2 (**Figure 3a-c**). One prediction of this paralog sub-



182 functionalization model is that more recently duplicated genes will have more similar expression patterns  
183 within each species.

184 We used a previously described metric to calculate the expression divergence ( $dT$ ) for each pair  
185 of paralogous genes within each species<sup>31</sup>. We found that paralog gene pairs that were generated  
186 through more recent duplication events showed lower  $dTs$  (empirical cumulative distribution function) than  
187 paralogs from more ancient duplication events (**Figure 3d-e, and Figure S6b**). To gain further insight  
188 into the divergence patterns of paralogous genes within species we determined putative ancestral gene  
189 expression patterns based on the minimal group of cell types that express either or both paralogs in both  
190 species (ancestral cell types) (Figure S12b'). This analysis revealed that divergence of paralog gene pairs  
191 was predominantly due to loss of co-expression rather than expansion of expression patterns into new  
192 cell types (Extended Data Figure 4 and Supplementary Information).

193 To test whether genes diverged differentially in zebrafish and Mexican tetra, we compared the  $dT$   
194 for each pair of paralogous genes which were conserved in both species. The expression patterns of  
195 most gene pairs have diverged in both species, but many have diverged in only zebrafish or Mexican  
196 tetra (**Figure 3f**). The most recently duplicated gene pairs also had  $dT$  patterns which were different  
197 between species. For example, 14 of 22 gene pairs that arose in the last common ancestor of zebrafish  
198 and Mexican tetra (*Otophysi*) have different  $dT$  levels and expression patterns in zebrafish and Mexican  
199 tetra (**Figure 3g**). For example, *etv5a* is expressed across several glutamatergic neuronal clusters in both  
200 Mexican tetra and zebrafish (**Figure 3h**). The expression of its paralog *etv5b* in the same subclusters was  
201 retained in zebrafish, but lost in Mexican tetra (**Figure 3h**).

202 For all gene pairs in both species, we determined the correlation of their expression patterns  
203 between zebrafish and Mexican tetra (**Figure 3i**). The oldest gene pairs examined had the most  
204 correlated expression patterns across species, including those gene pairs which arose prior to the last  
205 common ancestors of vertebrates (Bilateria, Chordata, and Vertebrata) (**Figure 3i**). No difference was  
206 observed between genes arising from either the 2R or 3R WGDs (**Figure 3i**). This is consistent with the  
207 functions and expression patterns of older genes having diverged long ago, and are therefore less likely  
208 to change further in more closely related species. These results indicate that the divergence in gene  
209 expression patterns for duplicated genes has occurred differently in zebrafish and Mexican tetra.

210

## 211 **The transcription factors associated with neuropeptide expression have diverged between** 212 **species**

213 Our comparison of gene expression signatures in cell types between species revealed  
214 maintenance of cellular function by either reducing transcriptional divergence and/or promoting the  
215 expression of functionally similar and paralogous genes (**Figures 2-3**). We therefore wondered whether  
216 the upstream regulatory mechanisms controlling the expression of functional terminal effector genes,  
217 such as neuropeptides, were also conserved between zebrafish and Mexican tetra. To test this  
218 hypothesis we identified putative gene regulatory networks using SCENIC/GENIE3 and compared these  
219 between species.

220 The putative gene regulatory networks (TF sets) for the same terminal effector genes were  
221 different between zebrafish and Mexican tetra, and many of the top TFs for each gene were not shared  
222 between species (Extended Data Figure 5). For example, the neuropeptide *vip* is highly associated with  
223 the TFs *gsx1*, *otpa*, *six3a*, *mllt11*, and *xbp1* in zebrafish, but with the TFs *klf17*, *id4*, *nr4a3*, *hes6*, and *ets2*  
224 in Mexican tetra (**Figure 4a**). To quantify how similar TF sets were between species, we calculated the  $S_I$   
225 between zebrafish and Mexican tetra for the TF sets of each terminal effector gene (**Figure 4b**). TF sets  
226 associated with neuropeptides had significantly higher  $S_I$  than TF sets associated with neurotransmitters,  
227 synaptic, or ion channel genes (**Figure 4b**). Thus, more of the same TFs remain associated with specific  
228 neuropeptides across species compared to other terminal effector genes.

229 We then quantified the similarity in the relative contribution of each TF (TF weights) for each  
230 terminal effector between species. No statistical difference was observed between the terminal effector

231 classes, and all gene sets had low or negative correlation between species (mean Pearson correlation  
232 between 0 and 0.15) (**Figure 4c**). For example, the TFs *prox1b* and *id4* appear in the TF sets for *vip* in  
233 both species, but are more predictive of *vip* expression in Mexican tetra than in zebrafish (**Figure 4a**).  
234 These results suggest that even in cases where the same TFs are associated with specific terminal  
235 effector genes across species, their relative contributions to putative target gene expression are not  
236 maintained.

237

### 238 **Transcription factors diverge more than target genes and can be replaced by non-paralogs**

239 Divergence in TF sets between species may be compensated by expression of paralogous TFs,  
240 similar to what we observed for subcluster transcriptomes. For example, the expression of the highly  
241 conserved neuropeptide *oxt* was correlated with the TFs *sim1b*, *otpa*, and *otpb* in zebrafish, while in  
242 Mexican tetra, *oxt* was highly correlated with *sim1a* and *otpb* (Extended Data Figure 5). *sim1b* is specific  
243 to the zebrafish genome, indicating that the paralogs *sim1a* and *sim1b* underwent sub-functionalisation in  
244 zebrafish, with *sim1a* losing its function and co-expression with *oxt* in zebrafish<sup>34</sup>. In contrast, in the case  
245 of the neuropeptide *vip*, there were no highly weighted TFs associated with its expression that were  
246 paralogs in the two species (**Figure 4a**). However, in general we did not observe compensation across  
247 species through the association of paralogous TFs with terminal effectors (Figure 4d).

248 The above results suggest that the expression patterns of TFs may be less conserved between  
249 species than the expression patterns of their target genes (non-TFs). To test this prediction, we  
250 calculated the *SI* across all subclusters using all marker genes, only those marker genes that were TFs,  
251 or only those marker genes that were neuropeptides and neurotransmitters (terminal effector genes). *SI*  
252 was significantly higher for transcription factors compared to all marker genes, and was highest for  
253 neuropeptides and neurotransmitters (**Figure 4e**). For example, the neuropeptidergic subclusters of the  
254 Neuronal\_01 parent cluster (including the *galn+*, *hcrt+*, and *oxt+* subclusters) have lower *SI* when  
255 considering only TFs then when considering NP/NT genes (**Figure 4f**). All together, both the SCENIC and  
256 *SI* results suggest that specific classes of genes in specific cell types may be more conserved between  
257 species, and may even be more conserved than the TFs which regulate them.

258

### 259 **Species-specific cellular novelty is associated with species-specific genetic novelty**

260 We next sought to define which genes were associated with species-specific subclusters. There  
261 were 43 species-specific subclusters in our integrated dataset, each composed of > 90% of cells from one  
262 species. Five of these species-specific subclusters were from the Mexican tetra-specific “Ciliated” cell  
263 type (**Figure 1** & Extended Data Figure 2b), 19 were hematopoietic subclusters (see below), and 6 were  
264 from the Oligodendrocyte, Endothelial, and Lymphatic parental clusters. There were 8 Mexican tetra-  
265 specific and 5 zebrafish-specific neuronal subclusters, which expressed a variety of genes indicating that  
266 they may represent different temporal or spatial cell states, captured in one species or the other (Figure  
267 5a, Supplementary Information and Supplementary Figure 7).

268 It has previously been reported that species- or lineage-specific genes may contribute to species-  
269 or lineage-specific morphological and cellular innovations<sup>35</sup>. Though they may be shared by other species  
270 of fish, for simplicity we refer here to genes that are only found in the genome of zebrafish or Mexican  
271 tetra as species-specific genes. In both zebrafish and Mexican tetra, a higher percentage of the genes  
272 expressed in species-specific subclusters were species-specific, compared to the genes expressed in  
273 shared subclusters (**Figure 5b-c**). Importantly, all of these subclusters were still identified in the absence  
274 of species-specific genes, suggesting that they have distinct expression patterns of orthologous genes as  
275 well (Extended Data Figure 6). Mexican tetra non-neuronal cells expressed significantly more species-  
276 specific genes as compared to neuronal subclusters, with the immune subclusters expressing the highest  
277 percentages (**Figure 5d**). In contrast, all zebrafish neuronal subclusters expressed more species-specific  
278 genes as compared to non-neuronal subclusters (**Figure 5e**).

279 Enrichment for species-specific genes was also apparent in the species-specific neuronal  
280 subclusters (**Figure 5f**). For example, the zebrafish-specific Neuronal\_04\_7 subcluster was distinguished  
281 by 5 members of the jacalin family of lectins (*jac2*, *jac3*, *jac6*, *jac8*, and *jac9*) (**Figure 5f**). This gene family  
282 has undergone an extensive species-specific gene expansion, resulting in 14 known genes in zebrafish,  
283 compared to only 2 genes in Mexican tetra<sup>36</sup>. Additionally, the Mexican tetra-specific subcluster  
284 Neuronal\_12\_5 expressed the neuronal calcium sensor (NCS) *HPCAL1*, which was generated by a  
285 Mexican tetra-specific duplication (**Figure 5f**).<sup>37</sup> One subcluster expressed the Mexican tetra-specific  
286 guanylate cyclase *ENSAMXG00000017498*, and *vipb*, which is a paralog of the neuropeptide *vip* that  
287 arose in a common ancestor of teleosts. *vip* is expressed in several subclusters shared by zebrafish and  
288 Mexican tetra, whereas *vipb* is expressed only in the Mexican tetra-specific Neuronal\_13\_2 (**Figure 5f**).  
289 This result suggests that the expression patterns of *vipb* and *ENSAMXG00000017498* have undergone  
290 neo-functionalization in Mexican tetra, but not zebrafish.

291 Three other species-specific subclusters were characterized by genes which were duplicated in a  
292 common ancestor of zebrafish and Mexican tetra, but subsequently lost in zebrafish (**Figure 5f**). These  
293 include Neuronal\_10\_5, which expresses the c-type lectin *COLEC12*, and Neuronal\_07\_0 and  
294 Neuronal\_07\_5 which both express *NPTX1* and *PPP3CA* (**Figure 5f**). Altogether, these results suggest  
295 that the main driver of cellular diversification may be species-specific expansion, retention, and neo-  
296 functionalization of the expression patterns of gene families.

297

## 298 **Comparisons of cell types between species-morphs**

### 299 **Transcriptional differences in shared neuropeptidergic cell types**

300 Zebrafish and Mexican tetra last shared a common ancestor roughly 200 million years ago, yet  
301 we found that the degree of cell type conservation between these two teleost species was extensive. We  
302 therefore wondered whether the surface- and cave-morphs of Mexican tetra, which shared a common  
303 ancestor 250-500 thousand years ago, had any detectable cell type differences. To identify the repertoire  
304 of cellular diversity in Mexican tetra, we performed subclustering on the Mexican tetra dataset alone,  
305 resulting in 166 subclusters, including 19 subclusters that were species-morph specific, and 147  
306 subclusters shared between species-morphs (Extended Data Figure 7a-c and Supplementary Figure 8;  
307 See Supplementary Data for the full list of subclusters and associated marker genes for both zebrafish  
308 and Mexican tetra).

309 Neuronal subclusters shared by cave- and surface-morphs were characterized by low  
310 dendrogram distance, high *SI*, and a lack of enrichment for genes associated with divergent genomic  
311 windows (**Figure 6a**). The most transcriptionally different subclusters between morphs were the *galn*<sup>+</sup> and  
312 *otpa*<sup>+</sup>/*oxt*<sup>+</sup> subclusters (**Figure 6a**). Further examination of these subclusters revealed that surface-morph  
313 *galn*<sup>+</sup> cells expressed *galn* at a significantly higher level than cave-morph cells (Extended Data Figure 8a).  
314 Similarly to what we observed between species (**Figure 5**), surface- and cave-morph *oxt*<sup>+</sup> cells were  
315 distinguished by the differential expression of gene duplications. Surface-morph *oxt*<sup>+</sup> cells co-expressed  
316 *oxt* and its paralogs *avp* and *ENSAMXG00000021172*, whereas cave-morph *oxt*<sup>+</sup> cells only expressed *oxt*  
317 (Extended Data Figure 8b). Co-expression of *oxt* and *avp* was not observed in zebrafish, with each  
318 neuropeptide expressed in its own subcluster. These results highlight transcriptional changes in  
319 conserved cell types which may be associated with cave-adaptation.

