

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

Emergence of neuronal diversity during vertebrate brain development

Bushra Raj^{1,2*}, Jeffrey A. Farrell^{1,3}, Jialin Liu^{2,4}, Jakob El Kholtej^{2,4}, Adam Carte^{1,4,5}, Joaquin Navajas Acedo^{2,4}, Lucia Y Du^{2,4}, Aaron McKenna⁶, Đorđe Relić^{3,7}, Jessica M. Leslie¹, and Alexander F. Schier^{1,2,4,8,9,10*}

¹ Department of Molecular and Cellular Biology, Harvard University, Cambridge, Massachusetts, USA

² Allen Discovery Center for Cell Lineage Tracing, Seattle, Washington, USA

³ Unit on Cell Specification and Differentiation, National Institute of Child Health and Human Development, NIH, Bethesda, Maryland, USA

⁴ Biozentrum, University of Basel, Switzerland

⁵ Systems, Synthetic, and Quantitative Biology Program, Harvard University, Cambridge, Massachusetts, USA

⁶ Department of Molecular and Systems Biology, Dartmouth Geisel School of Medicine, Lebanon, New Hampshire, USA

⁷ Swiss Institute of Bioinformatics (SIB), Basel, Switzerland

⁸ Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

⁹ Harvard Stem Cell Institute, Harvard University, Cambridge, Massachusetts, USA

¹⁰ Center for Brain Science, Harvard University, Cambridge, Massachusetts, USA

*Co-corresponding authors. Email: bushranraj@gmail.com (B.R.); alex.schier@unibas.ch (A.F.S.)

43 **ABSTRACT**

44
45 Neurogenesis comprises many steps from progenitor proliferation to neuronal differentiation and
46 maturation. These processes are highly regulated, but the landscape of transcriptional changes
47 underlying brain development are poorly characterized. Here, we describe a developmental
48 single-cell RNA-seq catalog of ~220,000 zebrafish brain cells encompassing 12 stages from 12
49 hours post-fertilization to 15 days post-fertilization. We characterize known and novel gene
50 markers for ~800 clusters and provide an overview of the diversification of neurons and
51 progenitors across these timepoints. We also introduce an optimized version of the GESTALT
52 lineage recorder that enables higher expression and recovery of Cas9-edited barcodes to query
53 lineage segregation. Cell type characterization indicates that most embryonic neural progenitor
54 states are transitory and transcriptionally distinct from neural progenitors of post-embryonic
55 stages. Reconstruction of cell specification trajectories reveals that late-stage retinal neural
56 progenitors transcriptionally overlap cell states observed in the embryo. The zebrafish brain
57 development atlas provides a resource to define and manipulate specific subsets of neurons and
58 to uncover the molecular mechanisms underlying vertebrate neurogenesis.

59

60

61 **INTRODUCTION**

62

63 The vertebrate brain develops from a limited pool of embryonic neural progenitor cells that cycle
64 through rounds of proliferation, diversification, and terminal differentiation into an extensive
65 catalogue of distinct neuronal and glial cell types. A central goal in developmental neurobiology
66 is to investigate how neuronal complexity arises through molecular specification and commitment
67 by studying the origins and fates of cells during development. Fundamental insights into these
68 processes have been gained via classic approaches using genetic markers, perturbations and
69 fate mapping (Cepko, 2014; Kretzschmar and Watt, 2012; Ma et al., 2017; Wamsley and Fishell,
70 2017; Wilson et al., 2002; Woo and Fraser, 1995; Woodworth et al., 2017). These approaches
71 have recently been complemented by single-cell genomics technologies in the developing
72 nervous system, including the spinal cord (Delile et al., 2019; Rosenberg et al., 2018); cortex
73 (Nowakowski et al., 2017; Zhong et al., 2018); olfactory system (H. Li et al., 2017); cerebellum
74 (Carter et al., 2018; Tambalo et al., 2020); retina (Clark et al., 2019; Hu et al., 2019; Xu et al.,
75 2020); and whole animal (Farnsworth et al., 2020). These studies have provided transcriptome-
76 level views of the rich heterogeneous states that cells progress through as they proliferate,
77 migrate and differentiate. Nevertheless, existing datasets are limited in their scope as they focus
78 on specific brain regions, survey limited timepoints or do not enrich for neural cell types, thereby
79 missing transitions and cellular diversity. Thus, there is a need for a large-scale
80 neurodevelopmental single-cell resource that profiles whole brain development across a range of
81 closely-spaced embryonic and post-embryonic stages. In addition, such an atlas would help
82 address fundamental questions about the dynamics of brain development. For example, it is
83 poorly understood how embryonic neural progenitors are molecularly related to post-embryonic
84 neural progenitors. Furthermore, the transcriptional programs that are activated or suppressed as
85 neural progenitors become fate-restricted and differentiate are largely unknown.

86

87 Here we present resources to obtain global views of neurogenesis, cell type heterogeneity,
88 specification trajectories and lineage relationships in the developing zebrafish brain. We
89 generated a single-cell RNA-seq (scRNA-seq) atlas consisting of ~220,000 cells from 12 hours
90 post fertilization (hpf) to 15 days post fertilization (dpf). We also created a new version of the
91 scGESTALT CRISPR-Cas9 lineage recorder (Raj et al., 2018b) with improved barcode capture
92 and used it to query early lineage decisions. Using the cell type atlas, we analyzed the expansion
93 of neuronal diversity, the loss of transitory embryonic progenitors, and the maintenance of distinct
94 larval progenitor states. We reconstructed cell specification trajectories of the zebrafish retina and
95 hypothalamus, revealing gene expression cascades and distinct specification programs.
96 Collectively, the zebrafish brain development atlas reveals molecular and cellular changes at an
97 unprecedented scale and resolution, and lays the foundation for the detailed analysis of neuronal
98 diversification.

99

100

101 RESULTS

102

103 Building a developmental atlas of the zebrafish brain with single-cell transcriptomics

104

105 To reveal the landscape of cell states and cell types during brain development, we profiled
106 223,037 cells across 12 stages of zebrafish embryonic and larval development using the 10X
107 Chromium scRNA-seq platform. Samples spanned from 12 hpf (shortly after gastrulation), when
108 the embryo is undergoing early developmental patterning, to 15 dpf, when larvae are mature,
109 exhibit complex behaviors, and are expected to exhibit substantial cell type diversity (Figure 1A).
110 To enrich for brain cell types, we dissected the heads of animals from 12 hpf to 3 dpf, and the
111 brains and eyes from 5 dpf to 15 dpf (Figure 1B). To determine cell type diversity in the head and
112 brain of zebrafish, data from each stage was analyzed individually using Louvain clustering
113 (Figure 1C and Sup Figure 1). This approach identified a total of 815 cell clusters across all 12
114 timepoints (Sup Table). To classify each cluster, we compared enriched gene markers with
115 existing gene expression annotations in the ZFIN database and literature, as described previously
116 (Raj et al., 2018b). Plotting expression of known cell type markers identified clusters
117 corresponding to neural progenitors (*sox19a*), dozens of neuron subtypes (*elavl3*, *gad2*,
118 *slc17a6b*), eye cells (*foxg1b*, *lim2.4*, *pmela*, *ca14*, *gnat1*, *opn1mw1*), radial glia (*mfg8a*, *s100b*),
119 neural crest (*sox10*), oligodendrocytes (*mbpa*), blood cells (*cahz*, *etv2*, *cd74a*), cartilage (*matn4*,
120 *col9a2*), pharyngeal arches (*pmp22a*, *prrx1b*, *barx1*), sensory placodes (*dlx3b*, *six1b*), and
121 epidermal cells (*epcam*, *cldni*), among others. As expected, cell type complexity increased with
122 developmental time. We validated new marker expression across several cell types identified in
123 our dataset, such as *sdpra* in the trigeminal placode, *sox1a* in the hypothalamus, and *ompa* in the
124 retina (Figure 1D-F). Our analysis also revealed groups of embryonic clusters that were absent
125 or transcriptionally distinct from larval clusters, suggesting that many embryonic cell states are
126 transitory. Several of these transitions are known developmental changes (e.g. loss of placodes
127 and rhombomeres), but changes in neural progenitor cell states are poorly understood (see
128 below).

129

130 To enable direct comparison of cell types across our time course, we subsetted the 12 hpf dataset
131 to only comprise neural populations and blood cells found in the brain, eliminating non-relevant
132 head cells from earlier stages, such as mesoderm, placodes, and periderm. This approach
133 resulted in an initial set of 21 clusters at 12 hpf (Figure 2A) that diversified into 98 clusters by 15
134 dpf (Figure 2C). Notably, most clusters could be uniquely identified using a minimal group of 2-3
135 enriched gene markers (Figure 2B, 2D). For example, at 12 hpf, the optic vesicle is identified by
136 expression of *rx2* and *rx3*; hindbrain rhombomeres 5/6 by *hoxb3a* and *eng2b*; and ventral
137 diencephalon by *nkx2.4a* and *dbx1a*. Similarly, at 15 dpf, the cerebellar granule cells are marked
138 by expression of *opr1b* and *zic2a*; optic tectum by *pax7a* and *tal1*; and a new retinal cell type by
139 *kidins220a*, *foxg1b* (exclusively detected in retinal cells) and *tbx3a*. We did not find unique gene
140 combinations for cycling progenitors, differentiating progenitors and newly born neurons, as many
141 of these subtypes had similar expression signatures of pan neuronal or pan progenitor marker
142 genes, such as *elavl3* and *tubb5* in neurons, and *rpl5a* and *npm1a* in progenitors (Figure 2D, grey
143 box).

144
145 At 12 hpf, the early demarcation of multiple brain regions is already apparent and by 15 dpf these
146 regions expand and diversify further. For example, the optic vesicle at 12 hpf is defined by one
147 cluster and is the origin of 18 retinal cell types at 15 dpf. Similarly, a single cluster of ventral
148 diencephalon cells (expressing *shha*, *nkx2.4a*, *nkx2.1*, *rx3*) at 12 hpf develops into 7 major
149 hypothalamus cell types at 15 dpf. An exception to this diversification is the loss of rhombomeres
150 (r1-r7) in the hindbrain (Moens and Prince, 2002).

151
152 To further explore brain neuronal subtypes at 15 dpf, we analyzed the expression of transcription
153 factors, neuropeptides and their receptors, and genes involved in neuronal physiology (e.g.
154 neurotransmitters, transporters, receptors, and channels) (R. Chen et al., 2017; Pandey et al.,
155 2018; Tiklová et al., 2019; Zeisel et al., 2018). Our results indicate that nearly all identified neuron
156 subtypes can be distinguished from one another via the expression of individual or combinations
157 of genes belonging to these categories (Figure 3A-C). For example, cluster 2 and 84 neurons are
158 GABAergic forebrain neurons that express *dlx2a* and *dlx5a*, while cluster 84 neurons additionally
159 express *six3b*, *gria1a* and *gria2b*.

160
161 We next asked if neuron clusters detected at 15 dpf are found in the earlier larval stages, when
162 most behavioral experiments are performed (Sup Figure 2). 68% (23/34 clusters) and 74% (25/34
163 clusters) of 15 dpf neuron clusters have a closely matching counterpart at 5 dpf and 8 dpf (based
164 on enriched marker gene expression), respectively (Figure 3D). Sampling issues might have
165 prevented the identification of additional overlapping clusters, but our data indicate a large overlap
166 between identified cell types from 5 to 15 dpf. These results suggest that the zebrafish brain
167 already has considerable cell type diversity at early larval stages. Furthermore, 97% (33/34) of 15
168 dpf clusters overlapped with clusters identified in our previously described 23-25 dpf juvenile brain
169 dataset (Raj et al., 2018b). Thus, by 15 dpf late larval stage, nearly all of the brain cell types that
170 persist into the early juvenile stage have already been established. Notably, among cell types that
171 are “missing” or under-represented at 15 dpf but readily detected at 23-25 dpf are cell types in
172 the optic tectum, cerebellum and the torus longitudinalis, suggesting that these structures undergo
173 further diversification after 15 dpf. In contrast, many cell types in the pallium, habenula (Pandey

174 et al., 2018), hypothalamus and preoptic area are detected across these stages, suggesting that
175 they develop earlier.

176

177 In summary, we generated a zebrafish brain development cell type atlas spanning 12 stages of
178 brain organogenesis. The complete dataset can be explored using the accompanying app:
179 https://github.com/brlauuu/zf_brain.

180

181 **Neurogenic expansion during brain development**

182

183 During development, cell composition shifts from predominantly progenitor populations to more
184 differentiated cell types (Schmidt et al., 2013). To better characterize how differentiation varies
185 during neuronal development, we first asked if our dataset captured the two neurogenic phases
186 (primary and secondary) before and after 2 dpf that have been traditionally defined through
187 histological analyses (Allende and Weinberg, 1994; Korzh et al., 1998; Mueller and Wullmann,
188 2003). We considered neural progenitors as non-differentiated neuronal precursor cells that may
189 or may not be proliferating, and express a subset of classical progenitor markers e.g. *sox19a*, *dla*,
190 *s100b*, and cell cycle genes. Since the brain is undergoing substantial molecular changes during
191 these developmental windows, we defined the transcriptional programs and cells that exhibit
192 these programs as progenitor cell states. We calculated the percentage of the dataset that
193 corresponds to neural progenitor cells, neurons (expressing markers such as *elav13*, *elav14*) or
194 other cell types across each timepoint in our dataset. Since the earlier stages (12 hpf to 3 dpf)
195 contained non-brain and non-eye cell types, while later stages were restricted to these tissues,
196 we subsetting the early timepoints to only brain and eye cells. With increasing developmental time,
197 we observed a progressive decrease in the fraction of the dataset comprising neural progenitor
198 cells (from 53.8% to 18.3%) with a concomitant increase in neurons (from 4.5% to 58%) (Figure
199 4A). For example, we observed an initial increase in the number of distinct progenitor clusters
200 from 12 hpf to 18 hpf (early embryo stages), while the number of neuron clusters remained low
201 (Figure 4A, right panels). From 20 hpf to 3 dpf (intermediate stages), the total progenitor clusters
202 decreased while neuron clusters started to increase. For example, neuronal clusters expanded
203 from 11 at 20 hpf to 23 at 36 hpf. This burst coincides with the presumed timing of late-stage
204 primary neurogenesis in zebrafish (Mueller and Wullmann, 2003). Notably, by 5 to 15 dpf (late
205 larva stages), a second expansion of neuronal populations, corresponding to the secondary
206 neurogenic phase (Mueller and Wullmann, 2003), had occurred (53 neuronal subtypes at 5 dpf).
207 At 5 dpf, we detected cell types identified as early as 36 hpf (e.g. *tal1*⁺, *gata3*⁺ neurons in the optic
208 tectum, and *tfap2e*⁺, *barhl2*⁺ neurons in the thalamus), as well as subtypes only observed during
209 the second phase, such as *nrgnb*⁺ *prkcda*⁺ neurons in the forebrain and cone bipolar cell
210 subtypes in the retina. Collectively, our dataset captures both phases of neurogenesis and reveals
211 the diversification of neurons in multiple brain structures.

212

213 **Dampening of spatial and developmental signatures during the transition from** 214 **embryonic to larval neural progenitors**

215

216 We next analyzed our dataset to determine how cell states change during the transition from the
217 embryonic to post-embryonic brain. The zebrafish brain undergoes lifetime constitutive

218 neurogenesis due to the persistence of neural progenitor pools distributed along the brain's axis
219 (Schmidt et al., 2013). However, the embryonic origins and transcriptional programs that underlie
220 their development are poorly understood. Furthermore, how the molecular identities of embryonic
221 and post-embryonic neural progenitor cell states compare have not been well characterized. To
222 address these questions, we asked how neural progenitor gene expression signatures globally
223 change from embryo to larva. Based on the results described above, we defined early embryonic
224 brain progenitors as neural cell transcriptional states from 12 hpf to 18 hpf, intermediate stage
225 brain progenitors as neural cell transcriptional states from 20 hpf to 3 dpf, and larval brain
226 progenitors as neural cell transcriptional states from 5 dpf to 15 dpf (Figure 4B, Sup Figure 3).
227 We determined the greatest sources of variation within these populations. For embryonic brain
228 progenitors we found that the top 3 principal components comprise genes implicated in spatial
229 and developmental patterning (Gibbs et al., 2017; Moens and Prince, 2002; Wilson et al., 2002;
230 Wilson and Rubenstein, 2000). Cells exhibit characteristic anteroposterior and dorsoventral axial
231 signatures (Figure 4C, top panel). For example, the telencephalon (anterior forebrain) is marked
232 by *foxg1a* and *emx3a* expression, the midbrain by *pax2a* and *eng2a*, and the hindbrain is
233 segmented into rhombomeres marked by distinct combinatorial patterns of *egr2b* and *hox* gene
234 expression. Furthermore, all cells are in a highly proliferative state with strong expression of cell
235 cycle genes such as *pcna*, *mki67* and *cdca7a*. Collectively, the expression signatures are
236 reflective of a developmental state during which the embryo is orchestrating a rapid expansion of
237 neural progenitor populations concurrent with their acquisition of positional information and overt
238 absence of differentiation (Schmidt et al., 2013; Stigloher et al., 2008).

239
240 In contrast, larval neural progenitors comprised two major groups: proliferating (expressing cell
241 cycle genes *pcna* and *top2a*) and non-proliferating (depleted expression of cell cycle markers)
242 (Figure 4D, bottom panel). Indeed, the top 3 principal components in the larval progenitors
243 comprised genes that mark stem cells (PC1, PC3) and differentiation (PC2). The non-proliferating
244 group is subdivided into radial glia (stem cells) and *her2*⁺ neural progenitors expressing proneural
245 genes *insm1b* and *scrt2*. The proliferating group is subdivided into *her2*⁺ and *scrt2*⁻ neural
246 progenitor cells, *her2*⁻ progenitors, *her2*⁺ and *neurod1*⁺ progenitor cells, and upper rhombic lip
247 progenitors (localized to cerebellum) expressing *atoh1c* and *oprd1b*.

248
249 Strikingly, most larval progenitors were characterized by a reduced spatial signature (except for
250 the cerebellar upper rhombic lip pool), such that cells were less enriched in region-specific
251 transcription factors relative to embryonic progenitors (Figure 4D, top panel). For example, radial
252 glia exist in multiple pools along the brain axis (Than-Trong and Bally-Cuif, 2015), but they formed
253 a single cluster in our dataset (marked by expression of *fabp7a*, *cx43*, *s100b* and *aqp1a.1*). This
254 result suggests that radial glia are largely transcriptionally similar. Although some expression of
255 region-specific transcription factors was detected in larval progenitor clusters, these signatures
256 were not sufficiently strong to resolve clusters as they were during embryonic stages.

257
258 To explore the apparent dearth of spatial signatures further, we calculated pairwise correlation
259 scores for 79 transcription factors and signaling proteins with known spatial expression patterns
260 in the forebrain and midbrain based on previously described histological analysis (ZFIN), and
261 which were identified as gene markers for neuronal clusters in our dataset. These genes showed

262 strongest correlations in embryonic progenitors, followed by intermediate stage progenitors, and
263 were weakly correlated in larval progenitors (Figure 4E).

264
265 Since spatial signatures are encoded by a combinatorial code of genes with overlapping
266 expression patterns, we asked whether the same subsets of genes co-varied with each of the 79
267 spatial markers across embryonic, intermediate, and larval neural progenitors. We found that
268 intermediate stage progenitors showed overlap in co-varying genes with both embryonic and
269 larval progenitors. For example, 44/79 genes had >40% overlap in their top 20 co-varying genes
270 between embryonic and intermediate stage progenitors, and 23/79 genes had >40% overlap
271 between intermediate and larval stage progenitors. In contrast, we found low overlap across
272 embryonic and larval stages (3/79 genes had >40% overlap in their top 20 co-varying genes).
273 Additionally, when we searched for genes that strongly co-varied with these 79 spatial markers
274 (Pearson correlation >0.4), we found 38 genes during embryonic stages, 17 genes during
275 intermediate stages, but only 4 genes during larval stages (Figure 4F).

276
277 Taken together, these results demonstrate that intermediate stage progenitors resemble a hybrid
278 of early embryonic and late larval progenitor signatures. Furthermore, the overall spatial code
279 between embryonic and larval progenitors are distinct, and the embryonic spatial code involves a
280 larger collection of genes. Notably, the signatures of larval progenitors resemble juvenile neural
281 progenitor pools (Raj et al., 2018b), indicating developmental switches in neural progenitor
282 identities from embryo to larva that are maintained to at least juvenile stages. Thus, embryonic
283 states that existed in early progenitors are largely altered in late-stage progenitors: while spatial
284 patterning signals are the greatest source of variation between embryonic neural progenitors,
285 these signals are dampened in post-embryonic neural progenitors.

286
287 **An optimized scGESTALT lineage recorder**

288
289 A long-term goal in developmental neurobiology is to understand the lineage relationships of
290 neurons. As a first step to derive lineage relationships of the cell types identified in the brain
291 development atlas, we performed lineage recording experiments with scGESTALT. This lineage
292 recorder enables simultaneous cell type and cell lineage identification by combining scRNA-seq
293 with CRISPR-Cas9 barcode editing (McKenna et al., 2016; Raj et al., 2018b). To enable higher
294 recovery of edited barcodes from single cells, we optimized the design and library preparation of
295 the lineage recording cassette, including barcode editing of a transgene coding region and
296 compatibility with the 10x platform (see Methods). To test the performance of this new recording
297 cassette, we barcoded early embryonic lineage relationships by injecting Cas9 protein and target
298 guide RNAs into 1-cell embryos (Figure 5A) and then isolated four 15 dpf larval brains. We
299 recovered barcodes and transcriptional profiles of 5,794 cells total (barcode recovery rate 30-75%
300 compared to 6-28% of our previous scGESTALT version (Raj et al., 2018b)). Edited barcodes
301 showed no overlap between animals, displayed a diverse spectrum of repair products that
302 spanned single and multiple sites, and were of varying clone sizes (Figure 5B-D, Sup Figure 4A).
303 These features closely resembled the editing patterns obtained with our previous recorders
304 (McKenna et al., 2016; Raj et al., 2018b). Using the recovered barcodes and associated
305 transcriptomes, we reconstructed lineage trees representing cell lineage segregations formed

306 during early embryogenesis (for one example see Sup Figure 4B). These lineage trees
307 accompany our transcriptional cell type atlas and are available to explore at
308 <https://scgestalt.mckennalab.org/>

309
310 Since the injection of editing reagents into 1-cell embryos saturates editing within 4-6 hours
311 (McKenna et al., 2016) , we expected early lineage divergences to be overrepresented in our
312 dataset. We first asked if our recorder captured diverse multi-lineage tissue origins of the eye,
313 which is derived from neuroectoderm, surface ectoderm and mesoderm (Figure 5E). Eye cell
314 types were identified as clusters that contained cells from scRNA-seq samples comprising eye
315 tissue exclusively. Retinal cell types were defined as clusters expressing the pan-retinal marker
316 *foxf1b* (Figure 1F), whereas non-retinal cell types were depleted in *foxf1b*. We performed
317 pairwise comparisons of all eye clusters with at least 4 independent barcodes (each with at least
318 2 cells). Since <1% of all barcoded cells were captured by scRNA-seq, we asked if there is cell
319 type-specific barcode enrichment greater than expected by chance (“lineage segregation” in
320 Figure 5E). For cluster pairs where we did not observe significant lineage segregation, we asked
321 if this was due to a lack of sampling (“lineage status undefined”) or true lack of cell type-specific
322 barcode enrichment (“no lineage segregation”). The latter case would indicate that two cell types
323 shared a more recent common ancestor than cell types that segregated earlier. We found that
324 multiple retinal and non-retinal cell types segregated from each other, as would be expected due
325 to early separation of their tissue origins. Interestingly, however, a few non-retinal cell types (e.g.
326 clusters 34, 44, 49) did not fully segregate from retinal cell types, suggesting that they shared a
327 common progenitor. Furthermore, there was extensive lineage segregation between various non-
328 retinal cell types (e.g. clusters 45, 47, 86). In contrast, we did not observe lineage segregation
329 between the different retinal cell types, likely due to the termination of barcode editing prior to
330 terminal divisions. The exception was cluster 28 (cones), which segregated from clusters 15 and
331 32 (cone bipolar cells) and 28 (retinal ganglion cells). Thus, lineage splits between retinal and
332 non-retinal cell types, and within non-retinal subtypes preceded most splits within retinal subtypes.

333
334 Next, we asked if our recorder captured lineage divergences between neurons across brain
335 regions and the retina. Although the hindbrain and retina formed distinct lineages early in
336 development, forebrain and midbrain neurons continued to share progenitors across the same
337 barcoding period (Figure 5F). Pairwise comparisons of all forebrain and midbrain clusters
338 revealed examples of emerging segregation along multiple spatial axes (Figure 5G). For example,
339 we saw evidence of dorsal-ventral split: cluster 9 pallium (dorsal) separated from cluster 25 sub-
340 pallium (ventral). Furthermore, barcode enrichments confirmed rostral-caudal splits: cluster 64
341 habenula separated from clusters 9 and 25 pallium (telencephalon, rostral) and clusters 0 and 13
342 optic tectum (caudal). Overall, the lineage segregations agreed with classic fate mapping
343 experiments (Woo and Fraser, 1995) and correlate with the anteroposterior and dorsoventral
344 gene expression signatures of early progenitors (Figure 4).

345
346 To query the lineage relationships of brain progenitor cell types, we performed pairwise
347 comparisons of progenitor clusters at 15 dpf (Figure 5H). Notably, the upper rhombic lip (URL)
348 progenitors (cluster 12) formed a separate lineage from all progenitor classes except cluster 74,
349 a cycling progenitor subtype expressing *pif1*. Since URL progenitors give rise to granule cells in

350 the cerebellum, we asked if the two cell types shared barcodes. We found that the proportion of
351 barcode overlap was highest between granule cells and URL progenitors (Figure 5I). The URL
352 progenitors formed a distinct cluster as early as 12 hpf (cluster 9) in our transcriptional dataset.
353 Thus, URL progenitors become discrete in both lineage and transcriptional signature relatively
354 early in development.

355
356 In summary, we present an optimized scGESTALT cassette with improved lineage barcode
357 expression and recovery by scRNA-seq. The barcodes display high sequence diversity, which is
358 important for generating large-scale distinct labels in a developing animal. The scGESTALT
359 transgenic line is available as a resource for the community and can be paired with other
360 transgenic lines for temporal, spatial or cell-type specific control of barcode editing (see
361 Discussion).

362

363 **Cell specification trajectories in the retina and hypothalamus**

364

365 With the exception of a few model systems (Clark et al., 2019; Delile et al., 2019; Guo and J. Y.
366 H. Li, 2019; Holguera and Desplan, 2018; Kim et al., 2019; Tambalo et al., 2019), little is known
367 about gene expression cascades that accompany the development of progenitors into terminally
368 differentiated neurons. To address how different neuronal populations become molecularly
369 specialized, we reconstructed gene expression trajectories from 12 hpf to 15 dpf. We first tested
370 our approach on the subsetted retina dataset in which cell types expand from a single cluster at
371 12 hpf to 18 clusters at 15 dpf (Figure 2). UMAP embedding of the subsetted dataset revealed
372 progressive paths from the embryonic state to defined cell types at 15 dpf (Figure 6A, Sup Figure
373 5A). One outlier cluster that expressed *kidins220a* and whose progenitor state may not have been
374 captured in our timepoints, was excluded from further analysis. Although UMAP represents
375 continuity in the data, it does not order individual cells according to their relative developmental
376 time (i.e. pseudotime). Therefore, we also used URD (Farrell et al., 2018) to construct a branching
377 specification tree that represents the developmental trajectories in the retina at a higher resolution
378 (Figure 6B, Sup Figure 5B, Sup Figure 6A-B). Many of the major branching features agreed with
379 the UMAP representation. For example, the trajectories revealed the early segregation of RPE,
380 shared branching of photoreceptor cells, a path towards multiple cone bipolar cell subtypes, and
381 a common branchpoint between amacrine and retinal ganglion cells (RGC).

382
383 Plotting gene expression of known early regulators of eye development and terminal cell type
384 markers on the URD tree supported the inferred specification branches (Figure 6C, Sup Figure
385 7). For example, *pax6a* was most enriched in the amacrine and RGC branches, and *vsx1* marked
386 cone bipolar cells with *fezf2* marking one specific subtype. Notably, our analysis also revealed
387 previously unknown markers and characteristics of horizontal and amacrine cells. Zebrafish
388 horizontal cells are GABAergic (*gad2⁺*, *gad1b⁺*), but unlike mammals where these cells do not
389 express GABA membrane uptake transporters (Deniz et al., 2011), zebrafish cells expressed
390 *slc6a1l* (likely a duplication of *slc6a1* involved in GABA uptake from the synaptic cleft), suggesting
391 that they may be capable of uptake. Additionally, whereas *slc32a1* GABA transporter is expressed
392 in mouse horizontal and amacrine cells (Cueva et al., 2002), we observed restriction of *slc32a1*

393 to amacrine cells and *slc6a11* to horizontal cells. Finally, we detected several novel horizontal cell
394 markers such as *ompa* and *prkaca* (Figure 1F).

395
396 To discover the gene expression trajectories from precursors to different retinal cell types, we
397 used differential gene expression approaches that characterize pseudotime-ordered molecular
398 trajectories. This analysis revealed known and novel regulatory steps (Figure 6D, Sup Figure 8).
399 For example, RGC specification trajectories confirmed several known differentiation regulators
400 including *sox11a*, *sox11b*, *sox6*, *irx4a*, and *pou4f2* (Rheaume et al., 2018). Similarly, known
401 regulators of photoreceptor differentiation such as *isl2a* (Fischer et al., 2011), *prdm1a* (Brzezinski
402 et al., 2010), *otx5* (Vicizian et al., 2003), and *crx* (Shen and Raymond, 2004) were expressed early
403 in our photoreceptor trajectories, while known regulators of cone versus rod fate, such as *six7*
404 (Ogawa et al., 2015), *nr2f1b* (Satoh et al., 2009), and *nr2e3* (J. Chen et al., 2005) were expressed
405 as those trajectories diverged. Furthermore, our analysis revealed novel transcription factors
406 within the gene expression cascades. For example, we detected *runx1t1*, *foxp1b*, *mef2aa* in the
407 RGC pathway; *tfap2a* in horizontal cell trajectory; and *tbx3a* and *tbx2a* in amacrine cell branches.
408 Interestingly, among signaling pathways, we found that both apelin receptors (*aplnra*, *aplnrb*)
409 were expressed in photoreceptor progenitors, while one of their ligands (*apln*) was expressed in
410 differentiating cones; this suggests a potential cell autonomous role for apelin signaling in
411 photoreceptor cells in addition to its role in preventing photoreceptor degeneration via vascular
412 remodeling (McKenzie et al., 2012).

413
414 A surprising result from this analysis was that a Muller glia pathway was detected earlier in
415 zebrafish than expected based on studies in mouse, where these cells are detected late (Centanin
416 and Wittbrodt, 2014; Clark et al., 2019). We found a cluster of cells as early as 20 hpf (cluster 50)
417 that expresses markers (e.g. *cahz*, *rlbp1a*) that are shared with the Muller glia cluster (cluster 33)
418 at 15 dpf (Sup Table). smFISH analysis of Muller glia markers validated their expression at 36 hpf
419 and 2 dpf (Sup Figure 9). Similarly, in our transcriptional trajectories (Figure 6B), the Muller glia
420 expression program is the earliest non-epithelial retinal program to diverge, commencing with the
421 expression of several *her*-family transcription factors (*her4*, *her12*, and *her15*), then proceeding
422 through a cascade of intermediate overlapping expression states such as onset of *fabp7a*,
423 *s100a10b*, and later connexin genes that are characteristic of Muller glia fate (Sup Figure 8). Cells
424 from all timepoints can already be found in the early part of the Muller glia branch. These
425 observations suggest that cells early in development transition from a naive progenitor state to a
426 Muller glia-like transcriptional state, and do so continually during larval development.

427
428 To extend our analysis to a central brain region, we reconstructed specification trajectories and
429 expression cascades for hypothalamic neurons. These cells expanded from a single ventral
430 diencephalon cluster at 12 hpf to 7 clusters at 15 dpf (Figure 6E-H, Sup Figure 6C-D, Sup Figure
431 10). The earliest branchpoint denoted segregation of *prdx1*⁺ and *prdx1*⁻ cells. Committed
432 hypothalamic progenitors in the *prdx1*⁻ trajectory gave rise to neuronal precursors expressing
433 proneural transcription factors such as *ascl1a*, *scrt2*, *insm1a* and *elavl3* (early neuron fate marker)
434 (Sup Figure 11). The specified cell types then matured over time and were characterized by
435 expression of neuronal maturation markers such as *tubb5*, *gap43*, *ywhag2*, *snap25a*, *scg2b* and
436 *elavl4*. The *prdx1*⁻ group further diverged into two major groups: *nrgna*⁺ and *nrgna*⁻ trajectories

437 (Figure 6F). The *nrgna*⁺ branch segregated into GABAergic *tac1*⁺, *synpr* subtype and GABAergic
438 *tac1*⁺, *synpr*⁺ positive subtype. The *nrgna*⁻ branch subdivided into glutamatergic *pdyn*⁺ neurons
439 and a GABAergic branch that further resolved to *sst1.1*⁺ and *tph2*⁺ neuron subtypes. We detected
440 expression of known regulators of hypothalamus development in the early branches such as
441 *shha*, *rx3*, *nkx2.4b*. We also identified new candidate regulators in later branches including *nrgna*
442 in the *synpr*⁺ and *synpr* trajectories, and *sox1a*, *sox1b* and *sox14* in the *pdyn*⁺ trajectory (Sup
443 Figure 12, Figure 1E). The results in the retina and hypothalamus demonstrate that the brain
444 development atlas can be used to reconstruct neuronal differentiation trajectories and define the
445 underlying gene expression cascades

446

447

448 **Differences in progenitor specification strategies between retina and hypothalamus**

449

450 Pseudotime analysis represents cell trajectories in relative but not absolute time (Bendall et al.,
451 2014; Trapnell et al., 2014). Therefore, comparing the developmental and pseudotime age of cells
452 can define whether molecular states are unique to a given developmental stage or persist through
453 development (Figure 6B, 6F). For example, mapping RGC and *pdyn*⁺ neurons from different
454 developmental stages onto the pseudotime trajectory showed the expected maturation of these
455 cell types with developmental age (Sup Figure 13). In addition, even at 15 dpf some RGC and
456 *pdyn*⁺ neurons were still in an immature state, consistent with the continuous growth and
457 differentiation in the zebrafish retina and brain (Centanin and Wittbrodt, 2014; Schmidt et al.,
458 2013).

459

460 To systematically analyze the relationships of pseudotime state and developmental stage, we
461 mapped differentiated cells, precursors and progenitors found in different pseudotime windows to
462 their origin in developmental time. We found that the proportion of differentiated cells increased,
463 whereas the number of early progenitors in both retina and hypothalamus decreased with
464 developmental age. In contrast, precursor cells from an intermediate pseudotime window were
465 present in embryo and larva. These precursor cells expressed genes that were an intermediate
466 of progenitor (e.g. *insm1a*, *her4.1* in hypothalamus (Xie and Dorsky, 2017); *hes2.2*, *rx2* in retina)
467 and early differentiation genes (e.g. *tubb5*, *gap43* in hypothalamus; *foxg1b* in retina). In addition,
468 a second class of retinal progenitors mapped to an earlier pseudotime trajectory but was also
469 present from embryonic to late larval stages (Figure 7, Sup Figure 14). Comparison of these
470 progenitors between 24-36 hpf and 15 dpf identified only 71 differentially expressed genes. The
471 majority of these genes (56/71) increased in all cells of the retina between these stages, while a
472 few (15/71) were only upregulated in the 15 dpf group. A similar population was not detected in
473 the hypothalamus. These observations suggest that as the retina grows, some progenitor cell
474 states observed in the embryo persist later in development without extensive maturation.

475

476

477 **DISCUSSION**

478

479 As the brain develops, embryonic neural progenitor pools transition through many cellular states
480 as they become more committed, diversify into post-embryonic neural progenitors, and undergo

481 terminal differentiation. Although regulators and transcriptional changes of this process have been
482 identified (e.g. using specific driver lines and *in situ* detection of select genes), the global
483 transcriptional networks mediating the sequential activation and maturation of neurogenic
484 programs from embryo to later stages are largely unknown. To help address this question, we
485 used scRNA-seq to generate a zebrafish brain development atlas. This resource supports the
486 identification of marker genes, the comparison of cell types, and the dissection of cell specification
487 and differentiation trajectories during vertebrate brain development.

488
489 Our data address how the transcriptional programs of neural progenitors vary and contribute to
490 fate-restriction during development. Different models to explain these processes have been
491 proposed. For example, neural progenitors of the medial and lateral mouse ganglionic eminence,
492 which give rise to cortical interneurons, have been found to converge to a shared mitotic signature
493 regardless of their region of origin, followed by expression of cardinal fate-specific transcription
494 factors post-mitotically (Mayer et al., 2018). In contrast, the spinal cord has dedicated pools of
495 domain-specific neural progenitors that retain domain-specific signatures (Delille et al., 2019;
496 Jessell, 2000; Lee and Pfaff, 2001; Sagner and Briscoe, 2019). Our results indicate that early
497 embryonic neural progenitors in the brain are transcriptionally distinct from late larval neural
498 progenitors. Gene expression profiles of neural progenitors switch from strong spatially
499 segregated signatures in early embryos to proliferative and non-proliferative states in late larvae.
500 These cell state changes might reflect developmental shifts from an establishment program during
501 gastrulation, where strong spatial patterning cues set up regional boundaries, to a maintenance
502 program at late stages, where progenitors are geographically confined and express dampened
503 regional restriction signatures. Although expression of some spatially-enriched transcription
504 factors (e.g. *pax6a*, *eng2a*, *nkx2.4a*) and signaling proteins detected in embryonic progenitors are
505 also detected in late progenitors, the overall signatures are different, as these factors co-vary with
506 different sets of genes in larva relative to embryo.

507
508 The expression of pan-progenitor markers at larval stages raises the question of how neural
509 progenitor pools remain or become fate restricted. There are several different scenarios that might
510 address this question. First, it is conceivable that embryo and larva share a minimal core set of
511 regionally-restricted transcription factors that are sufficient to ensure spatial restriction, despite
512 differences in their relative expression levels and downstream targets. Spatial genes that are
513 highly expressed in the embryo may be lowly expressed in the larva, and be sufficient to maintain
514 regionally-restricted cell states. Second, cell-type specific transcription factors rather than
515 spatially defined regulators might guide specification and differentiation at these stages,
516 independent of positional information. Such signatures would be difficult to analyze via scRNA-
517 seq, which is biased towards recovering highly expressed genes. Third, it is also possible that
518 restrictions at the genomic level, such as chromatin accessibility, may ensure that cells maintain
519 the signature of their spatial origin. Fate mapping experiments of early and late neural progenitors,
520 profiling open chromatin states of neural progenitors, and transcriptome analyses that recover
521 lowly expressed genes will provide further insight into these questions.