320 It was recently reported that the expression of the neuropeptide *hcrt* is upregulated in Pachon  
321 cave-morphs, and is associated with increased sleep/wake activity compared to surface-morphs<sup>38,39</sup>.  
322 Additionally, our genetic analysis suggested that genes associated with circadian rhythm and  
323 neuropeptidergic cell types were under selection in cave morphs (Supplementary Information and  
324 Supplementary Figure 9). We wondered if we could use our single-cell data to identify changes in the  
325 transcription factors or regulatory network underlying the expression of *hcrt* and other neuropeptides and  
326 terminal effector genes between morphs. The majority of terminal effector GRNs were more conserved  
327 between species-morphs than between species, including the TFs associated with *galn* (Extended Data

328 Figure 8c-d). The TFs associated with the *hcrt* were poorly correlated between species-morphs, with the  
329 TF *creb3l1* more highly associated with *hcrt* expression in surface-morph cells, compared to cave-morph  
330 cells (Extended Data Figure 8e). High association between *hcrt* and *creb3l1* was not observed in the  
331 zebrafish data, indicating that this association may be specific to Mexican tetra, and responsible for the  
332 increased *hcrt* expression previously observed. This analysis also provided a potential mechanism for the  
333 co-expression of *oxt* and *avp* in Mexican-tetra compared to zebrafish. Three of the top TFs associated  
334 with *oxt* (*creb3l1*, *otpb*, and *sim1a*), were also predictive of *avp* expression (Extended Data Figure 8f-g).  
335 Differential expression of neuropeptides within conserved cell types may therefore be a common cave  
336 adaptation strategy across morphs.

337

### 338 **Species-morph specific subclusters are species-specific and express cell-state transcriptomes**

339 Of the 19 species-morph specific subclusters, 4 were neuronal, 3 were from glial populations, and  
340 12 were from the hematopoietic lineage (**Figure 6a-b** and Extended Data Figure 7a-b). The majority  
341 (11/19) of these species-morph specific subclusters mapped to integrated identities that were also  
342 specific to Mexican tetra. This included 3 of the 4 neuronal subclusters and 7 of the 12 immune  
343 subclusters that were species morph-specific (Extended Data Figure 7a-b). This suggests that many of  
344 the cell types specific to Mexican tetra are associated with or were co-opted during adaptation to the cave  
345 environment. Similarly, expression of Mexican tetra-specific genes was enriched in cell types from the  
346 hematopoietic lineage, which represented the majority of the species and species-morph specific  
347 subclusters (Figure 5a). Pachon cave morphs have been reported to have a smaller and less active  
348 immune system than surface morphs<sup>40</sup>. Though concluding changes in cell type proportions is difficult in  
349 single-cell experiments, we consistently observed fewer immune cells across independent cave-morph  
350 samples than in surface-morph samples (Extended Data Figure 2b, Figure 6a, Extended Data 9).  
351 Furthermore, hematopoietic lineage cells from cave-morphs expressed high levels of *ccr9a* and *sat1b*,  
352 and low levels of *fabp11a*, conditions which have been linked to inflammation resistance (Extended Data  
353 Figure 9c-d)<sup>41-44</sup>. Altogether, these results suggest that cave-morphs have a reduced immune system that  
354 expresses a neuro-inflammation resistance cell state transcriptome.

355 Three of the four species-morph specific neuronal subclusters mapped to integrated subclusters  
356 that were also species-specific: surface-morph specific Neuronal\_09-4 mapped to Mexican tetra specific  
357 Neuronal\_03\_13 (*pou4f2* and *etv1* positive), cave-morph specific Neuronal\_00\_1 and Neuronal\_03\_6  
358 mapped to Mexican tetra specific Neuronal\_07\_6 (*ENSAMXG00000025407+*), and Neuronal\_12\_5  
359 (*HPCAL1+*) respectively (Extended Data Figure 7a-b). The identity of the cave-specific neuronal  
360 subcluster Neuronal\_03\_4 was less clear. Cells from this subcluster mapped to a cell type shared  
361 between species, and expressed a set of marker genes that was conserved across species (*rtn4rl2a*,  
362 *rtn4rl2b*, *cd9b*, and *penkb*) (Extended Data Figure 10a-d). However, Neuronal\_03\_4 cells also expressed  
363 an additional gene signature, which included the genes *rcan1a* and *prelid3b* (Extended Data Figure 10c).  
364 GO analysis of these differentially expressed genes between Neuronal\_03\_4 and the highly similar  
365 shared Neuronal\_03\_1 revealed enrichment for terms related to stress response, protein folding, and  
366 translation, including the heat-shock genes *hspb1*, *hspa4a*, and prolyl isomerase *fkbp4* (Extended Data  
367 Figure 10e-f). These results indicate that an ancestral cell type found in both zebrafish and Mexican tetra  
368 acquired a stress response transcriptional program in the cave lineage, resulting in a morph-specific cell  
369 state (Extended Data Figure 10g).

370

### 371 **DISCUSSION:**

372 How evolution generates and shapes cellular diversity is largely unknown. In this study we used  
373 single-cell transcriptomics, high resolution clustering, and cross-species integration to compare cell types  
374 of the teleost hypothalamus between two divergent teleosts, zebrafish and Mexican tetra. First, we  
375 observe extensive conservation of cell-types across roughly 150 million years of evolution between  
376 zebrafish and Mexican tetra (>75% of all subclusters were shared), providing a high resolution

377 quantification of the molecular similarity between cell types across such a large phylogenetic distance.  
378 Second, we show that cell types conserved between species are characterised by subfunctionalization of  
379 paralogous gene expression patterns and by gene regulatory divergence. Third, we find that species-  
380 specific cell types were associated with the evolution of gene families, linking genetic novelty with cellular  
381 novelty. Fourth, we identify transcriptomic, cellular and genomic changes associated with cave-adaptation  
382 in Mexican tetra.

### 383 384 **Shared cell types are characterized by regulatory divergence and shifts in paralog expression**

385 Hundreds of cell types have been cataloged in the brains of vertebrates, including fish, mice and  
386 humans, but their conservation between species is unclear<sup>8,11,23,45</sup>. We observed extensive conservation  
387 of 75% of cell-types between zebrafish and Mexican tetra, who last shared a common ancestor more than  
388 150 million years ago, before the break-up of Pangea<sup>46,47</sup>. In our analysis, shared cell types were even  
389 more similar when taking paralog expression into account. Up to 20% of the transcriptomic divergence of  
390 shared cell types between species was from preferential expression of functionally similar paralogous  
391 genes. These expression pattern differences, or paralog shifts, suggest that shared cell types often  
392 express paralogous genes. Similarity and shifts in paralog expression between species was highest for  
393 progenitor cell types, and for clusters compared to subclusters. Changes to cluster and progenitor  
394 populations would likely have pleiotropic effects that may have prevented transcriptional divergence. A  
395 comparison of single-cell atlases across animal phyla has also demonstrated shifts in paralog expression  
396 for homologous cell types<sup>48</sup>. In that study the authors argue that paralog shifts may be due to genetic  
397 compensation by paralog substitution. However, our analysis suggests that divergence patterns of  
398 paralogous genes were mostly due to loss of redundancy, differed between species, and scaled with  
399 evolutionary gene age. This observation suggests that following ancestral gene duplication, expression  
400 patterns of paralogous genes are shifted, caused by independent sub-functionalization of gene  
401 expression patterns in each species. Further work will be necessary to determine the exact evolutionary  
402 and molecular mechanisms that generate paralog shifts in shared and homologous cell types.

403 We found that the expression patterns of transcription factors and their putative associations with  
404 specific classes of terminal effectors were less conserved than the expression patterns of the terminal  
405 effectors themselves. This observation contrasts with the high inter-species conservation of 'core' TFs  
406 expressed during early lineage determination events<sup>7</sup>. It agrees, however, with the low inter-species  
407 conservation in the expression of TFs 're-used' multiple times in different tissues<sup>7</sup>. The TF code  
408 associated with specific cell types, such as hypothalamic neurons, may therefore not be highly conserved  
409 between species. Alternatively, the differences in GRNs we observe might be caused by convergence in  
410 the GRNs of non-orthologous cell types to regulate terminal effector genes, as has been postulated for  
411 neurotransmitters in the *Drosophila* brain<sup>52</sup>. We note, however, that neuropeptidergic cell types and  
412 effector gene expression patterns were highly conserved between zebrafish and Mexican Tetra. We  
413 therefore favor a process akin to developmental systems drift, where conserved homologous traits  
414 between species can have divergent gene regulatory underpinnings caused by neutral drift<sup>53</sup>. We  
415 speculate that there might be cellular systems drift, where selection acts to maintain the functional output  
416 of cell types, rather than the regulatory mechanisms which generate or maintain them.

417 Together, our results paint a picture of the evolutionary history of the hypothalamic cell types in  
418 two teleost species. Cell types are highly conserved between species, yet divergence in paralog  
419 expression and regulatory associations is common. These patterns suggest an interplay between dosage  
420 compensation and subfunctionalization of expression patterns after genome duplication, neutral evolution  
421 causing shifts in paralog expression and regulatory divergence, and stabilizing selection maintaining cell  
422 type functions.

### 423 424 **Species specific cellular novelty is associated with species-specific genetic novelty and paralog** 425 **neo-functionalization**

426 Cross-species comparisons using single-cell sequencing data typically only consider orthologous  
427 genes between the species of interest, limiting the identification of species-specific innovations<sup>5</sup>. Here we  
428 find that the majority of species-specific cell types between zebrafish and Mexican tetra were enriched for  
429 the expression of non-homologous genes between species. This observation extends previous studies  
430 that have linked the evolution and diversification of biological traits with genetic novelty<sup>35,54</sup>. For example,  
431 expression of human specific genes in radial glia has been linked to cortical evolution and the expansion  
432 of the neocortex in primates<sup>55</sup>. Indeed, we found that the expression of jacalin lectins, which are specific  
433 to the zebrafish lineage<sup>36</sup>, are associated with a zebrafish-specific neuronal cell type. These results  
434 illustrate how species-specific genetic novelty underlies species-specific cellular novelty.

435 Moreover, our results suggest that the generation of new cell types within teleosts may be driven  
436 by species-specific neo-functionalization of paralogous genes. Many of the species-specific cell types  
437 were associated with expression of genes generated by recent duplication events. Furthermore, the loss  
438 of ancestrally duplicated paralogs in zebrafish (*HPCAL1*, *COLEC12*, *NPTX1*, and *PPP3CA*) was also  
439 associated with Mexican tetra specific cell types. These genes had expression patterns that differed from  
440 their paralogs, suggesting they have gained new functions since their duplication. Previous studies have  
441 suggested that new cell types are generated first through the birth of similar or homologous sister cell  
442 types<sup>56</sup>. Genetic individuation of sister cell types through the generation of distinct core regulatory  
443 complexes would then allow subsequent divergence through acquisition of different pre-existing gene  
444 modules<sup>3</sup>. Our results suggest an alternative scenario wherein gene duplication may precede or even  
445 drive the partitioning of cellular functions into distinct cell types (sister cell types). For example, amino  
446 acid substitutions between *vip* and *vipb* may have endowed different functionality, promoting the  
447 generation of the *vipb* subcluster in Mexican tetra. This scenario is reminiscent of the evolution of rod and  
448 cone cells following opsin gene duplication<sup>57,58</sup>. These observations suggest paralog neo-functionalization  
449 as a basis for cell type diversification.

450 Similar to the relationships between homologous genes, homologous cell types (shared cell types  
451 or sister-cell types) could refer to populations separated by a speciation event (orthologous cell types), or  
452 through a cell type duplication event (paralogous cell types)<sup>56</sup>. The shared populations we observed  
453 between zebrafish and Mexican tetra may therefore represent orthologous cell types which were present  
454 in the last common ancestor of both species. Species-specific cell types derived from cell type duplication  
455 events within species may be paralogous to cell types shared between species. Future work will be  
456 necessary to unravel the complicated evolutionary history of gene and cell type diversification.

457

### 458 **Single-cell transcriptomic signatures associated with cave-adaptation**

459 Mexican tetra Pachon cave-morphs have previously been reported to have a smaller and  
460 differentially active immune system<sup>40</sup>. Our results extend these observations to the Tinaja and Molino  
461 cave-lineages. In addition, we observe expression of a neuro-inflammation resistance signature in the  
462 immune cells of all three cave-morphs. Inflammation and neurodegeneration are intricately connected,  
463 and associated with aging in many species, including humans<sup>59</sup>. Negligible senescence has been  
464 reported in cave-morphs compared to surface-morphs<sup>60</sup>. It is therefore intriguing to speculate that the lack  
465 of immune inflammation in the nervous system may contribute to the lack of age-related senescence in  
466 cave-morphs.

467 The species-morphs of the Mexican tetra have divergent behavioural phenotypes which have  
468 previously been linked to the hypothalamus<sup>38,39,61,62</sup>. We observed differences in the expression patterns  
469 of several neuropeptides associated with these behaviours. For example, decreased *galn* expression in  
470 cave-morphs versus surface-morphs could partially explain the loss of sleep or changes in appetite or  
471 aggression in cave-morphs<sup>63-65</sup>. Alterations in oxytocin cells might also be linked to changes in appetite,  
472 or the lack of social interactions (schooling) observed in cave-morphs<sup>63,66</sup>. In the future, the single-cell  
473 atlases presented here will be a powerful resource to explore the behavioural differences between both  
474 species and species-morphs.





476

## 477 **METHODS**

478

### 479 **Husbandry of zebrafish and Mexican tetra**

480 All animal work was performed at the facilities of Harvard University, Faculty of Arts & Sciences  
481 (HU/FAS). Mexican tetra husbandry was performed as previously described<sup>67</sup>. This study was approved  
482 by the Harvard University/Faculty of Arts & Sciences Standing Committee on the Use of Animals in  
483 Research & Teaching under Protocol No. 25–08. The HU/FAS animal care and use program maintains  
484 full AAALAC accreditation, is assured with OLAW (A3593-01), and is currently registered with the USDA.

485

### 486 **Processing of samples for scRNA-seq**

487 Wild type adult zebrafish, and wild type adult Mexican tetra surface- and cave-morphs were used  
488 for scRNA-seq analysis. All zebrafish used were approximately 2-3 months old, and all Mexican tetra  
489 were between 1-2 years of age. For all zebrafish samples, tissues from 4-6 individual zebrafish were  
490 pooled for downstream dissociation then split into 4 samples for single-cell encapsulation. As their brains  
491 are much larger, each sample for Mexican tetra was composed of a single individual fish. A total of 16  
492 samples of zebrafish (8 males, 8 females), and 16 individual Mexican tetra were used (8 male, and 8  
493 female), including 8 Choy surface-morphs, 4 Pachon cave-morphs, 2 Tinaja cave-morphs, and 2 Molino  
494 cave-morphs split evenly between males and females.

495 The same procedure was used to collect and dissociate single-cells from both zebrafish and  
496 Mexican tetra. Animals were sacrificed by first placing them on ice, followed by decapitation. Whole brains  
497 were removed and immediately placed in 4% low-melt agarose mixed 50:50 with Neurobasal media plus  
498 B27 supplement (2% agarose final solution) (ThermoFisher). Once solidified, 500  $\mu$ m sections were  
499 obtained from whole brains mounted in agarose using a vibratome (Leica VT1000S). The hypothalamus  
500 and pre-optic area were then dissected from vibratome sections and dissociated into single cells using  
501 the Papin Dissociation Kit (Worthington) as previously described<sup>23</sup>. Cells were counted using a  
502 hemocytometer and resuspended at a final concentration of 1000 cells/ $\mu$ l in Neurobasal media  
503 (ThermoFisher). Samples were run on the 10X Genomics scRNA-seq platform according to the  
504 manufacturer's instructions (Single Cell 3' v2 kit). Libraries were processed according to the  
505 manufacturer's instructions (Single Cell 3' v2 kit). Transcriptome libraries were sequenced using Nextera  
506 75 cycle kits at the Bauer Core Facility (Harvard). Protocol for cell dissociation is available at  
507 [https://github.com/maxshafer/Cavefish\\_Paper](https://github.com/maxshafer/Cavefish_Paper).

508 We recovered between 2998 and 5490 cells per sample for the zebrafish dataset, and 3029 and  
509 5919 cells per sample for the Mexican tetra dataset. Samples had a minimum of 18,347 reads per cell  
510 with averages of 31,656 and 29,536 reads per cell for the zebrafish and Mexican tetra datasets,  
511 respectively. Sequencing saturation was between 68%-90%, with means of 78% for zebrafish samples  
512 and 83% for Mexican tetra samples. Between 30% and 60% of reads per sample for both datasets were  
513 mapped confidently to their respective transcriptome (78% - 88% to the genome).

514

### 515 **Bioinformatic processing of raw sequencing data and independent cell type clustering and 516 subclustering analysis**

517 Transcriptome sequencing data were processed using Cell Ranger 2.1.0 according to the 10X  
518 guidelines to obtain cell by gene expression matrices for each sample. For zebrafish, reads were mapped  
519 to a transcriptome constructed using the GRCz10 genome assembly annotated using the RefSeq  
520 genome annotation for GRCz10 (NCBI). For Mexican tetra, reads were mapped to a transcriptome  
521 constructed using the AstMex102 genome assembly annotated using the Ensembl genome annotation for  
522 AstMex102 (Ensembl). Clustering analysis was performed using Seurat v3.2.0<sup>24</sup>. Due to the lack of a  
523 mitochondrial genome for *A. mexicanus*, we opted not to remove cells with high mitochondrial content  
524 from the Zebrafish dataset. The following options were used for PCA, knn graph construction, and

525 clustering for both zebrafish and Mexican tetra. Only cells with between 200 and 2500 expressed genes  
526 were used (*nFeature\_RNA*). Variable features were obtained using the *mean.var.plot (mvp)* selection  
527 method as in Seurat v2.3.4. The identified variable features were used for PCA, and the top 50 PCs were  
528 used for clustering, though similar results were obtained with variable PC numbers. A *k* of 30 (*k.param*),  
529 and an error bound of 0.5 (*nn.eps*) were used for constructing the Shared Nearest Neighbor (SNN) graph.  
530 Clusters were called using a resolution (*resolution*) of 0.6 using the original Louvain algorithm. Shared  
531 marker genes for each cluster were obtained using Seurat's FindConservedMarkers function, and  
532 species-specific marker genes were identified by first subsetting the Seurat object by species before  
533 running the FindMarkers function for each cell cluster and subcluster. These genes, as well as genes with  
534 known expression patterns in neuronal and hypothalamic cell types were used to annotate subclusters.  
535 Clusters were identified as GABAergic or Glutamatergic based on which marker genes they expressed  
536 most highly (*slc17a6a/slc17a6b* or *gad1b/gad2*). In many cases, clusters expressed markers of both, due  
537 to having both GABAergic and Glutamatergic cells and were therefore all annotated as Neuronal.

538 Independent subclustering analysis was performed by first subsetting the zebrafish or Mexican  
539 tetra data into individual clusters, then performing all of the steps of Seurat clustering on each cluster  
540 independently, including finding highly variable genes and principal component analysis. Parameters  
541 used for subclustering were the same as for clusters, except we used the resolution 0.4 (*resolution*) and  
542 15 PCs derived from the variable features of the cells in each cell cluster. Because different sets of  
543 variable genes were used for subclustering and construction of the tSNE projection for the full dataset,  
544 the positions of subcluster labels are not necessarily representative of the true differences between  
545 subclusters. Full analysis scripts for cell type clustering, R objects, and raw sequencing, including all  
546 variables used are available on GitHub ([https://github.com/maxshafer/Cavefish\\_Paper](https://github.com/maxshafer/Cavefish_Paper)). Raw count data  
547 is available in the **Supplemental Data**, and raw sequencing data is available on NCBI GEO.

548

#### 549 **Dataset integration and integrated clustering and subclustering analysis**

550 Datasets were initially integrated using Seurat's MergeObjects function. Given the large biological  
551 batch effects between the species, cells first clustered by species, then by cell type. To correct for  
552 species-specific batch effects, and identify shared and species-specific cell types we used Seurat v3.0.0  
553 to integrate the zebrafish and Mexican tetra datasets. Seurat uses Canonical Correlation Analysis (CCA)  
554 to identify correlated changes in the transcriptomes of cell types between species, and identifies the most  
555 similar clusters across species using Mutual Nearest Neighbour (MNN) analysis. This allows identification  
556 of cluster specific batch correction vectors, which are used to correct the expression values of a subset of  
557 genes between species. These genes and their corrected expression values are then used for  
558 dimensionality reduction and clustering analysis. All genes from both datasets were used in the  
559 integration process, and orthologous genes were identified by matching gene names. The best results  
560 were obtained by first clustering using 100 *dims*, a *k.param* of 20, a *res* of 0.15, and an *nn.eps* of 0, which  
561 segregated all non-neuronal cells into appropriate clusters. Following this, we used expression of the  
562 neuronal marker gene *gng3* to identify and combine all neuronal cell clusters together. To cluster the  
563 neuronal cells, we used 10 *dims*, a *k.param* of 20, and an *nn.eps* of 0 which generated the 14 neuronal  
564 populations used in this study. Integrated subclustering analysis was performed by first subsetting the  
565 integrated data into individual clusters, then performing all of the steps of Seurat integration and  
566 clustering on the zebrafish and Mexican tetra cells from each cluster independently. Integration, including  
567 CCA and MNN analysis, was performed independently on each integrated cluster to maximize the gene  
568 information used to identify shared and species-specific cellular heterogeneity. For subclustering we used  
569 between 5 and 50 *dims* for each cluster depending on the number of cells, and a *res* of 0.25. Following  
570 integration, several subclusters were identified as aberrant, and expressed marker genes from both non-  
571 neuronal and neuronal subclusters. These populations appeared to be created by the integration and  
572 batch correction process, and were derived mainly from Erythrocyte cells from both species. These  
573 subclusters were removed from downstream analysis.

574

### 575 **Trinarization of gene expression patterns**

576 To identify genes robustly expressed by each population (clusters and subclusters) we calculated  
577 trinarization scores for each gene per cluster and subcluster<sup>27</sup>. Trinarization scores represent the  
578 probability that each gene is actually expressed by each population, based on the detection frequency of  
579 each gene in each population, and the posterior distribution of the underlying population frequency of  
580 expression ( $\square$ ). We used a Bayesian beta-binomial model to trinarized the data and calculate the  
581 posterior distribution ( $\square$ ) using hyperparameter values of 1.5 for  $a$ , and 1 for  $b$  as previously described<sup>27</sup>.  
582 We called a gene expressed if  $P(\square > f) > (1 - \text{PEP})$ , where the fraction of cells expressing each gene ( $f$ ) is  
583 0.1, and the Posterior Error Probability (PEP) of 0.05. Similar results were obtained using values of 0.2  
584 and 0.35 for  $f$ . Trinarization scores were calculated for all clusters and subclusters, for the zebrafish,  
585 Mexican tetra, surface-morph, and cave-morph datasets (**Supplemental Data**). These scores were also  
586 used to determine expression patterns of duplicated genes for the calculation of expression divergence  
587 ( $dT$ ).

588

### 589 **Marker gene identification and calculation of the transcriptional similarity index (SI)**

590 To annotate subclusters and identify genes whose expression was enriched within clusters and  
591 subclusters we used Seurat to find marker genes for each population. This was done for all clusters and  
592 subclusters in both the zebrafish and Mexican tetra datasets, as well as for all of the integrated clusters  
593 and subclusters in the combined dataset. For the integrated clusters and subclusters, we used the  
594 uncorrected expression data (`DefaultAssay(object) <- "RNA"`), which allowed the detection of species-  
595 specific gene expression patterns. For each population, we identified marker genes independently for the  
596 zebrafish and Mexican tetra cells within that cluster or subcluster. To identify shared marker genes for  
597 each population, we used Seurat's `FindConservedMarkers` function, which uses meta analysis of  
598 statistical values for each gene in the marker genes for each species. For all cases we used the following  
599 variables for `FindMarkers` and `FindConservedMarkers`; `logfc.threshold` of 0.25 (default), `min.pct` of 0.1  
600 (default), `min.cells.per.ident` of 1000, and "wilcox" for `test.use`. Species-specific marker genes for each  
601 population were defined as the set difference between the marker genes for one species and conserved  
602 marker genes (Figure 2a). These lists were then used to calculate the Similarity Index ( $SI$ ) for each cluster  
603 and subcluster between zebrafish and Mexican tetra.  $SI$  was calculated with the following equation, where  
604  $G_T$  is the shared set of marker genes, and  $G_A$  and  $G_B$  are the total number of marker genes for species  $A$   
605 and  $B$ , including both species-specific and shared marker genes<sup>68</sup>.

606

$$SI = 1 - \sqrt{\left(1 - \frac{G_T}{G_B}\right) * \left(1 - \frac{G_T}{G_A}\right)}$$

607 The same procedure was used to identify species-morph specific marker genes, and marker  
608 genes conserved between species-morphs for each population within the Mexican tetra data for  
609 calculation of the  $SI$  between species-morphs. To calculate  $SI$  between across all sets of integrated  
610 clusters and subclusters, we used the conserved marker genes for each population, and compared their  
611 p-values using the same procedure as in Seurat's `FindConservedMarkers` function - using the `minimump`  
612 function from `metap` package - to determine shared marker genes for each pair of cluster or subclusters.  
613 We then calculated  $SI$  as above. All marker gene sets were filtered to contain only those genes which  
614 also passed the trinarization threshold for that population, and for each dataset, prior to calculation of  $SI$ .

615

### 616 **Paralog Identification and enrichment analysis across species**

617 Paralogous gene pairs and orthology confidence and gene order scores were identified using the  
618 Ensembl BioMart service, and accessed using the `biomaRt` R package<sup>69</sup>. For each gene that was

619 specifically expressed in one species, we identified all corresponding paralogous genes and determined if  
620 any of these genes were present in the conserved marker genes, or the marker genes specific to the  
621 other species (**Figure 2a**). This was done for all clusters and subclusters shared between zebrafish and  
622 Mexican tetra. Fisher's exact test was performed to calculate statistical enrichment for paralogous genes  
623 for each cluster and subcluster using the *fdrtool* R package<sup>70</sup>. The remaining species specific genes  
624 (those that were not paralogs of a conserved, or opposite species-specific gene) were then subjected to  
625 gene ontology analysis. Species-specific marker genes for each subcluster were pooled by cluster, and  
626 the RDAVIDWebService R package was used to submit each list for GO analysis by DAVID<sup>71</sup>.