522
523 Our reconstruction of specification trajectories for cell types in the retina and hypothalamus
524 revealed several findings. First, our data supports a multipotent progenitor model whereby

525 multiple differentiated cell types can be traced to common post-embryonic progenitors. For
526 example, all retinal neurons can be traced to an early pseudotime progenitor branch containing
527 cells from larval stages, consistent with multipotency and fate stochasticity of zebrafish retinal
528 progenitors (Boije et al., 2015; He et al., 2012). The early emergence of Muller glia observed in
529 both the time course atlas and eye trajectory reconstruction is particularly interesting in light of
530 clonal analyses. For example, single retinal progenitor cells in zebrafish give rise to clones
531 comprised of neurons and one Muller glia cell (Rulands et al., 2018). This observation has been
532 interpreted as evidence for a progenitor that first gives rise to neurons and then differentiates into
533 a Muller glia cell. However, it is also conceivable based on our data that an early common
534 progenitor divides, with one daughter expanding to give rise retinal neurons while the other
535 daughter forms Muller glia. Second, our results reveal that whereas progenitor cell types in the
536 rest of the brain appear molecularly distinct between the embryo and larva, there are progenitor
537 cell states in the eye that are maintained from the embryo to larva (Figure 4 and Figure 7). A
538 subset of 15 dpf retinal progenitors have similar transcriptional states as observed in the
539 embryonic eye. This observation raises the possibility that a subset of long-term retinal
540 progenitors may be “frozen” in an embryonic phase that could possibly underlie the multi-fate
541 potential of these cells. An independent study of zebrafish retinal stem cells has proposed a
542 similar conclusion (Xu et al., 2020). Collectively, these findings highlight differences in neurogenic
543 programs in the central nervous system, and underscore the power of investigating multiple
544 specification trajectories simultaneously.

545
546 Our results also highlight differences between zebrafish and mammalian neurogenesis. For
547 example, we detected pan-neuronal transcriptional signatures (e.g. *neurod1*, *ascl1a*, *insm1a*,
548 *neurog1*) in zebrafish radial glia and other progenitors at late stages of development, suggesting
549 that neurons remain the principal output of these cells. This is consistent with fate mapping studies
550 that have shown that zebrafish radial glia persist into adulthood and contribute to neurogenesis
551 (Schmidt et al., 2013). In contrast, radial glia progenitor cells in the developing embryonic mouse
552 brain shift from neurogenic to gliogenic programs (Mission et al., 1991; Schmechel and Rakic,
553 1979).

554
555 While developmental atlases and trajectories can help identify cellular differentiation paths, a full
556 understanding of cell type specification requires lineage tracing experiments. To catalyze such
557 approaches we introduced improvements to scGESTALT through a redesigned recorder cassette
558 for optimized mRNA expression and library compatibility with the 10X Chromium scRNA-seq
559 platform. The resulting higher recovery of barcodes allows more dense reconstruction of lineage
560 trees. Our analysis revealed differences between the timing of segregation between different brain
561 regions: neuronal lineages in the retina and hindbrain diverged earlier than the forebrain and
562 midbrain. These results complement classic zebrafish fate maps of brain compartmentalization
563 (Woo and Fraser, 1995) and recent analysis of clonal cells in forebrain and midbrain (Solek et al.,
564 2017). Furthermore, our findings support early transcriptional and lineage segregation of
565 cerebellar upper rhombic lip progenitors relative to other classes of progenitor cells. To query
566 additional lineage divergences and combine with cellular trajectories, our optimized recorder can
567 be readily adapted for barcoding lineages at developmental windows that correspond to different

568 branches of the specification trees (Raj et al., 2018b) or combined with cell- or tissue-specific
569 Cas9 driver lines to introduce lineage labels in populations of interest.

570
571 The resources presented here lay the groundwork for characterizing lineage histories and
572 transcriptional changes underlying the development and diversification of the vertebrate brain.
573 Future extensions include the generation of transgenic reporters to select populations of interest
574 and perform deeper analyses of cell type heterogeneity and differentiation (Pandey et al., 2018).
575 Cell specification trajectories can be extended to include additional subregions of the brain to
576 generate increasingly complex trees and combined with other zebrafish scRNA-seq datasets
577 (Cosacak et al., 2019; Farnsworth et al., 2020; Farrell et al., 2018; Lange et al., 2020; Pandey et
578 al., 2018; Tambalo et al., 2020; Wagner et al., 2018; Xu et al., 2020) to trace complete trajectories
579 from gastrulation to adulthood. Finally, it will be interesting to perform comparative studies by
580 using our atlas in conjunction with data described in a recent preprint (La Manno et al., 2020).

581

582

583 **METHODS**

584

585 **Zebrafish husbandry**

586 All vertebrate animal work was performed at the facilities of Harvard University, Faculty of Arts &
587 Sciences (HU/FAS). This study was approved by the Harvard University/Faculty of Arts &
588 Sciences Standing Committee on the Use of Animals in Research & Teaching under Protocol No.
589 25–08. The HU/FAS animal care and use program maintains full AAALAC accreditation, is
590 assured with OLAW (A3593-01), and is currently registered with the USDA.

591

592 **Chromogenic in situ hybridization**

593 Embryos were dechorionated with forceps and then fixed in 4% PFA in 1X PBS (pH 7.4) overnight
594 at 4°C. After fixation, embryos were dehydrated in methanol series (0%, 25%, 50%, 75% and
595 100% MetOH in PBSTween 0.3% (PBST)) and stored in 100% methanol at –20°C. Embryos were
596 rehydrated by reversing the methanol series for 10 min in each step at room temperature (RT)
597 and washed 2 × 5 min in PBST. To bleach pigment in 2 dpf fish, larvae were incubated for 10 min
598 in bleaching solution (3% H₂O₂/0.5% KOH in ddH₂O) at room temperature and washed 3 x 5 min
599 in PBST (Thisse et al., 2004). For permeabilization, 2 dpf larvae were incubated with Proteinase
600 K (10 µg/ml in PBST) for 2 min at RT and postfixed in 4% PFA in 1X PBS for 30 min at RT.
601 Afterwards, embryos were washed 3 x 5 min in PBST at RT, prehybridized in HYB⁺ solution (50%
602 Deionized Formamide (Amresco), 5X SSC (Ambion), 0.1% Tween-20, 5mg/ml Torula RNA
603 (Sigma) in ddH₂O) for 3 hours at 69°C, and hybridized overnight with the antisense probes diluted
604 in HYB⁺ at 69°C. The rest of the steps were performed as described previously, by hand (Navajas
605 Acedo et al., 2019). Before imaging, embryos were cleared using an increasing MetOH series.
606 For imaging of 12 hpf embryos, the yolk was dissected away, and the embryos were flat mounted
607 on a microscope slide and covered with a cover slip. Larvae were photographed on a Zeiss
608 AxioZoom.V16.

609 The antisense probes were synthesized from DNA fragments amplified from TLAB zebrafish
610 cDNA using the following primers: *klf17* (Fw GAAGGAAAGACTGCATCCTGAC; Rv
611 CTGCTGTCCCAAATAGGAGTT), *ptgs2a* (Fw CGAGGACTATGTTCAGCACTTG; Rv

612 TGCACATCGATCACAATACAAA), *tp63* (Fw TGCTTTGCTAAATTGTGCTGTC; Rv
613 ATTGCCGCTTATGAGAATCAAG), *cavin2a* (Fw GAGCCTTCTCGTGCTAACAAGT; Rv
614 CAGGCATTTTCAGTTCAATTTCA), *sox1a* (Fw AATCAAGACCGCGTAAAGAGAC; Rv
615 TTTGGTGGAGTGTCTGAATG), *pdyn* (Fw AAGAGAACGCCATACTGAAAGG; Rv
616 GCAGTTACGAATTGCCATGATA), *dlx1a* (Fw AAGGAGGAGAGGTTTCGTTTCA; Rv
617 AGTGTGTGTCAGCAGGTGTCTT).

618

619 smFISH staining and imaging

620

621 Single-molecule FISH probe sets were generated as previously described and coupled to either
622 Atto 647N NHS ester (Millipore Sigma #18373) (*foxg1b*, *cahz*) or Atto 550 NHS ester (Millipore
623 Sigma # 92835) (*ompa*, *rlbp1a*) (Lord et al., 2019). Sectioned larvae were affixed to polylysine-
624 coated #1.5 coverslips, and staining was carried out as previously described (Lord et al., 2019),
625 with each coverslip contained in a well of a plastic 6-well plate. During the probe hybridization
626 step, coverslips were placed upside-down onto a 100µl droplet of probe solution on Parafilm
627 (Farack and Itzkovitz, 2020). Sample mounting was performed as previously described (Lord et
628 al., 2019). Mounted samples were imaged on an Olympus spinSR spinning disk microscope
629 fitted with a UPLAPO 60X/1.5 oil immersion objective using 0.3µm slices.

630

631 smFISH image processing

632

633 All image processing was performed in Fiji (Schindelin et al., 2012). Rolling-ball background
634 subtraction (radius 25 pixels) was performed on smFISH channels before maximum intensity
635 projections were produced from 30 slices (Figure 1F) or 50 slices (Sup Figure 9) of processed z-
636 stacks. Channels were scaled individually, maximizing for visibility.

637

638 **Optimization of scGESTALT lineage cassette**

639 In our previous iteration of scGESTALT, the barcode capture rate by scRNA-seq was 6-28%. (Raj
640 et al., 2018b), thereby limiting the density of lineage tree reconstruction. To improve recovery we
641 adapted a different transgenic cassette (Yoshinari et al., 2012) for lineage recording. This cassette
642 has the following modifications compared to our previous recorder: (1) The heat-shock inducible
643 (*hsp70l*) promoter of the previous version is now replaced with a constitutive ubiquitous promoter
644 (*medaka beta-actin*) to drive strong widespread expression of the barcode mRNA. Expression of
645 the cassette was confirmed by fluorescence and the signal was more intense than that obtained
646 with the heat shock promoter. Furthermore, this version eliminates the requirement to heat shock
647 edited animals to express the barcode prior to scRNA-seq experiments. (2) We adapted the 3'
648 end of the DsRed open reading frame as a lineage recorder cassette with up to 8 sgRNA target
649 sites positioned next to each other. This vastly improved expression of the construct compared to
650 our previous version where the recording cassette was placed downstream of the DsRed open
651 reading frame. (3) We made library preparation compatible with the 10X Genomics platform.

652

653 To generate scGESTALT.2 barcode founder fish, one-cell embryos were injected with zebrafish
654 codon optimized Tol2 mRNA and pT2O_{lactb}:loxP-dsR2-loxP-EGFP vector (gift from Atsushi
655 Kawakami (Yoshinari et al., 2012)). Potential founder fish were screened for widespread DsRed

656 expression and grown to adulthood. Adult founder transgenic fish were identified by outcrossing
657 to wild type fish and screening clutches of embryos for ubiquitous DsRed expression. Single copy
658 scGESTALT.2 F1 transgenics were identified using qPCR, as described previously (McKenna et
659 al., 2016; Pan et al., 2013; Raj et al., 2018b).

660
661 SgRNAs specific to sites 1-8 of the scGESTALT.2 array were generated by in vitro transcription
662 as previously described (Raj et al., 2018a). To initiate early barcode editing, single copy
663 scGESTALT.2 F1 male transgenic adults were crossed to wildtype female adults and one-cell
664 embryos were injected with 1.5 nl of Cas9 protein (NEB) and sgRNAs 1-8 in salt solution (8 μ M
665 Cas9, 100 ng/ μ l pooled sgRNAs, 50 mM KCl, 3 mM MgCl₂, 5 mM Tris HCl pH 8.0, 0.05% phenol
666 red). Since editing results in loss of DsRed signal, transgenic animals were distinguished from
667 wild type animals by amplifying the scGESTALT.2 barcode by PCR using genomic DNA from the
668 tail fin at 15 dpf. In the experiments presented in this study, early lineage decisions were barcoded
669 by injecting reagents at the one-cell stage. It is worth noting that the scGESTALT.2 barcode can
670 be readily paired with a two-step barcoding protocol. This would require the establishment of a
671 second stable transgenic line for in vivo expression of Cas9 and a subset of sgRNAs matching
672 the target sequences of the new barcode cassette to enable sequential barcoding at early and
673 late stages. Such a line can be established using a similar step-by-step guidance that is detailed
674 in (Raj et al., 2018a).

675

676 **Processing of samples for scRNA-seq time course**

677 Wild type embryos (12 hpf, 14 hpf, 16 hpf, 18 hpf, 20 hpf, 24 hpf, 36 hpf) and larvae (2 dpf, 3 dpf,
678 5 dpf, 8 dpf) were used for scRNA-seq analysis. Samples for 15 dpf had a mix of wild type and
679 barcode edited larvae. Two of the 15 dpf samples consisted of only eye cells (no brain). Embryos
680 from 12 hpf to 36 hpf were first de-chorionated by incubating in 1 mg/ml pronase (Sigma-Aldrich)
681 at 28 C for 6-7 min until chorions began to blister, and then washed three times in ~200 ml of
682 zebrafish embryo medium (5 mM NaCl, 0.17 mM KCl, 0.33 mM CaCl₂, 0.33 mM MgSO₄, 0.1%
683 methylene blue) in a glass beaker. Embryos were de-yolked using two pairs of watchmaker
684 forceps, and the heads were chopped just anterior of the spinal cord. All processing steps were
685 done using 100 mm Petri dishes coated with Sylgard (Raj et al., 2018a). Samples from 2 and 3
686 dpf were processed similarly to the embryos, except they were not de-chorionated as they had
687 hatched out of the chorions. Larvae from 5 dpf to 15 dpf were dissected to remove whole brains
688 and eyes as described previously (Raj et al., 2018a). The following numbers of embryos and
689 larvae were used for each timepoint: 12 hpf – ~20 embryos; 14 hpf – ~20 embryos; 16 hpf – ~18
690 embryos; 18 hpf – ~18 embryos; 20 hpf – ~30 embryos; 24 hpf – ~30 embryos; 36 hpf – ~15
691 embryos; 2 dpf – ~30 larvae; 3 dpf – ~30 larvae; 5 dpf – ~25 larvae; 8 dpf – ~25 larvae; 15 dpf
692 – ~15 larvae. Tissues were dissociated into single cells using the Papin Dissociation Kit
693 (Worthington) as described previously (Raj et al., 2018a). Cells were resuspended in 50 μ l to 150
694 μ l of DPBS (Life Technologies) depending on anticipated amount of material, and counted using
695 a hemocytometer. Samples were run on the 10X Genomics scRNA-seq platform according to the
696 manufacturer's instructions (Single Cell 3' v2 kit). Libraries were processed according to the
697 manufacturer's instructions. Transcriptome libraries were sequenced using NextSeq 75 cycle kits.

698

699 **scGESTALT.2 library prep**

700 To generate scGESTALT.2 libraries, lineage edited 15 dpf samples post cDNA amplification and
701 prior to fragmentation were split into two halves. One half was processed for transcriptome
702 libraries as instructed by the manufacturer. The other half was processed for lineage libraries as
703 follows. To enrich for scGESTALT.2 lineage barcodes, 5 µl of the whole transcriptome cDNA
704 was PCR amplified using Phusion polymerase (NEB) and 10XPCR1_F (CTACACGACGCTCTT
705 CCGATCT) and GP10X2_R (GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT GCTGCTTC
706 ATCTACAAGGTGAAG). The reaction (98 C, 30 s; [98 C, 10 s; 67 C, 25 s; 72 C, 30 s] x 14-15
707 cycles; 72 C, 2 min) was cleaned up with 0.6X AMPure beads and eluted in 20 µl EB buffer
708 (Omega). Finally, adapters and sample indexes were incorporated in another PCR reaction
709 using Phusion polymerase and 10XP5Part1long (AATGATACGGCGACCACCGA
710 GATCTACTACTTTCC CTACACGACGCTCTTCCGATCT) and 10XP7Part2Ax
711 (CAAGCAGAAGACGGCATAACGAGAT-xxxxxxx-GTGACTGGAGTTCAGACGTGT), where x
712 represents index bases. These include A1: GGTTTACT; A2: TTTTCATGA; A3: CAGTACTG; A4:
713 TATGATTC. Thus, up to 4 scGESTALT.2 samples were multiplexed in a sequencing run.
714 Libraries were sequenced using MiSeq 300 cycle kits and 20% PhiX spike-in. Sequencing
715 parameters: Read1 250 cycles, Read2 14 cycles, Index1 8 cycles, Index2 8 cycles. Standard
716 sequencing primers were used.

717

718 **Bioinformatic processing of raw sequencing data and cell type clustering analysis**

719 Transcriptome sequencing data were processed using Cell Ranger 2.1.0 according to the
720 manufacturer's guidelines. scGESTALT.2 sequencing data were processed with a custom
721 pipeline (https://github.com/aaronmck/SC_GESTALT) as previously described (Raj et al., 2018b).
722 The scGESTALT.2 barcode for each cell was matched to its corresponding cell type (tSNE cluster
723 membership) assignment using the cell identifier introduced during transcriptome capture. Cells
724 with fewer than 500 expressed genes, greater than 9% mitochondrial content or very high
725 numbers of UMIs and gene counts that were outliers of a normal distribution (likely
726 doublets/multiplets) were removed from further analysis. Clustering analysis was performed using
727 the Seurat v2.3.4 package (Butler et al., 2018) as described previously (Raj et al., 2018b). For
728 Figure 3 and Sup Figure 2, we selected the list of transcription factors, neuropeptides and their
729 receptors, and genes involved in neuron electrophysiology from our enriched marker analysis and
730 previous literature (R. Chen et al., 2017; Pandey et al., 2018; Tiklová et al., 2019; Zeisel et al.,
731 2018).

732

733 **Construction of lineage trees from GESTALT barcodes.**

734 All unique barcodes were then encoded into an event matrix and weights file, as described
735 previously (McKenna et al., 2016; Raj et al., 2018b), and were processed using PHYLIP mix with
736 Camin-Sokal maximum parsimony (Felsenstein, 1989). Individual cells were then grafted onto the
737 leaves matching their barcode sequence. After the subtrees were attached, we repeatedly
738 eliminated unsupported internal branching by recursively pruning parent-child nodes that had
739 identical barcodes. Cell annotations are then added to the corresponding leaves. The resulting
740 tree was converted to a JSON object, annotated with cluster membership, and visualized with
741 custom tools using the D3 software framework.

742

743 **Lineage segregation analysis between cell types**

744 We combined all barcodes obtained from 4 fish. For our analysis, we only considered barcodes
745 with at least two cells, and we only analyzed cell types with at least 4 barcodes. To test
746 segregation between any two cell types/clusters, we first retrieved all barcodes that were present
747 in at least one of the two cell types. Then, we split these barcodes into two categories: “shared
748 barcode” or “specific barcode”. A shared barcode was defined as one that contains cells from
749 both cell types. In contrast, a specific barcode was defined as one that only contains cells from
750 one of the two cell types. Our null hypothesis is that the two cell types come from the same
751 ancestor at the time of Cas9 editing. Thus, we asked whether the number of observed specific
752 barcodes can be explained by chance under the null hypothesis. If it cannot be explained by
753 chance, it indicates that the two cell types have segregated.

754

755 To do so, we performed a randomization test as below:

- 756 1. We generated a pool of cells. The size of the pool is the total number of cells from the two
757 cell types. The ratio of the two cell types in the pool is equal to the ratio observed in the
758 real data. Under the null hypothesis, the pool of cells come from the same ancestor, so
759 they would share the same barcode.
- 760 2. For each barcode, we randomly sampled the same number of cells of this barcode from
761 the pool of cells.
- 762 3. We repeated this for all the barcodes, and then calculated the number of barcodes that
763 only contain one cell type (i.e. “specific barcode”).
- 764 4. We repeated steps 2 and 3 5000 times.
- 765 5. We calculated how many times (for example n times) the number of specific barcodes
766 from the random sampling process is greater than or equal to the number of specific
767 barcodes from the real data.
- 768 6. The probability that the number of specific barcodes can be explained by chance under
769 the null hypothesis is $n/5000$.
- 770 7. If the probability < 0.01 ($pvalue < 0.01$), we rejected the null hypothesis.

771

772 Next, for each cell type we split its corresponding pairwise comparison cell types into two
773 categories: “with segregation” or “other”. For the “other” category, we considered two
774 interpretations. First, it could signify that there is no segregation between the two cell types.
775 Second, it could suggest that we did not recover enough cells with barcode information, such that
776 there is not enough power to detect lineage segregation (low sampling). To distinguish between
777 the two scenarios, for each cell type in the two categories, we calculated the ratio between the
778 number of cells with barcodes and the number of all cells from scRNA-seq. If the ratio of one cell
779 type from the “other” category is greater than or equal to the smallest ratio from the first category
780 (“with segregation”), it indicates this cell type did not have low sampling issues. Thus, it supports
781 the interpretation that there is no segregation between the queried cell types. Otherwise, we
782 assign the cell type pair as “undefined” (i.e. insufficient sampling power to query lineage
783 segregation).

784

785 **Granule cell analysis**

786 For each progenitor cell type, we used barcodes that did not include any cells from the other nine
787 progenitor cell types. The Jaccard Index between each progenitor cell type and granule cell was
788 calculated as below:

$$789 \quad \text{Jaccard Index} = \frac{\text{the number of shared barcodes between the two cell types}}{\text{the number of barcodes in either cell type}}$$

790

791 **Analyzing dampened spatial correlations in progenitors**

792 Progenitors were isolated by subsetting the data to include clusters expressing markers such as
793 *sox19a*, *her* genes, *pcna*, *mki67*, *fabp7a*, *gfap*, *id1*, etc (Supplementary Table). Cells from 12 hpf
794 – 18 hpf were considered embryonic progenitors, cells from 20 hpf – 3 dpf were considered
795 intermediate progenitors, and cells from 5 dpf – 15 dpf were considered larval progenitors.
796 Variable genes were calculated for embryonic, intermediate and larval progenitors separately
797 using the `FindVariableGenes` function from Seurat v2.3.4 with parameters: *x.low.cutoff* = 0.015,
798 *x.high.cutoff* = 3, *y.cutoff* = 0.7. Then, a list of 79 transcription factors with known spatial signatures
799 was assembled by consulting previously described histological analysis (ZFIN) together with
800 those that were identified as gene markers for neuronal clusters in our dataset. Separately in the
801 three progenitor groups, the pairwise Pearson correlation was calculated pairwise between all
802 genes detected as variable in each progenitor group. For several thresholds between 0.2–0.8,
803 the number of genes that correlated more strongly than the threshold with any of the 79 spatial
804 transcription factors (excluding self-correlation) were determined. The strongest correlations were
805 observed in the embryonic population, followed by the intermediate population, and for any
806 threshold, more genes correlated with the spatial TFs in the embryonic progenitors than the larval
807 progenitors.

808

809 **Construction and analysis of branching transcriptional trajectories using URD**

810 We built branching transcriptional trajectories from cells of the retina and hypothalamus to
811 determine the molecular events that occur as cells diversify and differentiate in these tissues.
812 First, cells from the retina and hypothalamus were isolated from each stage by determining
813 clusters that belonged to these tissues by expression of marker genes.

814

815 Determination of variable genes

816 For URD trajectory analyses, a more restrictive set of variable genes was calculated on each
817 subset of the data, as previously described (Farrell et al., 2018; Pandey et al., 2018) using the
818 URD `findVariableGenes` function, with parameter *diffCV.cutoff* = 0.3. Briefly, a curve was fit that
819 related each gene's coefficient of variation to its mean expression level and represents the
820 expected coefficient of variation resulting from technical noise, given a gene's mean expression
821 value; genes with much higher coefficients of variation likely encode biological variability and were
822 used downstream.

823

824 Removal of outliers

825 Poorly connected outliers can disrupt diffusion map calculation and so were removed from the
826 data. A *k*-nearest neighbor network was calculated between cells (Euclidean distance in variable
827 genes) with 100 nearest neighbors. Cells were then removed based on either unusually high
828 distance to their nearest neighbor or unusually high distance to their 20th nearest neighbor, given

829 their distance to their nearest neighbor using the URD function *knnOutliers* (retina: *x.max* = 40,
830 *slope.r* = 1.05, *int.r* = 4.3, *slope.b* = 0.75, *int.b* = 11.5; hypothalamus: *x.max* = 40, *slope.r* = 1.1,
831 *int.r* = 3, *slope.b* = 0.66, *int.b* = 11.5).

832

833 Removal of doublets by NMF modules

834 To remove putative cell doublets (i.e. where two cells are encapsulated into a single droplet and
835 processed as one cell), which can disrupt trajectory relationships, we removed cells that
836 expressed multiple NMF (non-negative matrix factorization) modules characteristic of different
837 expression programs, as previously described (Siebert et al., 2019). NMF modules were
838 computed using a previously published NMF framework (<https://github.com/YiqunW/NMF>)
839 (Farrell et al., 2018). The analysis was performed on log-normalized read count data for a set of
840 variable genes using the *run_nmf.py* script with the following parameters: *-rep* 5 *-scl* "false" *-miter*
841 10000 *-perm* True *-run_perm* True *-tol* 1e-6 *-a* 2 *-init* "random" *-analyze* True. Several *k*
842 parameters were evaluated for each tissue, and *k* was chosen to maximize the number of
843 modules, while minimizing the proportion of modules defined primarily by a single gene (retina, *k*
844 = 45; hypothalamus, *k* =). Modules were used downstream that (a) had a ratio between their top-
845 weighted and second-highest weighted gene of < 5, and (b) exhibited a strong cell-type signature,
846 as determined by plotting on a UMAP representation and looking for spatial restriction. Pairs of
847 modules that were appropriate for using to remove doublets (and that did not define transition
848 states) were determined using the URD function *NMFDoubletsDefineModules* with parameters
849 *module.thresh.high* = 0.4, and *module.thresh.low* = 0.15. Putative doublets were identified using
850 the URD function *NMFDoubletsDetermineCells* with parameters *frac.overlap.max* = 0.03,
851 *frac.overlap.diff.max* = 0.1, *module.expressed.thresh* = 0.33 and were then removed.

852

853 Choice of root and tips

854 Branching transcriptional trajectories in the retina and hypothalamus were constructed using URD
855 1.1.1 (Farrell 2018). Briefly, cells from the first stage of the time course (12 hpf) were selected as
856 the 'root' or starting point for the tree. Terminal cell types comprised the clusters at 15 dpf from
857 these tissues, with the exception of clusters that were clearly progenitor or precursors based on
858 known gene expression (retina: 29, 39, 43). Additionally, in the retina, one cluster (96) was
859 excluded because it did not seem that any related cell types had been recovered in previous
860 stages.

861

862 Construction of branching transcriptional trajectories

863 A diffusion map was calculated using *destiny* (Haghverdi et al., 2015; 2016), using 140 (retina) or
864 100 (hypothalamus) nearest neighbors (approximately the square root of the number of cells in
865 the data), and with a globally-defined sigma of 14 (retina) or 8 (hypothalamus) — slightly smaller
866 than the suggested sigma from *destiny*. Pseudotime was then computed using the simulated
867 'flood' procedure previously described (Farrell et al., 2018), using the following parameters: *n* =
868 100, *minimum.cells.flooded* = 2. Biased random walks were performed to determine the cells
869 visited from each terminal population in the data as previously described (Farrell et al., 2018),
870 using the following parameters: *optimal.cells.forward* = 40, *max.cells.back* = 80, *n.per.tip* = 50000,
871 *end.visits* = 1. The branching tree was then constructed using URD's *buildTree* function with the
872 following parameters: *divergence.method* = "ks" (hypothalamus) or *divergence.method* =

873 "preference" (retina), *save.all.breakpoint.info* = TRUE, *cells.per.pseudotime.bin* = 40,
874 *bins.per.pseudotime.window* = 5, *p.thresh* = 0.0001 (hypothalamus) or , *p.thresh* = 0.01 (retina),
875 and *min.cells.per.segment* = 10. The resulting trees were then evaluated using known marker
876 genes and branch regulators.

877

878 Finding genes that vary during differentiation

879 Genes were selected for inclusion in gene cascades based on their differential expression relative
880 to other cell types in the tissue. See the Supplementary Analysis for the full set of commands
881 used. Within each tissue, cells were first compared in large populations that defined major cell
882 types (retina: cone bipolar cells, photoreceptors, amacrine cells, retinal ganglion cells, horizontal
883 cells, Muller glia, retinal pigmented epithelium; hypothalamus: *prdx1*+ neurons, *pdyn*+ neurons,
884 GABAergic *dlx*+ neurons, *nrgna*+ neurons). Comparisons were performed pairwise, and genes
885 were considered differential in a population if they were upregulated compared to at least 2
886 (hypothalamus) or 3 (retina) other groups. Genes were considered differentially expressed based
887 on their expression fold-change (retina: ≥ 1.32 -fold change, hypothalamus: ≥ 1.41 -fold change)
888 and their performance as a precision-recall classifier for the two cell populations compared (≥ 1.1 -
889 fold better than a random classifier). Additionally, the *aucprTestAlongTree* function from URD was
890 used to select additional genes by performing pairwise comparisons, starting from a terminal cell
891 type and comparing at each branchpoint along the way, back to the root (Farrell et al., 2018).
892 Genes were selected based on expression fold-change between branchpoints (hypothalamus:
893 ≥ 1.74 -fold upregulated; hypothalamus, populations with small cell numbers (GABAergic *dlx*+
894 cells): ≥ 1.51 -fold upregulated; retina: ≥ 1.32 -fold upregulated), their function as a precision-recall
895 classifier between branchpoints (hypothalamus: ≥ 1.2 -fold better than a random classifier;
896 hypothalamus, populations with small cell numbers (GABAergic *dlx*+ cells): ≥ 1.15 -fold better than
897 a random classifier; retina: ≥ 1.1 -fold better than a random classifier), their function as a precision
898 recall classifier globally (i.e. between the entire trajectory leading to a cell type and the rest of the
899 tissue): ≥ 1.03 -fold better than a random classifier, and their upregulation globally (i.e. between
900 the entire trajectory leading to a cell type and the rest of the tissue): ≥ 1.07 -fold upregulated.
901 Mitochondrial, ribosomal, and tandem duplicated genes were excluded. Cells were ordered
902 according to pseudotime, split into groups of at least 25 cells that differ at least 0.005 in
903 pseudotime, and the mean expression was determined with a 5-group moving window. A spline
904 curve was fit to the mean expression vs. pseudotime relationship of selected genes, using the
905 *smooth.spline* function from R's *stats* package, with the parameter *spar* = 0.5. Genes were then
906 sorted according to their peak expression in pseudotime, normalized to their max expression
907 observed in the tissue, and plotted on a heatmap.

908

909 Analyzing progenitor populations

910 To determine whether retinal progenitors mature transcriptionally over time, we looked for genes
911 that were differentially expressed between young and old progenitors. We chose cells that
912 occupied the same region of the URD tree from either early (24 / 36 hpf) or late (15 dpf) stages.
913 We looked for genes that were differentially expressed in 15 dpf progenitors that: (1) were 1.1-
914 fold better as a precision-recall classifier than random, (2) changed ≥ 1.32 -fold in expression, (3)
915 were expressed in at least 20% of progenitors, (4) had a mean expression value ≥ 0.8 , and (5)

916 were more differentially expressed than equally sized cell populations chosen at random at least
917 99% of the time.

918
919 To determine whether cells were found in progenitor or precursor states long-term, we first defined
920 progenitor and precursor states by cells' assignment in the URD tree, cross-referenced with the
921 expression of progenitor / precursor markers. We then determined how many cells from different
922 stages fell into each of these different states.

923
924

925 **ACKNOWLEDGEMENTS**

926
927 We thank members of the Schier lab for discussion and advice, the Bauer Core Facility (Harvard)
928 and the Molecular Biology Core Facility (Dana Farber Cancer Institute) for sequencing services,
929 the Harvard zebrafish facility staff for technical support, and the Imaging Core Facility of the
930 Biozentrum for microscopy resources. We thank M. Shafer and J. Gagnon for comments on the
931 manuscript, and H. Boije for comments on retinal lineages. This work was supported by a
932 postdoctoral fellowship from the Canadian Institutes of Health Research and 1K99HD098298 to
933 B.R., 1K99HD091291 to J.A.F., R00HG010152 to A.M., R01HD85905, DP1HD094764, ERC
934 834788, an Allen Discovery Center grant, and a McKnight Foundation Technological Innovations
935 in Neuroscience Award to A.F.S.

936
937

938 **AUTHOR CONTRIBUTIONS**

939
940 B.R. and A.F.S. conceived and designed the study. B.R., J.A.F., J.L., J.E.K, and A.F.S. interpreted
941 the data. B.R., J.A.F. and A.F.S. wrote the manuscript. B.R. and J.L.L. generated transgenic lines.
942 B.R. performed scRNA-seq and scGESTALT experiments and data processing. J.L. analyzed
943 scGESTALT data with assistance from B.R. and J.E.K. L.Y.D. generated violin plots of neuron
944 subtype diversity. J.N.A. performed chromogenic in situs. A.N.C. and J.E.K. performed smFISH
945 experiments. J.A.F. performed URD trajectory analysis with assistance from B.R. A.M. generated
946 lineage trees. D.R. developed the R Shiny app for scRNA-seq data exploration.

947
948

949 **REFERENCES**

950
951 Allende, M.L., Weinberg, E.S., 1994. The expression pattern of two zebrafish achaete-scute
952 homolog (ash) genes is altered in the embryonic brain of the cyclops mutant. *Dev. Biol.* 166,
953 509–530. doi:10.1006/dbio.1994.1334
954 Bendall, S.C., Davis, K.L., Amir, E.-A.D., Tadmor, M.D., Simonds, E.F., Chen, T.J., Shenfeld,
955 D.K., Nolan, G.P., Pe'er, D., 2014. Single-cell trajectory detection uncovers progression and
956 regulatory coordination in human B cell development. *Cell* 157, 714–725.
957 doi:10.1016/j.cell.2014.04.005
958 Boije, H., Rulands, S., Dudczig, S., Simons, B.D., Harris, W.A., 2015. The Independent
959 Probabilistic Firing of Transcription Factors: A Paradigm for Clonal Variability in the
960 Zebrafish Retina. *Developmental Cell* 34, 532–543. doi:10.1016/j.devcel.2015.08.011

961 Brzezinski, J.A., Lamba, D.A., Reh, T.A., 2010. Blimp1 controls photoreceptor versus bipolar
962 cell fate choice during retinal development. *Development* 137, 619–629.
963 doi:10.1242/dev.043968

964 Butler, A., Hoffman, P., Smibert, P., Papalexi, E., Satija, R., 2018. Integrating single-cell
965 transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol*
966 36, 411–420. doi:10.1038/nbt.4096

967 Carter, R.A., Bihannic, L., Rosencrance, C., Hadley, J.L., Tong, Y., Phoenix, T.N., Natarajan, S.,
968 Easton, J., Northcott, P.A., Gawad, C., 2018. A Single-Cell Transcriptional Atlas of the
969 Developing Murine Cerebellum. *Current Biology* 28, 2910–2920.e2.
970 doi:10.1016/j.cub.2018.07.062

971 Centanin, L., Wittbrodt, J., 2014. Retinal neurogenesis. *Development* 141, 241–244.
972 doi:10.1242/dev.083642

973 Cepko, C., 2014. Intrinsically different retinal progenitor cells produce specific types of progeny.
974 *Nat. Rev. Neurosci.* 15, 615–627. doi:10.1038/nrn3767

975 Chen, J., Rattner, A., Nathans, J., 2005. The rod photoreceptor-specific nuclear receptor Nr2e3
976 represses transcription of multiple cone-specific genes. *J. Neurosci.* 25, 118–129.
977 doi:10.1523/JNEUROSCI.3571-04.2005

978 Chen, R., Wu, X., Jiang, L., Zhang, Y., 2017. Single-Cell RNA-Seq Reveals Hypothalamic Cell
979 Diversity. *Cell Rep* 18, 3227–3241. doi:10.1016/j.celrep.2017.03.004

980 Clark, B.S., Stein-O'Brien, G.L., Shiau, F., Cannon, G.H., Davis-Marcisak, E., Sherman, T.,
981 Santiago, C.P., Hoang, T.V., Rajaii, F., James-Esposito, R.E., Gronostajski, R.M., Fertig,
982 E.J., Goff, L.A., Blackshaw, S., 2019. Single-Cell RNA-Seq Analysis of Retinal Development
983 Identifies NFI Factors as Regulating Mitotic Exit and Late-Born Cell Specification. *Neuron*
984 102, 1111–1126.e5. doi:10.1016/j.neuron.2019.04.010

985 Cosacak, M.I., Bhattarai, P., Reinhardt, S., Petzold, A., Dahl, A., Zhang, Y., Kizil, C., 2019.
986 Single-Cell Transcriptomics Analyses of Neural Stem Cell Heterogeneity and Contextual
987 Plasticity in a Zebrafish Brain Model of Amyloid Toxicity. *Cell Rep* 27, 1307–1318.e3.
988 doi:10.1016/j.celrep.2019.03.090

989 Cueva, J.G., Haverkamp, S., Reimer, R.J., Edwards, R., Wässle, H., Brecha, N.C., 2002.
990 Vesicular gamma-aminobutyric acid transporter expression in amacrine and horizontal cells.
991 *J. Comp. Neurol.* 445, 227–237. doi:10.1002/cne.10166

992 Delile, J., Rayon, T., Melchionda, M., Edwards, A., Briscoe, J., Sagner, A., 2019. Single cell
993 transcriptomics reveals spatial and temporal dynamics of gene expression in the developing
994 mouse spinal cord. *Development* 146, dev173807. doi:10.1242/dev.173807

995 Deniz, S., Wersinger, E., Schwab, Y., Mura, C., Erdelyi, F., Szabó, G., Rendon, A., Sahel, J.-A.,
996 Picaud, S., Roux, M.J., 2011. Mammalian retinal horizontal cells are unconventional
997 GABAergic neurons. *J. Neurochem.* 116, 350–362. doi:10.1111/j.1471-4159.2010.07114.x

998 Farack, L., Itzkovitz, S., 2020. Protocol for Single-Molecule Fluorescence In Situ Hybridization
999 for Intact Pancreatic Tissue. *STAR Protocols* 1, 100007. doi:10.1016/j.xpro.2019.100007

1000 Farnsworth, D.R., Saunders, L.M., Miller, A.C., 2020. A single-cell transcriptome atlas for
1001 zebrafish development. *Dev. Biol.* 459, 100–108. doi:10.1016/j.ydbio.2019.11.008

1002 Farrell, J.A., Wang, Y., Riesenfeld, S.J., Shekhar, K., Regev, A., Schier, A.F., 2018. Single-cell
1003 reconstruction of developmental trajectories during zebrafish embryogenesis. *Science* 360,
1004 eaar3131. doi:10.1126/science.aar3131

1005 Felsenstein, J., 1989. PHYLIP - Phylogeny Inference Package (Version 3.2). *Cladistics*, Vol. 5
1006 (1989), pp. 164-166 5, 164–166.