627

### 628 **Calculation of expression divergence ( $dT$ ), redundancy, and gene expression expansion**

629 For both zebrafish and Mexican tetra, we used Ensembl's Biomart tool to identify paralogous  
630 gene pairs in both species, and the last common ancestor which shares each gene duplication  
631 (**Supplemental Data**). To calculate the expression divergence for paralogous gene pairs, we calculated  
632 the following for each paralogous gene pair in each species. The subclusters which expressed each gene  
633 above the trinarization threshold were used as input for the calculation of expression divergence ( $dT$ ) as  
634 previously described<sup>31</sup>. Expression divergence was calculated for both zebrafish and Mexican tetra  
635 separately by comparing the number of subclusters that express either paralog ( $n_{tu}$ ) to the number of  
636 subclusters that express both paralogs ( $n_{ti}$ ), with the following equation.

637

$$dT = \frac{(n_{tu} - n_{ti})}{n_{ti}}$$

638 Gene pairs where neither gene was expressed in our datasets were not included in this analysis  
639 To determine the expression pattern of the ancestral gene prior to duplication in the common ancestor of  
640 zebrafish and Mexican tetra (putative ancestral gene expression patterns), we used the intersection of the  
641 subclusters that expressed either or both paralogs in both species (**Figure 3c**) Redundancy of  
642 expression for paralogous gene pairs was calculated as  $1 - dT$  within ancestral cell types, representing  
643 how much of the ancestral gene expression pattern was conserved between paralog gene pairs. The  
644 number of non-ancestral subclusters which only expressed either or both paralogs in only one species  
645 was used to determine the amount of expansion of paralogous gene expression patterns. Importantly,  
646 these metrics cannot account for the possibility that ancestral expression in one or more cell types was  
647 lost for both paralogs in one species, but retained for at least one paralog in the other. Paralogous genes  
648 generated by the vertebrate 2R or teleost 3R whole genome duplication events (Ohnologs) were  
649 identified from the OHNOLOG repository using the "Strict 2R" and "Strict 3R" datasets for zebrafish (*D.*  
650 *rerio*) (<https://ohnologs.curie.fr>)<sup>72</sup>.

651

### 652 **Identification and analysis of putative gene regulatory networks using SCENIC/GENIE3**

653 We used SCENIC/GENIE3 to identify transcription factors (TFs) that were predictive of the  
654 expression of terminal effector genes associated with the functions of the hypothalamus, including  
655 neuropeptides, neurotransmitter or synapse associated genes, and ion channel genes<sup>33,34</sup>. This  
656 analysis outputs numerical weights for the association between each TF and each terminal effector gene  
657 in the two species, which are used to determine TF sets for each terminal effector gene. We  
658 downsampled the zebrafish and Mexican tetra datasets to ensure equal cell numbers across subclusters.  
659 This analysis uses single-cell information, but is independent of cell cluster and subcluster identities. Lists  
660 of transcription factors, neuropeptides, neurotransmitter related, synaptic, and ion channel genes were  
661 identified using ZebrafishMine, and used to identify orthologous genes in Mexican tetra with the same  
662 ontology<sup>73</sup>. Gene lists used in this study are available at ([https://github.com/maxshafer/Cavefish\\_Paper](https://github.com/maxshafer/Cavefish_Paper)).  
663 For the list of transcription factors, we used search terms "transcription" and "transcription factor activity",  
664 which resulted in a combined list of 3141 unique gene names. Datasets were first downsampled such that

665 the same number of cells from each subcluster were included to reduce the effects of differential  
666 subcluster abundances between species on the GRN analysis. Cutoff values for *minSamples* and  
667 *minCountsPerGene* for the *geneFiltering* function used to filter out lowly detected or expressed genes  
668 were determined such that *hcrt* was included in all analyses. SCENIC then uses the GENIE3 algorithm to  
669 generate random forest weights for each transcription factor and target gene, based on the predictive  
670 power of each transcription factor in determining the expression level for each target gene. The lists of  
671 transcription factors and their corresponding weights for each target gene were used in downstream  
672 analysis. We used the “top50” cutoff from SCENIC to determine transcription factors to calculate the *S'*  
673 between species or species-morphs. The same procedure was used for analysis of Mexican tetra  
674 surface- and cave-morphs. We used customized versions of some SCENIC functions, including  
675 *geneFiltering*, *runGenie3*, and *runSCENIC\_1\_coexNetwork2modules*, to allow use of our gene lists, and  
676 to allow easier implementation on a laptop (specifically the ability to stop and restart the analysis).

677 To test whether divergence in all TF sets were mitigated by association with paralogous TFs, we  
678 calculated paralog enrichment in the species-specific TFs for each terminal effector gene as done  
679 previously (Figure 2). The majority of TF sets (260 out of 435, 60%) were composed of roughly the  
680 number of paralogs expected by random chance (odds ratio ~ 1), and 26 terminal effector genes had  
681 significantly fewer paralogs than expected (Figure 4d). Therefore, divergence in the gene regulatory  
682 networks of neuropeptides and other terminal effector classes is not compensated through the expression  
683 of paralogous TFs.

684

#### 685 **Analysis of species-specific cell types and identification of species-specific genes**

686 Clusters which were specific to either a species or a species-morph were identified by calculating  
687 the proportion of cells which came from each species or species-morph for each cluster or subcluster.  
688 Clusters which were composed of 90% or more cells from one species or species-morph were considered  
689 specific to that species. Cells belonging to other species or species-morphs in those clusters were likely  
690 incorrectly assigned due to the limited gene information used for integrated clustering and subclustering  
691 analysis. These subclusters were not enriched for technical artifacts due to the encapsulation or  
692 clustering methods used (Supplementary Figure 7). We examined the genes that differentiated each  
693 species-specific subcluster from the other cells in the same parent cluster for clues to their origins or  
694 functions (Figure 5). Two of these subclusters were distinguished by the expression of unique TFs  
695 (Neuronal\_03\_13 and Neuronal\_04\_5). For example, the zebrafish-specific Neuronal\_04\_5 subcluster  
696 expressed *meis1b*, *six6a*, and *six6b* in addition to the parent cluster marker genes *cbln1* and *adcyap1b*.  
697 Other species-specific neuronal subclusters were characterized by cell cycle genes (Neuronal\_02\_3),  
698 genes related to axonal guidance or remodelling (Neuronal\_04\_6), or expressed different  
699 neurotransmitters (Neuronal\_05\_0 and Neuronal\_05\_1) (Figure 5). These subclusters may therefore  
700 reflect different temporal or spatial cell states, captured in one species or the other. The presence,  
701 absence, and orthology of the specific duplicated genes discussed in the current report, including *vipb*,  
702 *HPCAL1*, and the jacalin lectin genes, was confirmed using the most recent Ensembl release (Ensembl  
703 Release 101), which includes newer versions of both the zebrafish (GRCz11) and Mexican tetra  
704 (*Astyanax mexicanus*-2.0) genome assemblies<sup>34</sup>. For analysis of subclustering in the absence of non-  
705 homologous genes, all non-homologous genes were removed from the Variable Features for each cluster  
706 prior to subclustering using the same parameters as above.

707

708

#### 709 **Construction of Sankey diagrams and other plots**

710 Sankey diagrams were constructed using the *networkD3* R package. We used the Seurat  
711 wrappers for *ggplot2* functions to construct tSNE graphs and DotPlots of expression values across  
712 clusters or subclusters. Custom R scripts were used to construct the rest of the plots using *ggplot2*,  
713 including the gene ontology analysis, and the multi-layered circular plots (Figure 6a). Graphical tables



714 were constructed using the formatable R package. All scripts used to construct figures are available on  
715 GitHub (<https://github.com/maxshafer/Cavefish Paper>). Final figures were assembled using Affinity  
716 Designer (Serif Europe).

717 Species and species-morph dendrograms, as well as subcluster dendrograms were constructed  
718 using both the *Seurat*, *ggtree*, and *phylogram* R packages<sup>24,77,78</sup>. Pseudo-bulk expression data for each  
719 cluster and subcluster were used to calculate the dendrogram dissimilarity values for **Figure 6a**. For  
720 calculation of the similarity between species and species-morph, pseudo-bulk expression data was  
721 generated for zebrafish, Choy surface, and Pachon, Tinaja, and Molino cave-morph samples by  
722 averaging the expression of each gene across all cells within each cluster and subcluster. Species and  
723 species-morph dendrograms were then generated for each population, based on the similarity in whole  
724 transcriptomes using the *BuildClusterTree* function in *Seurat*. The *ggtree* package was used to construct  
725 the density dendrogram, where the colour of the edges corresponds to the number of subclusters which  
726 support each arrangement of the dendrogram. The distance between surface- and cave-morph versions  
727 of each subcluster on the dendrogram was used for plotting.

728

### 729 **DATA AVAILABILITY**

730 Processed single-cell RNAseq counts and metadata, marker gene lists, trinarized gene lists, SI  
731 results, SCENIC results, results from genetic analysis, and GO lists are available as supplementary data.  
732 Raw sequencing results are available at the Sequence Read Archive (SRA) under BioProject ID  
733 PRJNA754013.

734

### 735 **ACKNOWLEDGEMENTS**

736 We thank Clifford J. Tabin for providing *Astyanax mexicanus* samples, and advice on  
737 experimental design. We thank members of the Schier lab for discussion and advice, including Drs. B.  
738 Raj, J. Liu, P. Abitua, and A. Nichols, and the Harvard zebrafish and cavefish facilities staff, including  
739 Brian Martineau, for technical support. We thank Drs. Gray Camp, Walter Salzburger and Nathalie  
740 Jurisch-Yaksi for helpful comments on the manuscript. This work was supported by a postdoctoral  
741 fellowship from the Canadian Institutes of Health Research (CIHR) to M.E.R.S., a grant from the Swiss  
742 National Science Foundation (SNSF) to M.E.R.S. (SPARK 196313), grants from SNSF (SPARK 195955)  
743 and the University of Basel to A.N.S., and an NIH grant (DP1HD094764), an ERC Advanced grant  
744 (834788), an Allen Discovery Center grant, and a McKnight Foundation Technological Innovations in  
745 Neuroscience Award to A.F.S..

746

### 747 **AUTHOR CONTRIBUTIONS STATEMENT**

748 M.E.R.S. and A.F.S. conceived and designed the study. A.N.S. and M.E.R.S. conceived and  
749 performed Similarity Index analysis. M.E.R.S. performed all other experiments and analysis, including  
750 scRNA-seq experiments, and all bioinformatic analysis. M.E.R.S., A.N.S., and A.F.S. wrote the  
751 manuscript. All authors read and approved of the manuscript.

752

### 753 **COMPETING INTERESTS STATEMENT**

754 The authors declare no competing interests.

755

### 756 **FIGURE LEGENDS**

757

758 **Figure 1. Integration of zebrafish and Mexican tetra single-cell data reveals extensive conservation**  
759 **of cell types**

760 (a) UMAP reduction of integrated zebrafish and Mexican tetra cells coloured by species. Datasets were  
761 integrated with Mutual Nearest Neighbour (MNN) and Canonical Correlation Analysis (CCA) using *Seurat*.

762 (b) UMAP reduction of integrated zebrafish and Mexican tetra cells coloured by annotated cell type. (c)

763 Sankey diagram of relationships between zebrafish, integrated, and Mexican tetra annotated clusters  
764 from Figure S3. Heights of squares and thickness of connecting lines are relative to the number of cells  
765 per identity or connection, respectively. (d) Circular heatmap of the proportion of zebrafish (dark blue), or  
766 Mexican tetra (yellow) cells per integrated subcluster. Subclusters are grouped first by cluster, and  
767 clusters are arranged by the dendrogram of cluster similarity, shown in the center of the circular heatmap.  
768 Red outlines indicate subclusters with > 90% of cells from one species (species-specific). (e) Density  
769 dendrogram for all shared subclusters across species and species-morphs. The density dendrogram was  
770 constructed using dendrograms for the similarity between species and species-morph versions of each  
771 subcluster identity shared between zebrafish and Mexican tetra. Darkness of lines indicate the level of  
772 support for each branch.