1007 Fischer, A.J., Bongini, R., Bastaki, N., Sherwood, P., 2011. The maturation of photoreceptors in
1008 the avian retina is stimulated by thyroid hormone. *Neuroscience* 178, 250–260.
1009 doi:10.1016/j.neuroscience.2011.01.022

1010 Gibbs, H.C., Chang-Gonzalez, A., Hwang, W., Yeh, A.T., Lekven, A.C., 2017. Midbrain-
1011 Hindbrain Boundary Morphogenesis: At the Intersection of Wnt and Fgf Signaling. *Front.*
1012 *Neuroanat.* 11, 64. doi:10.3389/fnana.2017.00064

1013 Guo, Q., Li, J.Y.H., 2019. Defining developmental diversification of diencephalon neurons
1014 through single cell gene expression profiling. *Development* 146, dev174284.
1015 doi:10.1242/dev.174284

1016 Haghverdi, L., Büttner, F., Theis, F.J., 2015. Diffusion maps for high-dimensional single-cell
1017 analysis of differentiation data. *Bioinformatics* 31, 2989–2998.
1018 doi:10.1093/bioinformatics/btv325

1019 Haghverdi, L., Büttner, M., Wolf, F.A., Büttner, F., Theis, F.J., 2016. Diffusion pseudotime
1020 robustly reconstructs lineage branching. *Nat. Methods* 13, 845–848.
1021 doi:10.1038/nmeth.3971

1022 He, J., Zhang, G., Almeida, A.D., Cayouette, M., Simons, B.D., Harris, W.A., 2012. How variable
1023 clones build an invariant retina. *Neuron* 75, 786–798. doi:10.1016/j.neuron.2012.06.033

1024 Holguera, I., Desplan, C., 2018. Neuronal specification in space and time. *Science* 362, 176–
1025 180. doi:10.1126/science.aas9435

1026 Hu, Y., Wang, X., Hu, B., Mao, Y., Chen, Y., Yan, L., Yong, J., Dong, J., Wei, Y., Wang, W.,
1027 Wen, L., Qiao, J., Tang, F., 2019. Dissecting the transcriptome landscape of the human
1028 fetal neural retina and retinal pigment epithelium by single-cell RNA-seq analysis. *PLoS*
1029 *Biol.* 17, e3000365. doi:10.1371/journal.pbio.3000365

1030 Jessell, T.M., 2000. Neuronal specification in the spinal cord: inductive signals and
1031 transcriptional codes. *Nat. Rev. Genet.* 1, 20–29. doi:10.1038/35049541

1032 Kim, D.W., Washington, P.W., Wang, Z.Q., Lin, S., Sun, C., Jiang, L., Blackshaw, S., 2019.
1033 Single cell RNA-Seq analysis identifies molecular mechanisms controlling hypothalamic
1034 patterning and differentiation. *bioRxiv* 83, 657148. doi:10.1101/657148

1035 Korzh, V., Sleptsova, I., Liao, J., He, J., Gong, Z., 1998. Expression of zebrafish bHLH genes
1036 *ngn1* and *nrd* defines distinct stages of neural differentiation. *Dev. Dyn.* 213, 92–104.
1037 doi:10.1002/(SICI)1097-0177(199809)213:1<92::AID-AJA9>3.0.CO;2-T

1038 Kretzschmar, K., Watt, F.M., 2012. Lineage tracing. *Cell* 148, 33–45.
1039 doi:10.1016/j.cell.2012.01.002

1040 La Manno, G., Siletti, K., Furlan, A., Gyllborg, D., Vinsland, E., Langseth, C.M., Khven, I.,
1041 Johnsson, A., Nilsson, M., Lönnerberg, P., Linnarsson, S., 2020. Molecular architecture of
1042 the developing mouse brain. *bioRxiv* 135C, 2020.07.02.184051.
1043 doi:10.1101/2020.07.02.184051

1044 Lange, C., Rost, F., Machate, A., Reinhardt, S., Lesche, M., Weber, A., Kuscha, V., Dahl, A.,
1045 Rulands, S., Brand, M., 2020. Single cell sequencing of radial glia progeny reveals the
1046 diversity of newborn neurons in the adult zebrafish brain. *Development* 147, dev185595.
1047 doi:10.1242/dev.185595

1048 Lee, S.-K., Pfaff, S.L., 2001. Transcriptional networks regulating neuronal identity in the
1049 developing spinal cord. *Nat. Neurosci.* 4, 1183–1191. doi:10.1038/nn750

1050 Li, H., Horns, F., Wu, B., Xie, Q., Li, J., Li, T., Luginbuhl, D.J., Quake, S.R., Luo, L., 2017.
1051 Classifying *Drosophila* Olfactory Projection Neuron Subtypes by Single-Cell RNA
1052 Sequencing. *Cell* 171, 1206–1220.e22. doi:10.1016/j.cell.2017.10.019

1053 Lord, N.D., Carte, A.N., Abitua, P.B., Schier, A.F., 2019. The pattern of Nodal morphogen
1054 signaling is shaped by co-receptor expression. *bioRxiv* 125, 2019.12.30.891101.
1055 doi:10.1101/2019.12.30.891101

1056 Ma, J., Shen, Z., Yu, Y.-C., Shi, S.-H., 2017. Neural lineage tracing in the mammalian brain.
1057 *Curr. Opin. Neurobiol.* 50, 7–16. doi:10.1016/j.conb.2017.10.013

1058 Mayer, C., Hafemeister, C., Bandler, R.C., Machold, R., Batista Brito, R., Jaglin, X., Allaway, K.,
1059 Butler, A., Fishell, G., Satija, R., 2018. Developmental diversification of cortical inhibitory
1060 interneurons. *Nature* 555, 457–462. doi:10.1038/nature25999

1061 McKenna, A., Findlay, G.M., Gagnon, J.A., Horwitz, M.S., Schier, A.F., Shendure, J., 2016.
1062 Whole-organism lineage tracing by combinatorial and cumulative genome editing. *Science*
1063 353, aaf7907. doi:10.1126/science.aaf7907

1064 McKenzie, J.A.G., Fruttiger, M., Abraham, S., Lange, C.A.K., Stone, J., Gandhi, P., Wang, X.,
1065 Bainbridge, J., Moss, S.E., Greenwood, J., 2012. Apelin is required for non-neovascular
1066 remodeling in the retina. *Am. J. Pathol.* 180, 399–409. doi:10.1016/j.ajpath.2011.09.035

1067 Mission, J.P., Takahashi, T., Caviness, V.S., 1991. Ontogeny of radial and other astroglial cells
1068 in murine cerebral cortex. *Glia* 4, 138–148. doi:10.1002/glia.440040205

1069 Moens, C.B., Prince, V.E., 2002. Constructing the hindbrain: insights from the zebrafish. *Dev.*
1070 *Dyn.* 224, 1–17. doi:10.1002/dvdy.10086

1071 Mueller, T., Wullmann, M.F., 2003. Anatomy of neurogenesis in the early zebrafish brain. *Brain*
1072 *Res. Dev. Brain Res.* 140, 137–155. doi:10.1016/s0165-3806(02)00583-7

1073 Navajas Acedo, J., Voas, M.G., Alexander, R., Woolley, T., Unruh, J.R., Li, H., Moens, C.,
1074 Piotrowski, T., 2019. PCP and Wnt pathway components act in parallel during zebrafish
1075 mechanosensory hair cell orientation. *Nat Comms* 10, 3993–17. doi:10.1038/s41467-019-
1076 12005-y

1077 Nowakowski, T.J., Bhaduri, A., Pollen, A.A., Alvarado, B., Mostajo-Radji, M.A., Di Lullo, E.,
1078 Haeussler, M., Sandoval-Espinosa, C., Liu, S.J., Velmeshev, D., Ounadjela, J.R., Shuga, J.,
1079 Wang, X., Lim, D.A., West, J.A., Leyrat, A.A., Kent, W.J., Kriegstein, A.R., 2017.
1080 Spatiotemporal gene expression trajectories reveal developmental hierarchies of the human
1081 cortex. *Science* 358, 1318–1323. doi:10.1126/science.aap8809

1082 Ogawa, Y., Shiraki, T., Kojima, D., Fukada, Y., 2015. Homeobox transcription factor Six7
1083 governs expression of green opsin genes in zebrafish. *Proceedings of the Royal Society B:*
1084 *Biological Sciences* 282, 20150659. doi:10.1098/rspb.2015.0659

1085 Pan, Y.A., Freundlich, T., Weissman, T.A., Schoppik, D., Wang, X.C., Zimmerman, S., Ciruna,
1086 B., Sanes, J.R., Lichtman, J.W., Schier, A.F., 2013. Zebrow: multispectral cell labeling for
1087 cell tracing and lineage analysis in zebrafish. *Development* 140, 2835–2846.
1088 doi:10.1242/dev.094631

1089 Pandey, S., Shekhar, K., Regev, A., Schier, A.F., 2018. Comprehensive Identification and
1090 Spatial Mapping of Habenular Neuronal Types Using Single-Cell RNA-Seq. 28, 1052–
1091 1065.e7. doi:10.1016/j.cub.2018.02.040

1092 Raj, B., Gagnon, J.A., Schier, A.F., 2018a. Large-scale reconstruction of cell lineages using
1093 single-cell readout of transcriptomes and CRISPR-Cas9 barcodes by scGESTALT. *Nat*
1094 *Protoc* 13, 2685–2713. doi:10.1038/s41596-018-0058-x

1095 Raj, B., Wagner, D.E., McKenna, A., Pandey, S., Klein, A.M., Shendure, J., Gagnon, J.A.,
1096 Schier, A.F., 2018b. Simultaneous single-cell profiling of lineages and cell types in the
1097 vertebrate brain. *Nat Biotechnol* 36, 442–450. doi:10.1038/nbt.4103

1098 Rheaume, B.A., Jereen, A., Bolisetty, M., Sajid, M.S., Yang, Y., Renna, K., Sun, L., Robson, P.,
1099 Trakhtenberg, E.F., 2018. Single cell transcriptome profiling of retinal ganglion cells
1100 identifies cellular subtypes. *Nat Comms* 9, 1–17. doi:10.1038/s41467-018-05134-3

1101 Rosenberg, A.B., Roco, C.M., Muscat, R.A., Kuchina, A., Sample, P., Yao, Z., Graybuck, L.T.,
1102 Peeler, D.J., Mukherjee, S., Chen, W., Pun, S.H., Sellers, D.L., Tasic, B., Seelig, G., 2018.
1103 Single-cell profiling of the developing mouse brain and spinal cord with split-pool barcoding.
1104 *Science* 360, 176–182. doi:10.1126/science.aam8999

1105 Rulands, S., Iglesias-Gonzalez, A.B., Boije, H., 2018. Deterministic fate assignment of Müller
1106 glia cells in the zebrafish retina suggests a clonal backbone during development. *Eur. J.*
1107 *Neurosci.* 48, 3597–3605. doi:10.1111/ejn.14257

1108 Sagner, A., Briscoe, J., 2019. Establishing neuronal diversity in the spinal cord: a time and a
1109 place. *Development* 146, dev182154. doi:10.1242/dev.182154

1110 Satoh, S., Tang, K., Iida, A., Inoue, M., Kodama, T., Tsai, S.Y., Tsai, M.-J., Furuta, Y.,
1111 Watanabe, S., 2009. The spatial patterning of mouse cone opsin expression is regulated by

1112 bone morphogenetic protein signaling through downstream effector COUP-TF nuclear
1113 receptors. *J. Neurosci.* 29, 12401–12411. doi:10.1523/JNEUROSCI.0951-09.2009

1114 Schindelin, J., Arganda-Carreras, I., Frise, E., Kaynig, V., Longair, M., Pietzsch, T., Preibisch,
1115 S., Rueden, C., Saalfeld, S., Schmid, B., Tinevez, J.-Y., White, D.J., Hartenstein, V., Eliceiri,
1116 K., Tomancak, P., Cardona, A., 2012. Fiji: an open-source platform for biological-image
1117 analysis. *Nat. Methods* 9, 676–682. doi:10.1038/nmeth.2019

1118 Schmechel, D.E., Rakic, P., 1979. A Golgi study of radial glial cells in developing monkey
1119 telencephalon: morphogenesis and transformation into astrocytes. *Anat. Embryol.* 156,
1120 115–152. doi:10.1007/BF00300010

1121 Schmidt, R., Strähle, U., Scholpp, S., 2013. Neurogenesis in zebrafish - from embryo to adult.
1122 *Neural Development* 8, 3. doi:10.1186/1749-8104-8-3

1123 Shen, Y.-C., Raymond, P.A., 2004. Zebrafish cone-rod (crx) homeobox gene promotes
1124 retinogenesis. *Dev. Biol.* 269, 237–251. doi:10.1016/j.ydbio.2004.01.037

1125 Siebert, S., Farrell, J.A., Cazet, J.F., Abeykoon, Y., Primack, A.S., Schnitzler, C.E., Juliano,
1126 C.E., 2019. Stem cell differentiation trajectories in Hydra resolved at single-cell resolution.
1127 *Science* 365, eaav9314. doi:10.1126/science.aav9314

1128 Solek, C.M., Feng, S., Perin, S., Weinschutz Mendes, H., Ekker, M., 2017. Lineage tracing of
1129 *dlx1a/2a* and *dlx5a/6a* expressing cells in the developing zebrafish brain. *Dev. Biol.* 427,
1130 131–147. doi:10.1016/j.ydbio.2017.04.019

1131 Stigloher, C., Chapouton, P., Adolf, B., Bally-Cuif, L., 2008. Identification of neural progenitor
1132 pools by E(Spl) factors in the embryonic and adult brain. *Brain Res. Bull.* 75, 266–273.
1133 doi:10.1016/j.brainresbull.2007.10.032

1134 Tambalo, M., Mitter, R., Wilkinson, D.G., 2020. A single cell transcriptome atlas of the
1135 developing zebrafish hindbrain. *Development* 147, dev184143. doi:10.1242/dev.184143

1136 Tambalo, M., Mitter, R., Wilkinson, D.G., 2019. A single cell transcriptome atlas of the
1137 developing zebrafish hindbrain. *bioRxiv* 124, 745141. doi:10.1101/745141

1138 Than-Trong, E., Bally-Cuif, L., 2015. Radial glia and neural progenitors in the adult zebrafish
1139 central nervous system. *Glia* 63, 1406–1428. doi:10.1002/glia.22856

1140 Thisse, B., Heyer, V., Lux, A., Alunni, V., Degraeve, A., Seilliez, I., Kirchner, J., Parkhill, J.-P.,
1141 Thisse, C., 2004. Spatial and temporal expression of the zebrafish genome by large-scale in
1142 situ hybridization screening. *Methods Cell Biol.* 77, 505–519. doi:10.1016/s0091-
1143 679x(04)77027-2

1144 Tiklová, K., Björklund, Å.K., Lahti, L., Fiorenzano, A., Nolbrant, S., Gillberg, L., Volakakis, N.,
1145 Yokota, C., Hilscher, M.M., Hauling, T., Holmström, F., Joodmardi, E., Nilsson, M., Parmar,
1146 M., Perlmann, T., 2019. Single-cell RNA sequencing reveals midbrain dopamine neuron
1147 diversity emerging during mouse brain development. *Nat Comms* 10, 581–12.
1148 doi:10.1038/s41467-019-08453-1

1149 Trapnell, C., Cacchiarelli, D., Grimsby, J., Pokharel, P., Li, S., Morse, M., Lennon, N.J., Livak,
1150 K.J., Mikkelsen, T.S., Rinn, J.L., 2014. The dynamics and regulators of cell fate decisions
1151 are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* 32, 381–386.
1152 doi:10.1038/nbt.2859

1153 Viczian, A.S., Vignali, R., Zuber, M.E., Barsacchi, G., Harris, W.A., 2003. *XOtx5b* and *XOtx2*
1154 regulate photoreceptor and bipolar fates in the *Xenopus* retina. *Development* 130, 1281–
1155 1294. doi:10.1242/dev.00343

1156 Wagner, D.E., Weinreb, C., Collins, Z.M., Briggs, J.A., Megason, S.G., Klein, A.M., 2018.
1157 Single-cell mapping of gene expression landscapes and lineage in the zebrafish embryo.
1158 *Science* 360, 981–987. doi:10.1126/science.aar4362

1159 Wamsley, B., Fishell, G., 2017. Genetic and activity-dependent mechanisms underlying
1160 interneuron diversity. *Nat. Rev. Neurosci.* 18, 299–309. doi:10.1038/nrn.2017.30

1161 Wilson, S.W., Brand, M., Eisen, J.S., 2002. Patterning the zebrafish central nervous system.
1162 *Results Probl Cell Differ* 40, 181–215.

1163 Wilson, S.W., Rubenstein, J.L., 2000. Induction and dorsoventral patterning of the
 1164 telencephalon. *Neuron* 28, 641–651. doi:10.1016/s0896-6273(00)00171-9
 1165 Woo, K., Fraser, S.E., 1995. Order and coherence in the fate map of the zebrafish nervous
 1166 system. *Development* 121, 2595–2609.
 1167 Woodworth, M.B., Girsakis, K.M., Walsh, C.A., 2017. Building a lineage from single cells: genetic
 1168 techniques for cell lineage tracking. *Nat. Rev. Genet.* 18, 230–244.
 1169 doi:10.1038/nrg.2016.159
 1170 Xie, Y., Dorsky, R.I., 2017. Development of the hypothalamus: conservation, modification and
 1171 innovation. *Development* 144, 1588–1599. doi:10.1242/dev.139055
 1172 Xu, B., Tang, X., Jin, M., Zhang, H., Du, L., Yu, S., He, J., 2020. Unifying developmental
 1173 programs for embryonic and postembryonic neurogenesis in the zebrafish retina.
 1174 *Development* 147, dev185660. doi:10.1242/dev.185660
 1175 Yoshinari, N., Ando, K., Kudo, A., Kinoshita, M., Kawakami, A., 2012. Colored medaka and
 1176 zebrafish: Transgenics with ubiquitous and strong transgene expression driven by the
 1177 medaka β -actin promoter. *Develop. Growth Differ.* 54, 818–828. doi:10.1073/pnas.94.8.3789
 1178 Zeisel, A., Hochgerner, H., Lönnerberg, P., Johnsson, A., Memic, F., van der Zwan, J., Häring,
 1179 M., Braun, E., Borm, L.E., La Manno, G., Codeluppi, S., Furlan, A., Lee, K., Skene, N.,
 1180 Harris, K.D., Hjerling-Leffler, J., Arenas, E., Ernfors, P., Marklund, U., Linnarsson, S., 2018.
 1181 Molecular Architecture of the Mouse Nervous System. *Cell* 174, 999–1014.e22.
 1182 doi:10.1016/j.cell.2018.06.021
 1183 Zhong, S., Zhang, S., Fan, X., Wu, Q., Yan, L., Dong, J., Zhang, H., Li, L., Le Sun, Pan, N., Xu,
 1184 X., Tang, F., Zhang, J., Qiao, J., Wang, X., 2018. A single-cell RNA-seq survey of the
 1185 developmental landscape of the human prefrontal cortex. *Nature* 555, 524–528.
 1186 doi:10.1038/nature25980
 1187
 1188

1189 **Figure 1. Developmental compendium of zebrafish head and brain cell types**

1190 **A.** Schematic of the developmental stages profiled. Red hatched line represents head regions that were
 1191 selected for enrichment of brain cells in early development. Samples from 5 to 15 dpf were dissected to obtain
 1192 brain and eye specifically. h, hours post fertilization; d, days post fertilization
 1193 **B.** Schematic of scRNA-seq using 10X Genomics platform.
 1194 **C.** Cell type heterogeneity within each stage. Clusters at each stage were assigned to a region or tissue type
 1195 based on known markers and color coded to reflect their classification. tSNE implementations: Barnes-Hut (12h
 1196 to 3d), Fourier transform (5d and 15d).
 1197 **D.** In situ hybridization for novel markers in the trigeminal placode at 12 hpf. *kif17* is expressed on the anterior
 1198 polster and ventral mesoderm, delineating the border of the embryo. Trigeminal ganglia markers *ptgs2a*, *tp63*
 1199 and *sdpra* (*cavin2a*) are expressed bilaterally (asterisks) posterior to the eye. Eyes are delineated by dotted
 1200 lines. A: Anterior; P: Posterior. Scale bar equals 100 μ m.
 1201 **E.** In situ hybridization validation of novel marker *sox1a* in the hypothalamus at 2 dpf. Top panels, lateral view
 1202 of brain; Bottom panels, ventral view of brain. *dlx1a* and *pdyn* are known hypothalamus. Eyes are delineated by
 1203 dotted lines. VHyp: Ventral Hypothalamus; TVZ: Telencephalic Ventricular Zone; ADi: Anterior Diencephalon;
 1204 AFb: Anterior Forebrain; VDi: Ventral Diencephalon; Le: Lens. Scale bar equals 200 μ m.
 1205 **F.** smFISH validation of novel marker *ompa* in horizontal cells of the retina at 5 dpf. Left panel, retina section
 1206 stained with DAPI (grey), pan-retinal *foxg1b* (cyan) and *ompa* (yellow). Strong yellow signal in photoreceptors
 1207 represent autofluorescence. White box indicates area that was zoomed in for the right panels. Dotted lines
 1208 indicate the horizontal cell layer. PR, photoreceptor cells; HC, horizontal cells; BC, bipolar cells; AC, amacrine
 1209 cells; RGC, retinal ganglion cells
 1210
 1211

1211 **Figure 2. Brain cell type diversification from 12 hpf to 15 dpf**

1212 **A.** tSNE plot of 12 hpf dataset. Only clusters corresponding to neural and blood cell types are shown. Inferred
1213 identities of each cluster are described.
1214 **B.** Dot plot of gene expression pattern of select marker genes (columns) for each cluster (row). Dot size
1215 indicates the percentage of cells expressing the marker; color represents the average scaled expression level.
1216 **C.** tSNE plot of 15 dpf dataset. Inferred identities of each cluster are described.
1217 **D.** Dot plot of gene expression patterns of select marker genes for each cluster. Layout is same as **(B)**. Grey
1218 box represents generic neuronal and progenitor genes.
1219 tSNE implementations: Barnes-Hut **(A)**, Fourier transform **(C)**

1220

1221 **Figure 3. Neuron subtype diversity at 15 dpf**

1222 **A-C.** Violin plots of select marker gene expression in identified brain neuron subtypes from 15 dpf. Retina
1223 neurons and nascent neurons are omitted from the analysis. Cluster numbers are indicated at the bottom along
1224 with their inferred spatial location in the brain. Cluster 76 has unknown spatial location. Detailed cluster
1225 descriptions are in Supplementary Table 1 and can be explored interactively in the accompanying app at
1226 https://github.com/brlauuu/zf_brain.

1227 **A.** Expression of transcription factors.

1228 **B.** Expression of neuropeptides and their receptors.

1229 **C.** Expression of genes involved in neuron electrophysiology including neurotransmitters, transporters,
1230 receptors, and channels.

1231 **D.** Matrix showing whether neuron subtypes identified at 15 dpf are also detected in earlier larval (5 and 8 dpf)
1232 and later juvenile (25 dpf (Raj et al., 2018b)) stages. Clusters were matched across stages by comparing
1233 marker gene expression. The cluster number at 15 dpf is shown and an orange circle indicates that the subtype
1234 is detected in another stage.

1235

1236 **Figure 4. Developmental diversification of neurons and progenitors**

1237 **A.** Area plot of the percentage of dataset at each timepoint corresponding to neural progenitors, neurons, and
1238 other cell types. Right panels, Total number of clusters of progenitors and neurons at each stage of brain
1239 development.

1240 **B.** tSNE plot of embryonic, intermediate and larval neural progenitors. All progenitor cells were analyzed
1241 together after subsetting from the whole dataset.

1242 **C-D.** Heatmaps of select gene expression in early embryonic **(C)** and late larval **(D)** brain neural progenitors.

1243 Top panel, genes enriched in embryonic progenitors. Bottom panel, genes enriched in larval progenitors.

1244 Embryonic progenitors have a strong spatial signature (forebrain, midbrain, hindbrain) and are depleted in
1245 genes that distinguish larval progenitor subtypes **(C)**. Larval progenitors segregate into non-proliferative and
1246 proliferative groups that can be resolved into additional subtypes characterized by expression of various gene
1247 combinations **(D)**. TF, transcription factor. *pax6a is expressed in multiple regions

1248 **E.** Heatmap of Pearson correlation values of 79 spatial markers in embryonic, intermediate and larval neural
1249 progenitors. Spatial markers were selected based on existing literature. Groups of co-varying genes in the
1250 midbrain and forebrain are highlighted with dashed boxes.

1251 **F.** Plot showing number of highly variable genes that co-vary with any of the selected 79 spatial markers in
1252 embryonic and larval progenitors. Co-variation was determined by Pearson correlation, with several thresholds
1253 (from stringent to relaxed) displayed along the x-axis.

1254

1255 **Figure 5. Optimization of scGESTALT lineage recorder for better barcode recovery**

1256 **A.** Schematic overview of CRISPR-Cas9 lineage recording. Optimized scGESTALT comprises a barcode
1257 cassette in the 3' end of DsRed transgene (single copy) and the medaka beta-actin promoter. Embryos are
1258 injected with Cas9 protein and DsRed sgRNAs and animals are profiled at 15 dpf by scRNA-seq.

1259 **B.** Pairwise comparisons using cosine dissimilarity of barcode edit patterns from four (ZF1-4) edited 15 dpf
1260 larval brains.

- 1261 **C.** Chord diagram of the nature and frequency of deletions within and between target sites. Each colored sector
 1262 represents a target site. Links between target sites represent inter-site deletions; self-links represent intra-site
 1263 deletions. Link widths are proportional to the edit frequencies.
- 1264 **D.** Type of edit at each target site within the barcode from edited ZF1-4 larval brains.
- 1265 **E.** Heat map of lineage relationships between non-retinal and retinal cell types in the eye. All clusters with >3
 1266 cells and all barcodes with >1 cell were used to determine if there is enrichment of cell type-specific barcodes
 1267 across each cluster pair. Blue indicates significant enrichment and lineage segregation. Purple indicates no
 1268 significant enrichment and no lineage segregation. Grey indicates insufficient sampling power and undefined
 1269 lineage status. Cluster numbers are indicated (e.g. C45) and either cell type gene markers (e.g. *cldna*⁺) or the
 1270 exact name of the cell type (e.g. cone bipolar cells) are indicated along the rows. Along the columns, the
 1271 numbers within the brackets indicate the number of barcodes and number of cells, respectively, for that cluster.
- 1272 **F.** Heat map of lineage relationships between brain regions and the retina. Neuron clusters that could be
 1273 pseudospatially assigned to the each region were used (see Supplementary Table). Analysis, layout and color
 1274 code are same as in E.
- 1275 **G.** Heat map of lineage relationships between neuronal cell types in the forebrain and midbrain. Analysis,
 1276 layout and color code are same as in E. The brain region each cluster belongs to is indicated (e.g. pallium,
 1277 hypothalamus), and for clusters where a more precise location could be inferred a gene marker is indicated
 1278 (e.g. *pitx2*⁺).
- 1279 **H.** Heat map of lineage relationships between brain progenitor clusters. Analysis, layout and color code is same
 1280 as in E. Cell type marker genes are indicated along with the cluster number. URL, upper rhombic lip
- 1281 **I.** Bar plot of the proportion (based on Jaccard Index) of granule cell (cerebellum neurons) barcodes that are
 1282 shared with each brain progenitor cluster. Cluster numbers are the same as in H.

1283

1284 **Figure 6. Cell specification trajectories in the retina and hypothalamus**

- 1285 **A.** UMAP visualization of retinal cell types. Retinal cells (based on clustering analysis) from 12 hpf to 15 dpf
 1286 were subsetted from the full dataset and analyzed together. Cells are color coded by stage.
- 1287 **B.** Cell specification tree of zebrafish retinal development. Trajectories were generated by URD and visualized
 1288 as a branching tree. Cells are color coded by stage. 12 hpf cells were assigned as the root and 15 dpf
 1289 differentiated cells were assigned as tips. CBP, cone bipolar cells (6 subtypes are numbered); RGC, retinal
 1290 ganglion cells; RPE, retinal pigment epithelium
- 1291 **C.** Expression of select genes are shown on the retina specification tree.
- 1292 **D.** Heat maps of gene expression cascades of photoreceptor cell trajectories and retinal ganglion cell
 1293 trajectories. Cells were selected based on high expression along trajectories leading to these cell types,
 1294 compared to expression along opposing branchpoints. Red, high expression. Yellow, low expression
- 1295 **E.** UMAP visualization of hypothalamus cell types. Hypothalamus cells (based on clustering analysis) from 12
 1296 hpf to 15 dpf were subsetted from the full dataset and analyzed together. Cells are color coded by stage.
- 1297 **F.** Cell specification tree of zebrafish hypothalamus development. Trajectories were generated by URD and
 1298 visualized as a branching tree. Cells are color coded by stage. 12 hpf cells were assigned as the root and 15
 1299 dpf differentiated cells were assigned as tips.
- 1300 **G.** Expression of select genes are shown on the hypothalamus specification tree.
- 1301 **H.** Heat map of gene expression cascade of *nrgna*⁺ cell trajectories. Red, high expression. Yellow, low
 1302 expression

1303

1304 **Figure 7. Progenitor differences between retina and hypothalamus**

- 1305 Retinal and hypothalamus cells were divided into progenitor (purple), precursor (orange), and differentiated
 1306 (blue) cells, as shown on the URD tree. The fraction of cells in each of these transcriptional states was then
 1307 determined for three developmental periods (12–24 hpf, 36 hpf – 3 dpf, and 5–15 dpf). In the retina, cells can
 1308 be found in a progenitor state (light purple) that persists post-embryonically.

Figure 1
FIGURE 1

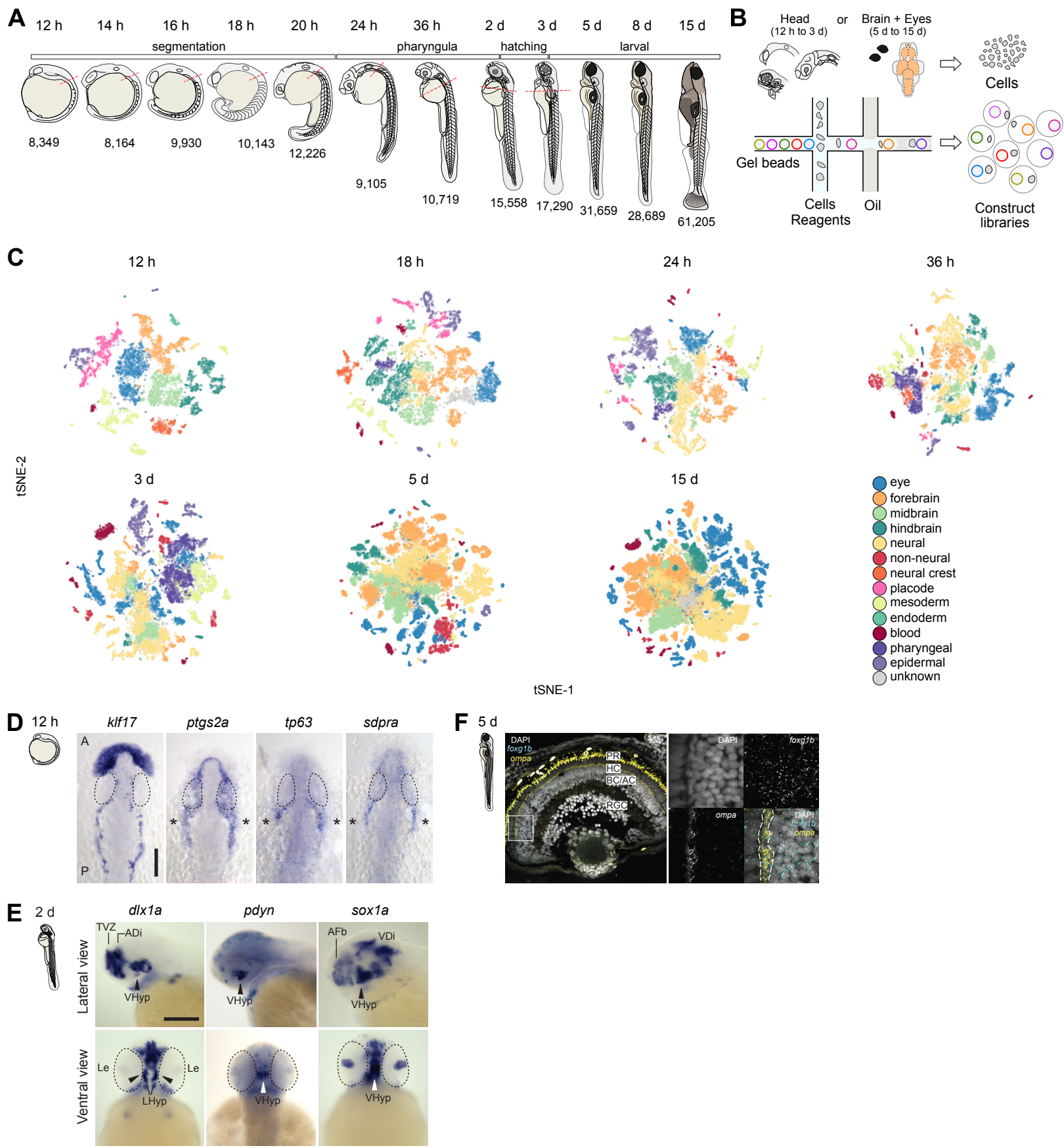
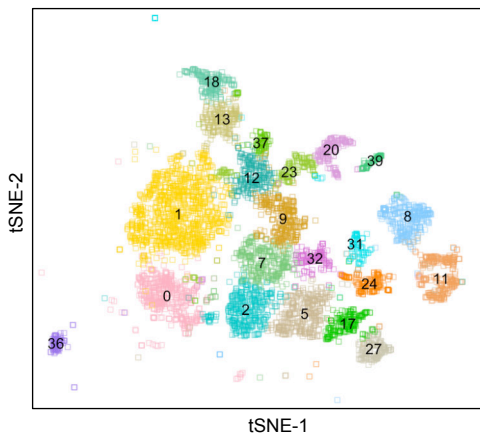


Figure 2

A

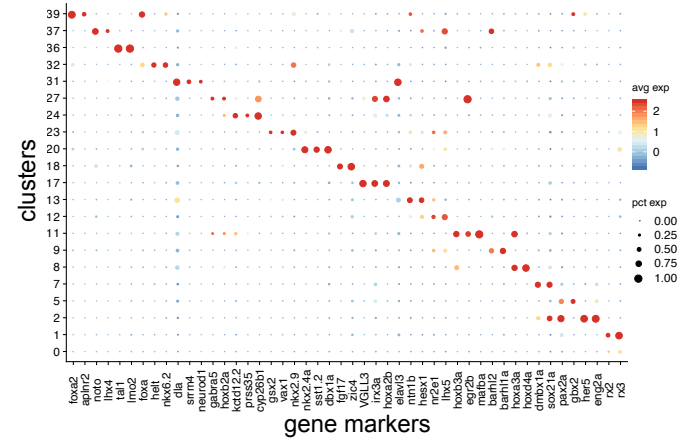


12 h



- 0 - optic vesicle
- 1 - optic vesicle
- 2 - midbrain
- 5 - mid/hindbrain boundary
- 7 - midbrain
- 8 - hindbrain r7
- 9 - diencephalon
- 11 - hindbrain r5/6
- 12 - dorsal diencephalon
- 13 - telencephalon
- 17 - hindbrain r1/2
- 18 - telencephalon
- 20 - ventral diencephalon
- 23 - ventral diencephalon, optic stalk
- 24 - hindbrain r4
- 27 - hindbrain r3
- 31 - neurons
- 32 - ventral midbrain
- 36 - rostral blood island
- 37 - diencephalon
- 39 - floor plate

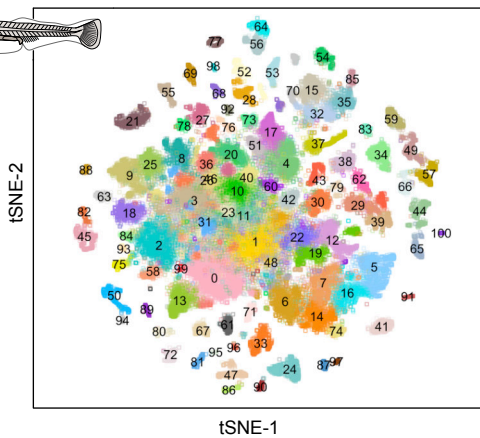
B



C

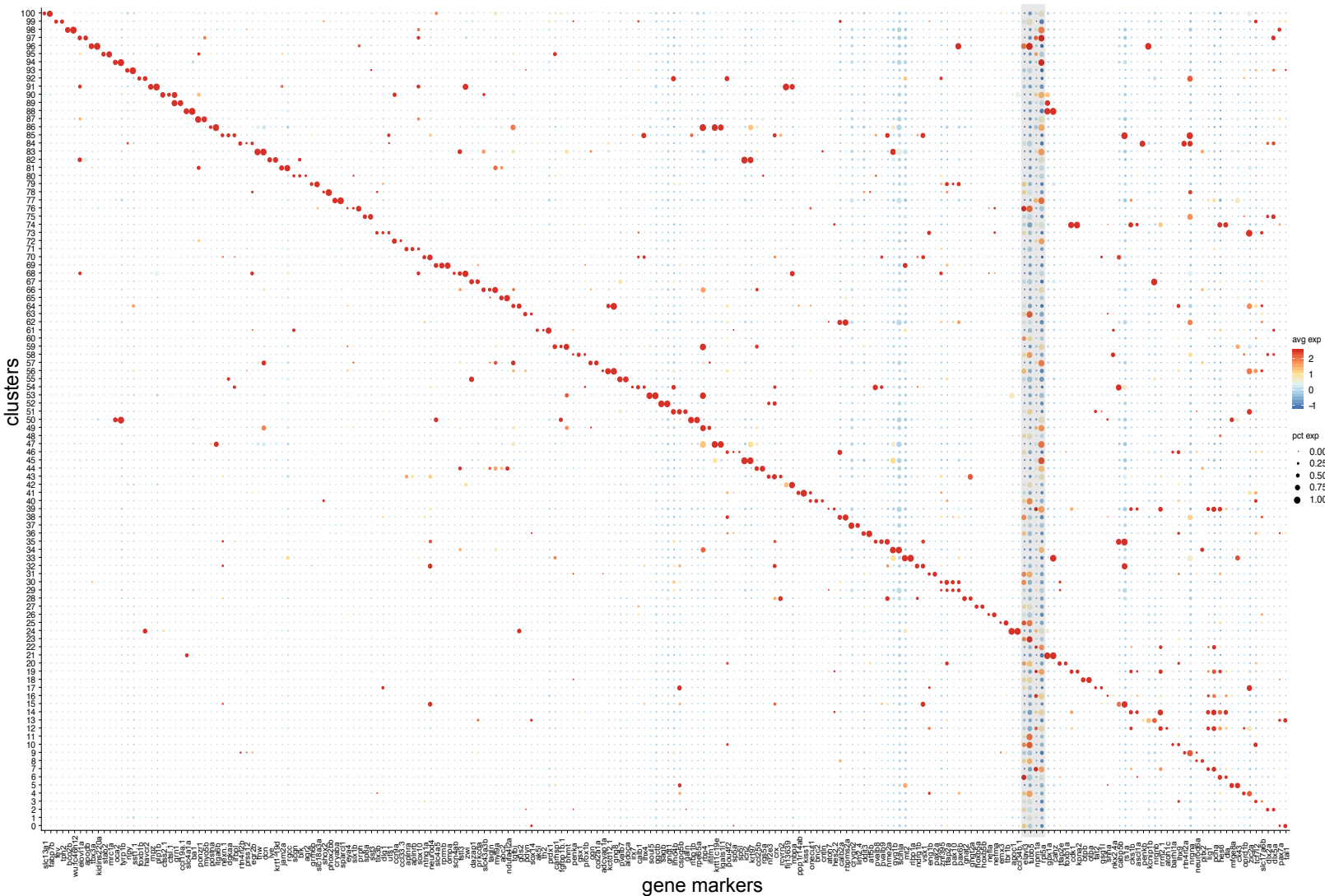


15 d



- 0,13 - optic tectum
- 2,58 - ventral forebrain
- 4,17,73 - granule cells
- 5 - radial glia
- 6,7,14,48 - progenitors
- 8,40,46,67 - diencephalon
- 9,25 - pallium
- 10,20,92 - thalamus
- 12 - URL progenitors
- 13,84 - telencephalon
- 15,32,35,54,70,85 - cone bipolar cells
- 16 - ventral progenitors
- 18 - tuberculum
- 19,74 - cycling progenitors
- 21,88 - erythrocytes
- 24 - microglia
- 27,31 - hindbrain/cranial nerve
- 28 - cones
- 29 - retina progenitors
- 30,79 - amacrine cells
- 26,36 - midbrain
- 33 - muller glia
- 34 - cornea
- 37 - lens
- 38,62 - retinal ganglion cells
- 39 - retina progenitors
- 41 - ventral habenula
- 42,68,91 - oligodendrocytes
- 43 - photoreceptor precursors
- 44,49,66 - cornea epithelium
- 45,82,83 - eye cells
- 47,86 - eye epidermis
- 50 - RPE
- 51 - torus longitudinalis
- 52 - rods
- 53 - glia
- 55 - purkinje cells
- 56,64 - dorsal habenula
- 57 - eye cartilage
- 59 - eye glia
- 61,63,75,93,99 - hypothalamus
- 65 - cornea epithelial progenitors
- 69 - horizontal cells
- 71 - OPC
- 72 - neutrophils
- 73 - lens epithelium
- 77 - lens epithelium
- 78 - hindbrain
- 80 - metaphocytes
- 81 - myeloid cells
- 87,94 - pigment cells
- 89,90 - blood/immune cells
- 95 - perivascular FGP cells
- 96 - retina neurons
- 97,98 - neural crest pigment cells
- 100 - glia progenitors

D



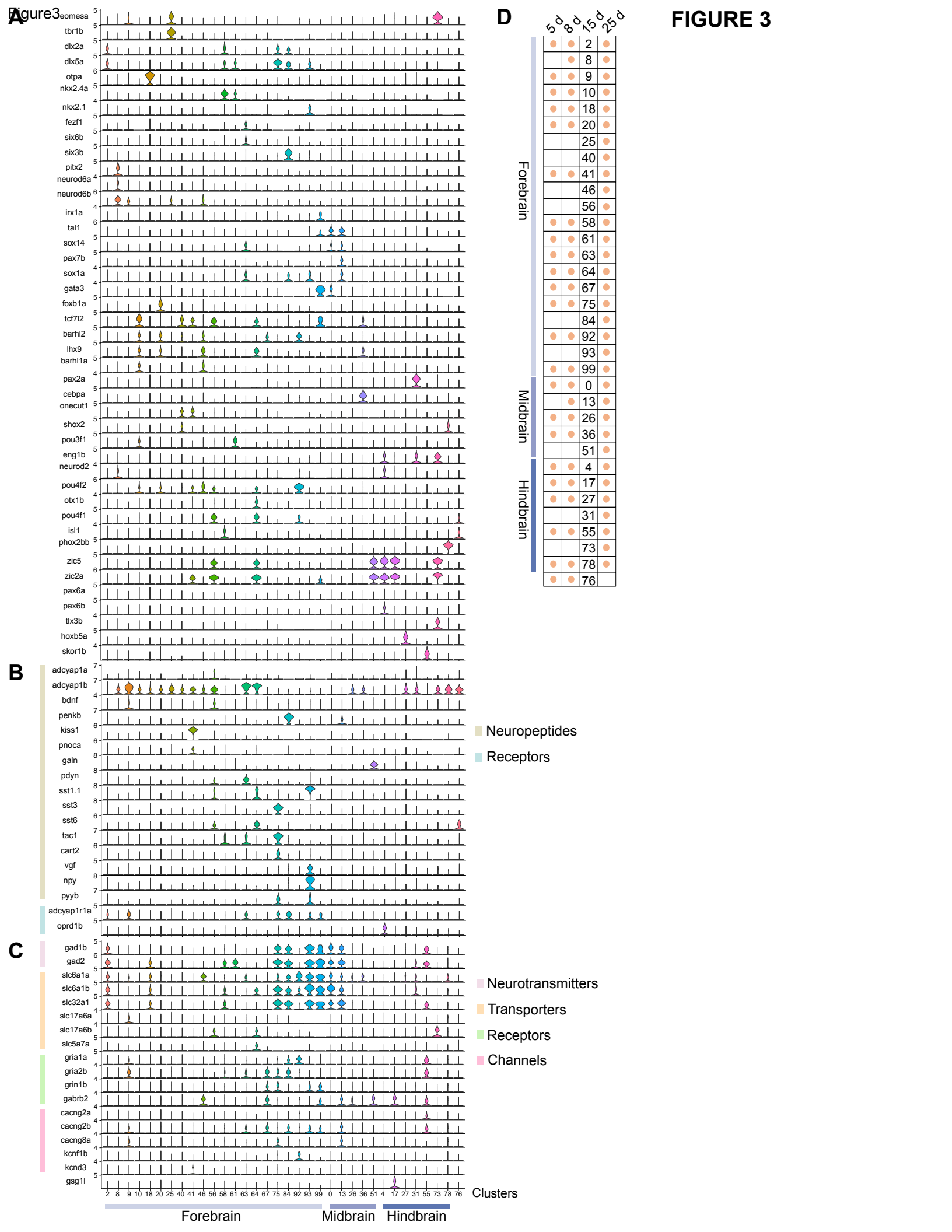
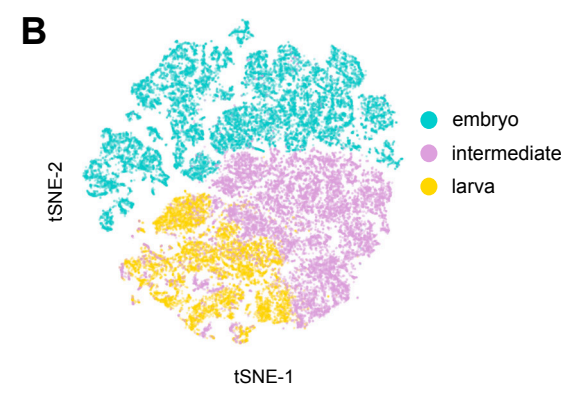
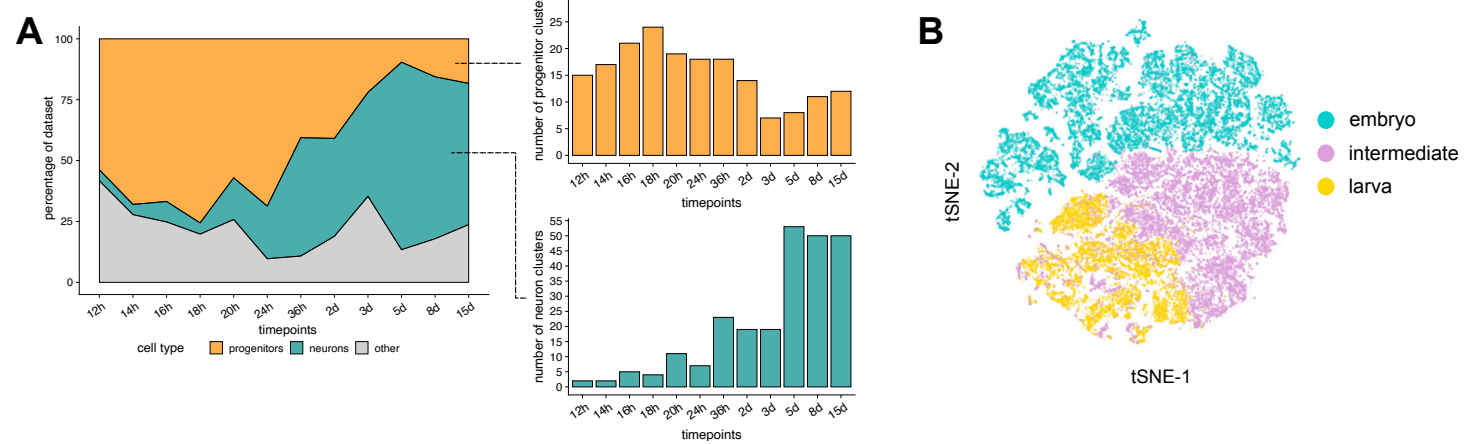
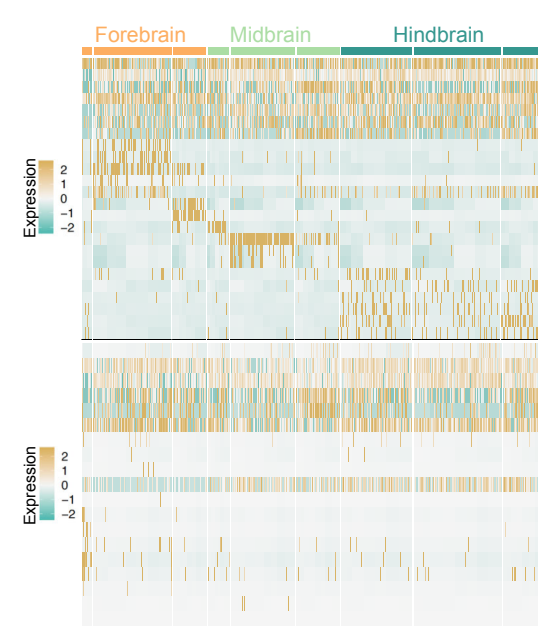


Figure 4

FIGURE 4



C Early embryonic brain neural progenitors



D Late larval brain neural progenitors

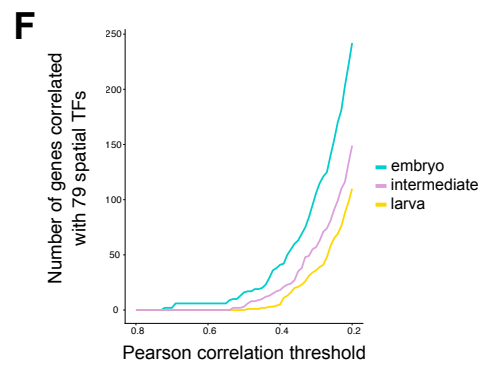
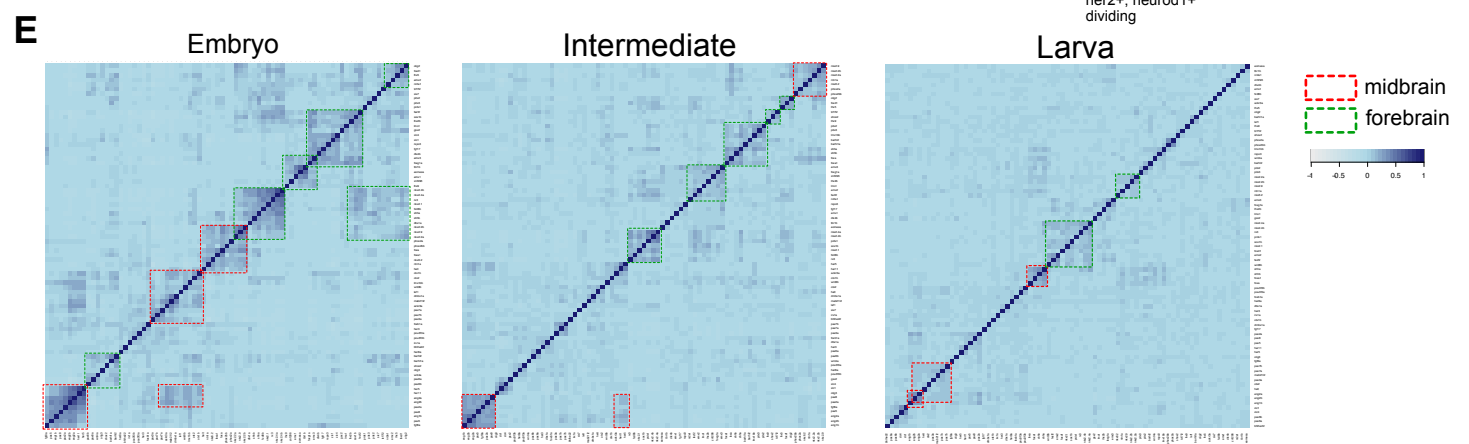
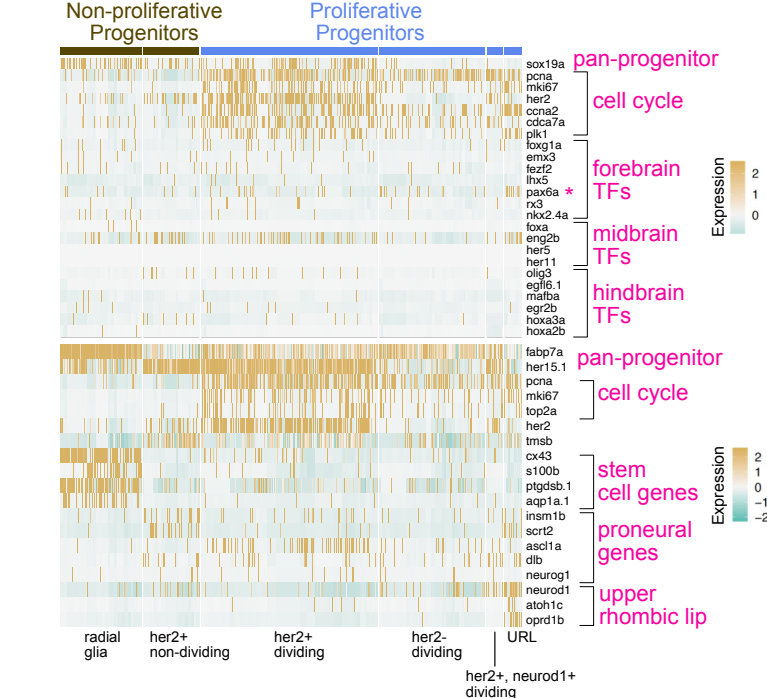


FIGURE 5

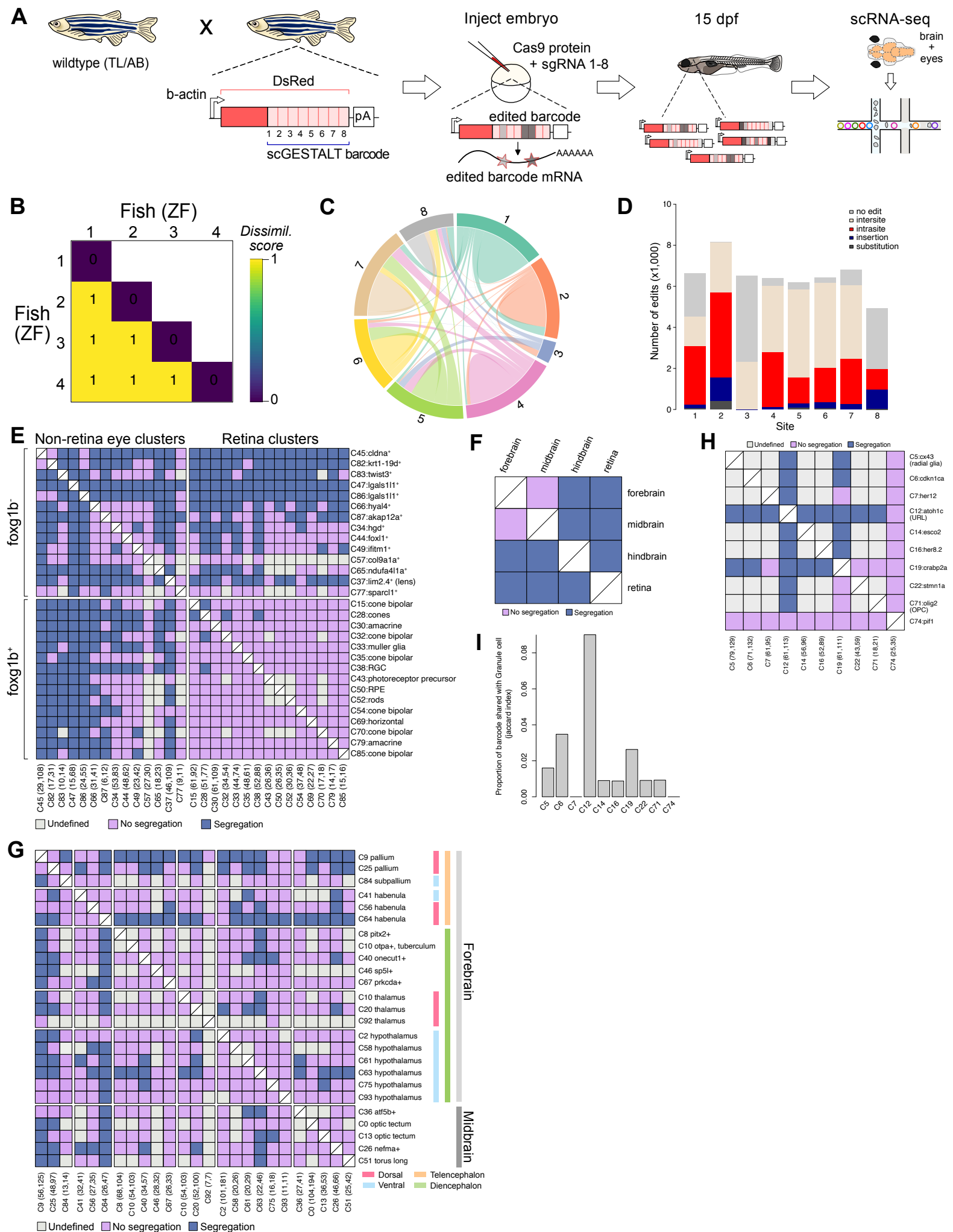


Figure6
FIGURE 6

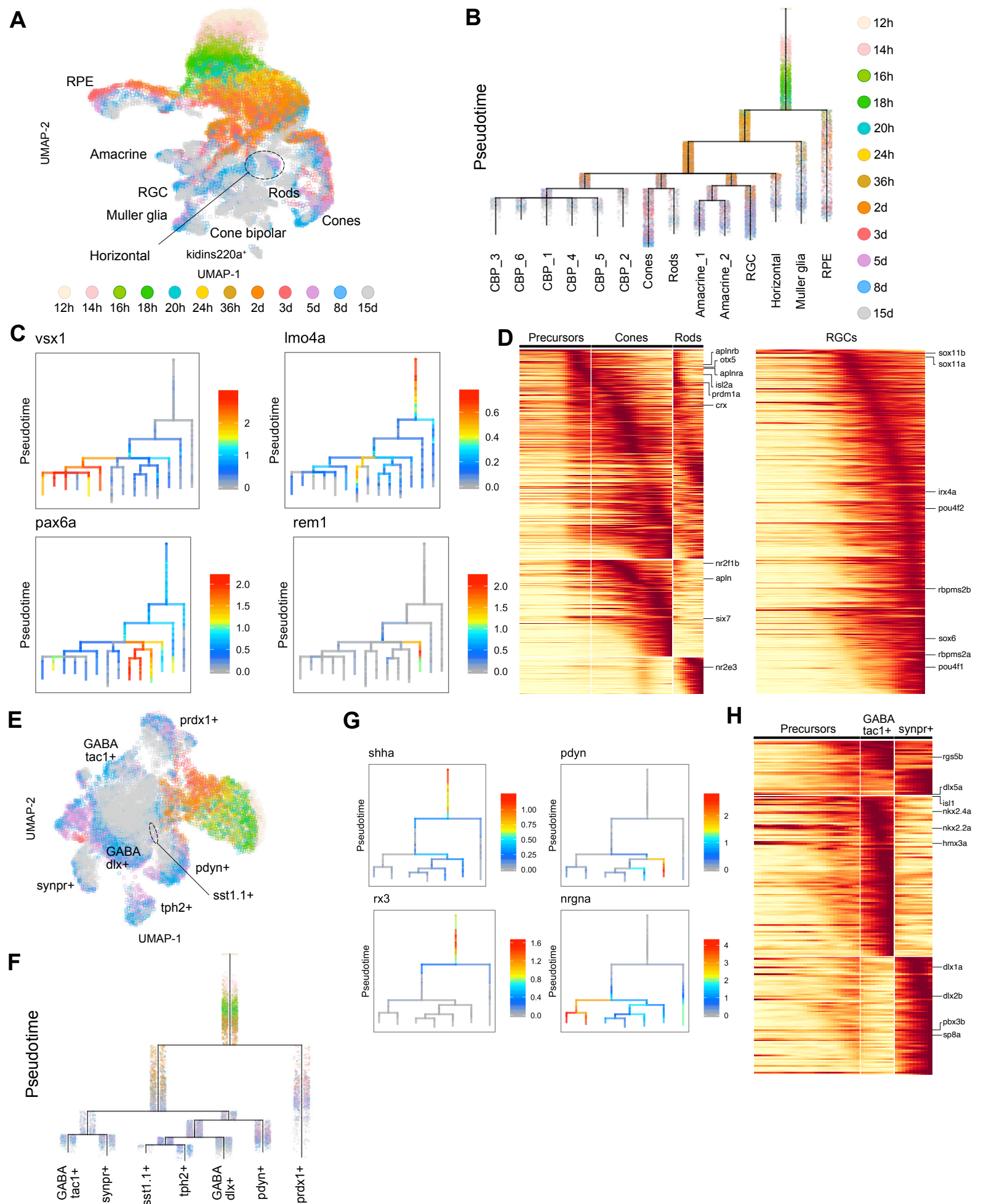
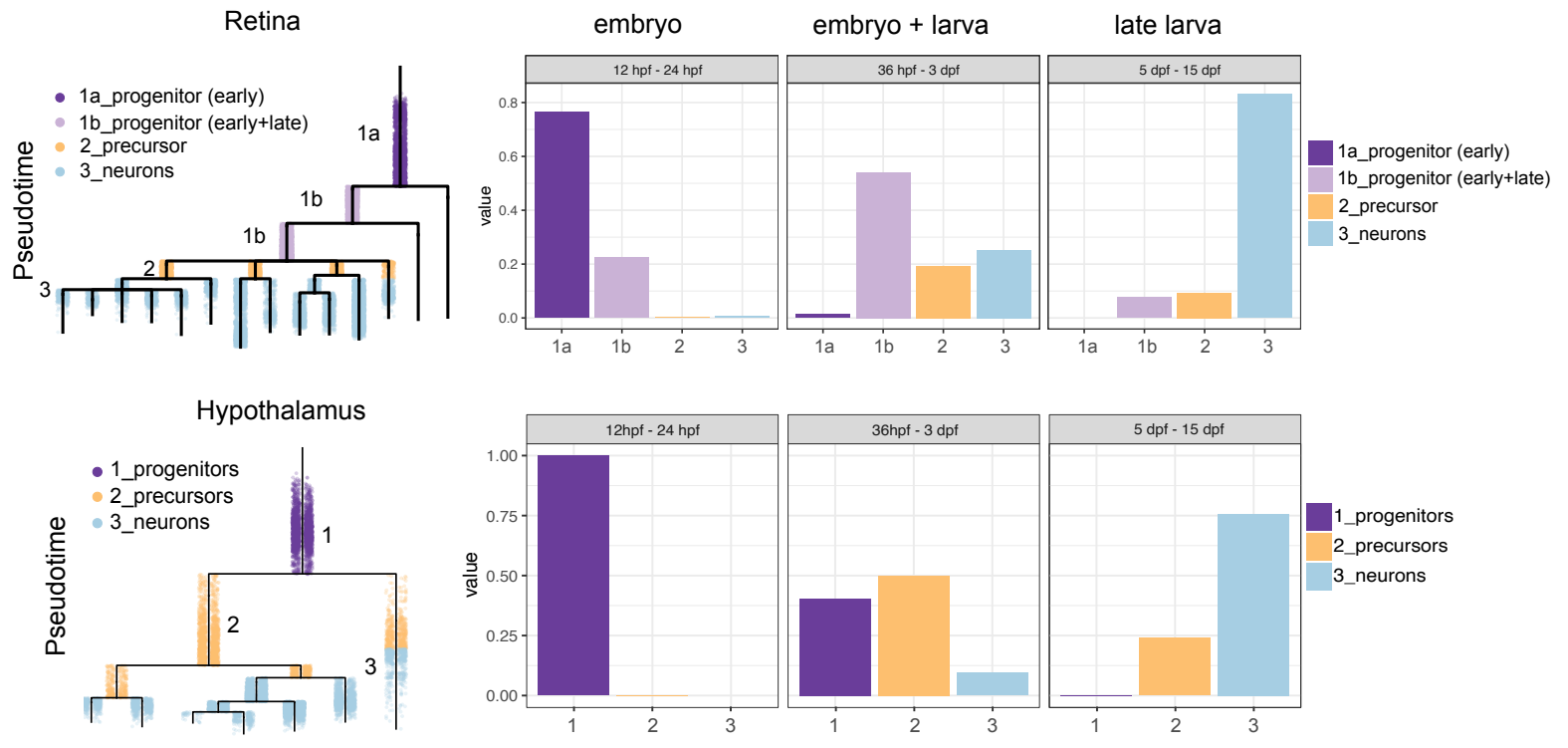


Figure 7
FIGURE 7



Supplemental Information

Emergence of neuronal diversity during vertebrate brain development

Bushra Raj^{1,2*}, Jeffrey A. Farrell^{1,3}, Jialin Liu^{2,4}, Jakob El Kholtei^{2,4}, Adam Carte^{1,4,5}, Joaquin Navajas Acedo^{2,4}, Lucia Y Du^{2,4}, Aaron McKenna⁶, Đorđe Relić^{3,7}, Jessica M. Leslie¹, and Alexander F. Schier^{1,2,4,8,9,10*}

¹ Department of Molecular and Cellular Biology, Harvard University, Cambridge, Massachusetts, USA

² Allen Discovery Center for Cell Lineage Tracing, Seattle, Washington, USA

³ Unit on Cell Specification and Differentiation, National Institute of Child Health and Human Development, NIH, Bethesda, Maryland, USA

⁴ Biozentrum, University of Basel, Switzerland

⁵ Systems, Synthetic, and Quantitative Biology Program, Harvard University, Cambridge, Massachusetts, USA

⁶ Department of Molecular and Systems Biology, Dartmouth Geisel School of Medicine, Lebanon, New Hampshire, USA

⁷ Swiss Institute of Bioinformatics (SIB), Basel, Switzerland

⁸ Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

⁹ Harvard Stem Cell Institute, Harvard University, Cambridge, Massachusetts, USA

¹⁰ Center for Brain Science, Harvard University, Cambridge, Massachusetts, USA

*Co-corresponding authors. Email: bushranraj@gmail.com (B.R.); alex.schier@unibas.ch (A.F.S.)

Sup Figure 1. Zebrafish brain cell types identified at each stage of time course

tSNE plots of cell types at each stage of the time course. Cells are color coded by stage.

tSNE implementations: Barnes-Hut (12 hpf to 3 dpf, 8 dpf), Fourier transform (5 dpf and 15 dpf). Cluster numbers are indicated on each plot. Detailed cluster descriptions are in Supplementary Table 1 and can be explored interactively in the accompanying app at https://github.com/brlauuu/zf_brain.

Sup Figure 2. Neuron subtype diversity at 5 and 8 dpf

Violin plots of select marker gene expression in identified brain neuron subtypes from 5 dpf (left) and 8 dpf (right). Retina neurons and nascent (immature) neurons are omitted from the analysis. Cluster numbers are indicated at the bottom along with their inferred spatial location in the brain. Clusters 69 (5 dpf) and 51 (8 dpf) have unknown spatial location. Detailed cluster descriptions are in Supplementary Table 1 and can be explored interactively in the accompanying app at https://github.com/brlauuu/zf_brain.

Sup Figure 3. Embryonic, intermediate and larval stage neural progenitor populations

tSNE plots showing embryonic (12 hpf – 18 hpf), intermediate (20 hpf – 3 dpf) and larval (5 dpf – 15 dpf) stage progenitor clusters that were subsetted from the dataset. Cluster numbers match plots shown in Sup Figure 1. Detailed cluster descriptions are in Supplementary Table 1 and can be explored interactively in the accompanying app at https://github.com/brlauuu/zf_brain.

Sup Figure 4. Optimized scGESTALT lineage recorder enables reconstruction of more dense lineage trees

A. Size and diversity of clones from edited ZF1-4 larval brains. Each colored rectangle represents a unique clone and the area represents the size of the clone.

B. A reconstructed scGESTALT brain lineage tree from one zebrafish 15 dpf brain. 302 barcodes recovered from ZF1 were assembled into a lineage tree. Barcode edits are represented as red (deletions), blue (insertions), and black (substitutions). Associated cells are color coded by cell type and region. Interactive trees are presented at: <https://scgestalt.mckennalab.org/>. Each tip on the tree has an associated cell type assignment (color coded), a lineage barcode schematic, and a cluster number. For reasons of space, the tree is split into multiple columns and dashed lines connect subsections of the tree together.

Sup Figure 5.

Retinal cell type marker expression and trajectory analysis

A. UMAP plots highlighting expression of select genes enriched in retinal cell types. *rx3*, *vsx2*, *hes2.2* are enriched in early embryonic retinal progenitors; *foxg1b* is enriched in differentiated cells; *pax6a* is enriched in progenitors, retinal ganglion cells (RGC) and amacrine cells; *crx* is enriched in photoreceptor cells and cone bipolar cells; *gngt2a* is enriched in cones; *gnat1* is enriched in rods; *ompa* is enriched in horizontal cells; *tfap2a* is enriched in RGCs and horizontal cells; *apoeb* is enriched in early progenitors and muller glia; *rbpms2a* is enriched in amacrine cells; *vsx1* is enriched in cone bipolar cells; *cabp5b* is enriched in cone bipolar cells; *rpe65a* is enriched in retinal pigment epithelium; *kidins220a* is enriched in new retinal subtype.

B. tSNE plot of 15 dpf brain and eye cell types. Retinal cell types used as the endpoint cell types (tips) for URD analysis are color coded. Cluster number and cell type description are indicated on the legend. Cluster 96 was discarded from all analysis, see Results.

Sup Figure 6.

Retina and hypothalamus URD trajectory analysis

A. Cell specification tree of zebrafish retinal development generated with URD, reproduction of Figure 6B for comparison. CBP, cone bipolar cells (6 subtypes are numbered); RGC, retinal ganglion cells; RPE, retinal pigment epithelium

B. Assessing robustness of cell assignment on URD retina trees. Trees were recalculated with random subsets of 50% of the cells from the original retinal dataset (sampled per stage so that proportions of cells from each stage remained

constant). The parameters used were the same, except the number of nearest neighbors used was reduced to reflect the smaller dataset. The position of cells on trees was compared to the original tree. Cells assigned the same segment are colored blue (“Same”), cells that moved to a parent or child segment (“Parent/Child”) are colored green, cells that changed assignment to a different location are colored red (“Changed”), and cells that were not assigned a location in the original tree are colored grey. 82.1% cells retained their original assignment, 12.5% cells shifted to either parent or child segments (reflects small shifts in pseudotime of branchpoint assignment), and 5.5% cells shifted to a different location.

C. Cell specification trees of zebrafish hypothalamus development generated with URD, reproduction of Figure 6F for comparison.

D. Assessing robustness of cell assignment on URD hypothalamus trees. Trees were recalculated with random subsets of 50% of the cells from the original hypothalamus dataset (sampled per stage so that proportions of cells from each stage remained constant). The parameters used were the same, except the number of nearest neighbors used was reduced to reflect the smaller dataset. The position of cells on trees was compared to the original tree. Cells assigned the same segment are colored blue (“Same”), cells that moved to a parent or child segment (“Parent/Child”) are colored green, cells that changed assignment to a different location are colored red (“Changed”), and cells that were not assigned a location in the original tree are colored grey. 80.1% cells retained their original assignment, 13.6% cells shifted to either parent or child segments (reflects small shifts in pseudotime of branchpoint assignment or small changes in tree structure), and 6.3% cells shifted to a different location.

Sup Figure 7. Gene expression along retina cell trajectories

Expression of select genes are shown on the retina URD specification tree. Cell types: 1. Cone bipolar cell (CBP_3); 2. Cone bipolar cell (CBP_6); 3. Cone bipolar cell (CBP_1); 4. Cone bipolar cell (CBP_4); 5. Cone bipolar cell (CBP_5); 6. Cone bipolar cell (CBP_2); 7. Cones; 8. Rods; 9. Amacrine cells (Amacrine_1); 10. Amacrine cells (Amacrine_2); 11. Retinal ganglion cells (RGC); 12. Horizontal cells; 13. Muller glia; 14. Retinal pigment epithelium (RPE)

Sup Figure 8. Gene expression cascades of retinal cell trajectories

Heat maps of gene expression cascades of photoreceptor cell, amacrine cell, retinal ganglion cell, muller glia, horizontal cell and retinal pigment epithelium cell trajectories. Cells were selected based on high expression along trajectories leading to these cell types, compared to expression along opposing branchpoints. Red, high expression. Yellow, low expression. X-axis represents cell states along the cascade progression.

Sup Figure 9. Müller glia-like cells are detected early in zebrafish retina development

Detection of Muller glia markers *cahz* and *rlbp1a* in the retina at **A.** 36 hpf and **B.** 2 dpf

A. Left panel, retina section stained with DAPI (grey), *cahz* (cyan) and *rlbp1a* (yellow). White box indicates area that was zoomed in for the right panels.

B. Left panel, retina section stained with DAPI (grey), *cahz* (cyan) and *rlbp1a* (yellow). White and red boxes indicate area that were zoomed in for the middle and right panels, respectively. Middle panels denote cells that co-express *rlbp1a* and *cahz*. Right panels denote cells that are *rlbp1a*⁺ and *cahz*⁻

Sup Figure 10. Hypothalamus cell type marker expression and trajectory analysis

A. UMAP plots highlighting expression of select genes enriched in hypothalamus cell types. *dbx1a*, *fezf2*, *rx3* are enriched in early embryonic hypothalamus progenitors; *fezf1* is enriched in *pdyn*⁺ subtype; *nrgna* is enriched in two subtypes (*synpr*⁺ and *synpr*), *dlx2a* is expressed in several subtypes; *nkx2.4a* is expressed in early progenitors and *prdx1*⁺ subtype; *synpr*, *npv*, *tph2*, *pdyn* and *prdx1* are enriched in specific subtypes

B. tSNE plot of 15 dpf brain and eye cell types. Hypothalamus cell types used as the endpoint cell types (tips) for URD analysis are color coded. Cluster number and cell type description are indicated on the legend.

Sup Figure 11. Gene expression along hypothalamus cell trajectories

Expression of select genes are shown on the hypothalamus specification tree. Cell types: 1. GABA *tac1*⁺, *nrgna*⁺; 2. *synpr*⁺; *nrgna*⁺; 3. *sst1.1*⁺; 4. *tph2*⁺; 5. GABA *dlx*⁺; 6. *pdyn*⁺; 7. *prdx1*⁺

Sup Figure 12. Gene expression cascades of hypothalamus cell trajectories

Heat maps of gene expression cascades of profiled hypothalamus cell trajectories. Cells were selected based on high expression along trajectories leading to these cell types, compared to expression along opposing branchpoints. Red, high expression. Yellow, low expression. X-axis represents cell states along the cascade progression.

Sup Figure 13. Cell type maturation along URD trajectories

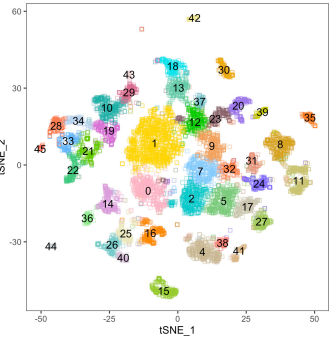
Retinal ganglion cells (**A**) and *pdyn*⁺ hypothalamic neurons (**B**) cells were plotted (red circles) on URD cell specification trajectories across the stages indicated. The cell types matured with developmental age, as expected. Additionally, later stages also contained immature cell states (early pseudotime) consistent with continuous neurogenesis.

Sup Figure 14. Gene markers of embryonic and larval progenitors in retina and hypothalamus

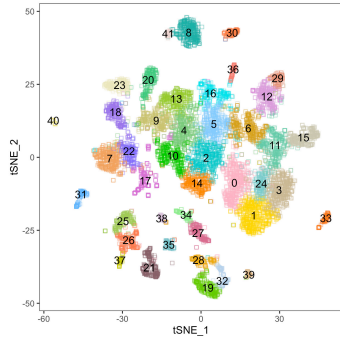
Dot plot of gene expression pattern of select marker genes that were used to define progenitor and precursor states (rows) for segments (columns) of the retina (left) or hypothalamus (right) cell specification trees (see Figure 7E). Dot size indicates the percentage of cells expressing the marker; color represents the average scaled expression level.