773

## 774 **Figure 2. Shared subclusters are highly similar between species and express paralogous genes**

775 (a) Diagram illustrating the relationships between gene sets for each subcluster used in the current study  
776 for calculation of the Similarity Index ( $SI$ ).  $G_A$  and  $G_B$  represent the total marker gene sets in the two  
777 species examined.  $G_T$  are genes that are found in both  $G_A$  and  $G_B$ . Species-specific marker genes are  
778 those which do not overlap with the other species ( $G_A$  or  $G_B$  minus  $G_T$ ), though may be paralogs of a  
779 marker gene from the other species. (b) The row-scaled  $SI$  for all subclusters between zebrafish and  
780 Mexican tetra based on marker genes filtered for genes which pass the trinarization threshold. Yellow  
781 indicates the highest  $SI$  value among Mexican tetra subclusters for each Zebrafish subcluster. (c) The  
782 Similarity Index ( $SI$ ) for progenitor and differentiated neuronal subclusters between zebrafish and Mexican  
783 tetra based on marker genes filtered for genes which pass the trinarization threshold. Two sample t-test  
784 p-value = 0.007661. (d) The  $SI$  for clusters and the mean of the  $SI$  for subclusters grouped by cluster  
785 between zebrafish and Mexican tetra coloured by cluster based on marker genes filtered for genes which  
786 pass the trinarization threshold. Paired t-test p-value = 0.003012. (e) Comparison of the  $SI$  for the same  
787 subclusters between species (purple), and between species-morphs (yellow), calculated using marker  
788 gene sets. (f) The percentage of species-specific marker genes for each subcluster which were paralogs  
789 of either the conserved or opposite species-specific marker gene for zebrafish and Mexican tetra. (g) The  
790 odds ratio for the enrichment of paralogs in the species-specific genes for each subcluster for zebrafish  
791 and Mexican tetra. (h) The percentage of paralog shifting events shared by both surface- and cave-  
792 morphs of Mexican tetra. Paralog shifts are separated by whether they were from the zebrafish or  
793 Mexican tetra species-specific marker genes. (i) Relationship between the Similarity Index and the mean  
794 of the percentage paralogs (for zebrafish and Mexican tetra) for each subcluster. (j) Comparison of the  $SI$   
795 (blue) and corrected- $SI$  (yellow) for subclusters between zebrafish and Mexican tetra. For all boxplots,  
796 box bounds represent the first and third quartiles and whiskers 1.5 times the interquartile range, thicker  
797 line represents the median.

798

## 799 **Figure 3. Paralog shifts are due to differential divergence after duplication between species**

800 (a) Example gene tree for a gene that is duplicated once in the common ancestor of two extant species  
801 (Species 1 and Species 2). (b) Model for paralog gene expression pattern divergence after gene  
802 duplication in a common ancestor. Over time (c), expression of Gene Xa is retained in the cell types of  
803 interest (filled rectangles) in Species 1 and Gene Xb is lost (empty rectangles), whereas expression of  
804 only Gene Xb is retained in the same cell types in Species 2. Ancestral cell types are defined as any cell  
805 type that expresses either paralog in both species (yellow highlight). Gene expression pattern expansion  
806 is represented by cell types that express only 1 paralog in 1 species (green highlight). (d) Empirical  
807 cumulative distribution function (ECDF) for expression divergence ( $dT$ ) for gene pairs grouped by their  
808 last common ancestor in zebrafish. From the oldest (Opisthokonta, yellow), to the most recent common  
809 ancestor (Otophysi, red), and to those gene duplicates which are only found in Danio rerio (dark red). (e)  
810 ECDF for expression divergence ( $dT$ ) for gene pairs grouped by their last common ancestor in Mexican  
811 tetra. From the oldest (Opisthokonta, yellow), to the most recent common ancestor (Otophysi, red), and to

812 those gene duplicates which are only found in *Astyanax mexicanus* (dark red). (f) Relationship between  
813 the expression divergence ( $dT$ ) for gene pairs in zebrafish and the expression divergence ( $dT$ ) for gene  
814 pairs in Mexican tetra for all gene pairs (black dots). (g) Relationship between the expression divergence  
815 ( $dT$ ) for gene pairs which arose in the last common ancestor of zebrafish and Mexican tetra (Otophysi,  
816 pink). (h) DotPlot of the expression of the paralog pairs *etv5a* / *etv5b* in zebrafish (top) and Mexican tetra  
817 (bottom) across cell clusters. (i) Ridge plots of the Pearson correlation of the binarized expression  
818 patterns across subclusters for gene pairs shared by zebrafish and Mexican tetra grouped by their last  
819 common ancestor or by their originating WGD event.

820

#### 821 **Figure 4. Divergence of gene regulatory networks underlying neuronal genes**

822 (a) Random forest weight for orthologous transcription factors in the zebrafish (y-axis) and Mexican tetra  
823 (x-axis) data for the neuropeptide *vip*. Colours indicate whether those transcription factors are in the top  
824 2% of transcription factors for each gene in either zebrafish (blue) and Mexican tetra (red), both (yellow),  
825 or none (black). (b) Similarity Index between the transcription factor sets for zebrafish and Mexican tetra  
826 for neuropeptides, neurotransmitters, synaptic genes, and ion channels. (c) Correlation between the  
827 random forest weights for transcription factor sets associated with each for neuropeptide,  
828 neurotransmitter, synaptic, or ion channel genes between zebrafish and Mexican tetra. (d) Odds ratio  
829 from Fisher's exact test for the enrichment for paralogous genes in the transcription factors associated  
830 with each gene, red dots indicated significant enrichment. (e) Similarity Index for all subclusters shared  
831 between zebrafish and Mexican tetra using either only neuropeptides and neurotransmitter related genes  
832 (purple), only transcription factors (green), or all marker genes (yellow). (f) Similarity Index for individual  
833 neuropeptidergic GABA\_1 subclusters between zebrafish and Mexican tetra using either only  
834 neuropeptides and neurotransmitter related genes (purple), only transcription factors (yellow), or all  
835 marker genes. For all boxplots, box bounds represent the first and third quartiles and whiskers 1.5 times  
836 the interquartile range, thicker line represents the median.

837

#### 838 **Figure 5. Species-specific subclusters are associated with species-specific genes**

839 (a) Sankey diagram of shared and species-specific subclusters, indicating the species (zebrafish or  
840 Mexican tetra) and the cellular lineage they belong to (Ciliated, Glial, Hematopoietic, or Neuronal). (b)  
841 The percentage of expressed genes (counts > 10) in each zebrafish subcluster which are non-  
842 homologous genes between species in shared or species specific cell subclusters. (c) The percentage of  
843 expressed genes (counts > 10) in each Mexican tetra subcluster which are non-homologous genes  
844 between species in shared or species specific cell subclusters. (d) The percentage of expressed genes  
845 (counts > 10) in each zebrafish subcluster which are non-homologous genes between species in neuronal  
846 or non-neuronal subclusters. (e) The percentage of expressed genes (counts > 10) in each Mexican tetra  
847 subcluster which are non-homologous genes between species in neuronal or non-neuronal subclusters.  
848 (c) DotPlot of the species-specific subcluster marker genes (y-axis) across subclusters (x-axis). Blue  
849 boxes highlight expression of specific paralogous genes in different subclusters. Gene expression is  
850 quantified by both the percentage of cells which express each gene (dot size) and the average  
851 expression in those cells (colour scale). For all boxplots, box bounds represent the first and third quartiles  
852 and whiskers 1.5 times the interquartile range, thicker line represents the median.

853

#### 854 **Figure 6. Divergence in subcluster repertoires and transcriptomes across Pachon, Tinaja, and 855 Molino cave-morphs**

856 (a) Graphical summary of subcluster and transcriptional differences between Mexican tetra surface- and  
857 cave-morphs, and between Pachon, Tinaja, and Molino cave-morphs. The first layer indicates the cluster  
858 identity (from Figure S2b), and the text label indicates the subcluster (from Figure S8). The second layer  
859 indicates the proportion of cells in each subcluster that come from a surface-morph (green) or a cave-  
860 morph (yellow). Red outlines indicate morph-specific subclusters (> 90% of cells from either surface- or

861 cave-morphs). Third layer indicates the proportion of cave-morph cells from each subcluster that come  
862 from the Pachon (orange), Tinaja (blue), or Molino (green) cave-morph samples. The fourth layer displays  
863 the Similarity Index between the surface-morph, cave-morph for shared marker genes for each  
864 subcluster. The fifth layer displays the percentage of marker genes for each subcluster that is also  
865 associated with a divergent genomic window ( $F_{ST}$  genes). Finally, the sixth layer displays the Dendrogram  
866 Distance, which is the distance between the surface- and cave-morph versions of each subcluster on a  
867 dendrogram based on the subcluster transcriptomes. **(b)** Dendrogram of the relationships between  
868 species and species-morphs in this study coloured by the number of cell type changes normalized by the  
869 evolutionary time determined from<sup>18</sup>.

870 **REFERENCES**

- 871 1. S. Ramón y Cajal. Histology of the human nervous system & vertebrates. *Maloine Paris* (1911).
- 872 2. Zeng, H. & Sanes, J. R. Neuronal cell-type classification: challenges, opportunities and the path forward. *Nat. Rev.*  
873 *Neurosci.* **18**, 530–546 (2017).
- 874 3. Arendt, D., Bertucci, P. Y., Achim, K. & Musser, J. M. Evolution of neuronal types and families. *Curr. Opin. Neurobiol.*  
875 **56**, 144–152 (2019).
- 876 4. Ohno, S. *Evolution by Gene Duplication*. (Springer-Verlag, 1970). doi:10.1007/978-3-642-86659-3.
- 877 5. Shafer, M. E. R. Cross-Species Analysis of Single-Cell Transcriptomic Data. *Front. Cell Dev. Biol.* **7**, (2019).
- 878 6. Tosches, M. A. *et al.* Evolution of pallium, hippocampus, and cortical cell types revealed by single-cell transcriptomics  
879 in reptiles. *Science* **360**, 881–888 (2018).
- 880 7. Briggs, J. A. *et al.* The dynamics of gene expression in vertebrate embryogenesis at single-cell resolution. *Science* **360**,  
881 (2018).
- 882 8. La Manno, G. *et al.* Molecular Diversity of Midbrain Development in Mouse, Human, and Stem Cells. *Cell* **167**, 566-  
883 580.e19 (2016).
- 884 9. Pollen, A. A. *et al.* Establishing Cerebral Organoids as Models of Human-Specific Brain Evolution. *Cell* **176**, 743-  
885 756.e17 (2019).
- 886 10. Kanton, S. *et al.* Organoid single-cell genomic atlas uncovers human-specific features of brain development. *Nature*  
887 **574**, 418–422 (2019).
- 888 11. Krienen, F. M. *et al.* Innovations present in the primate interneuron repertoire. *Nature* 1–8 (2020) doi:10.1038/s41586-  
889 020-2781-z.
- 890 12. Pasquier, J. *et al.* Gene evolution and gene expression after whole genome duplication in fish: the PhyloFish database.  
891 *BMC Genomics* **17**, 368 (2016).
- 892 13. Ravi, V. & Venkatesh, B. The Divergent Genomes of Teleosts. *Annu. Rev. Anim. Biosci.* **6**, 47–68 (2018).
- 893 14. Voldoire, E., Brunet, F., Naville, M., Volff, J.-N. & Galiana, D. Expansion by whole genome duplication and evolution of  
894 the sox gene family in teleost fish. *PLOS ONE* **12**, e0180936 (2017).
- 895 15. Gross, J. B. The complex origin of *Astyanax* cavefish. *BMC Evol. Biol.* **12**, 105 (2012).
- 896 16. Keene, A., Yoshizawa, M. & McGaugh, S. E. *Biology and Evolution of the Mexican Cavefish*. (Academic Press, 2015).
- 897 17. Mitchell, R. W., Russell, W. H. & Elliott, W. R. *Mexican Eyeless Characin Fishes, Genus Astyanax: Environment,*  
898 *Distribution, and Evolution*. (Texas Tech Press, 1977).
- 899 18. Herman, A. *et al.* The role of gene flow in rapid and repeated evolution of cave-related traits in Mexican tetra, *Astyanax*  
900 *mexicanus*. *Mol. Ecol.* **27**, 4397–4416 (2018).
- 901 19. Prober, D. A., Rihel, J., Onah, A. A., Sung, R.-J. & Schier, A. F. Hypocretin/Orexin Overexpression Induces An  
902 Insomnia-Like Phenotype in Zebrafish. *J. Neurosci.* **26**, 13400–13410 (2006).