Sup Figure 1

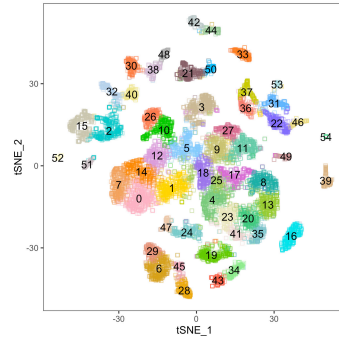
12 hpf



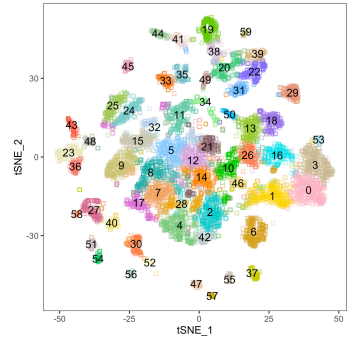
14 hpf



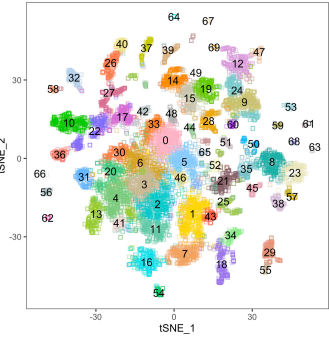
16 hpf



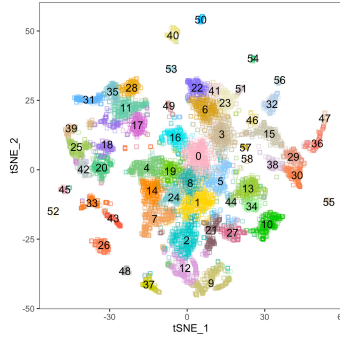
18 hpf



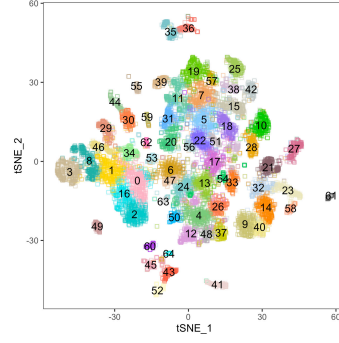
20 hpf



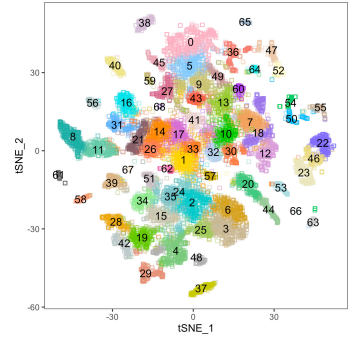
24 hpf



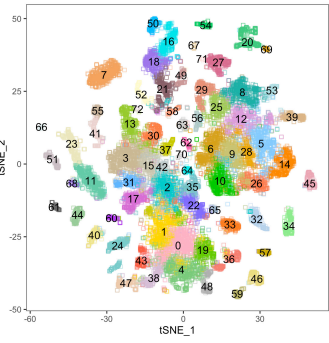
36 hpf



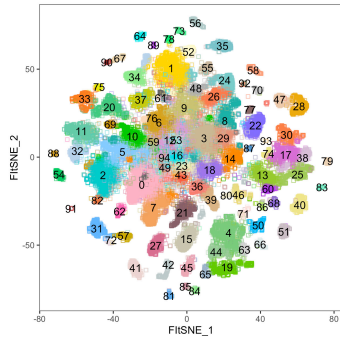
2 dpf



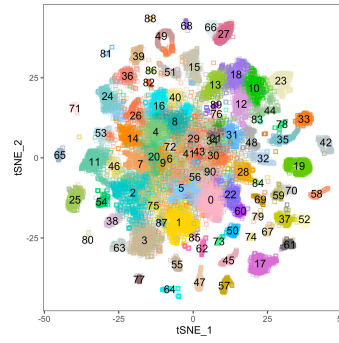
3 dpf



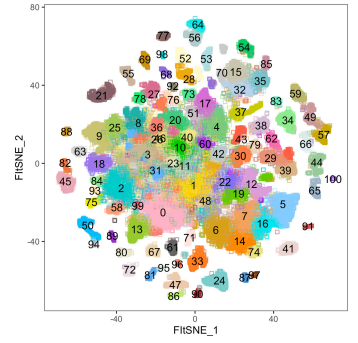
5 dpf



8 dpf



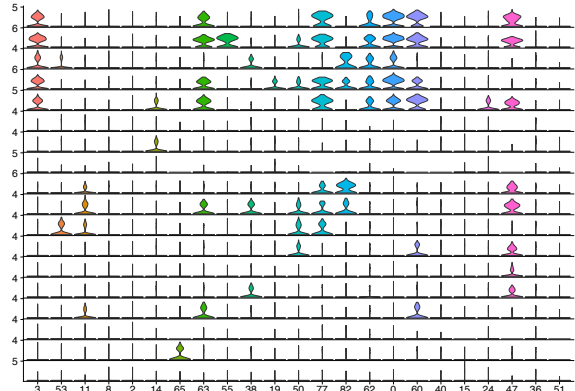
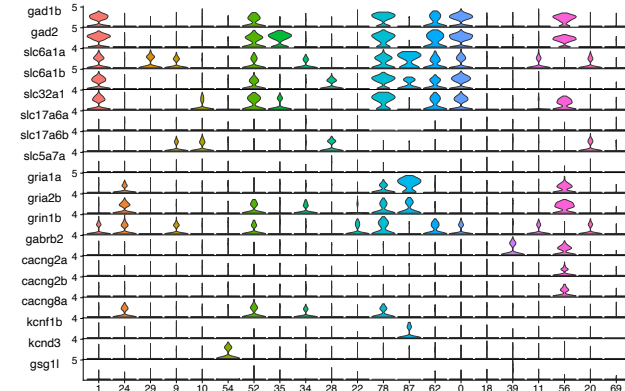
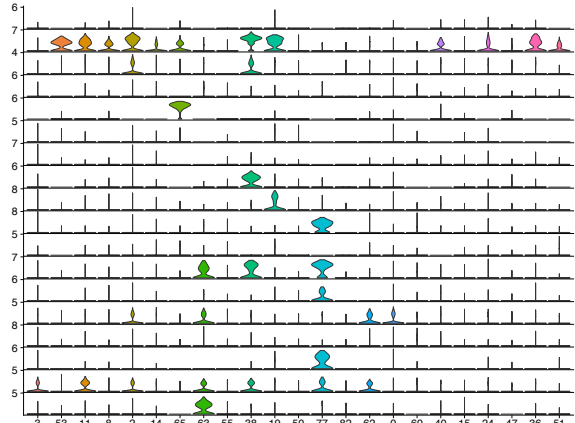
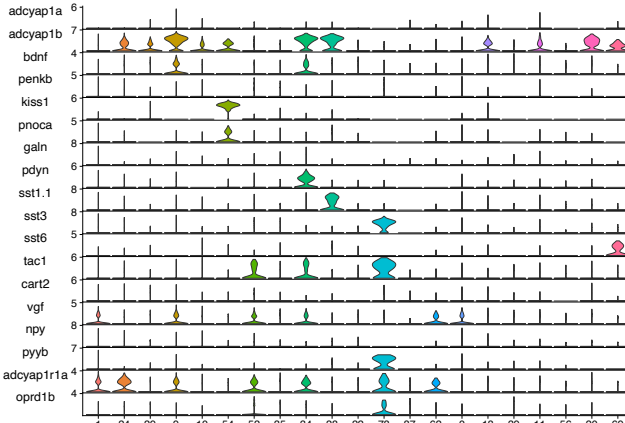
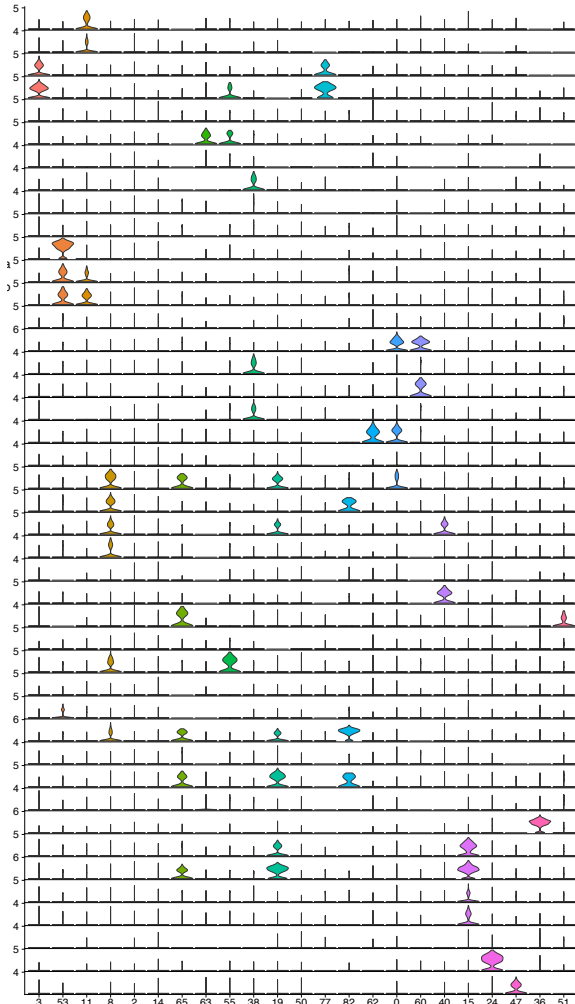
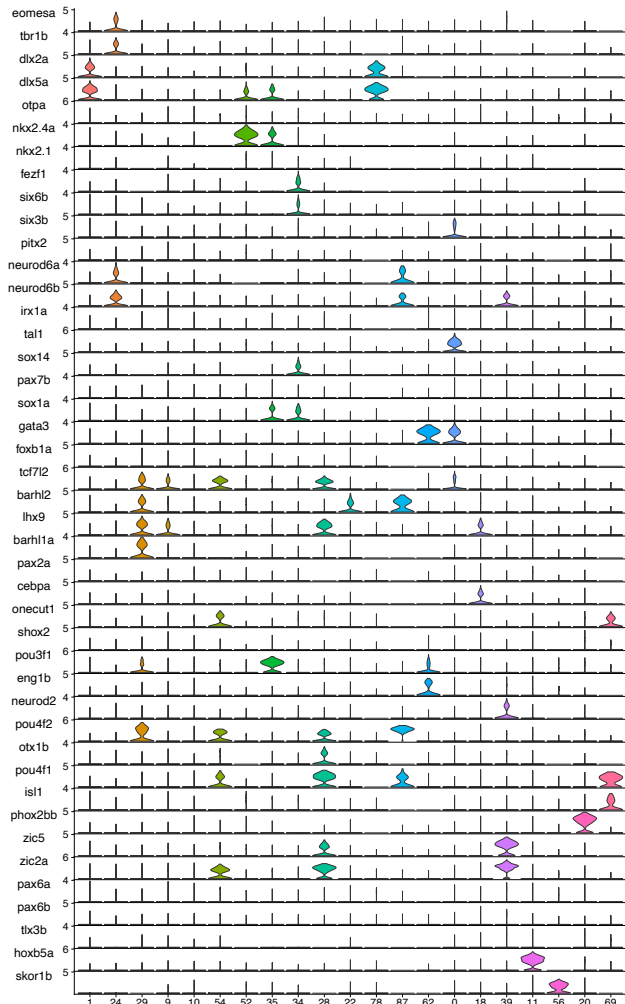
15 dpf



Sup Figure 2

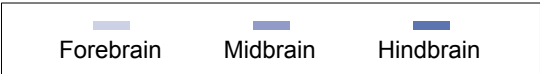
5 dpf

8 dpf

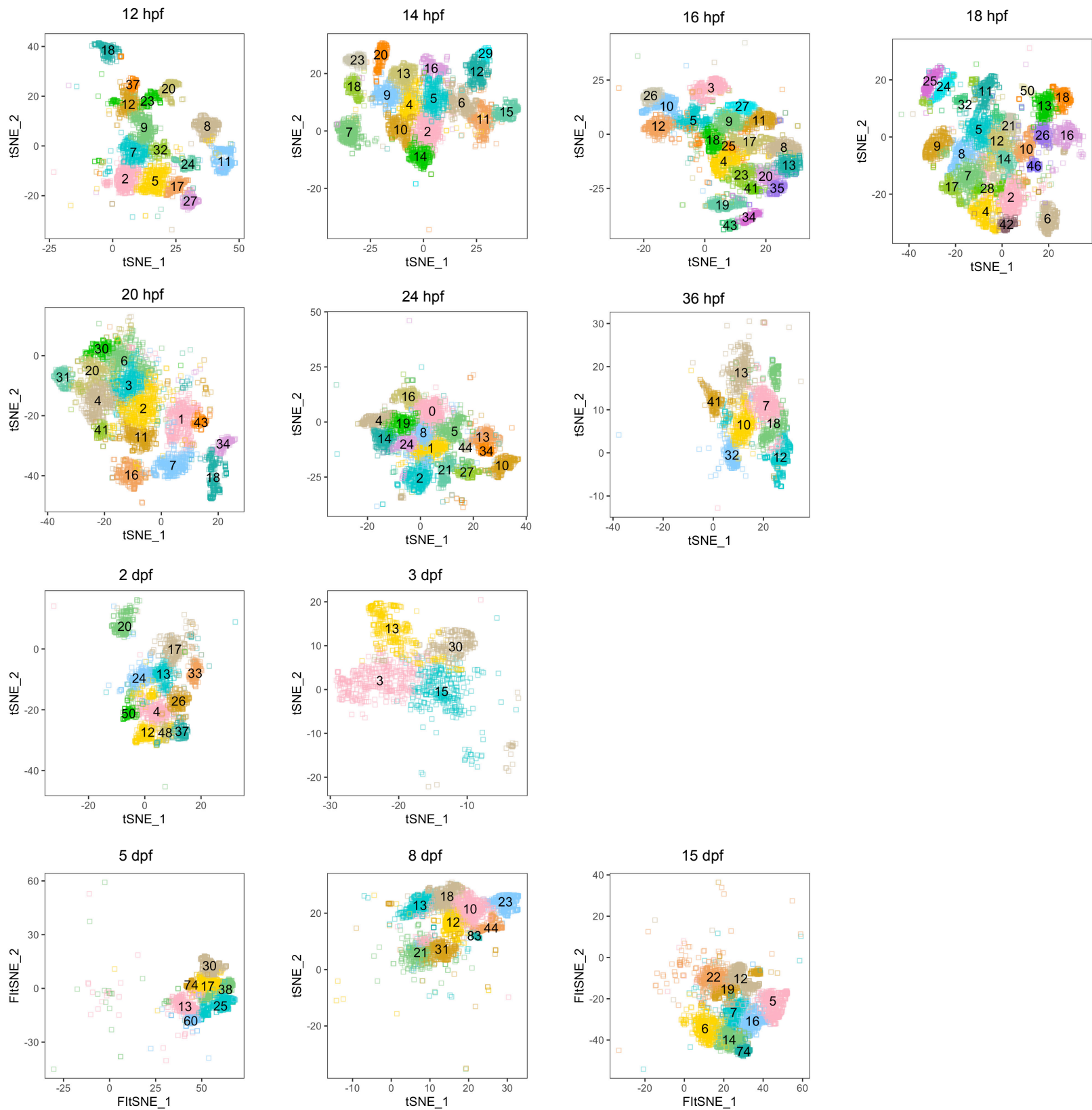


- Neuropeptides
- Receptors
- Neurotransmitters
- Transporters
- Receptors
- Channels

Clusters

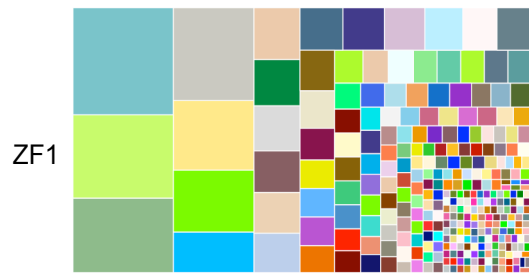


Sup Figure 3

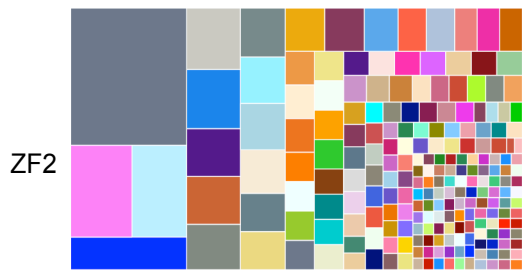


Sup Figure 4

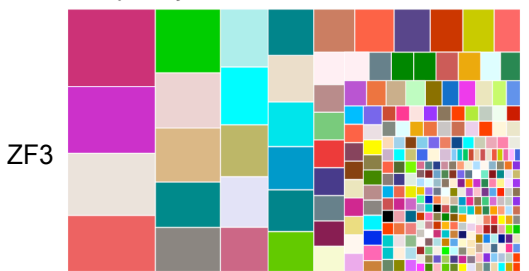
A



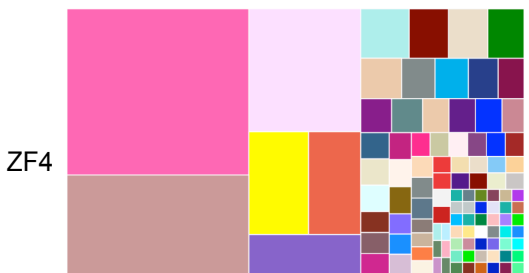
302 clones
2,331 cells
largest clone = 207 cells



216 clones
968 cells
largest clone = 130 cells

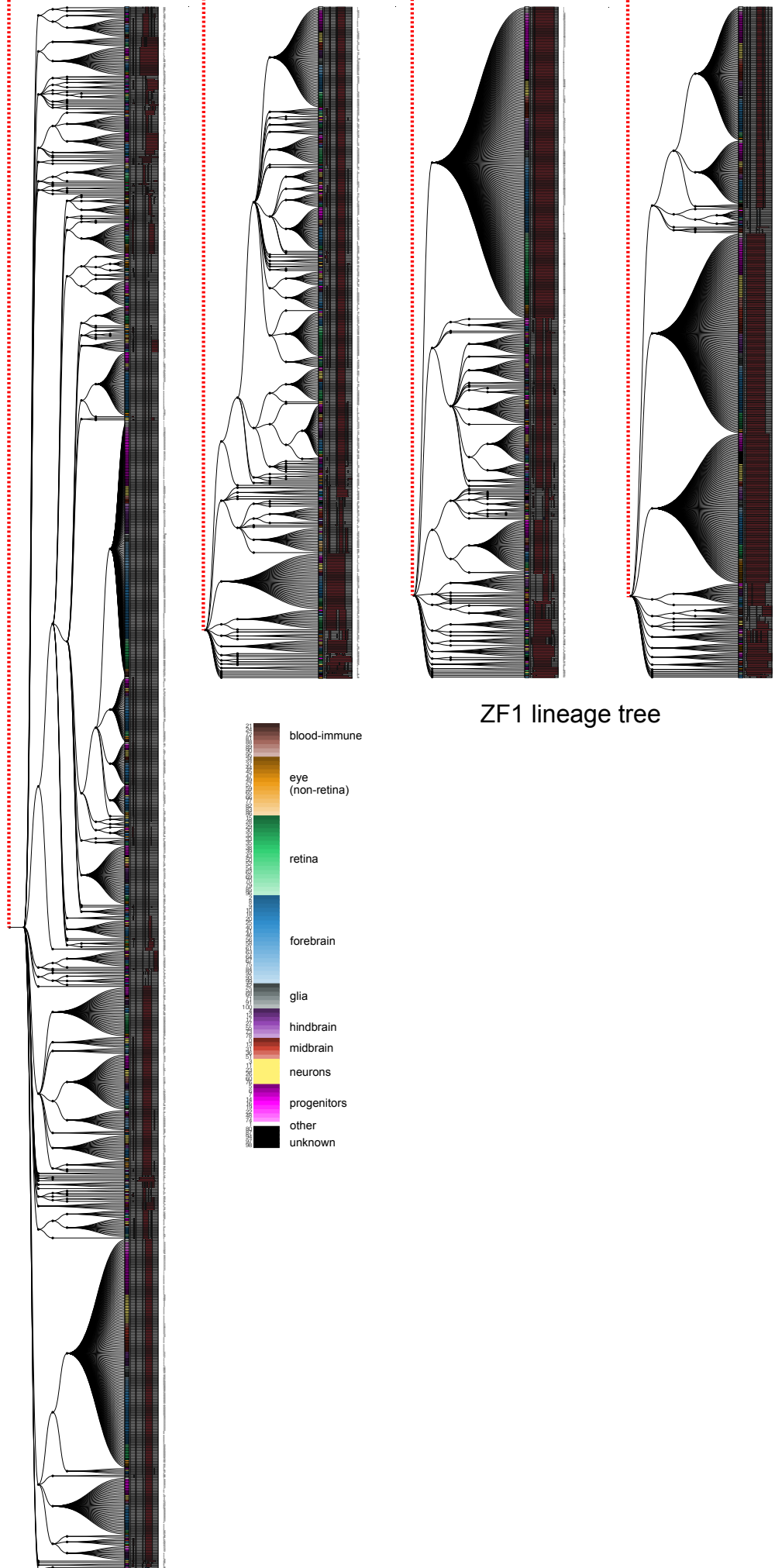


312 clones
1,685 cells
largest clone = 96 cells

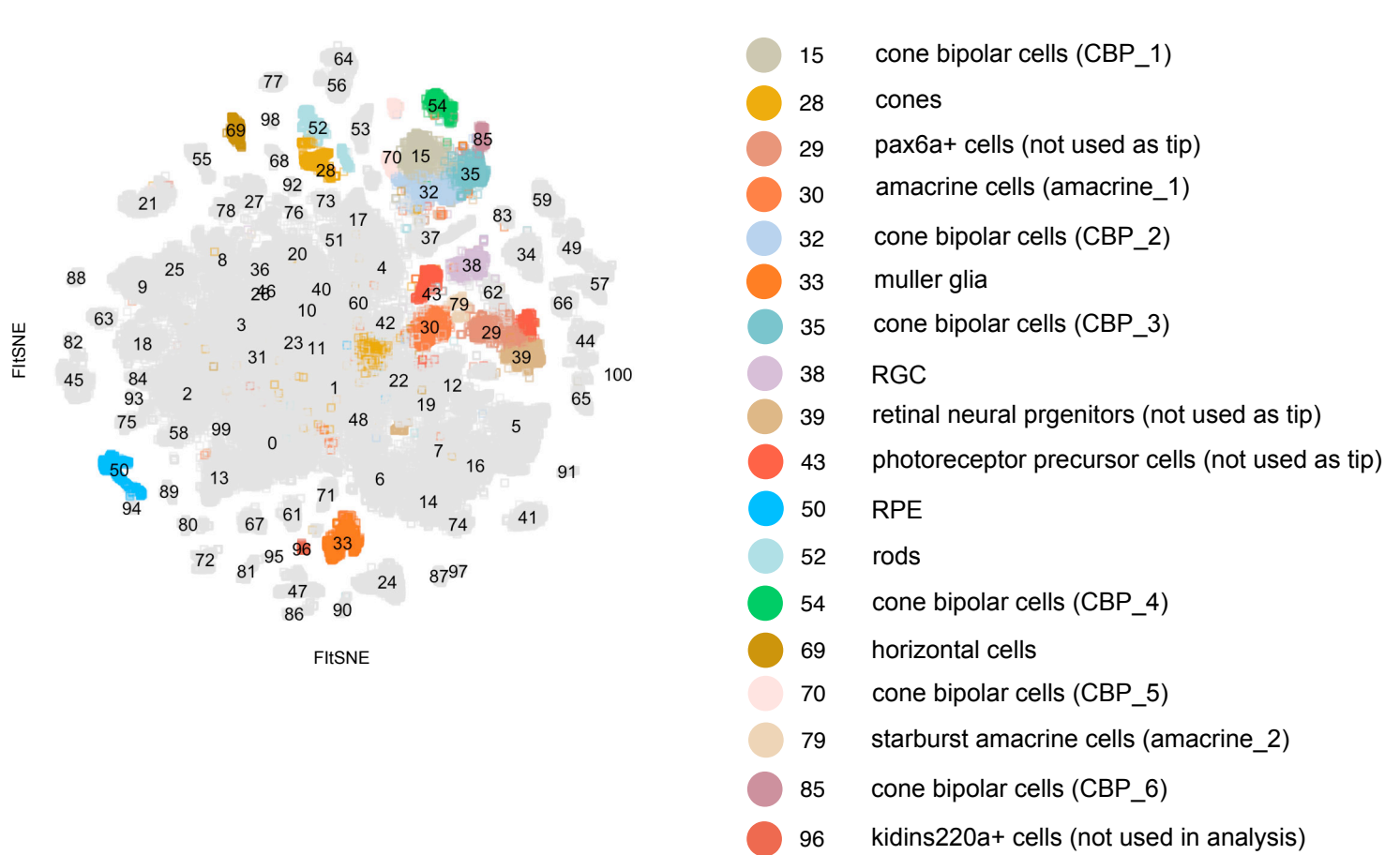
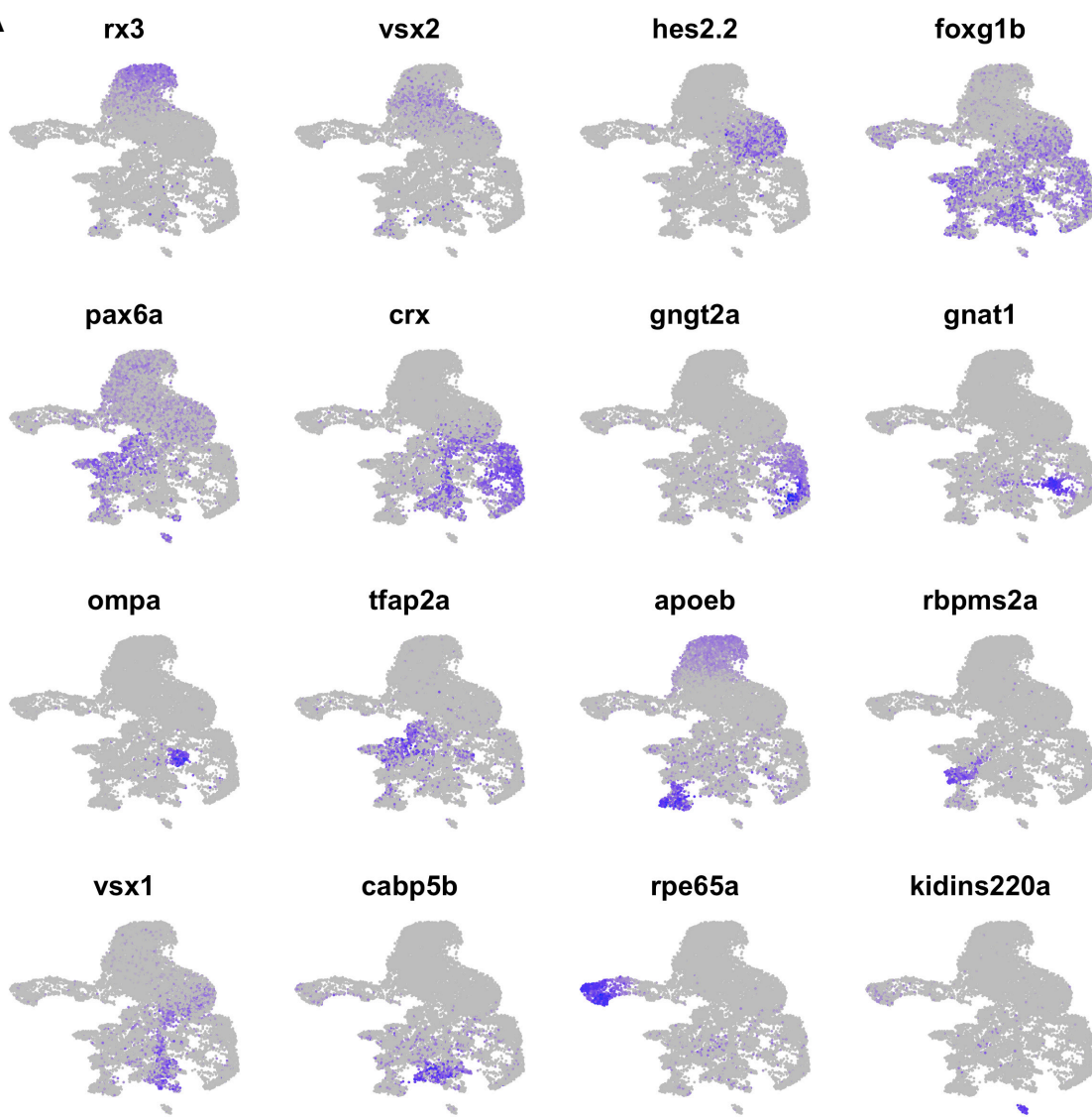


106 clones
812 cells
largest clone = 202 cells

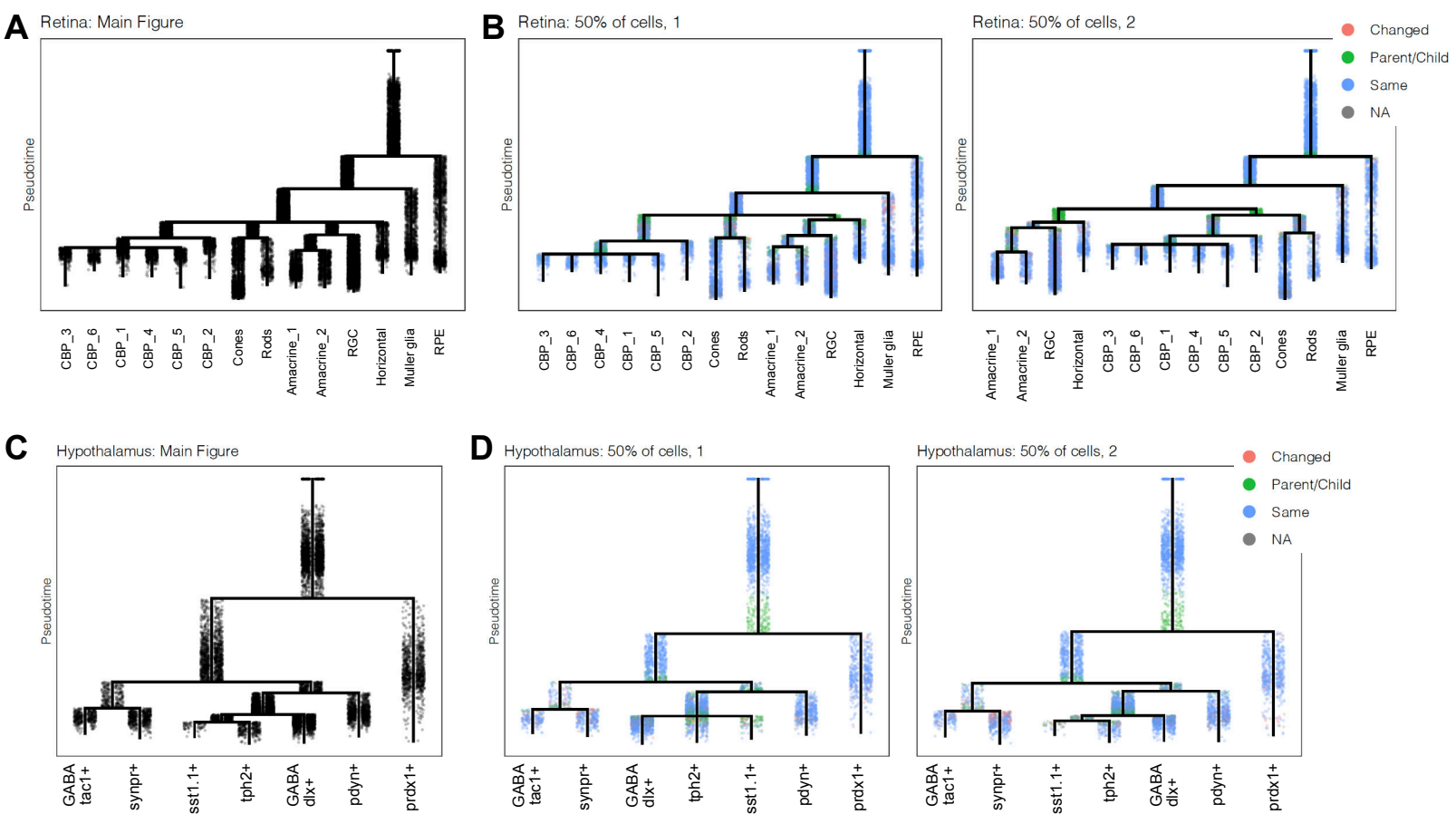
B



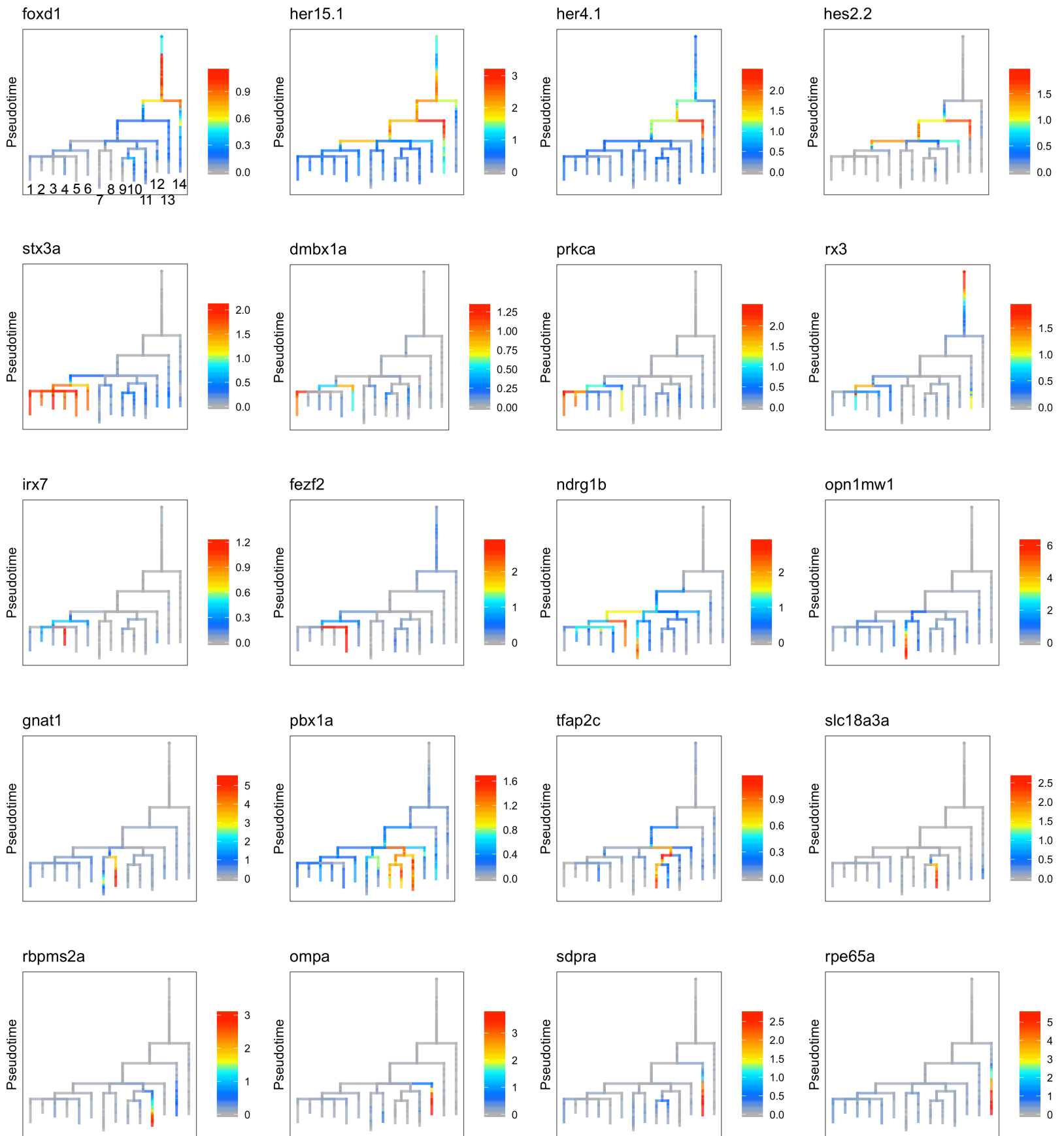
ZF1 lineage tree



Sup Figure 6

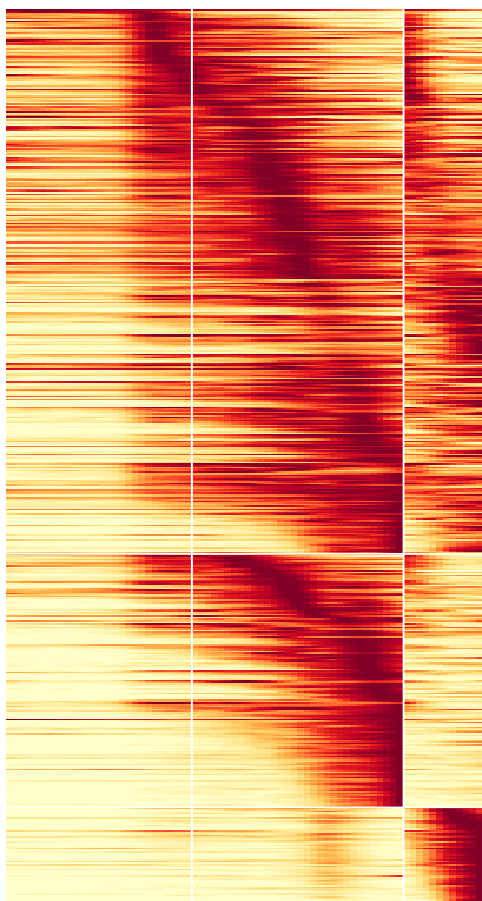


Sup Figure 7



Sup Figure 8

Photoreceptors

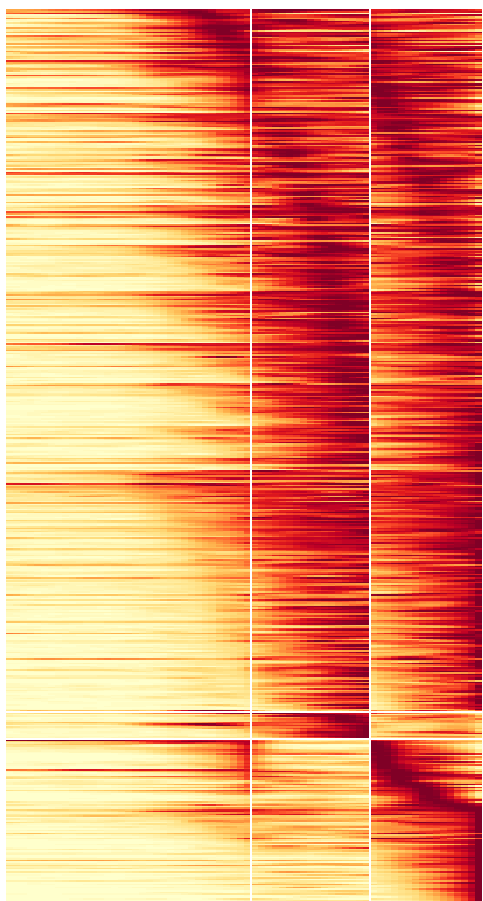


Precursors

Cones

Rods

Amacrine Cells

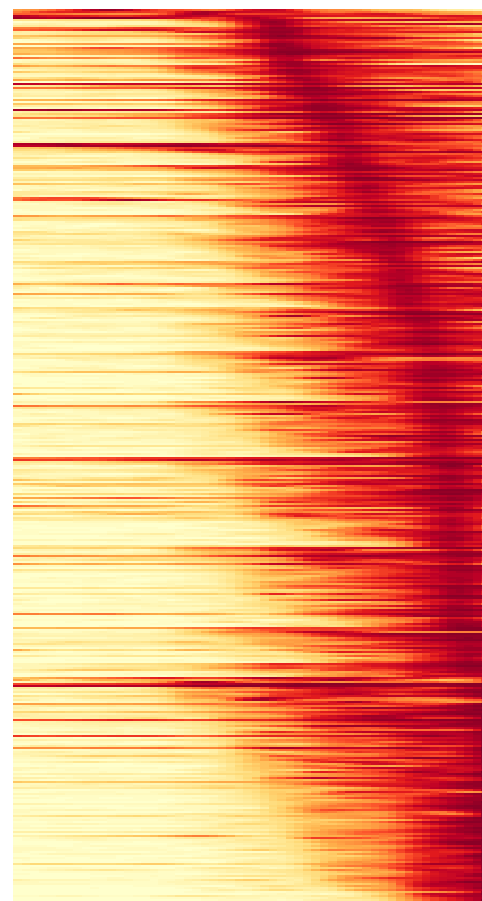


Precursors

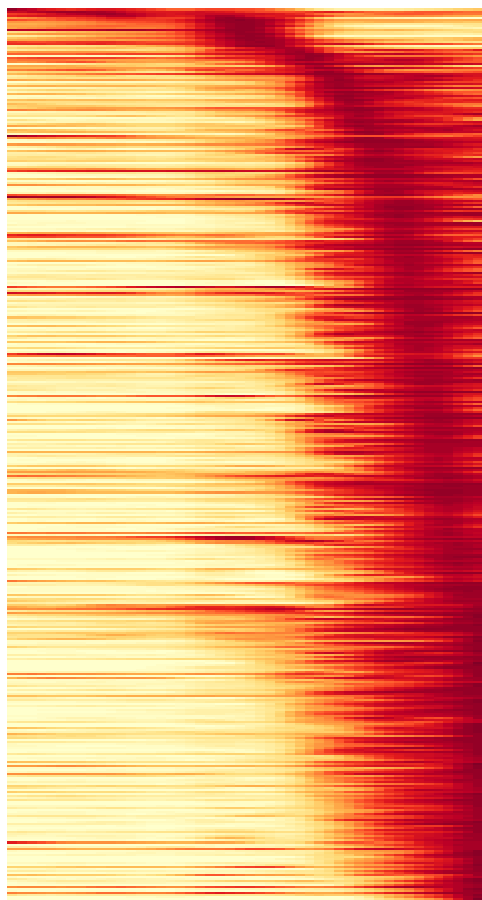
Amacrine

Starburst

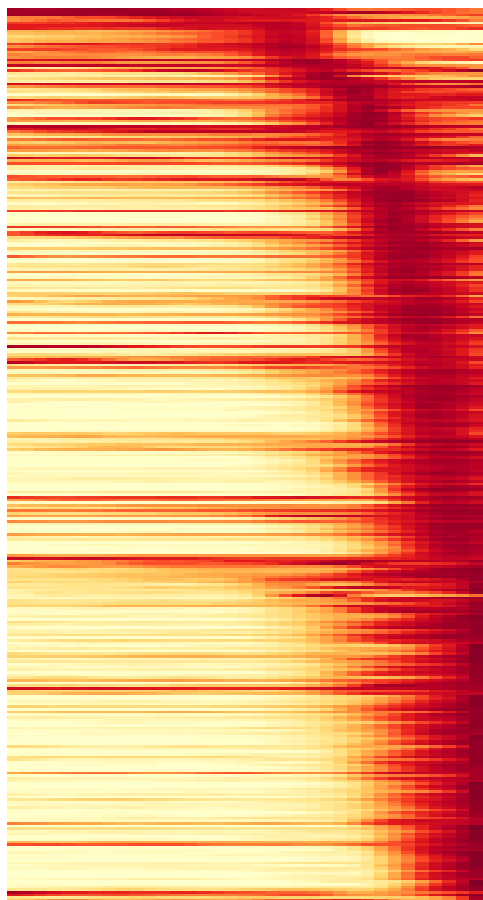
Retinal Ganglion Cells



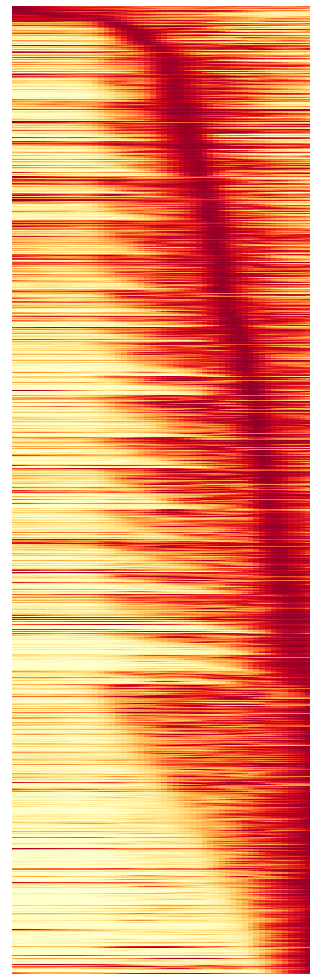
Muller Glia



Horizontal Cells



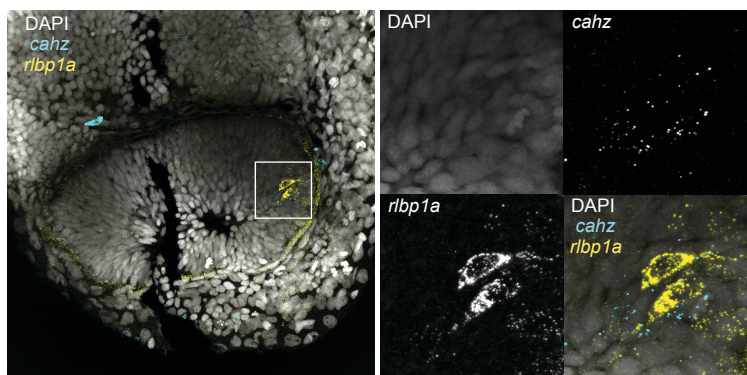
Retinal Pigmented Epithelium



Sup Figure 9

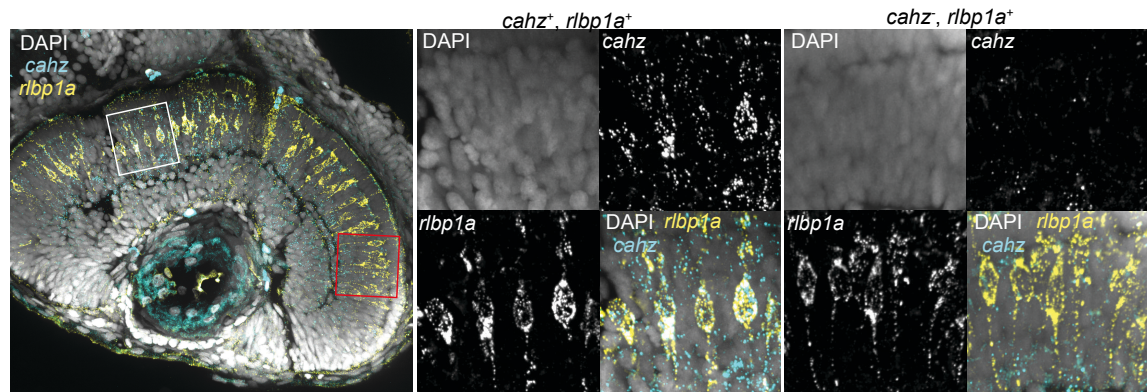
A

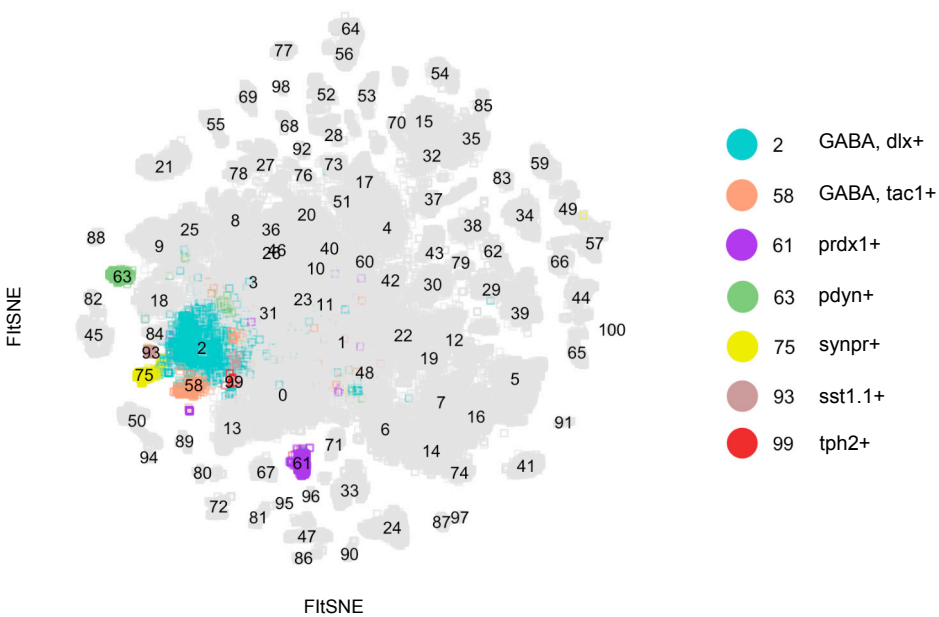
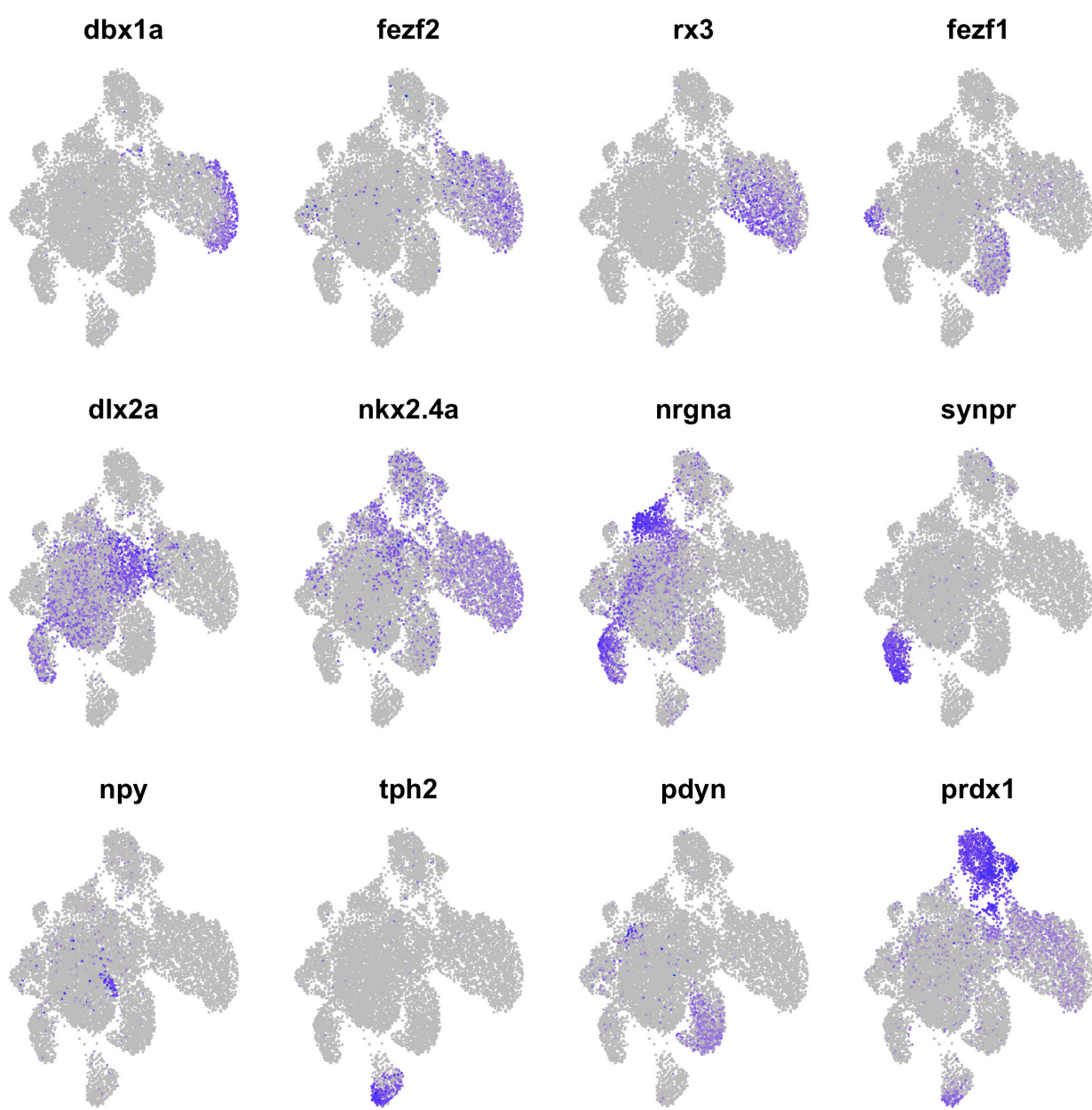
36 h



B

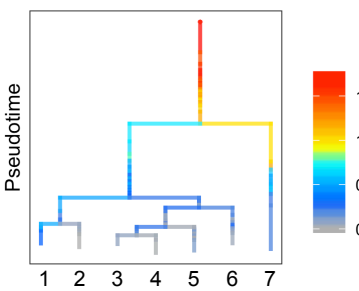
2 d



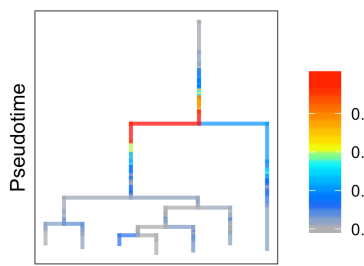


Sup Figure 11

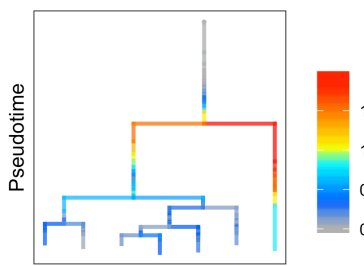
nkx2.4b



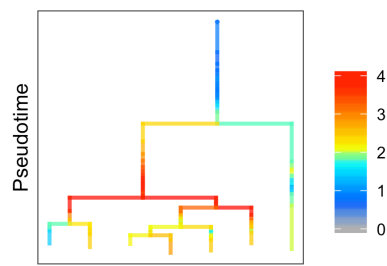
ascl1a



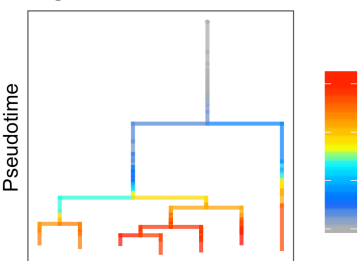
insm1a



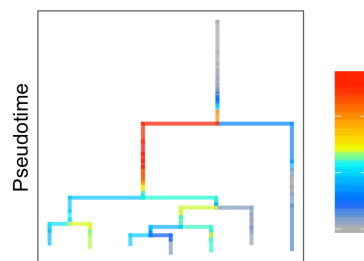
tubb5



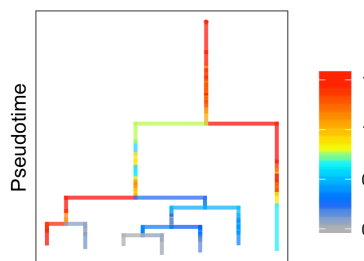
scg2b



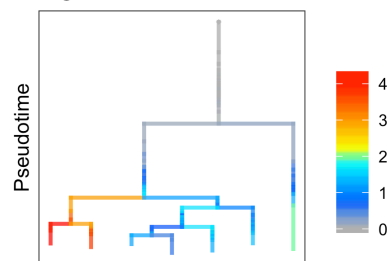
dlx2a



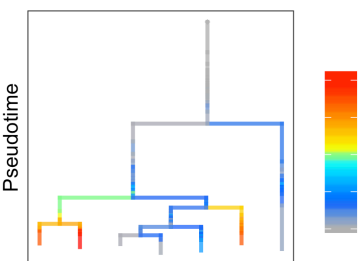
nkx2.4a



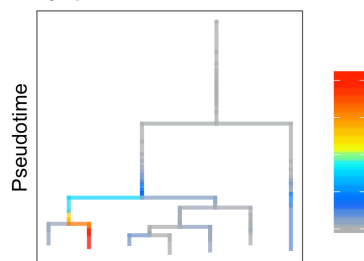
nrgna



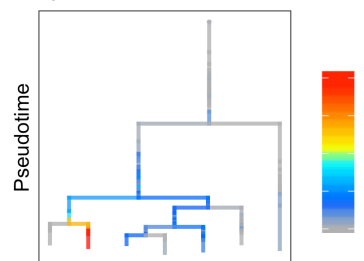
tac1



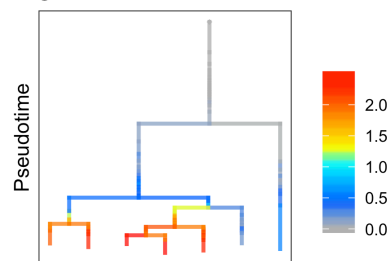
synpr



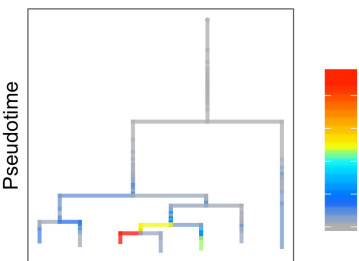
sp8a



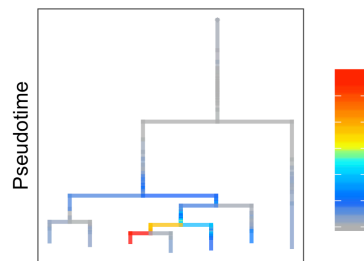
gad1b



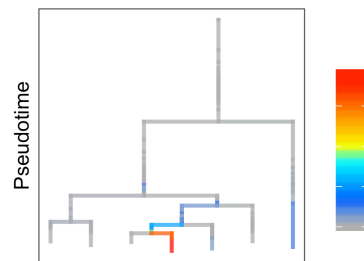
npy



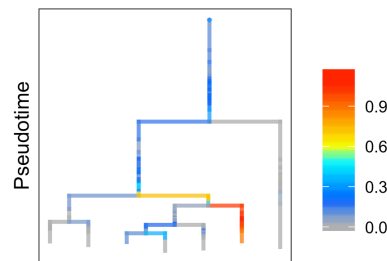
sst1.1



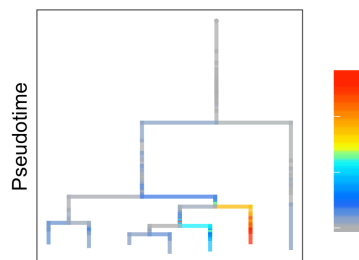
tph2



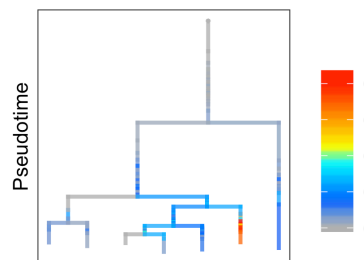
fezf1



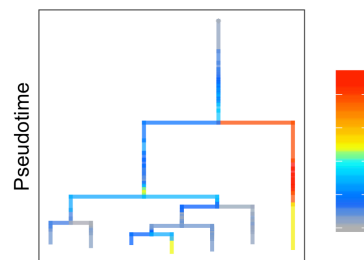
pdyn



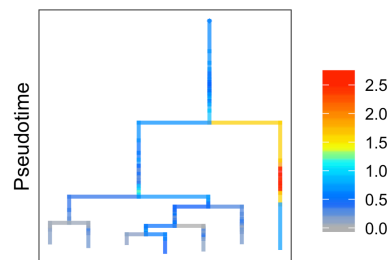
slc17a6b



prdx1

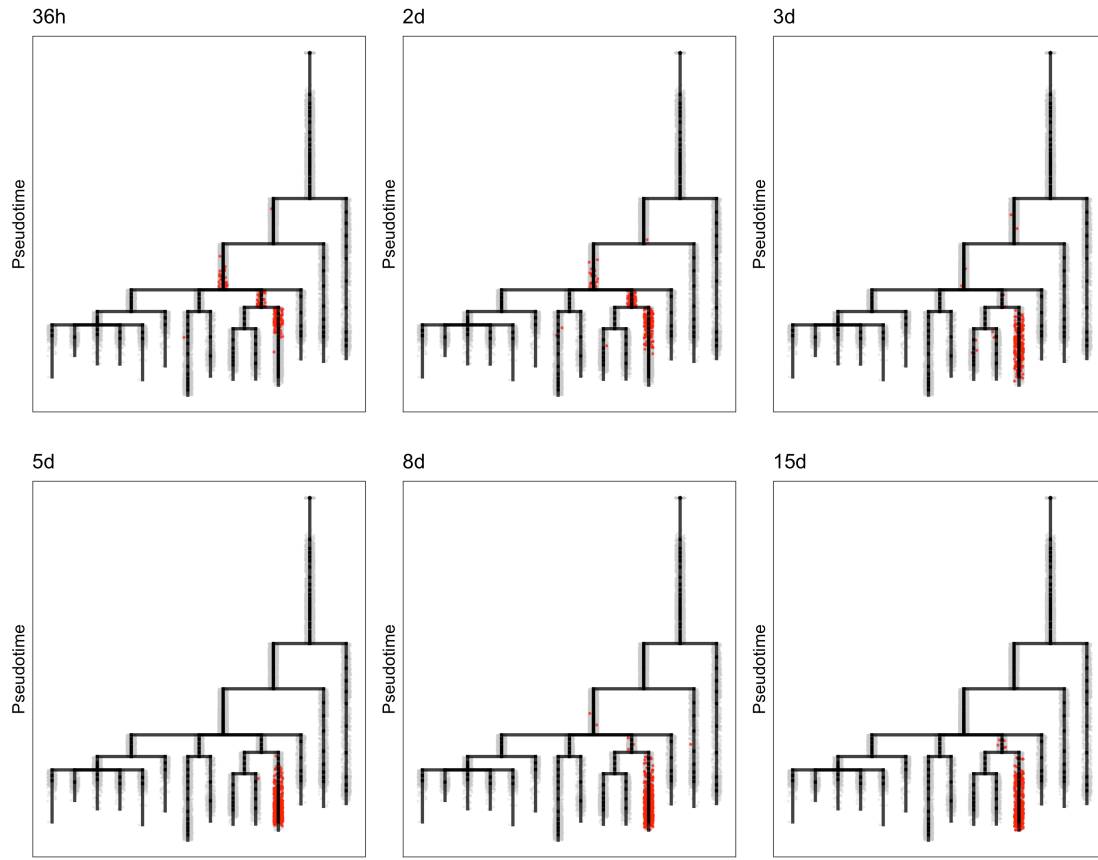


pou3f1

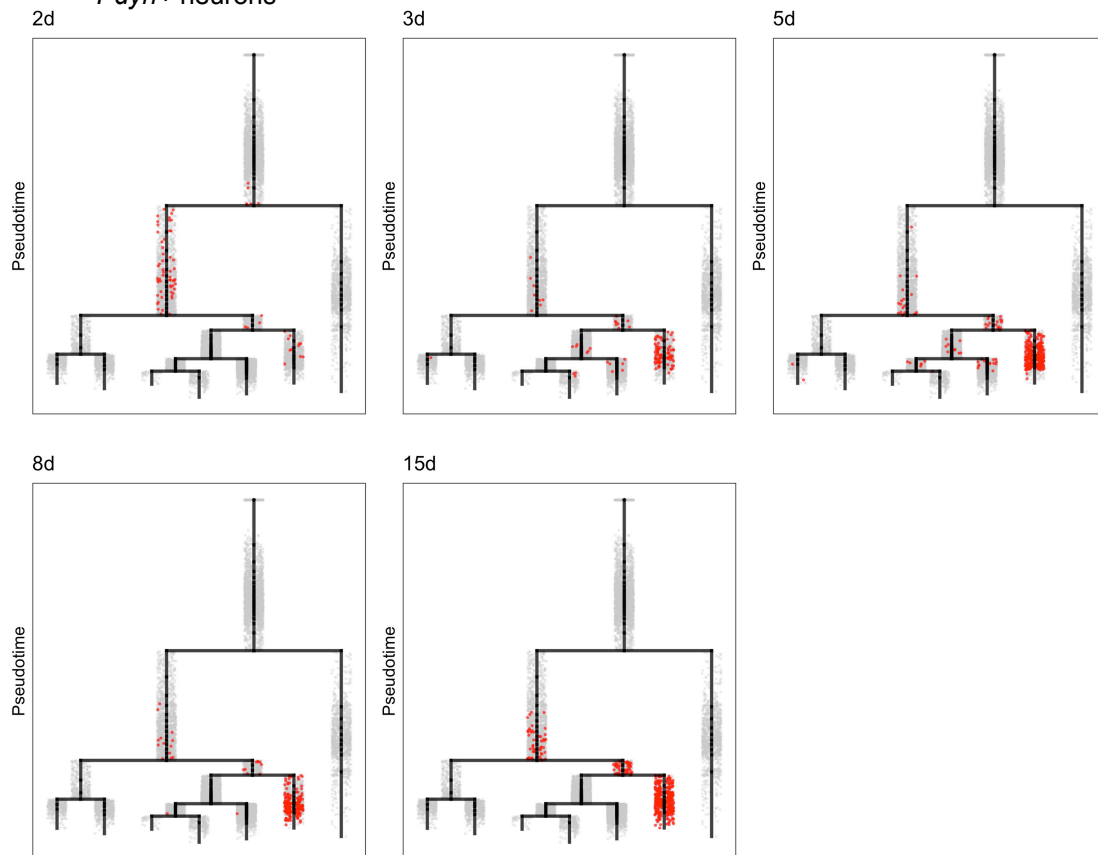


Sup Figure 13

A Retinal ganglion cells

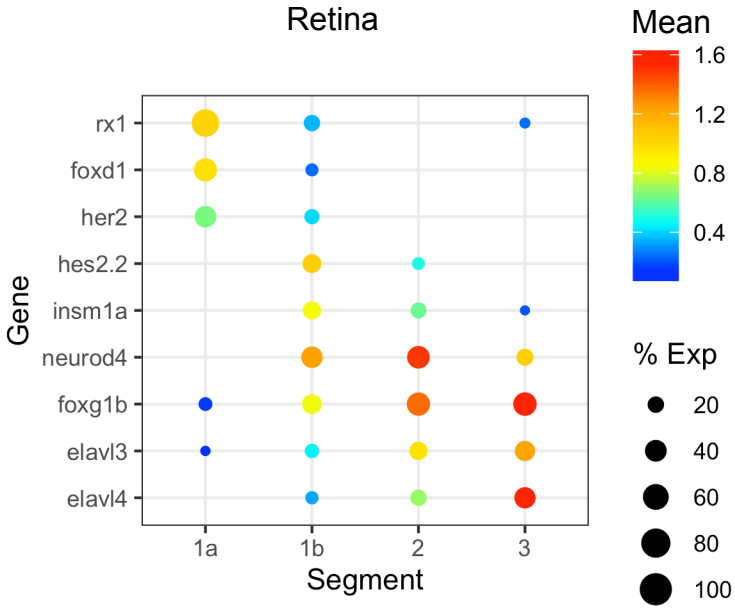


B *Pdyn+* neurons

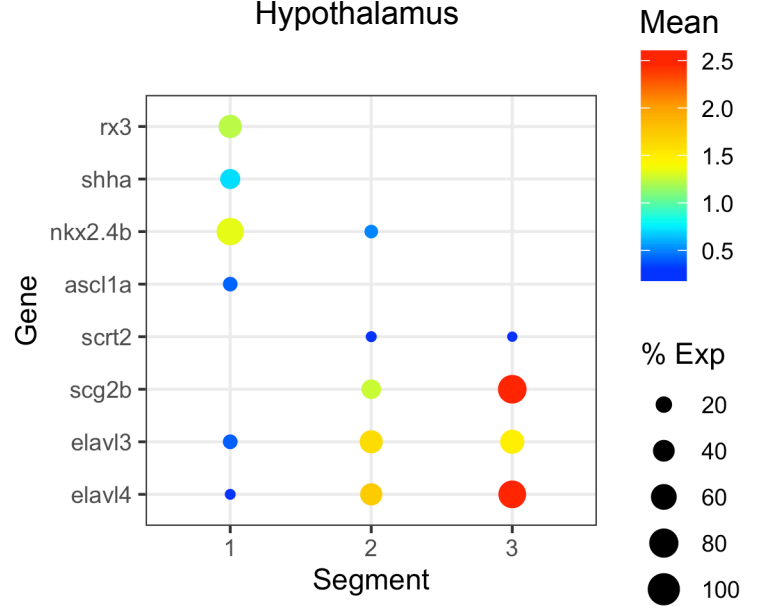


Sup Figure 14

Retina



Hypothalamus



Retina: 1 - URD object & doublet removal

Jeff Farrell

8/22/2019

Contents

Import data into URD	1
Convert Seurat object to URD	1
Combined individual stage clustering	1
Calculate highly variable genes	2
Calculate KNN graph and remove outliers	7
Remove <i>kidins220a+</i> population	8
Remove cell type doublets	8
Add UMAP projection	8
Load NMF results and import into object	9
Select cell-type specific modules	9
Determine which module pairs to use for doublet removal	12

Import data into URD

Convert Seurat object to URD

We first loaded a Seurat object that contained just cells from the clusters that belonged to the hypothalamus from each stage.

```
suppressPackageStartupMessages(library(URD))
suppressPackageStartupMessages(library(Seurat))

base.path <- "~/urd-cluster-bushra/"

# Load Seurat object that has been cropped to hypothalamus cells
object.seurat <- readRDS(paste0(base.path, "obj/retina.new_seurat.rds"))

# Convert to URD object
suburd <- seuratToURD(object.seurat)
```

Combined individual stage clustering

Bushra had performed individual clusterings across each stage with different resolutions. Here, it was better to create a single identifier that included stage + cluster information to combine all those clusterings (while preventing any overlap).

```

stages <- sort(unique(suburd@meta$stage))
clust.res.used <- paste0("res.", c("4.5", "4", "5", "5", "4.5", "5", "6",
  "6", "6", "5.5", "6", "5"))
names(clust.res.used) <- stages
suburd@group.ids$cluster <- NA
for (stage in stages) {
  suburd@group.ids[cellsInCluster(suburd, "stage", stage), "cluster"] <- paste0(stage,
    "-", suburd@group.ids[cellsInCluster(suburd, "stage", stage), clust.res.used[stage]])
}

```

Calculate highly variable genes

We calculated highly variable genes for each stage, used genes that were found as highly in at least two stages, but were not mitochondrial, ribosomal, heat-shock protein, or tandem duplicated genes.

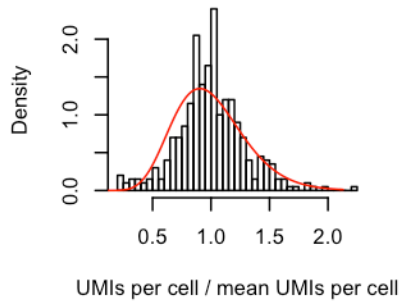
```

# Calculated on each stage separaely, final gene list was all genes
# that were 'variable' in at least two stages NB: For a couple of
# stages, the gamma fit was poor -- the library size distribution
# seemed bimodal. Have seen this before in 10X data, but not sure what
# it means.
var.genes.by.stage <- lapply(stages, function(stage) {
  findVariableGenes(suburd, cells.fit = cellsInCluster(suburd, "stage",
    stage), set.object.var.genes = F, diffCV.cutoff = 0.3, main.use = stage,
    do.plot = T)
})

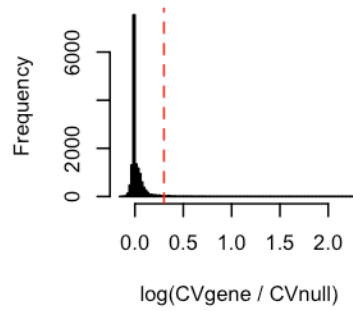
```


01-12h

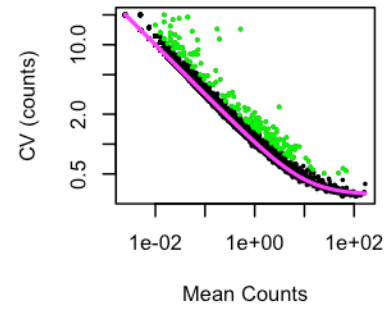
Size Factors & Gamma Fit ($\alpha=10$)



Diff CV

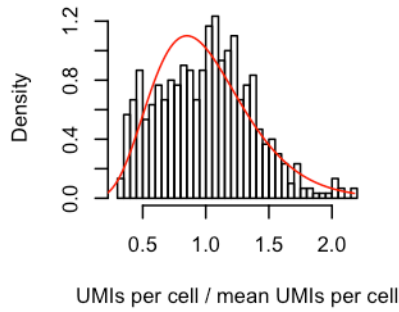


Selection of Variable Genes

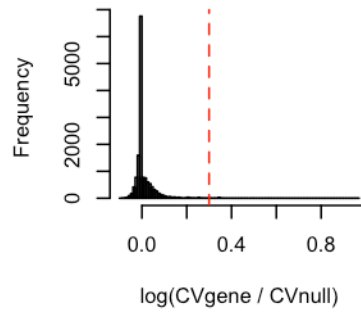


02-14h

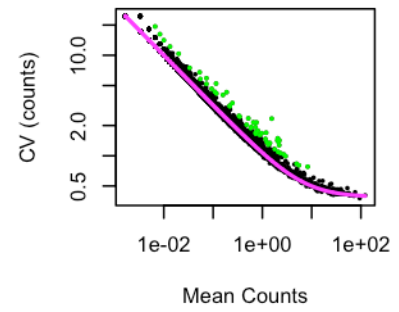
Size Factors & Gamma Fit ($\alpha=6.6$)



Diff CV

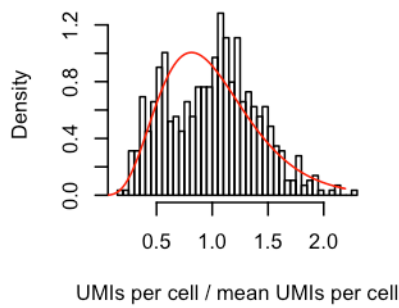


Selection of Variable Genes

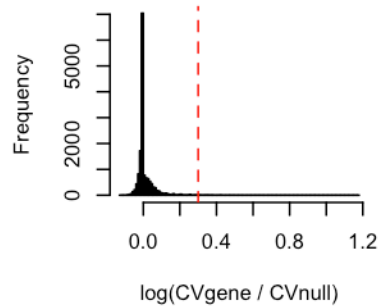


03-16h

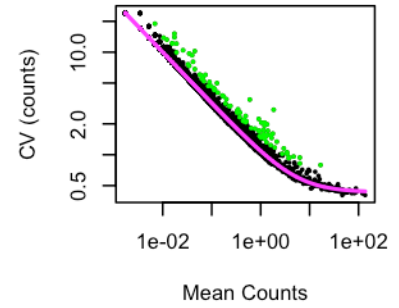
Size Factors & Gamma Fit ($\alpha=5.3$)



Diff CV

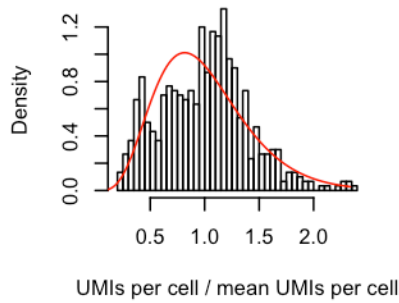


Selection of Variable Genes

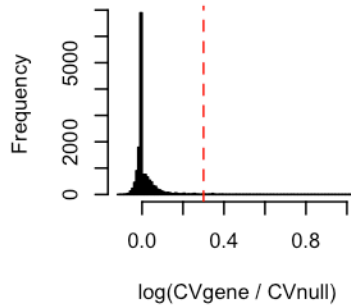


04-18h

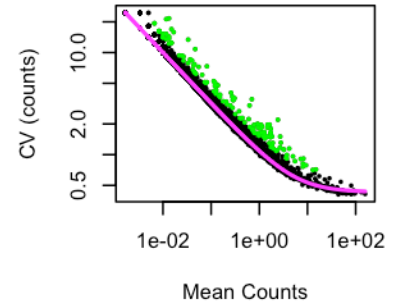
Size Factors & Gamma Fit ($\alpha=5.4$)



Diff CV

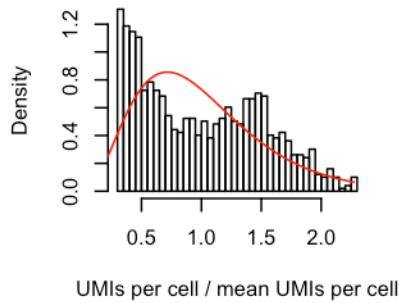


Selection of Variable Genes

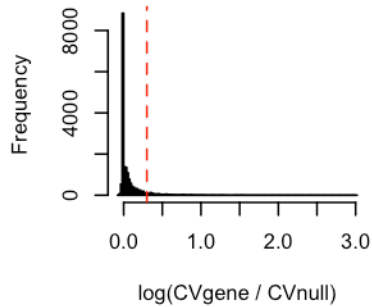


05-20h

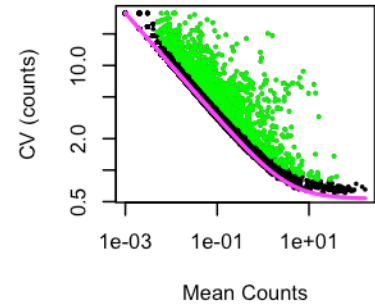
Size Factors & Gamma Fit ($\alpha=3.5$)



Diff CV

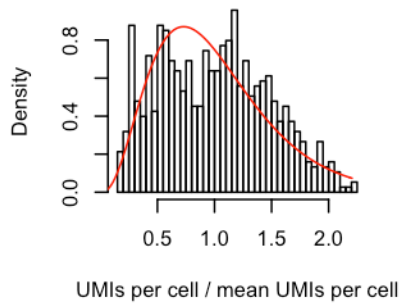


Selection of Variable Genes

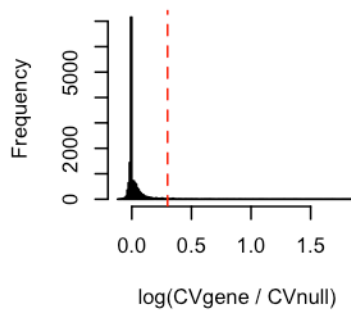


06-24h

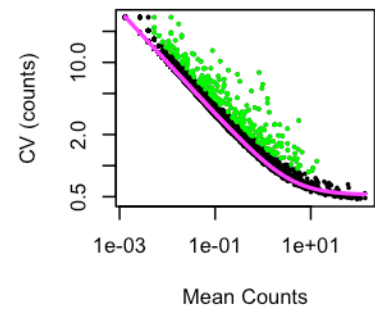
Size Factors & Gamma Fit ($\alpha=3.7$)



Diff CV

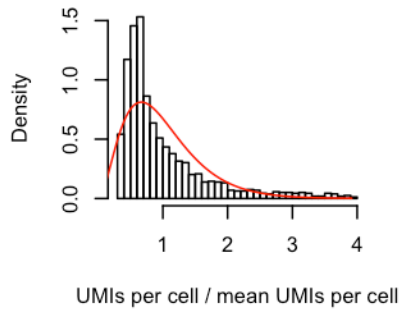


Selection of Variable Genes

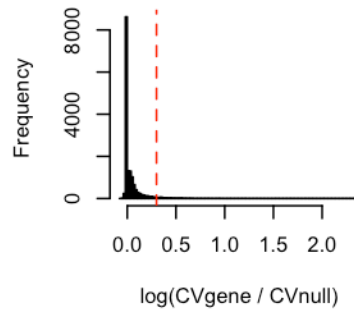


07-36h

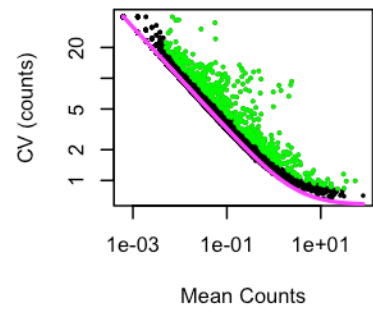
Size Factors & Gamma Fit ($\alpha=3.0$)