- 903 20. Richter, C., Woods, I. G. & Schier, A. F. Neuropeptidergic Control of Sleep and Wakefulness. *Annu. Rev. Neurosci.* **37**,  
904 503–531 (2014).
- 905 21. Xie, Y. & Dorsky, R. I. Development of the hypothalamus: conservation, modification and innovation. *Development* **144**,  
906 1588–1599 (2017).
- 907 22. Denes, A. S. *et al.* Molecular Architecture of Annelid Nerve Cord Supports Common Origin of Nervous System  
908 Centralization in Bilateria. *Cell* **129**, 277–288 (2007).
- 909 23. Raj, B. *et al.* Simultaneous single-cell profiling of lineages and cell types in the vertebrate brain. *Nat. Biotechnol.* **36**,  
910 442–450 (2018).
- 911 24. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888–1902.e21 (2019).
- 912 25. D’Gama, P. P. *et al.* Diversity and Function of Motile Ciliated Cell Types within Ependymal Lineages of the Zebrafish  
913 Brain. *bioRxiv* 2021.02.17.431442 (2021) doi:10.1101/2021.02.17.431442.
- 914 26. Liu, J. *et al.* Evolutionarily conserved regulation of hypocretin neuron specification by Lhx9. *Development* **142**, 1113–  
915 1124 (2015).
- 916 27. Zeisel, A. *et al.* Molecular Architecture of the Mouse Nervous System. *Cell* **174**, 999–1014.e22 (2018).
- 917 28. Liang, C., Musser, J. M., Cloutier, A., Prum, R. O. & Wagner, G. P. Pervasive Correlated Evolution in Gene Expression  
918 Shapes Cell and Tissue Type Transcriptomes. *Genome Biol. Evol.* **10**, 538–552 (2018).
- 919 29. Gu, X. Understanding tissue expression evolution: from expression phylogeny to phylogenetic network. *Brief.*  
920 *Bioinform.* **17**, 249–254 (2016).
- 921 30. Dunn, C. W., Zapata, F., Munro, C., Siebert, S. & Hejnl, A. Pairwise comparisons across species are problematic  
922 when analyzing functional genomic data. *Proc. Natl. Acad. Sci.* **115**, E409–E417 (2018).
- 923 31. Farré, D. & Albà, M. M. Heterogeneous Patterns of Gene-Expression Diversification in Mammalian Gene Duplicates.  
924 *Mol. Biol. Evol.* **27**, 325–335 (2010).
- 925 32. Pratapa, A., Jalihal, A. P., Law, J. N., Bharadwaj, A. & Murali, T. M. Benchmarking algorithms for gene regulatory  
926 network inference from single-cell transcriptomic data. *Nat. Methods* **17**, 147–154 (2020).
- 927 33. Aibar, S. *et al.* SCENIC: single-cell regulatory network inference and clustering. *Nat. Methods* **14**, 1083–1086 (2017).
- 928 34. Yates, A. D. *et al.* Ensembl 2020. *Nucleic Acids Res.* **48**, D682–D688 (2020).
- 929 35. Santos, M. E., Bouquin, A. L., Crumière, A. J. J. & Khila, A. Taxon-restricted genes at the origin of a novel trait allowing  
930 access to a new environment. *Science* **358**, 386–390 (2017).
- 931 36. Cao, J. & Lv, Y. Evolutionary analysis of the jacalin-related lectin family genes in 11 fishes. *Fish Shellfish Immunol.* **56**,  
932 543–553 (2016).
- 933 37. Burgoyne, R. D. Neuronal calcium sensor proteins: generating diversity in neuronal Ca<sup>2+</sup> signalling. *Nat. Rev.*  
934 *Neurosci.* **8**, 182–193 (2007).
- 935 38. Jaggard, J. B. *et al.* Hypocretin underlies the evolution of sleep loss in the Mexican cavefish. *eLife* **7**, e32637 (2018).



- 936 39. Alié, A. *et al.* Developmental evolution of the forebrain in cavefish, from natural variations in neuropeptides to behavior.  
937 *eLife* **7**, e32808 (2018).
- 938 40. Peuß, R. *et al.* Adaptation to low parasite abundance affects immune investment and immunopathological responses of  
939 cavefish. *Nat. Ecol. Evol.* **4**, 1416–1430 (2020).
- 940 41. Eisenberg, T. *et al.* Induction of autophagy by spermidine promotes longevity. *Nat. Cell Biol.* **11**, 1305–1314 (2009).
- 941 42. Duffy, C. M., Xu, H., Nixon, J. P., Bernlohr, D. A. & Butterick, T. A. Identification of a fatty acid binding protein4-UCP2  
942 axis regulating microglial mediated neuroinflammation. *Mol. Cell. Neurosci.* **80**, 52–57 (2017).
- 943 43. Zhang, H. *et al.* Polyamines Control eIF5A Hypusination, TFEB Translation, and Autophagy to Reverse B Cell  
944 Senescence. *Mol. Cell* **76**, 110-125.e9 (2019).
- 945 44. Rodrigues, F. T. S. *et al.* Major depression model induced by repeated and intermittent lipopolysaccharide  
946 administration: Long-lasting behavioral, neuroimmune and neuroprogressive alterations. *J. Psychiatr. Res.* **107**, 57–67  
947 (2018).
- 948 45. Saunders, A. *et al.* Molecular Diversity and Specializations among the Cells of the Adult Mouse Brain. *Cell* **174**, 1015-  
949 1030.e16 (2018).
- 950 46. Hughes, L. C. *et al.* Comprehensive phylogeny of ray-finned fishes (Actinopterygii) based on transcriptomic and  
951 genomic data. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 6249–6254 (2018).
- 952 47. Nakatani, M., Miya, M., Mabuchi, K., Saitoh, K. & Nishida, M. Evolutionary history of Otophysi (Teleostei), a major clade  
953 of the modern freshwater fishes: Pangaeen origin and Mesozoic radiation. *BMC Evol. Biol.* **11**, 177 (2011).
- 954 48. Tarashansky, A. J. *et al.* Mapping single-cell atlases throughout Metazoa unravels cell type evolution. *eLife* **10**, e66747  
955 (2021).
- 956 49. Kondrashov, F. A. & Koonin, E. V. A common framework for understanding the origin of genetic dominance and  
957 evolutionary fates of gene duplications. *Trends Genet.* **20**, 287–290 (2004).
- 958 50. Thomson, G. J. *et al.* Metabolism-induced oxidative stress and DNA damage selectively trigger genome instability in  
959 polyploid fungal cells. *EMBO J.* **38**, e101597 (2019).
- 960 51. Gillard, G. B. *et al.* Comparative regulomics supports pervasive selection on gene dosage following whole genome  
961 duplication. *Genome Biol.* **22**, 103 (2021).
- 962 52. Konstantinides, N. *et al.* Phenotypic Convergence: Distinct Transcription Factors Regulate Common Terminal  
963 Features. *Cell* **174**, 622-635.e13 (2018).
- 964 53. True, J. R. & Haag, E. S. Developmental system drift and flexibility in evolutionary trajectories. *Evol. Dev.* **3**, 109–119  
965 (2001).
- 966 54. Hilgers, L., Hartmann, S., Hofreiter, M. & von Rintelen, T. Novel Genes, Ancient Genes, and Gene Co-Option  
967 Contributed to the Genetic Basis of the Radula, a Molluscan Innovation. *Mol. Biol. Evol.* **35**, 1638–1652 (2018).
- 968 55. Florio, M. *et al.* Evolution and cell-type specificity of human-specific genes preferentially expressed in progenitors of

969 fetal neocortex. *eLife* **7**, e32332 (2018).

970 56. Arendt, D. *et al.* The origin and evolution of cell types. *Nat. Rev. Genet.* **17**, 744–757 (2016).

971 57. Viets, K., Eldred, K. C. & Johnston, R. J. Mechanisms of Photoreceptor Patterning in Vertebrates and Invertebrates.

972 *Trends Genet.* **32**, 638–659 (2016).

973 58. Bowmaker, J. K. Evolution of vertebrate visual pigments. *Vision Res.* **48**, 2022–2041 (2008).

974 59. Chitnis, T. & Weiner, H. L. CNS inflammation and neurodegeneration. *J. Clin. Invest.* **127**, 3577–3587 (2017).

975 60. Riddle, M. R. *et al.* Insulin resistance in cavefish as an adaptation to a nutrient-limited environment. *Nature* **555**, 647–

976 651 (2018).

977 61. Elipot, Y. *et al.* A mutation in the enzyme monoamine oxidase explains part of the *Astyanax* cavefish behavioural

978 syndrome. *Nat. Commun.* **5**, 3647 (2014).

979 62. Yoshizawa, M. *et al.* Distinct genetic architecture underlies the emergence of sleep loss and prey-seeking behavior in

980 the Mexican cavefish. *BMC Biol.* **13**, 15 (2015).

981 63. Fischer, E. K. & O'Connell, L. A. Modification of feeding circuits in the evolution of social behavior. *J. Exp. Biol.* **220**,

982 92–102 (2017).

983 64. Kroeger, D. *et al.* Galanin neurons in the ventrolateral preoptic area promote sleep and heat loss in mice. *Nat.*

984 *Commun.* **9**, 4129 (2018).

985 65. Yamashita, J. *et al.* Male-predominant galanin mediates androgen-dependent aggressive chases in medaka. *eLife* **9**,

986 66. Wee, C. L. *et al.* Social isolation modulates appetite and defensive behavior via a common oxytocinergic circuit in larval

987 zebrafish. *bioRxiv* 2020.02.19.956854 (2020) doi:10.1101/2020.02.19.956854.

988 67. Elipot, Y., Legendre, L., Père, S., Sohm, F. & Rétaux, S. *Astyanax* Transgenesis and Husbandry: How Cavefish Enters

989 the Laboratory. *Zebrafish* **11**, 291–299 (2014).

990 68. Sawh, A. N. *et al.* Lamina-Dependent Stretching and Unconventional Chromosome Compartments in Early *C. elegans*

991 Embryos. *Mol. Cell* **78**, 96-111.e6 (2020).

992 69. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the integration of genomic datasets with the

993 R/Bioconductor package biomaRt. *Nat. Protoc.* **4**, 1184–1191 (2009).

994 70. Strimmer, K. *fdrtool*: a versatile R package for estimating local and tail area-based false discovery rates. *Bioinformatics*

995 **24**, 1461–1462 (2008).

996 71. Fresno, C. & Fernández, E. A. RDAVIDWebService: a versatile R interface to DAVID. *Bioinformatics* **29**, 2810–2811

997 (2013).