Diff CV

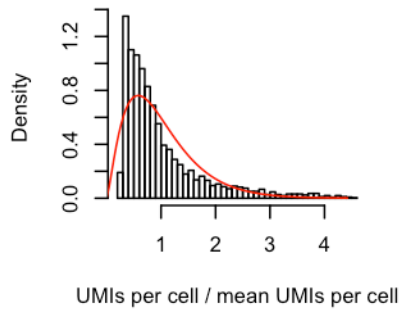


Selection of Variable Genes

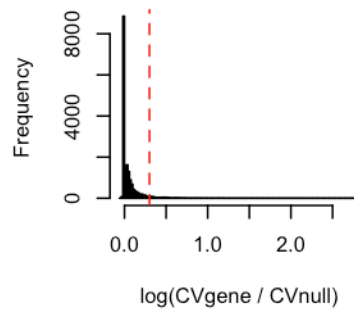


08-2d

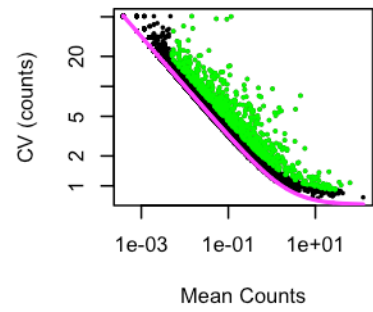
Size Factors & Gamma Fit ($\alpha=2.3$)



Diff CV

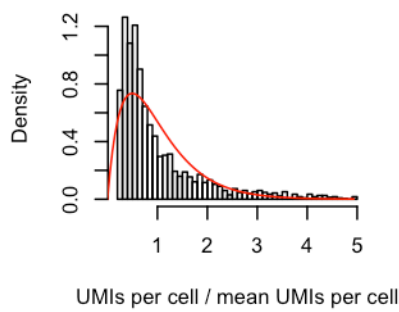


Selection of Variable Genes

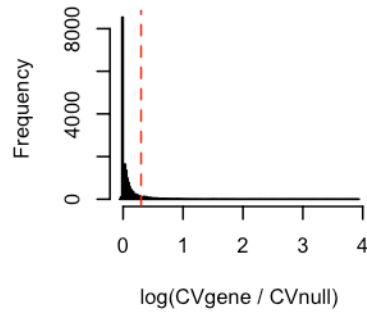


09-3d

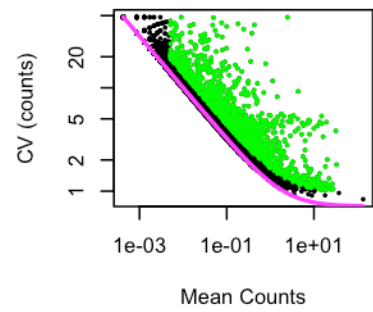
Size Factors & Gamma Fit ($\alpha=1.9$)



Diff CV

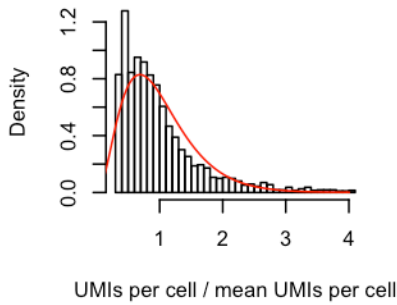


Selection of Variable Genes

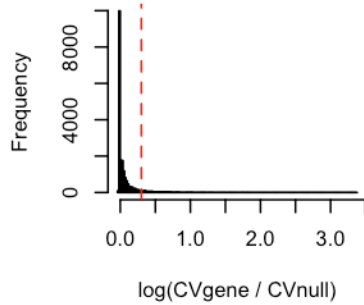


10-5d

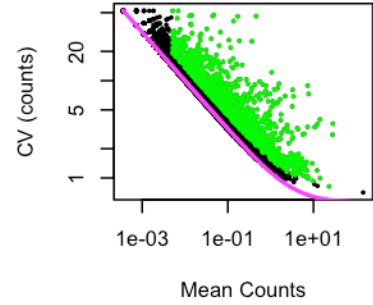
Size Factors & Gamma Fit ($\alpha=3.1$)



Diff CV

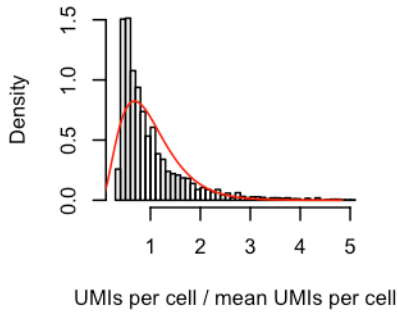


Selection of Variable Genes

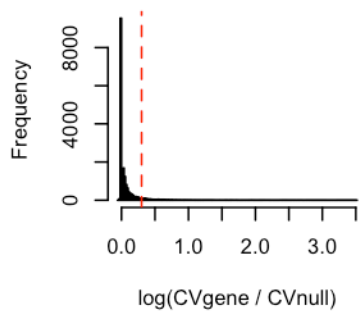


11-8d

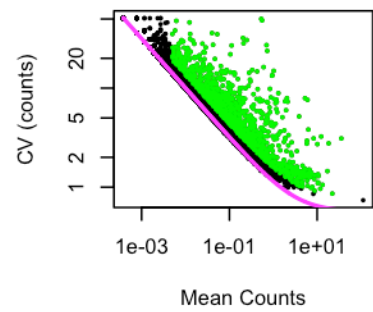
Size Factors & Gamma Fit ($\alpha=3.1$)



Diff CV

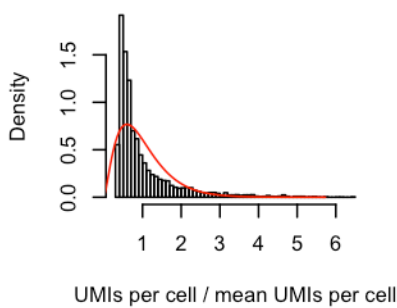


Selection of Variable Genes

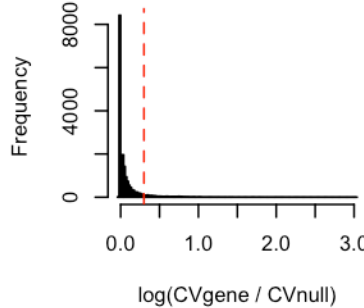


12-15d

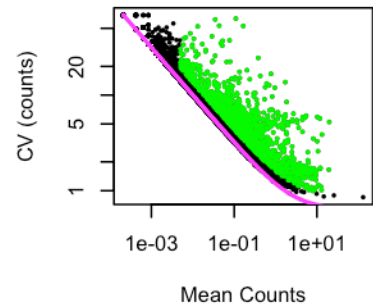
Size Factors & Gamma Fit ($\alpha=2.4$)



Diff CV



Selection of Variable Genes



```
names(var.gen.es.by.stage) <- stages
var.gen.es <- sort(unique(unlist(var.gen.es.by.stage)))
print(paste0("Length of variable genes is ", length(var.gen.es)))
```

```
## [1] "Length of variable genes is 2636"
```

```
var.gen.es.twice <- names(which(table(unlist(var.gen.es.by.stage)) >= 2))
print(paste0("Length of variable genes shared across at least 2 stages is ",
length(var.gen.es.twice)))
```

```
## [1] "Length of variable genes shared across at least 2 stages is 1724"
```

```
# Remove mitochondrial genes
var.mito <- grep("^mt-|^AC0", var.genes.twice, value = T)
# Remove ribosomal genes
var.ribo <- grep("^rps|^rpl", var.genes.twice, value = T)
# Remove hsp genes
var.hsp <- grep("^hsp", var.genes.twice, value = T)
# Remove genes with duplicates
var.dups <- grep("of many", var.genes.twice, value = T)
suburd@var.genes <- setdiff(var.genes.twice, c(var.mito, var.ribo, var.hsp,
var.dups))
print(paste0("Length of final variable genes list (after removing mito, ribo, hsp genes) is ",
length(suburd@var.genes)))
```

```
## [1] "Length of final variable genes list (after removing mito, ribo, hsp genes) is 1595"
```

To prevent downstream problems, we also removed any cells from the data that had the exact same expression of the variable genes (i.e. cells with completely duplicated coordinates in the high-dimensional space we would use for analysis downstream).

```
# Check for duplicate data points - cells with exact same expression of
# variable genes
vg.dups <- duplicated(as.data.frame(as.matrix(t(suburd@logupx.data[suburd@var.genes,
])))
if (length(which(vg.dups)) > 0) {
  print(paste("Removing", length(which(vg.dups)), "cell(s) with duplicated variable gene expression."))
  not.dup.cells <- colnames(suburd@logupx.data)[!vg.dups]
  suburd <- urdSubset(suburd, not.dup.cells)
}
```

```
## [1] "Removing 6 cell(s) with duplicated variable gene expression."
```

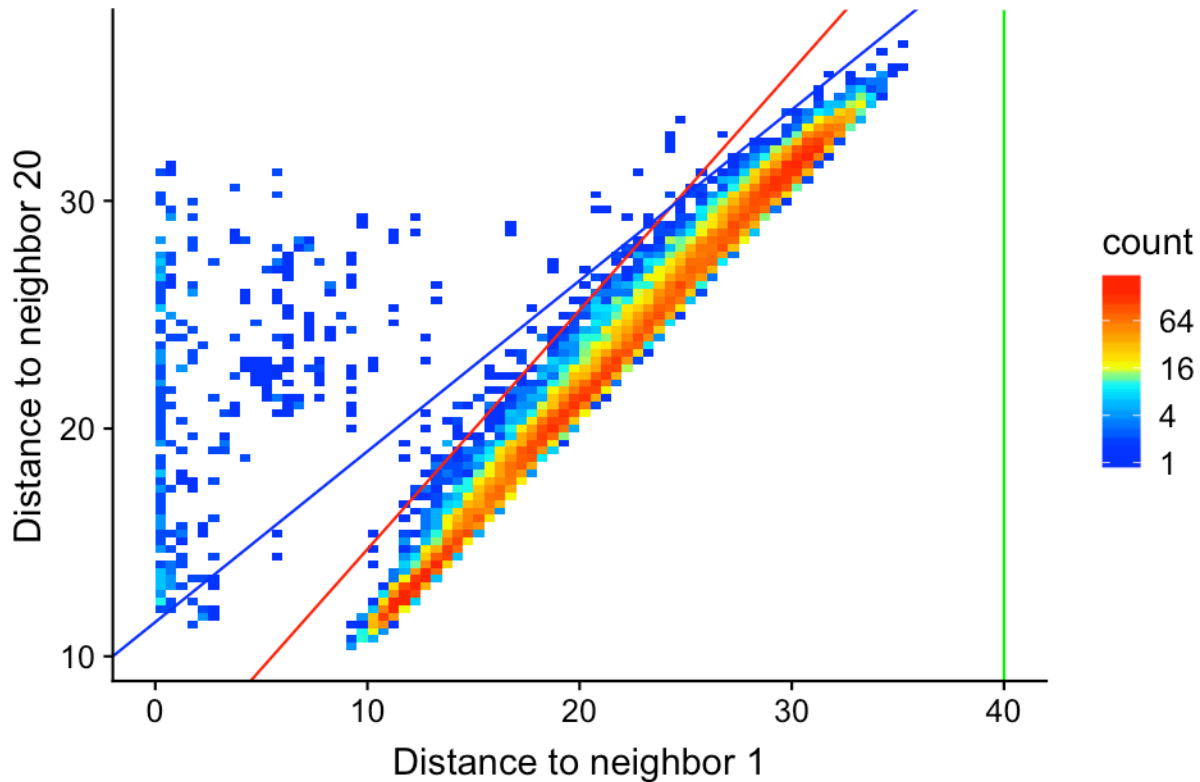
Calculate KNN graph and remove outliers

We then calculated a k-nearest neighbor graph and removed cells that had unusual distance to their nearest neighbor, or unusual distance to their 20th nearest neighbor (given their distance to their nearest neighbor). These sorts of outliers often cause problems or skew diffusion maps (used downstream).

```
# Calculate k-nn
suburd <- calcKNN(suburd)

# Check what the outliers are
outliers <- knnOutliers(suburd, nn.1 = 1, nn.2 = 20, x.max = 40, slope.r = 1.05,
  int.r = 4.2, slope.b = 0.75, int.b = 11.5, title = "Identifying Outliers by k-NN Distance.")
```

Identifying Outliers by k-NN Distance.



```
length(outliers)
```

```
## [1] 521
```

```
suburd <- urdSubset(suburd, cells.keep = setdiff(colnames(suburd@logupx.data),  
outliers))
```

Remove *kidins220a+* population

A cell cluster was observed in 15 dpf that was positive for expression of *kidins220a*, and *foxcg1b* (which is exclusive to the retina). However, no similar clusters were observed in other stages, suggesting that we did not recover the progenitors of this population, so we excluded it from the URD analysis.

```
suburd <- urdSubset(suburd, cells.keep = setdiff(colnames(suburd@logupx.data),  
cellsInCluster(suburd, "cluster", "12-15d-96")))
```

Remove cell type doublets

Add UMAP projection

While not strictly required, a UMAP projection can make it easier to assess the expression of NMF modules and whether thresholds for overlap are set correctly.

```
## add UMAP command

# Load pre-calculated UMAP
umap <- readRDS(paste0(base.path, "/umap/umap_retina.rds"))

# Add projection to URD object
suburd@tsne.y <- umap[colnames(suburd@logupx.data), ]
```

Load NMF results and import into object

NMF results were calculated by providing `suburd@logupx.data` to an external NMF pipeline written in Python. The output results are imported here, scaled, and added to the URD object.

```
# Load the NMF results
load(paste0(base.path, "/NMF/retina/result_tbls.Robj"))

# The results object contains NMF runs for several K values. k=45 was
# chosen for this tissue, so this extracts the results for that
# particular parameter
k.use <- "45"
nmf.cells <- result_obj[[paste0("K=", k.use)]][[1]]$C
rownames(nmf.cells) <- paste0("nmf", 1:nrow(nmf.cells))
colnames(nmf.cells) <- gsub("\\.", "-", colnames(nmf.cells))
nmf.genes <- result_obj[[paste0("K=", k.use)]][[1]]$G
colnames(nmf.genes) <- paste0("nmf", 1:nrow(nmf.cells))

# Some stages were subsampled in the original object and accidentally
# cropped out cells that had scGESTALT barcodes. Those were added back
# in, and their expression was decomposed with original NMF gene matrix
# to give an additional NMF cell matrix for those cells.
new.nmf.c <- read.csv(paste0(base.path, "/NMF/retina/retina_new_nmfC_k45.csv"),
  row.names = 1)
rownames(new.nmf.c) <- paste0("nmf", 1:nrow(new.nmf.c))
colnames(new.nmf.c) <- gsub("\\.", "-", colnames(new.nmf.c))

# Combine old and new NMF results
nmf.cells <- cbind(nmf.cells, new.nmf.c)

# Trim NMF results to match cells in current object
nmf.cells <- nmf.cells[, colnames(suburd@logupx.data)]

# Scale NMF results 0-1
nmf.cells.scaled <- sweep(nmf.cells, 1, apply(nmf.cells, 1, max), "/")

# Add scaled NMF results to the URD object
suburd@nmf.c1 <- as(t(as.matrix(nmf.cells.scaled)), "dgCMatrix")
```

Select cell-type specific modules

Several NMF modules will be poor markers of cell types — these are often modules driven mostly by the expression of 1-2 genes (where the gene loading of the first gene is much greater than that

of the fourth gene, for instance), or modules that don't exhibit any restriction in a tSNE or UMAP projection.

```
# Plot size parameters
plot.height = 8
plot.width = 8
dpi = 150

# Plot every module to determine which exhibit cell-type specificity
# This saves directly to the hard drive: two example plots are shown
# below.

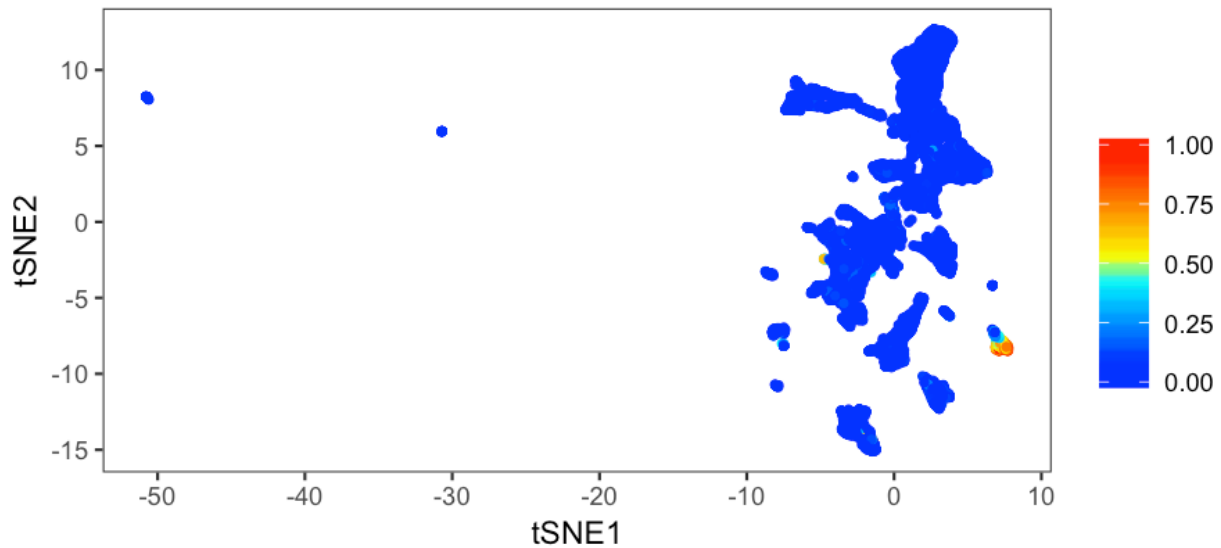
# for (n in colnames(suburd@nmf.c1)) { png(paste0(path, '/doublets/',
# subset, '-plots/', n, '.png'), width=dpi*plot.width,
# height=dpi*plot.height) plot(plotDim(suburd, n)) dev.off() }

gridExtra::grid.arrange(grobs = list(plotDim(suburd, "nmf4", plot.title = "nmf4: strong cell-type restriction"),
  plotDim(suburd, "nmf2", plot.title = "nmf2: poor cell-type restriction")),
  ncol = 1)

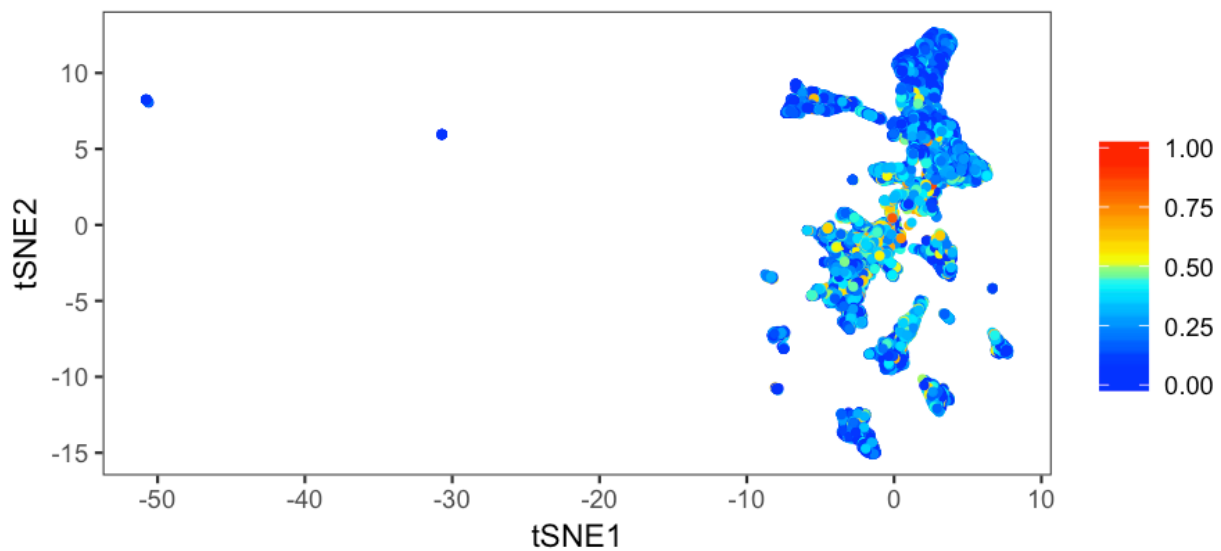
## Warning: Removed 520 rows containing missing values (geom_point).

## Warning: Removed 520 rows containing missing values (geom_point).
```


nmf4: strong cell-type restriction



nmf2: poor cell-type restriction



```
# Module Gene 1 : Gene 4 Ratios
top.genes <- result_obj[[paste0("K=", k.use)]][[1]]$top30genes
top.weights <- top.genes[, grep("Weights", colnames(top.genes), value = T)]
colnames(top.weights) <- paste0("nmf", 1:nrow(nmf.cells))
top.weights.ratio <- top.weights[1, ]/top.weights[4, ]

# Which modules exhibit cell-type restriction?
modules.ok.ratio <- names(top.weights.ratio)[which(top.weights.ratio <
5)]
restricted.modules <- paste0("nmf", c(4:5, 8:13, 15:18, 21:29, 31:35, 37:39,
41:44))
good.modules <- intersect(modules.ok.ratio, restricted.modules)
```

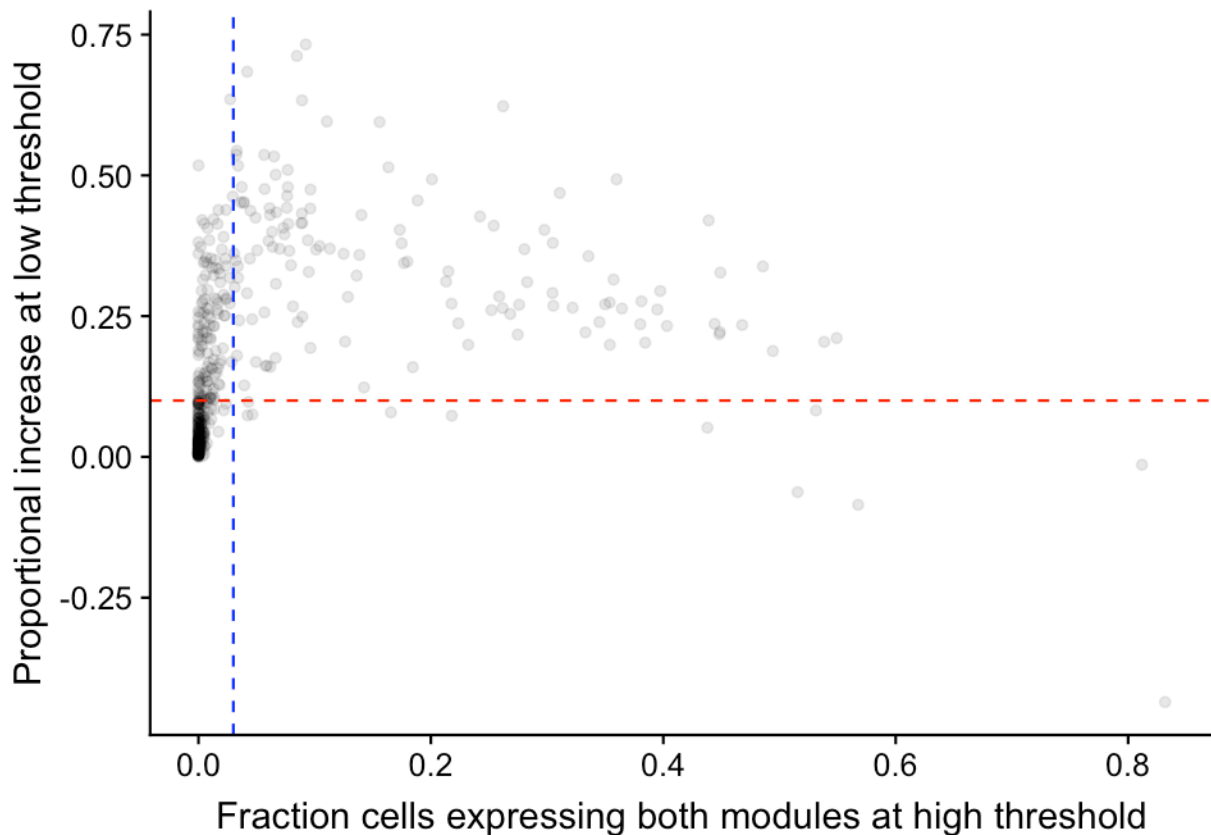
Determine which module pairs to use for doublet removal

We consider NMF modules pairwise and only use those pairs that don't are non-overlapping in the data. (In other words, NMF modules that are mutually exclusive in the majority of the data.) Here, we determine thresholds for selecting those module pairs.

```
# Determine overlaps between module pairs
nmf.doublet.combos <- NMFDoubletsDefineModules(suburd, modules.use = good.modules,
  module.thresh.high = 0.4, module.thresh.low = 0.15)

# Determine thresholds for NMF modules
frac.overlap.max = 0.03
frac.overlap.diff.max = 0.1
module.expressed.thresh = 0.33

# Determine which module pairs to use for doublets
NMFDoubletsPlotModuleThresholds(nmf.doublet.combos, frac.overlap.max = frac.overlap.max,
  frac.overlap.diff.max = frac.overlap.diff.max)
```



```
# These commands save plots directly to the hard-drive.

# Make plots to see how your thresholds are
NMFDoubletsPlotModuleCombos(suburd, path = paste0(path, "/doublets/"), subset,
  "-doublet-combos/"), module.combos = nmf.doublet.combos, module.expressed.thresh = module.expressed.thresh,
  frac.overlap.max = frac.overlap.max, frac.overlap.diff.max = frac.overlap.diff.max,
  boundary = "pass", sort = "near", n.plots = 25)
NMFDoubletsPlotModuleCombos(suburd, path = paste0(path, "/doublets/"), subset,
```

```

"--ok-combos/"), module.combos = nmf.doublet.combos, module.expressed.thresh = module.expressed.thresh,
frac.overlap.max = frac.overlap.max, frac.overlap.diff.max = frac.overlap.diff.max,
boundary = "discarded", sort = "near", n.plots = 25)

```

```

# Define doublet cells

```

```

nmf.doublets <- NMFDoubletsDetermineCells(suburd, nmf.doublet.combos, module.expressed.thresh = module.expressed.thresh,
frac.overlap.max = frac.overlap.max, frac.overlap.diff.max = frac.overlap.diff.max) # 376 cells / 19837 cells

```

```

# Plot doublet cells on the UMAP

```

```

suburd <- groupFromCells(suburd, "nmf.doublets", cells = nmf.doublets)
plot(plotDimHighlight(suburd, clustering = "nmf.doublets", cluster = "TRUE",
plot.title = paste0("NMF doublets: ", length(nmf.doublets), " cells"),
point.size = 2, highlight.color = "blue"))

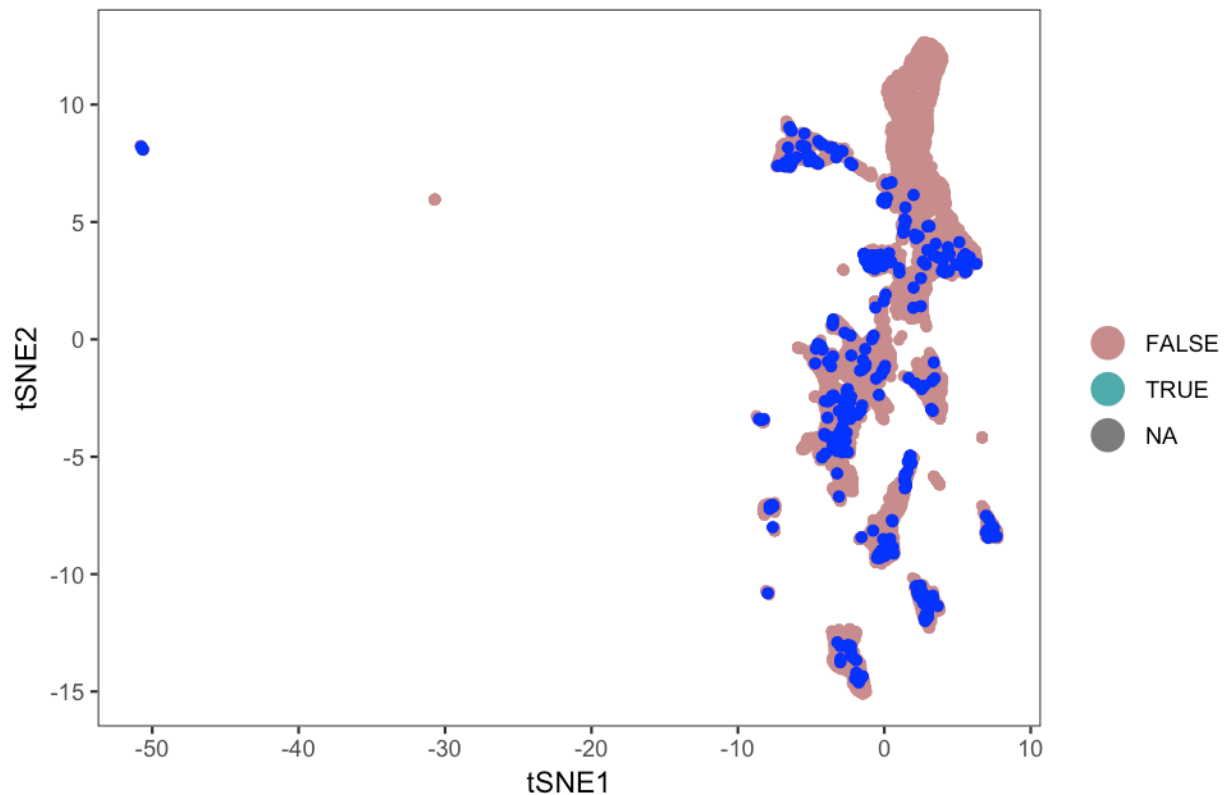
```

```

## Warning: Removed 520 rows containing missing values (geom_point).

```

NMF doublets: 376 cells (Highlight TRUE)



```

# Crop object to exclude doublets

```

```

suburd.cropped <- urdSubset(suburd, cells.keep = setdiff(colnames(suburd@logupx.data),
nmf.doublets))

```

And then save the completed object for use downstream in building a tree using URD.

```

saveRDS(suburd.cropped, file = paste0(base.path, "/obj/URD_retina_ND.rds"))

```

Retina: 2 - URD tree

Jeff Farrell

9/07/2019

Contents

Load data	1
Processed on the cluster	1
Calculate diffusion map and pseudotime	2
Calculate diffusion map	2
Calculate pseudotime	4
Calculate biased transition matrix	9
Perform biased random walks	9
Determine tips	9
Perform the biased random walks	11
Process the random walks	12
Build the URD tree	12
Save the URD tree	13

Load data

```
suppressPackageStartupMessages(library(URD))
suppressPackageStartupMessages(library(Seurat))

base.path <- "~/urd-cluster-bushra/"

# Load processed URD object
object <- readRDS(paste0(base.path, "obj/URD_retina_ND.rds"))
```

Processed on the cluster

Most of the following steps were run on a computing cluster. These individual tissue subsets can be run on a modern, well-equipped laptop. The use of a computing cluster allows multiple parameter choices to be tried in parallel, and also allows further parallelization of the random walk procedure, speeding it up. Below, we should the commands that one would run on their laptop, and then generally load the pre-processed results from the cluster that were used in the paper. If you want to parallelize your own processing on a compute cluster, the scripts we used will be available at <http://github.com/farrellja/URD/cluster/>

Calculate diffusion map and pseudotime

These two steps are run in the cluster script URD-DM-PT.R.

Calculate diffusion map

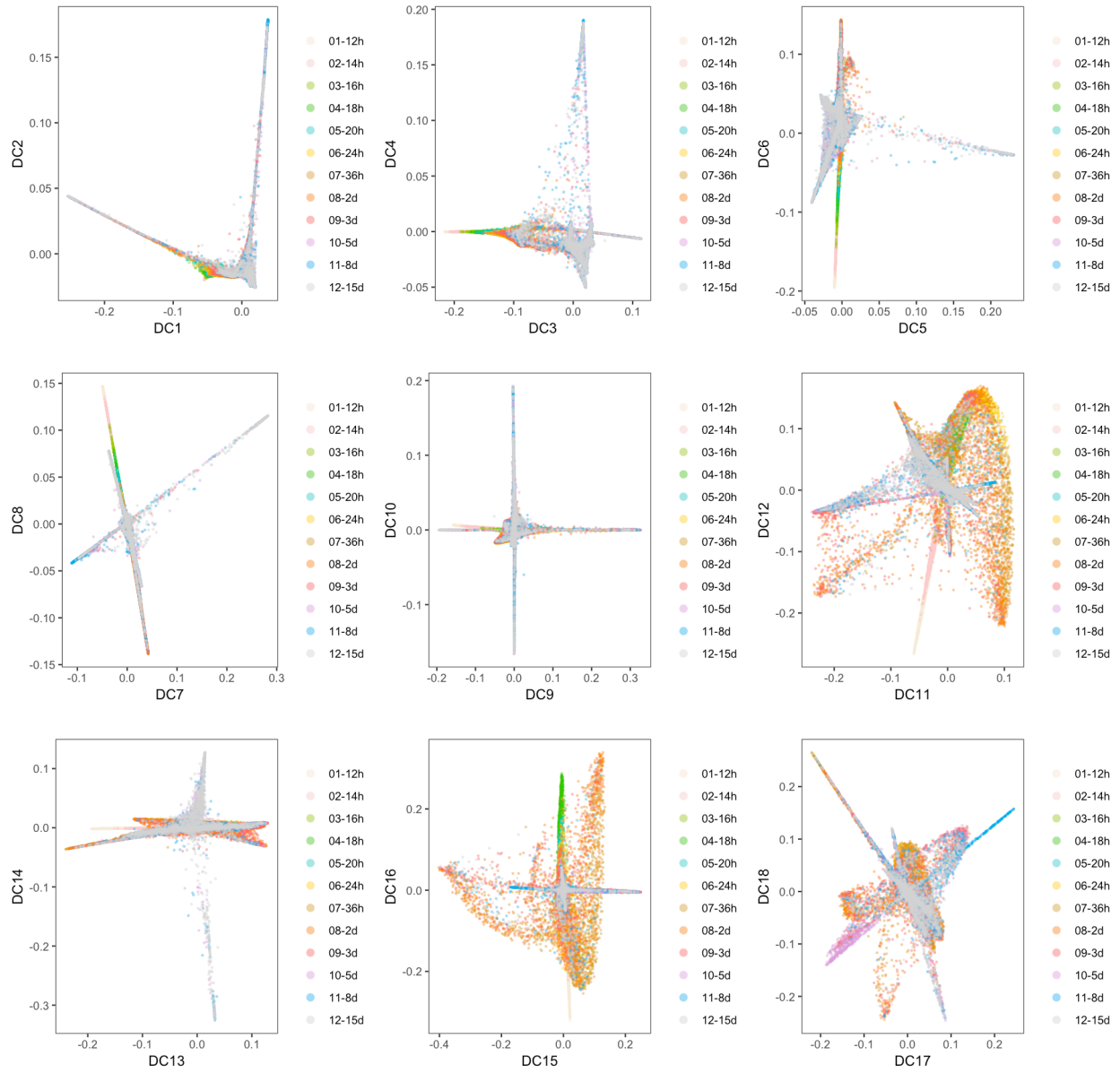
```
# To run locally: Calculate a diffusion map projection
object <- calcDM(object, knn = 140, sigma.use = 14)

# Or: Load a pre-computed diffusion map projection
dm <- readRDS(paste0(base.path, "dm/dm_retinawokND_knn-140_sigma-14.rds"))
object <- importDM(object, dm)

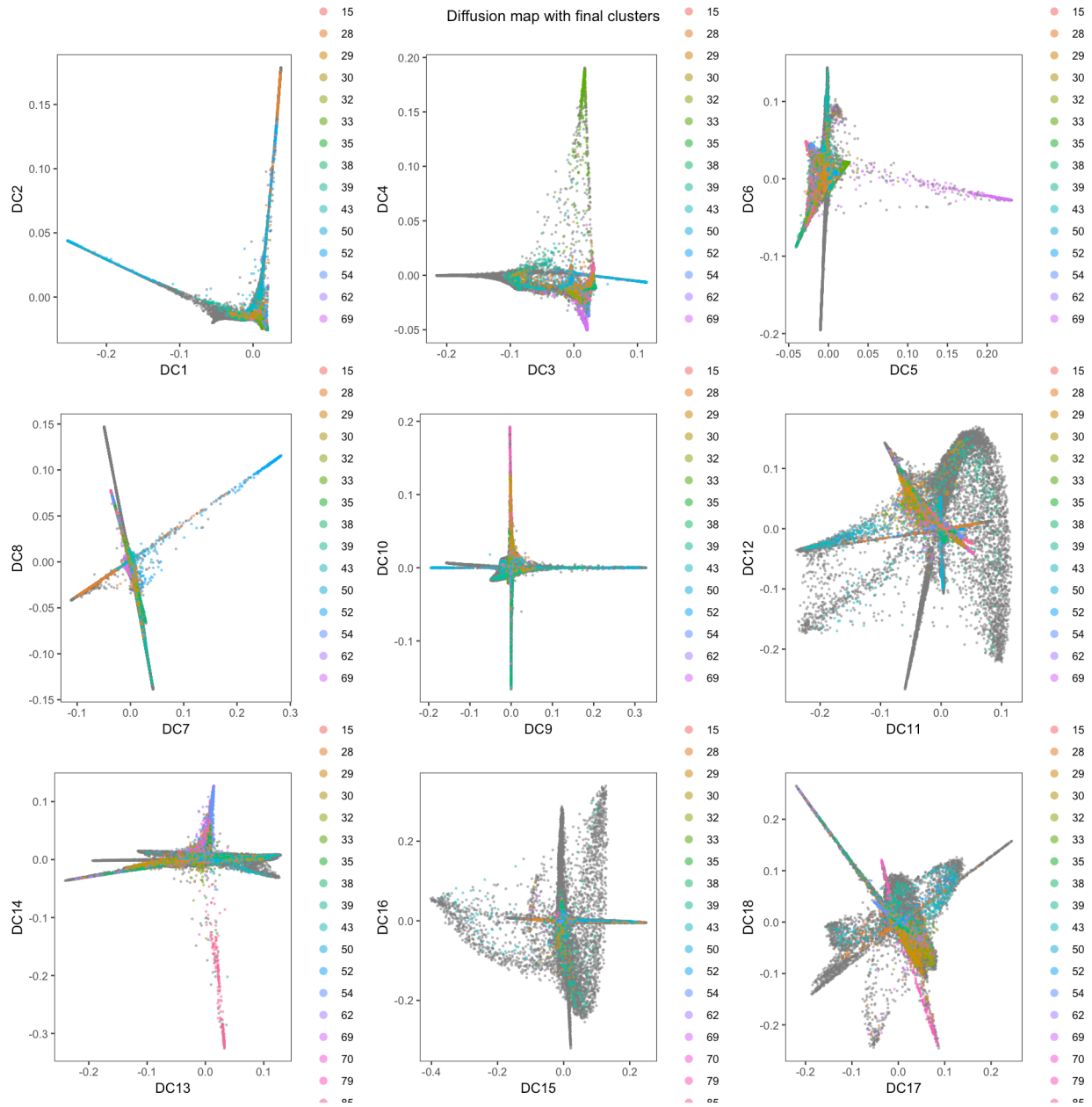
# Plot diffusion maps
stage.colors <- c("antiquewhite", "#FFCCCC", "#99CC00", "#33CC00", "cyan3",
  "gold", "goldenrod", "darkorange", "indianred1", "plum", "deepskyblue2",
  "lightgrey")

# Plot by stage
plotDimArray(object = object, reduction.use = "dm", dims.to.plot = 1:18,
  label = "stage", plot.title = "", outer.title = "Diffusion map labeled by Stage",
  legend = T, alpha = 0.45, discrete.colors = stage.colors)
```