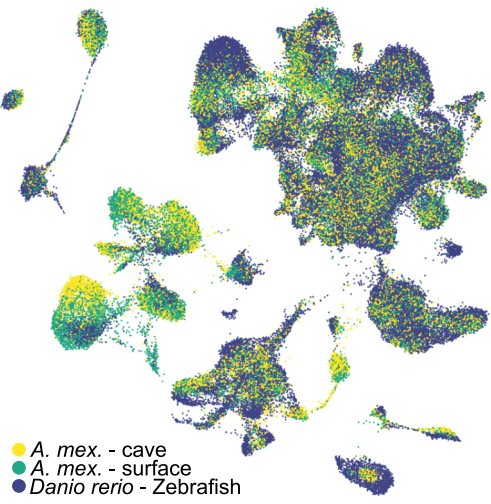
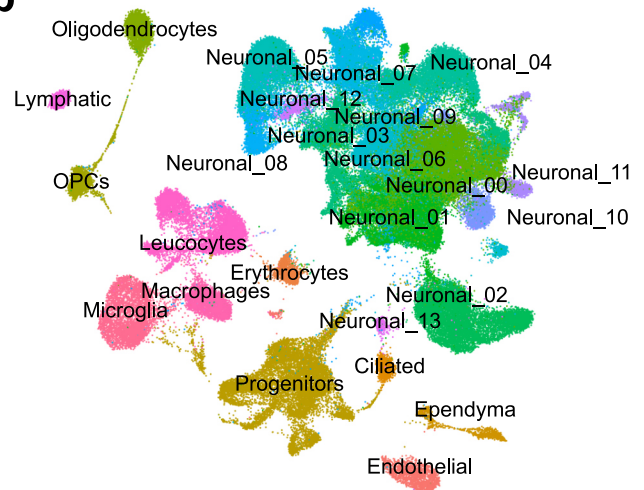
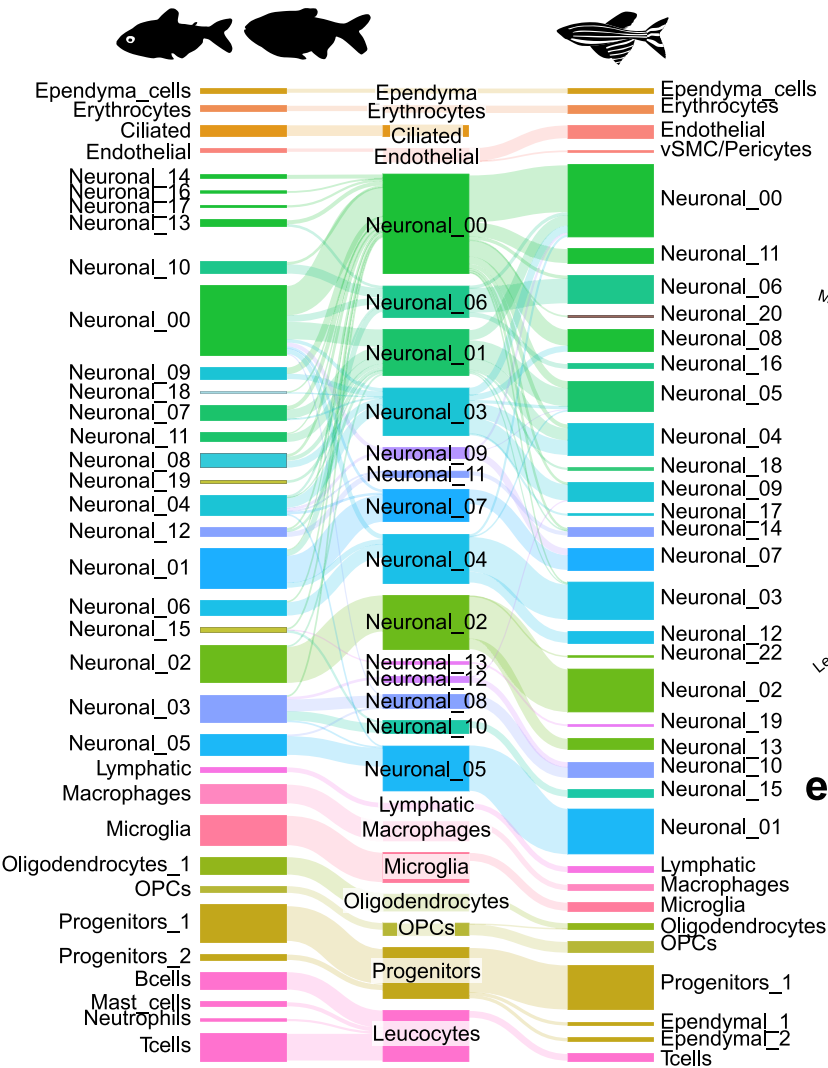
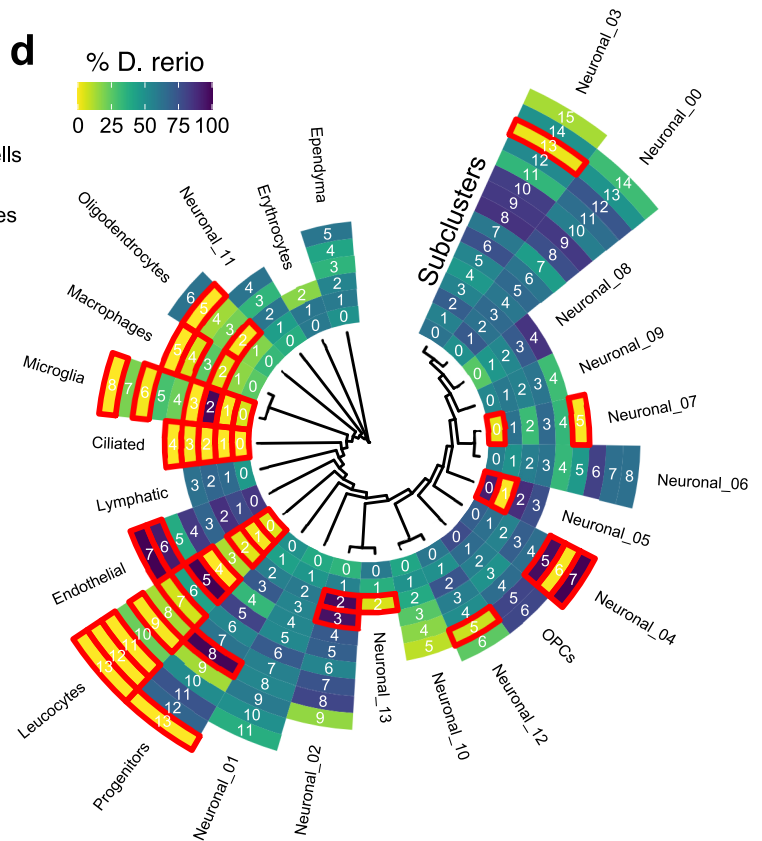
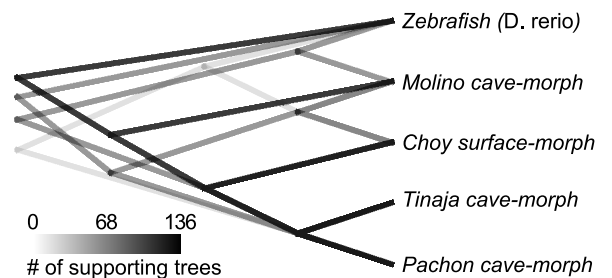
998 72. Singh, P. P. & Isambert, H. OHNOLOGS v2: a comprehensive resource for the genes retained from whole genome

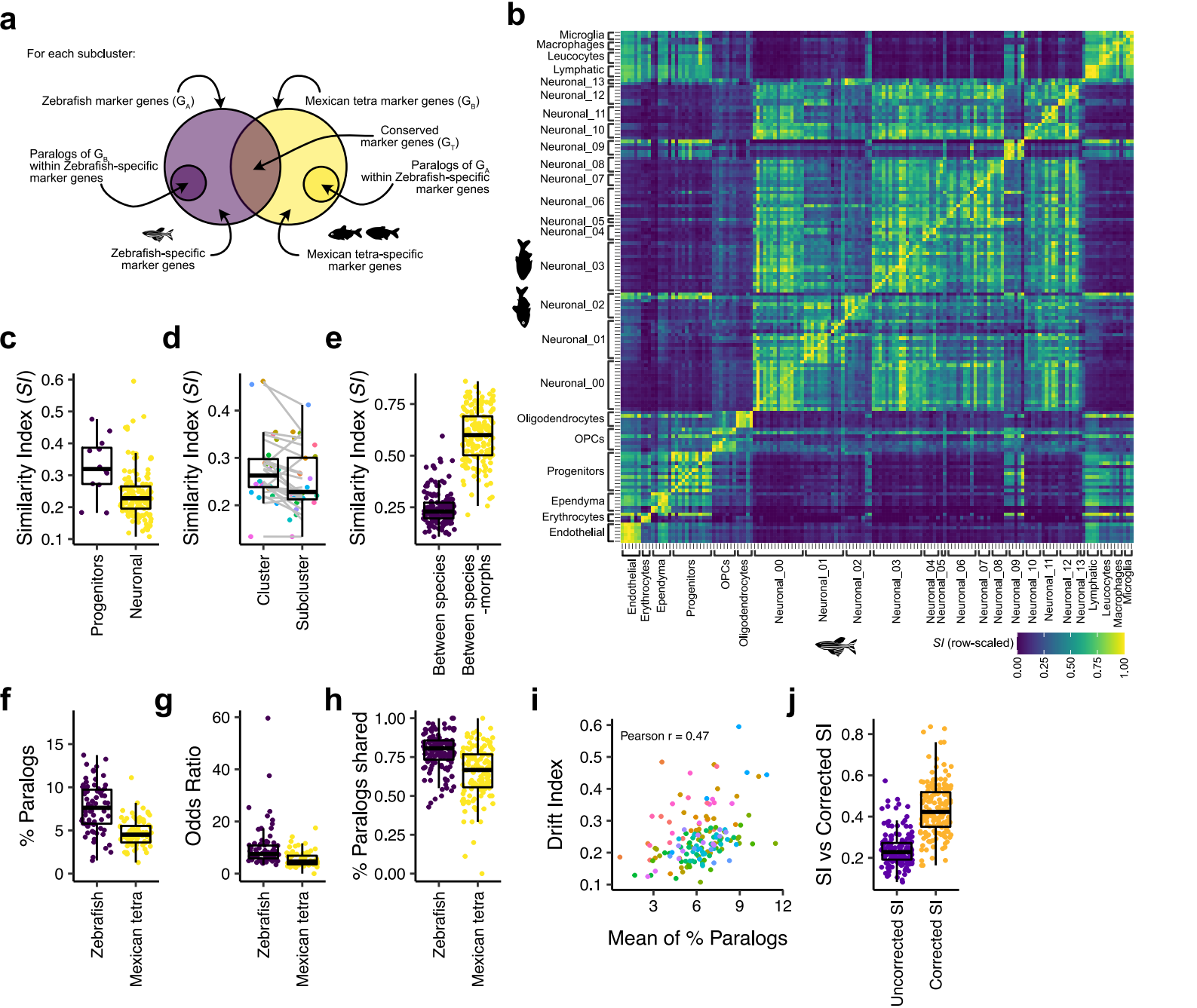
999 duplication in vertebrates. *Nucleic Acids Res.* **48**, D724–D730 (2020).

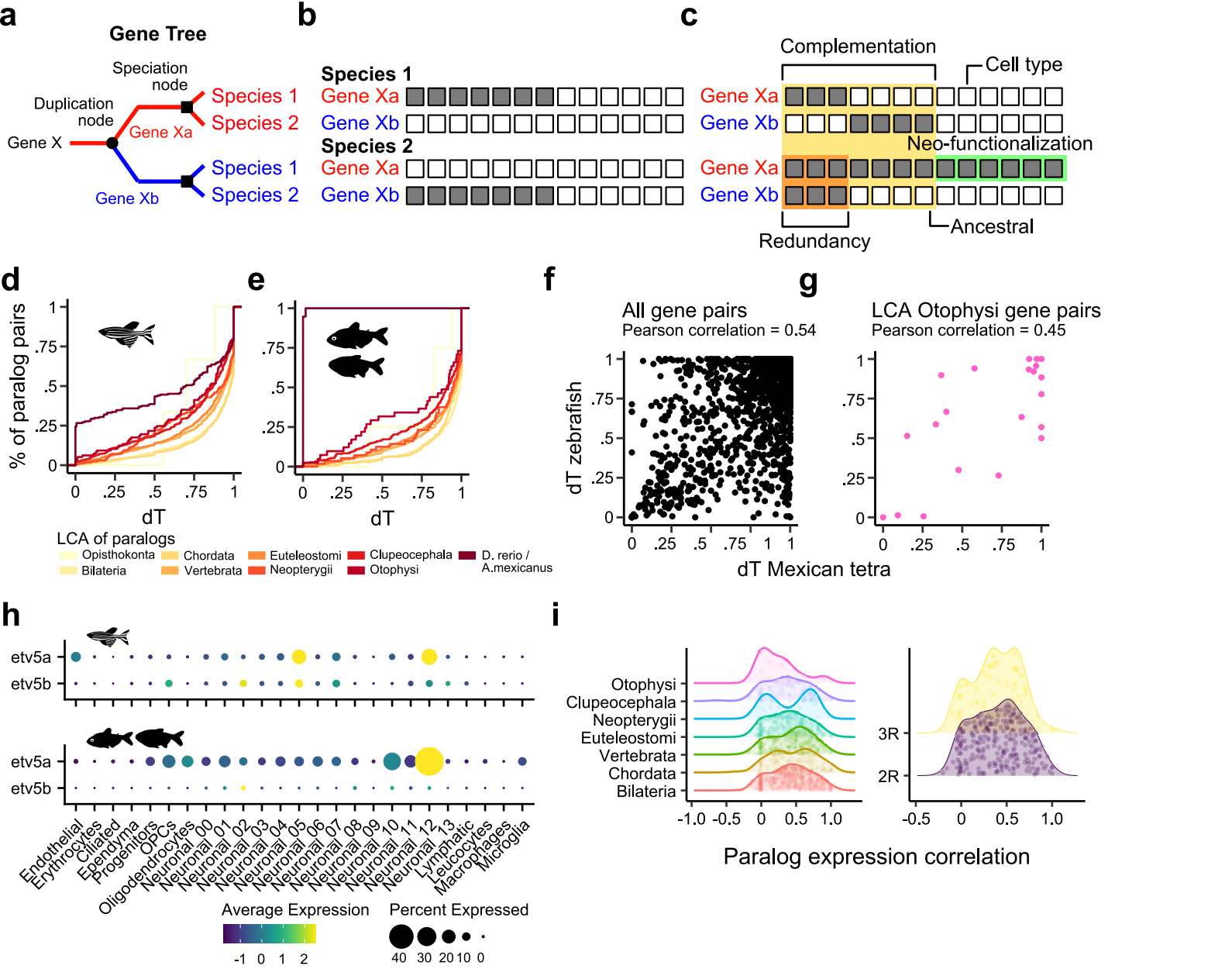
1000 73. Smith, R. N. *et al.* InterMine: a flexible data warehouse system for the integration and analysis of heterogeneous

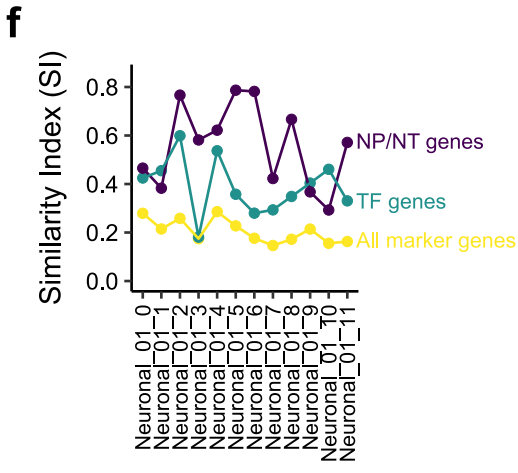
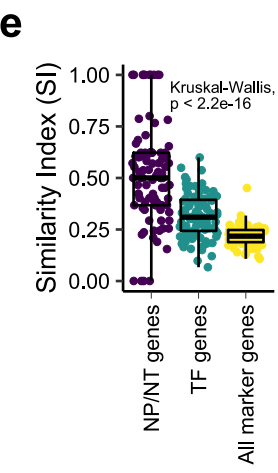
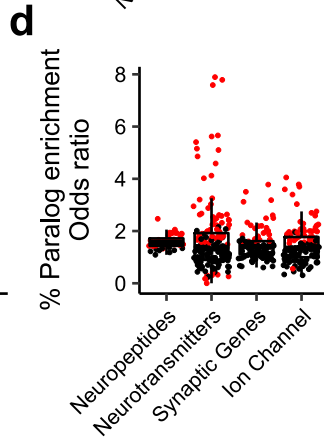
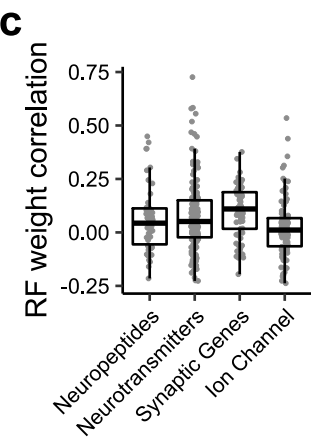
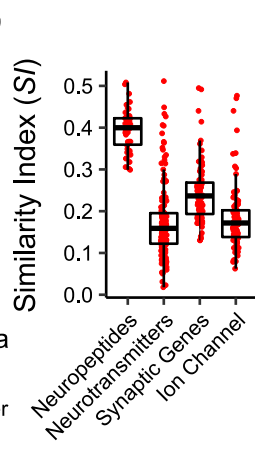
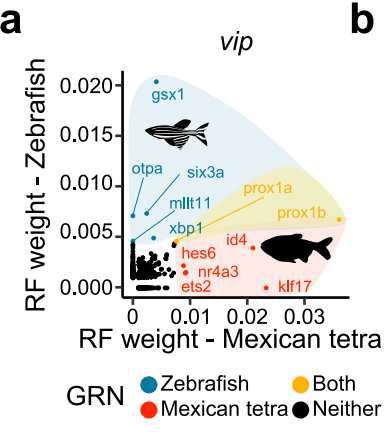
1001 biological data. *Bioinformatics* **28**, 3163–3165 (2012).

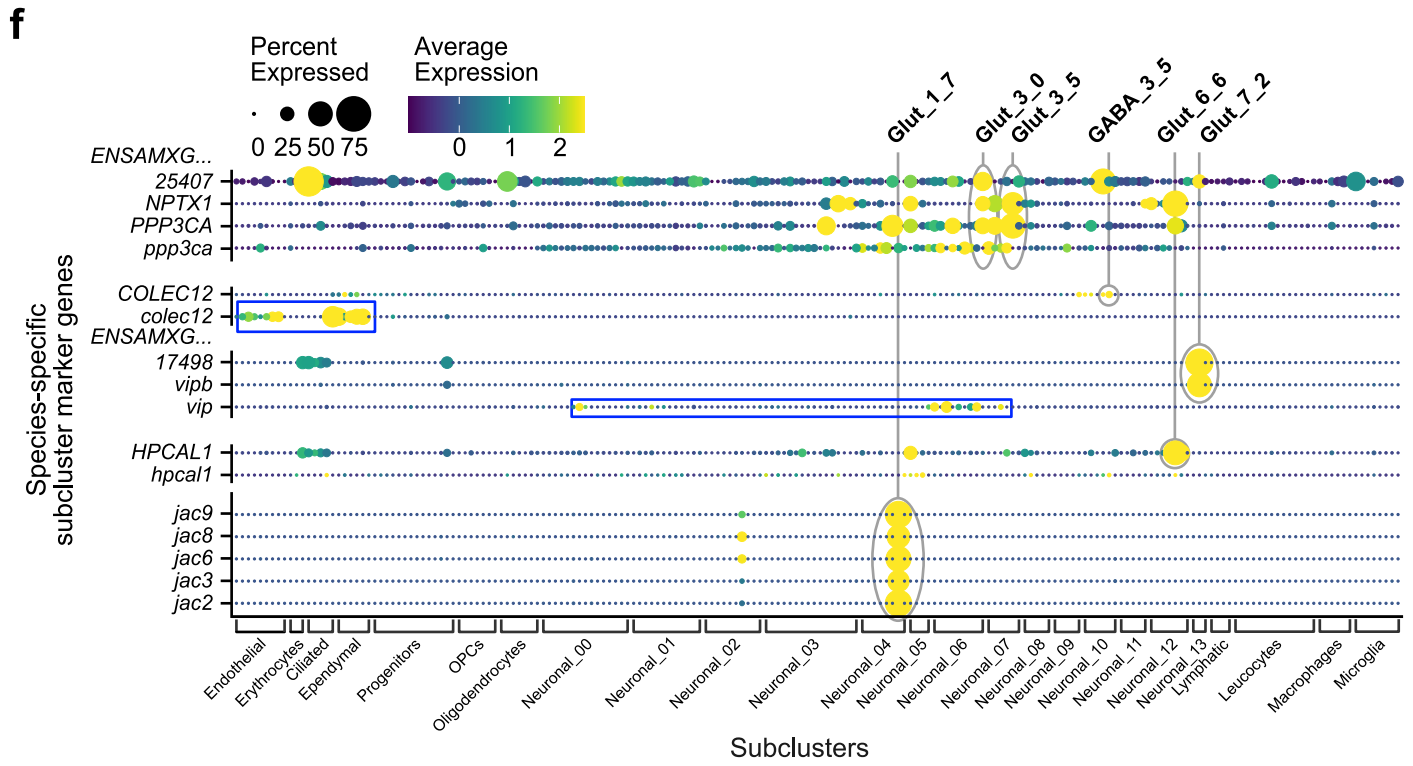
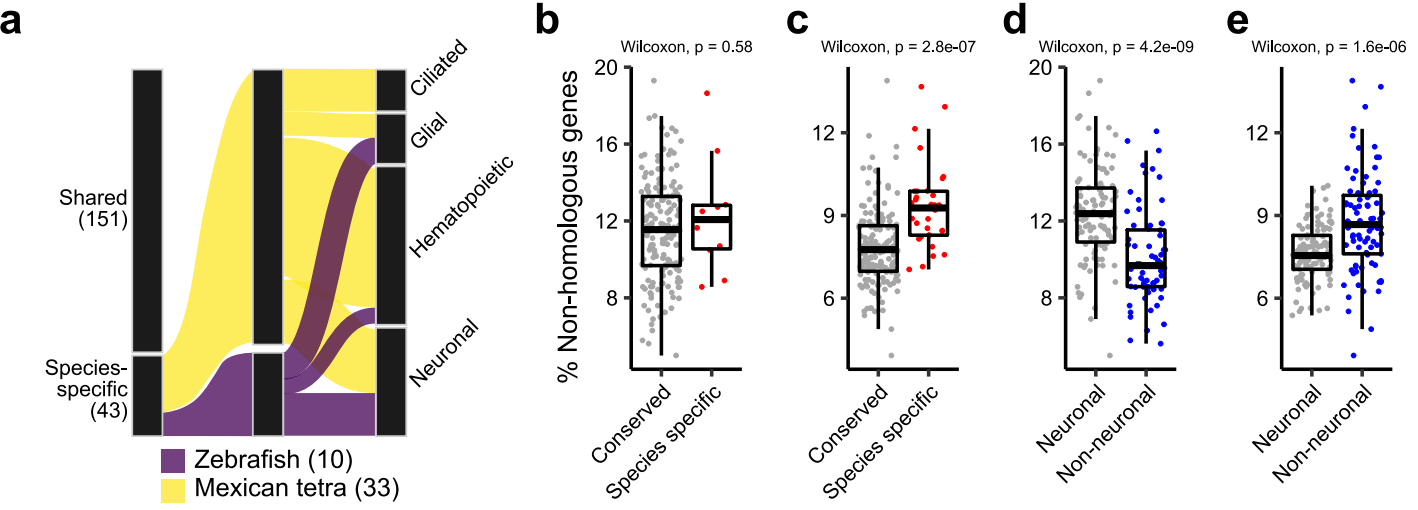
- 1002 74. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data.  
1003 *Nat. Genet.* **43**, 491–498 (2011).
- 1004 75. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler transform. *Bioinformatics* **26**, 589–  
1005 595 (2010).
- 1006 76. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
- 1007 77. Wilkinson, S. P. & Davy, S. K. phylogram: an R package for phylogenetic analysis with nested lists. *J. Open Source*  
1008 *Softw.* **3**, 790 (2018).
- 1009 78. Yu, G. Using ggtree to Visualize Data on Tree-Like Structures. *Curr. Protoc. Bioinforma.* **69**, e96 (2020).
- 1010

**a****b****c****d****e**

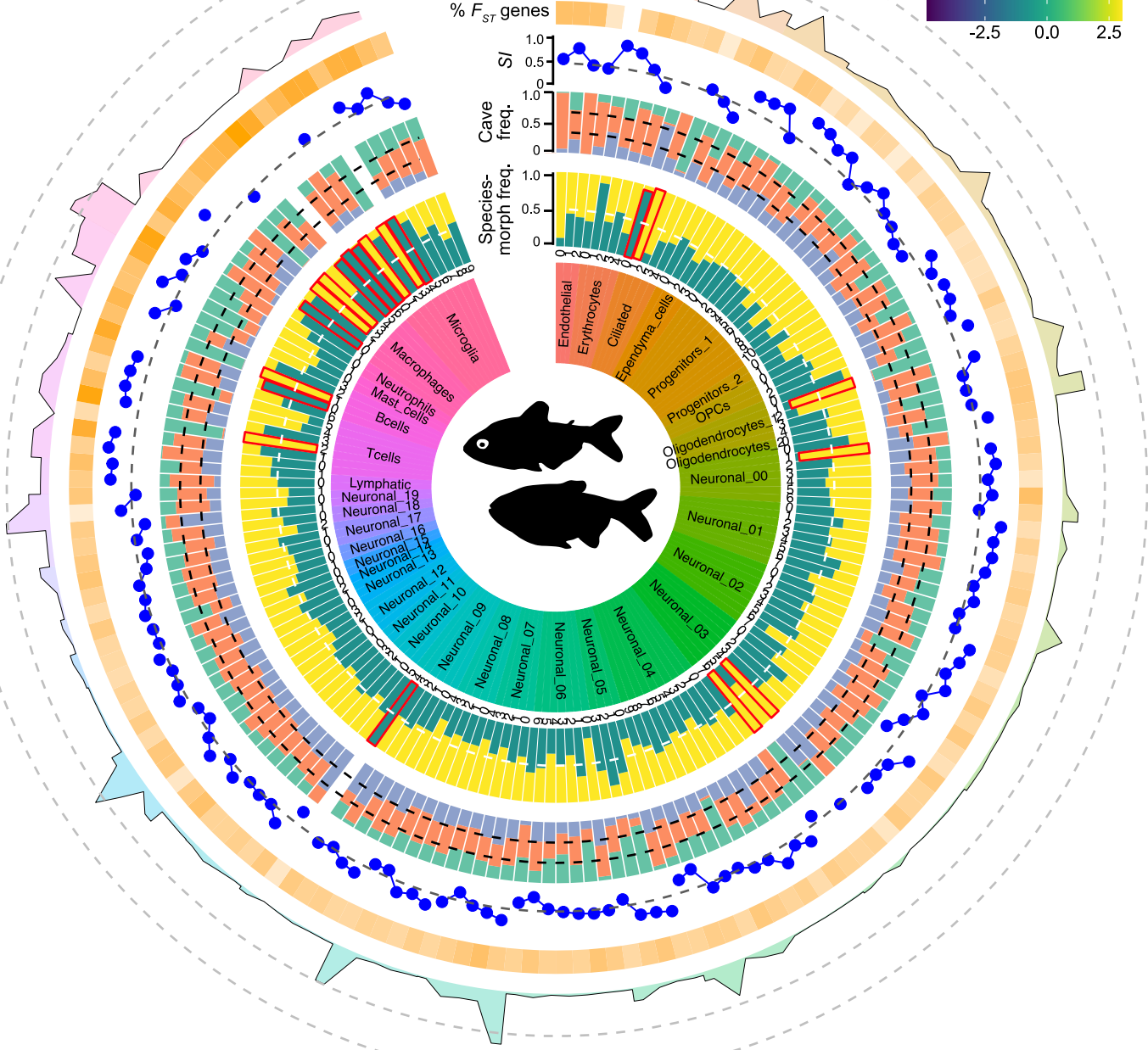
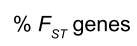
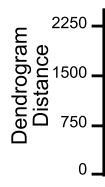
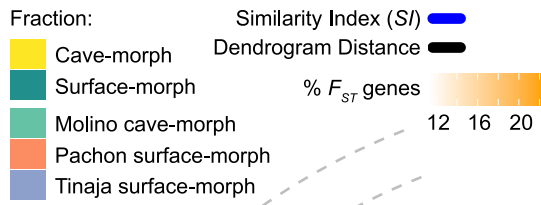










**a****b**