Diffusion map labeled by Stage



```
# Plot with final cell types labeled
object@group.ids$final.cluster <- NA
object@group.ids[cellsInCluster(object, "stage", "12-15d"), "final.cluster"] <- object@group.ids[cellsInCluster(
  "stage", "12-15d"), "res.5"]
plotDimArray(object = object, reduction.use = "dm", dims.to.plot = 1:18,
  label = "final.cluster", plot.title = "", outer.title = "Diffusion map with final clusters",
  legend = T, alpha = 0.6)
```



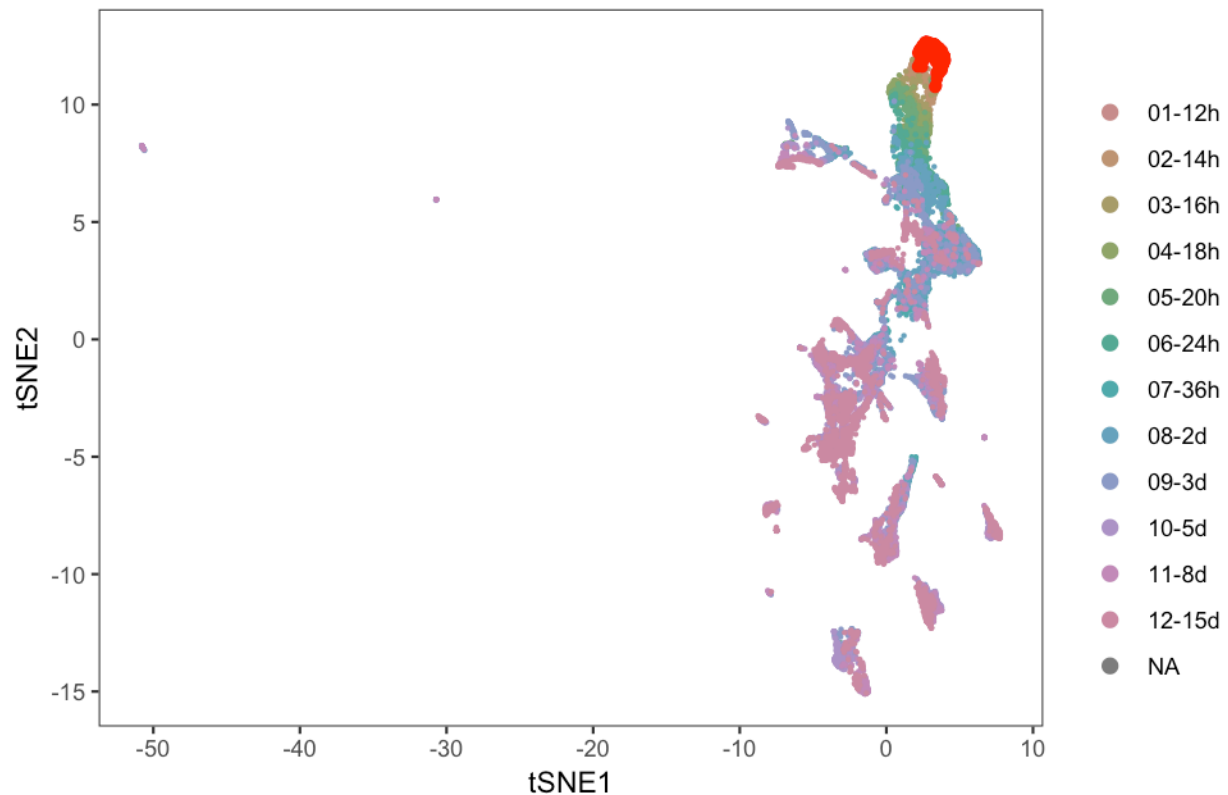
Calculate pseudotime

URD requires a starting point or 'root' for determining pseudotime. Here, we used all cells from the first timepoint (i.e. 12 hpf) as the root.

```
# Here, we used all cells from the first timepoint (i.e. 12 hours) as
# the root.
root.cells <- cellsInCluster(object, "stage", "01-12h")
plotDimHighlight(object, "stage", "01-12h", plot.title = "Root is 12 hpf cells")
```

```
## Warning: Removed 500 rows containing missing values (geom_point).
```

Root is 12 hpf cells (Highlight 01-12h)



```
# To run locally: Run graph-search simulations to determine pseudotime  
flood.result <- floodPseudotime(object, root.cells = root.cells, n = 100,  
  minimum.cells.flooded = 2, verbose = T)
```

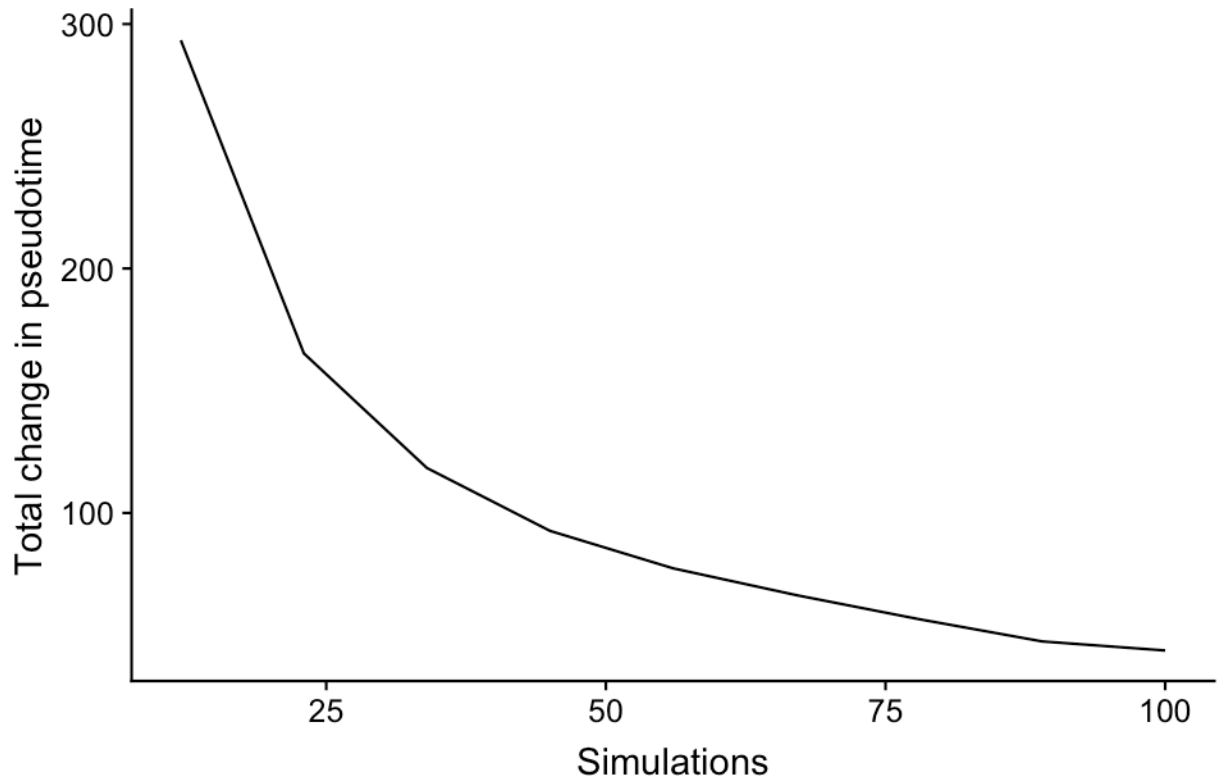
```
# Or load a pre-computed graph-search simulation result  
flood.result <- readRDS(paste0(base.path, "flood/flood_retinanewnokND_knn-140_sigma-14.rds"))
```

```
# Process the graph-search simulations to determine the pseudotime of  
# each cell  
object <- floodPseudotimeProcess(object, flood.result, floods.name = "pseudotime",  
  max.frac.NA = 0.4, pseudotime.fun = mean, stability.div = 10)
```

```
# If enough simulations have been run, then as additional simulations  
# are added, the overall change in pseudotime of cells should reach an  
# asymptote. If it does not, then floodPseudotime should be run with a  
# higher n.
```

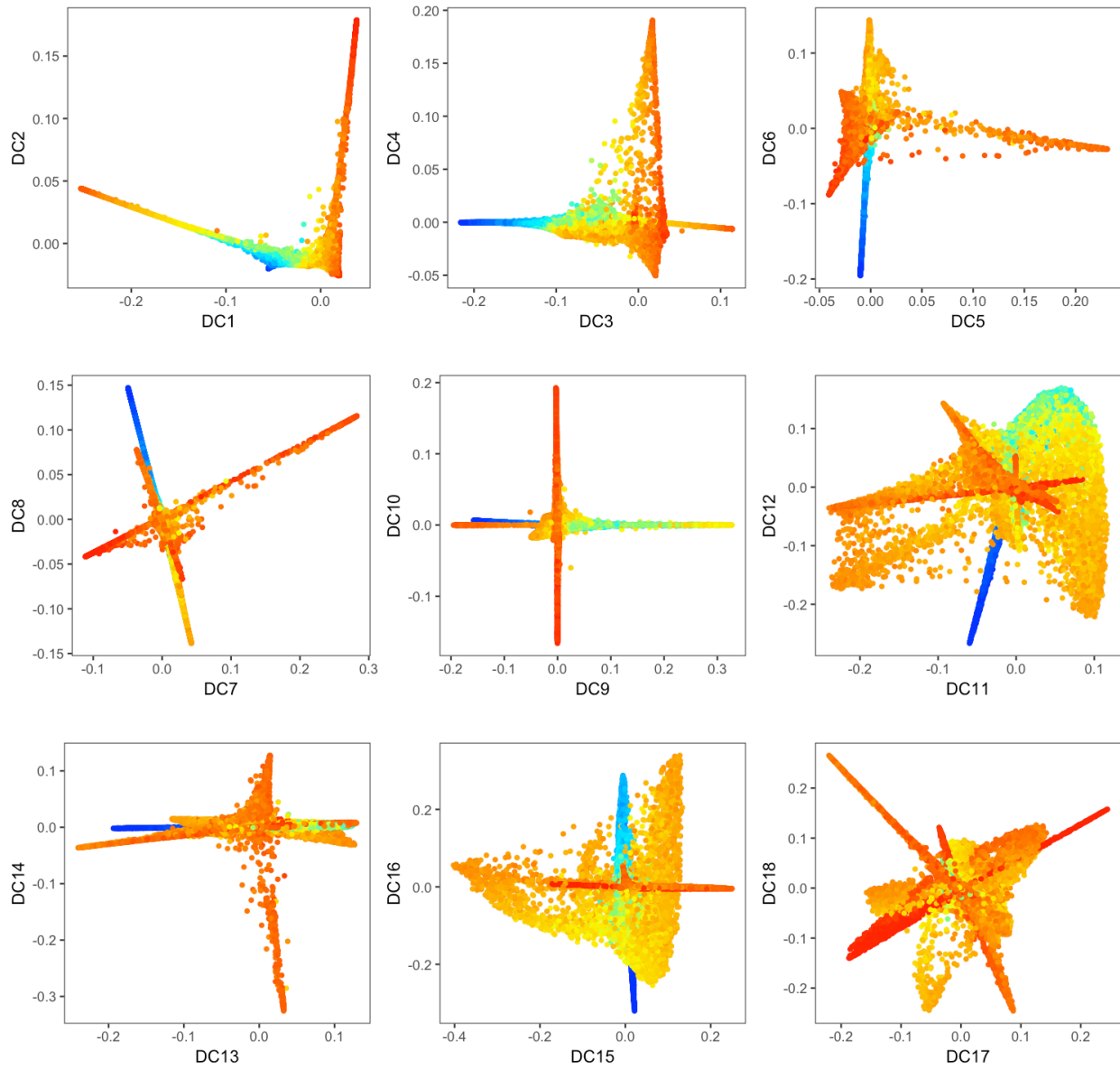
```
pseudotimePlotStabilityOverall(object)
```


Overall Pseudotime Stability



```
plotDimArray(object = object, reduction.use = "dm", dims.to.plot = 1:18,  
  label = "pseudotime", plot.title = "", outer.title = "Diffusion Map labeled by pseudotime",  
  legend = F, alpha = 0.4)
```

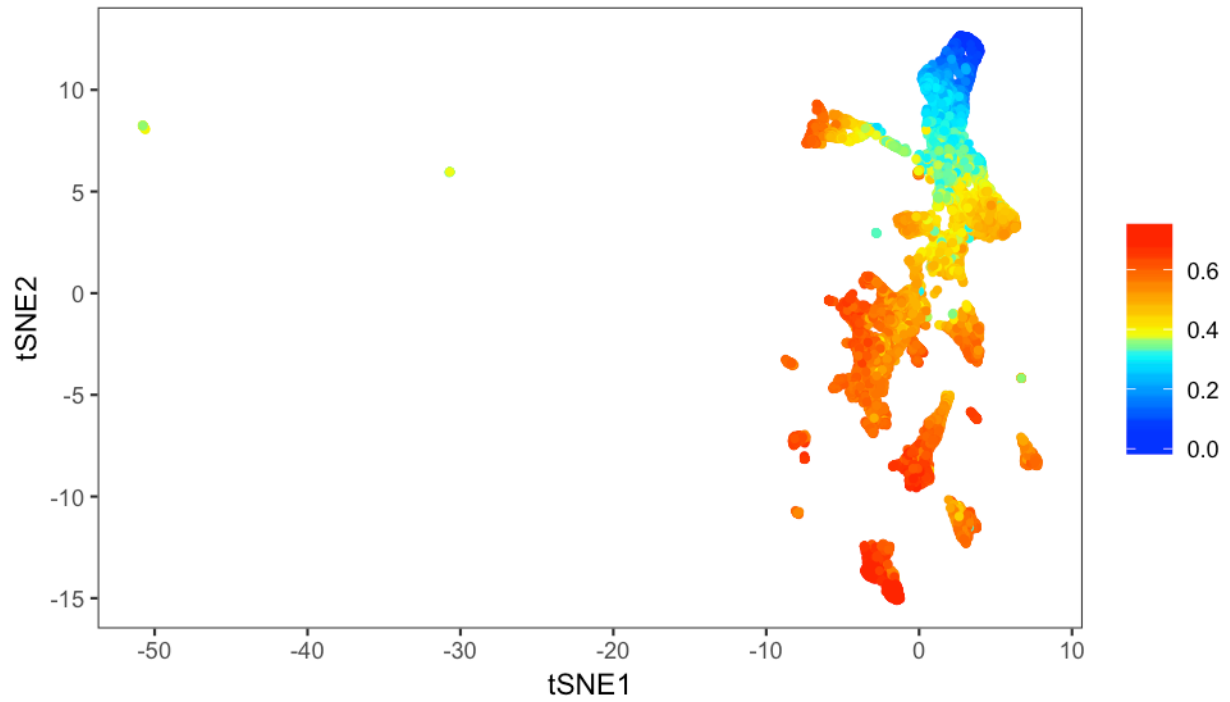
Diffusion Map labeled by pseudotime



```
plotDim(object, "pseudotime", plot.title = "UMAP projection colored by pseudotime")
```

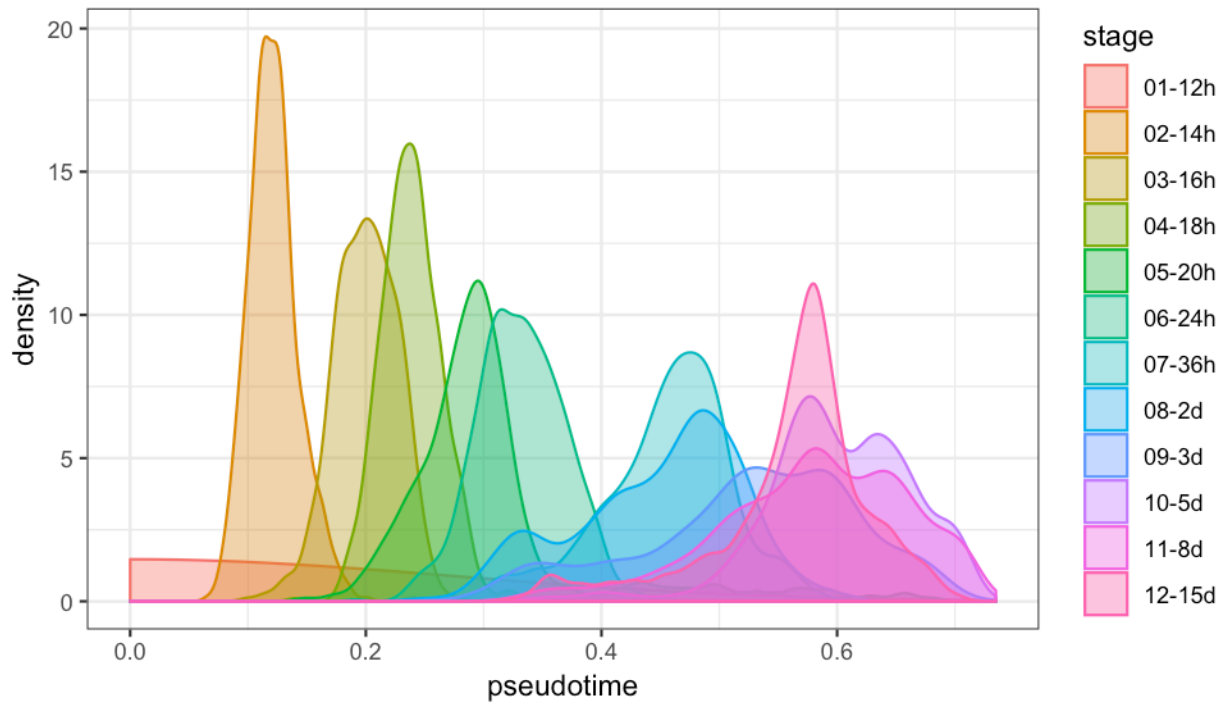
```
## Warning: Removed 500 rows containing missing values (geom_point).
```

UMAP projection colored by pseudotime



```
plotDists(object, "pseudotime", "stage", plot.title = "Pseudotime by stage")
```

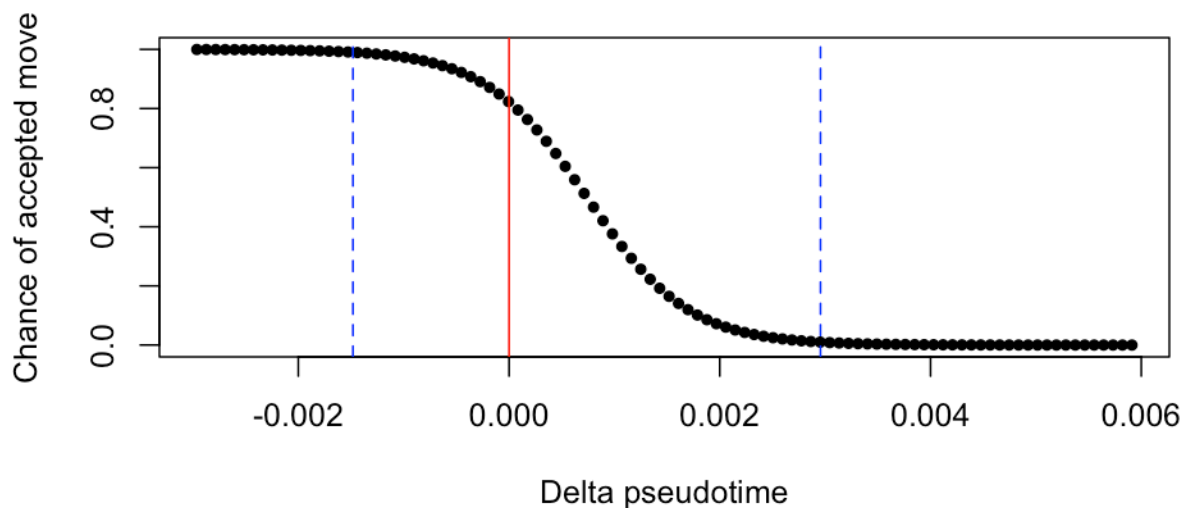
Pseudotime by stage



Calculate biased transition matrix

In order to perform biased random walks, we must first bias the transition matrix to ensure that walks proceed towards the root and do not turn into other differentiated cell types. This is performed in the cluster script `URD-TM.R`.

```
# Calculate parameters for biasing the transition matrix.
diffusion.logistic <- pseudotimeDetermineLogistic(object, "pseudotime",
  optimal.cells.forward = 40, max.cells.back = 80, pseudotime.direction = "<",
  do.plot = T, print.values = T)
```



```
## [1] "Mean pseudotime back (~80 cells) 0.00295569803650678"
## [1] "Chance of accepted move to equal pseudotime is 0.822024945232085"
## [1] "Mean pseudotime forward (~40 cells) -0.00148141113789574"
```

```
# Calculate the biased matrix.
biased.tm <- pseudotimeWeightTransitionMatrix(object, pseudotime = "pseudotime",
  logistic.params = diffusion.logistic, pseudotime.direction = "<")
```

Perform biased random walks

Then, we perform biased walks starting from each tip. Visited cells are inferred to lie along the trajectory that connects the root to each cell type. This is performed in the cluster script `URD-Walk.R`.

Determine tips

We used clusters from 15 dpf as the tips for performing biased random walks. Here we define the cells belonging to each of those clusters.

```

# All clusters at 15 days
clusters.15day <- unique(object@group.ids[grep("15d", object@group.ids$stage),
  "res.5"])
# All cells at 15 days
cells.15day <- rownames(object@group.ids)[grep("15d", object@group.ids$stage)]
# Cell lists of each cluster at 15dpf
cells.15dpf.clusters <- lapply(clusters.15day, function(clust) intersect(cells.15day,
  cellsInCluster(object, "res.5", clust)))
names(cells.15dpf.clusters) <- paste0("15d-", clusters.15day)

```

We also load a .csv file that contains information about the tips. It has four columns:

- id: Cluster ID for the tip
- use: Whether this cluster should be used when building the tree
- name: The name for this tip, which will be used on 2D plots
- short.name: The 'short' name for this tip, which would be used on 3D plots (though we did not use that feature in this study).

```

# Load CSV
tip.names <- read.csv(paste0(base.path, "tips/tip_names_retinanewnokND.csv"),
  header = F, stringsAsFactors = F, colClasses = c("character", "logical",
  "character", "character"))

# Name columns and rows
names(tip.names) <- c("id", "use", "name", "short.name")
rownames(tip.names) <- gsub("_", "-", tip.names$id)

# Sort alphabetically
tip.names <- tip.names[order(rownames(tip.names)), ]

```

These are the tips that were considered during the construction of the retina URD tree (some were excluded during tree construction later in the `buildTree` command).

```

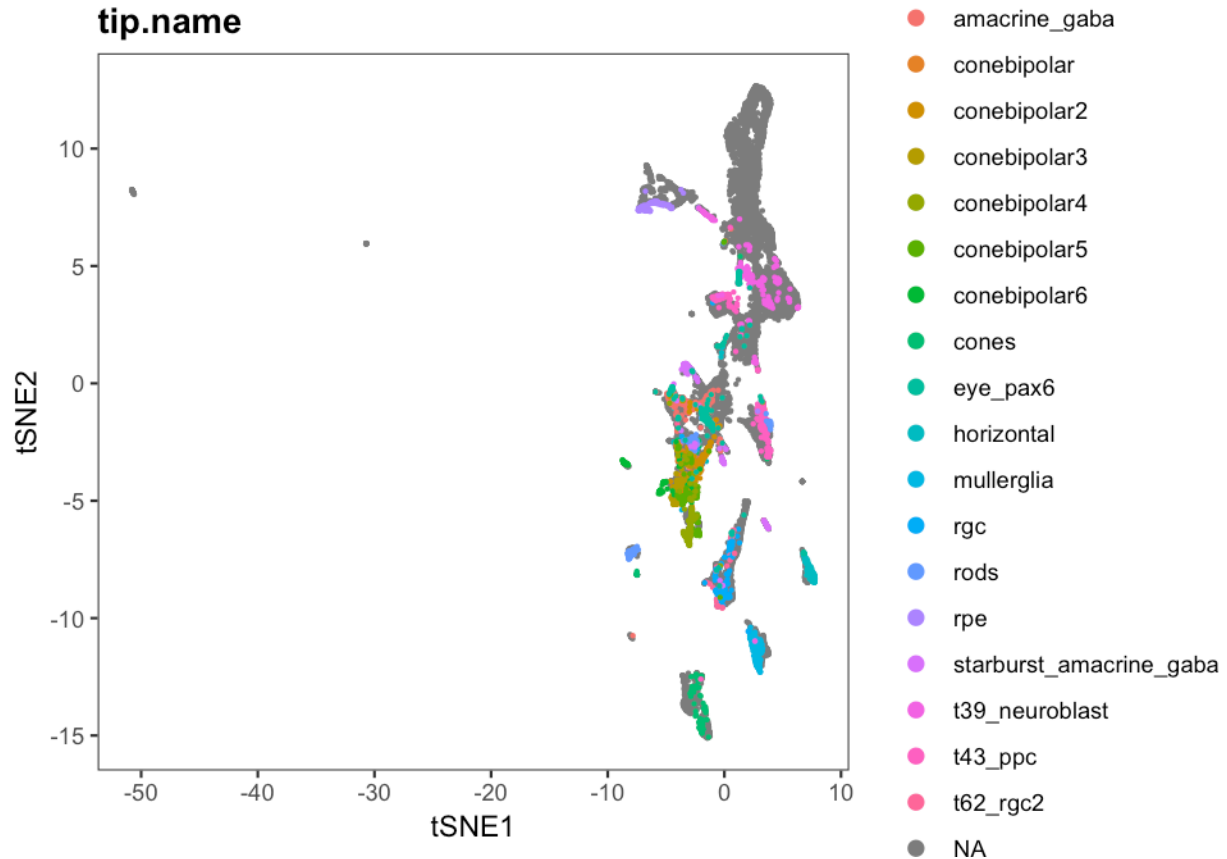
# Define a 'tips' clustering
object@group.ids$tip <- NA
object@group.ids$tip.id <- NA
object@group.ids$tip.name <- NA

# If the tip will be used in the tree, define its cells in the
# clustering
for (i in 1:nrow(tip.names)) {
  tip.cells <- cells.15dpf.clusters[[rownames(tip.names)[i]]]
  object@group.ids[tip.cells, "tip"] <- as.character(i)
  object@group.ids[tip.cells, "tip.id"] <- rownames(tip.names)[i]
  object@group.ids[tip.cells, "tip.name"] <- as.character(tip.names[i,
    "name"])
}

# Plot the tips
plotDim(object, "tip.name")

```

```
## Warning: Removed 500 rows containing missing values (geom_point).
```



Perform the biased random walks

Biased random walks then need to be run starting from each tip. This can be performed on a laptop, but is an ideal candidate for parallelization on a cluster. (The walks from each tip can be run as a separate job.)

```
## IF RUNNING LOCALLY

# Loop through each cluster
walks <- lapply(rownames(tip.names), function(c) {
  # Exclude any tip cells that for whatever reason didn't end up in the
  # biased TM (e.g. maybe not assigned a pseudotime).
  tip.cells <- intersect(cells.15dpf.clusters[[c]], rownames(biased.tm))
  # Perform the random walk simulation
  this.walk <- simulateRandomWalk(start.cells = tip.cells, transition.matrix = biased.tm,
    end.cells = root.cells, n = 50000, end.visits = 1, verbose.freq = 1000,
    max.steps = 5000)
  return(this.walk)
})
names(walks) <- rownames(tip.names)

# Alternatively, this loop is automated by the function
# simulateRandomWalksFromTips
```

Alternatively, a set of pre-calculated walks can be loaded. Since the walks are a simulation (and

therefore not deterministic), this is particularly crucial for reproducing results.

```
## IF LOADING PRE-CALCULATED WALKS

# Get list of files in the walks directory
walks.files <- list.files(paste0(base.path, "/walks/retinawalkND/"),
  pattern = ".rds")

# Load the walks previously performed for each cluster
walks <- lapply(rownames(tip.names), function(c) {
  walk.file <- grep(pattern = paste0("_tip-", c, "_"), x = walks.files,
    value = T)[1]
  return(readRDS(paste0(base.path, "/walks/retinawalkND/", walk.file)))
})
names(walks) <- rownames(tip.names)
```

Process the random walks

The walks are then converted to visitation frequency by importing them into the URD object.

```
for (i in 1:nrow(tip.names)) {
  # Load the individual walk visitation frequencies into the object
  object <- processRandomWalks(object, walks = walks[[i]], walks.name = i,
    n.subsample = 1, verbose = F)
}
```

Build the URD tree

Then, a branching tree is constructed, by joining trajectories in an agglomerative fashion when cells are highly visited by walks from multiple tips. The following steps were performed in the cluster script URD-Tree.R.

```
# Tree building is destructive, so create a copy of the object
object.tree <- object
```

```
# Load tip cells
object.tree <- loadTipCells(object.tree, "tip")
```

```
# Determine tips to use
tips.to.use <- which(tip.names$use)
```

```
# Build the tree
object.tree <- buildTree(object.tree, pseudotime = "pseudotime", divergence.method = "preference",
  cells.per.pseudotime.bin = 40, bins.per.pseudotime.window = 5, save.all.breakpoint.info = T,
  p.thresh = 0.01, verbose = F, tips.use = as.character(tips.to.use))
```

```
# Name the tips of the tree
object.tree <- nameSegments(object.tree, segments = tips.to.use, segment.names = as.character(tip.names[tips.to.use,
  "name"]), short.names = as.character(tip.names[tips.to.use, "short.name"]))
```

```
plotTree(object.tree, "stage", discrete.colors = stage.colors, label.segments = T)
```


Retina: 3 - URD Cascades and Figures

Jeff Farrell

10/08/2019, updated 07/30/2020

Contents

Load data	2
Plot gene expression on the tree	2
Plot tree by stage	2
Plot tree with gene expression: main figures	3
Plot tree with gene expression: supplemental figures	4
Determine genes enriched in trajectories to particular cell types	5
Comparison between major cell types	5
AUCPR along tree	7
Functions for curating differential expression results	7
threshold.tree.markers	7
threshold.clade.markers	8
divide.branches	8
Functions for heatmap generation	9
Color scale	9
determine.timing	9
filter.heatmap.genes	10
Heatmaps of gene cascades	10
Photoreceptors	11
Prepare cascade	11
Generate heatmap: all genes	12
Generate heatmap: main figure	15
Amacrine cells	17
Prepare cascade	17
Generate heatmap: all genes	18
Retinal ganglion cells	20
Prepare cascade	20
Generate heatmap: all genes	20
Generate heatmap: main figure	22
Horizontal Cells	24
Prepare cascade	24
Generate heatmap: all genes	24
Muller Glia	27
Prepare cascade	27
Generate heatmap: all genes	27
Retinal Pigmented Epithelium	29

Prepare cascade	29
Generate heatmap: all genes	29
Continuous differentiation	31
RGC cells	31
Progenitor cells	31
Progenitors over time	32
Identify populations	32
Differential expression between neural progenitor populations	36
boot.fc	36
Empirical p-value	37
Tissue-specific changes	38
Limit to well-expressed	39
Result	39
Preservation of embryonic molecular profiles in larval progenitors	40
Identify populations to compare	40
Determine proportion of cells in each state	42

Load data

```
# Load URD
library(URD)

## Loading required package: ggplot2

## Loading required package: Matrix

## Registered S3 method overwritten by 'xts':
##   method      from
##   as.zoo.xts  zoo

# Basic location
base.path <- "~/Documents/R sessions/urd-cluster-bushra/"

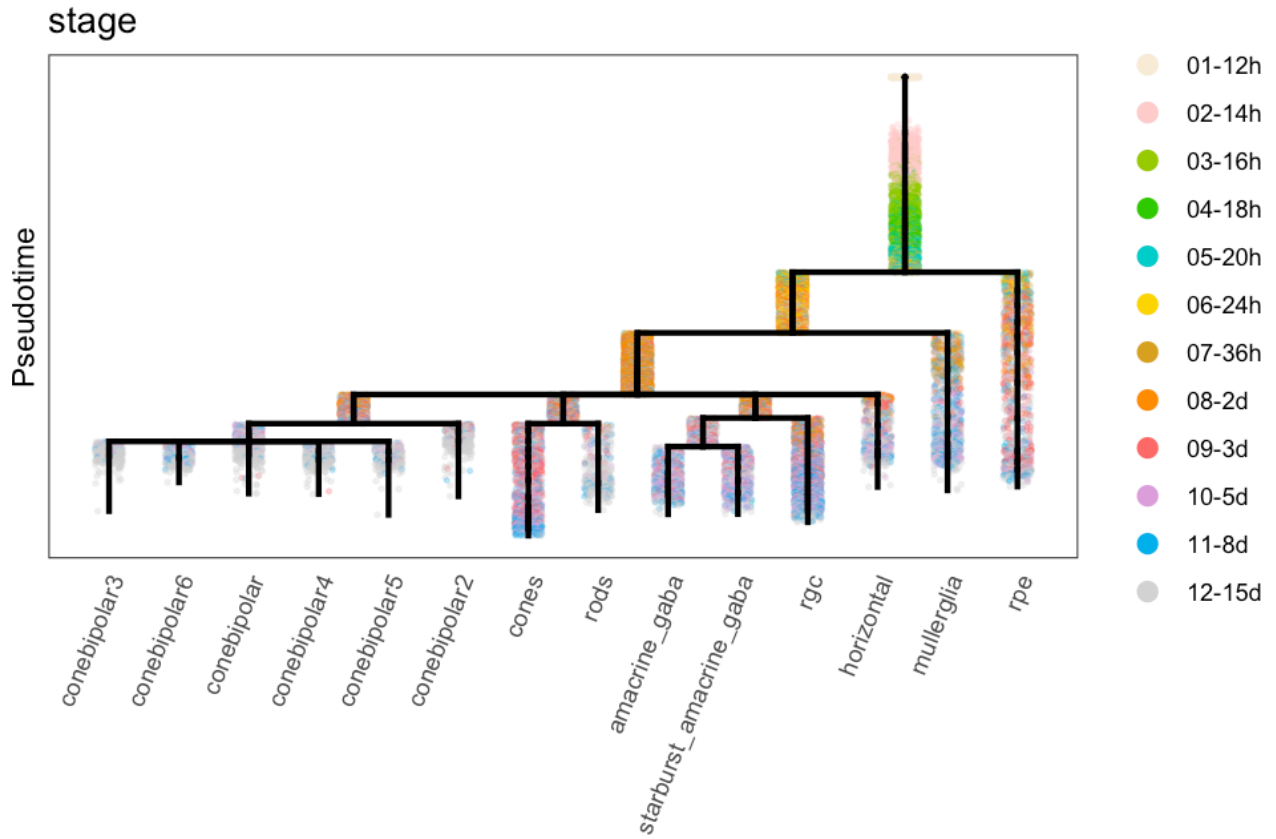
# Load completed retina tree object
obj.path <- paste0(base.path, "tree/retinaneokND/tree-retinaneokND_knn-140_sigma-14_40F-80B_NO-15d-29-15d-39-")
obj <- readRDS(obj.path)
```

Plot gene expression on the tree

Plot tree by stage

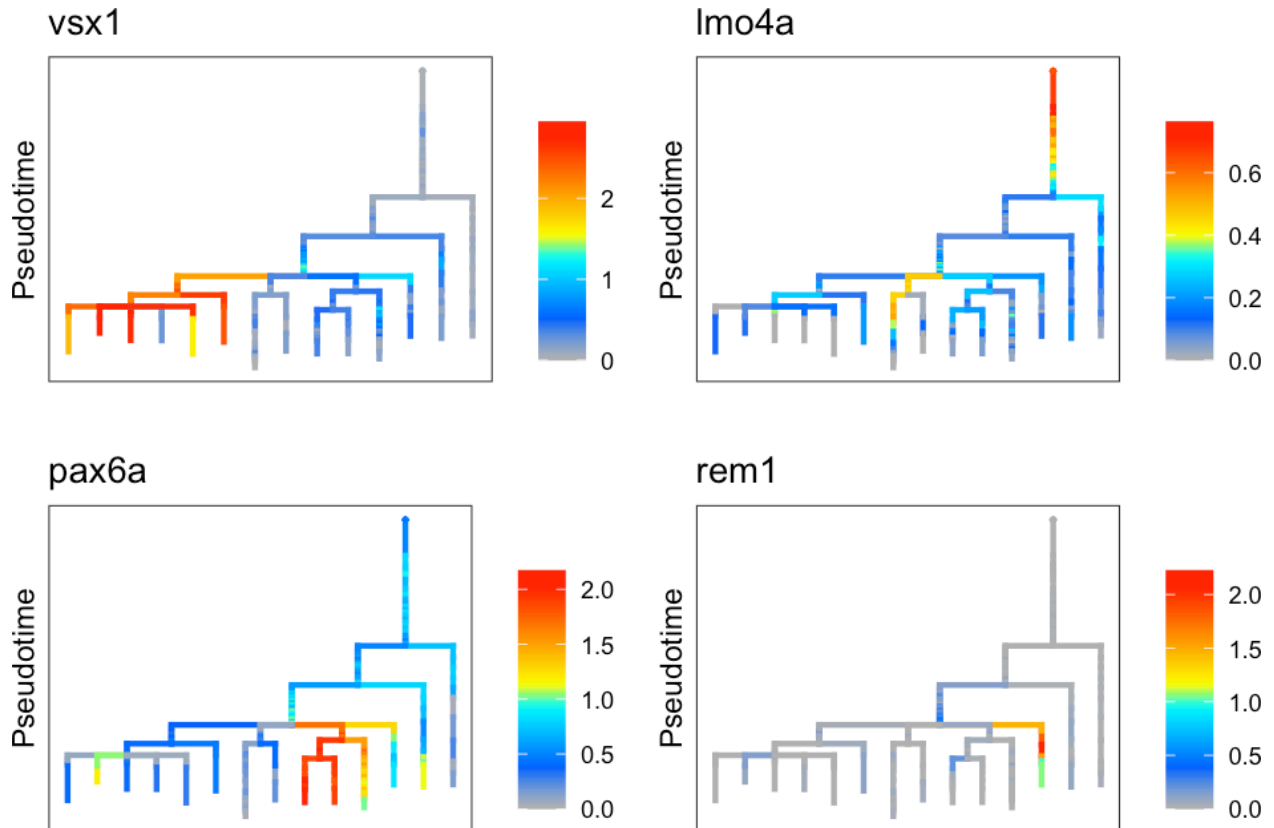
```
stage.colors <- c("antiquewhite", "#FFCCCC", "#99CC00", "#33CC00", "cyan3", "gold",
  "goldenrod", "darkorange", "indianred1", "plum", "deepskyblue2", "lightgrey")

plotTree(obj, "stage", label.type = "group", discrete.colors = stage.colors)
```



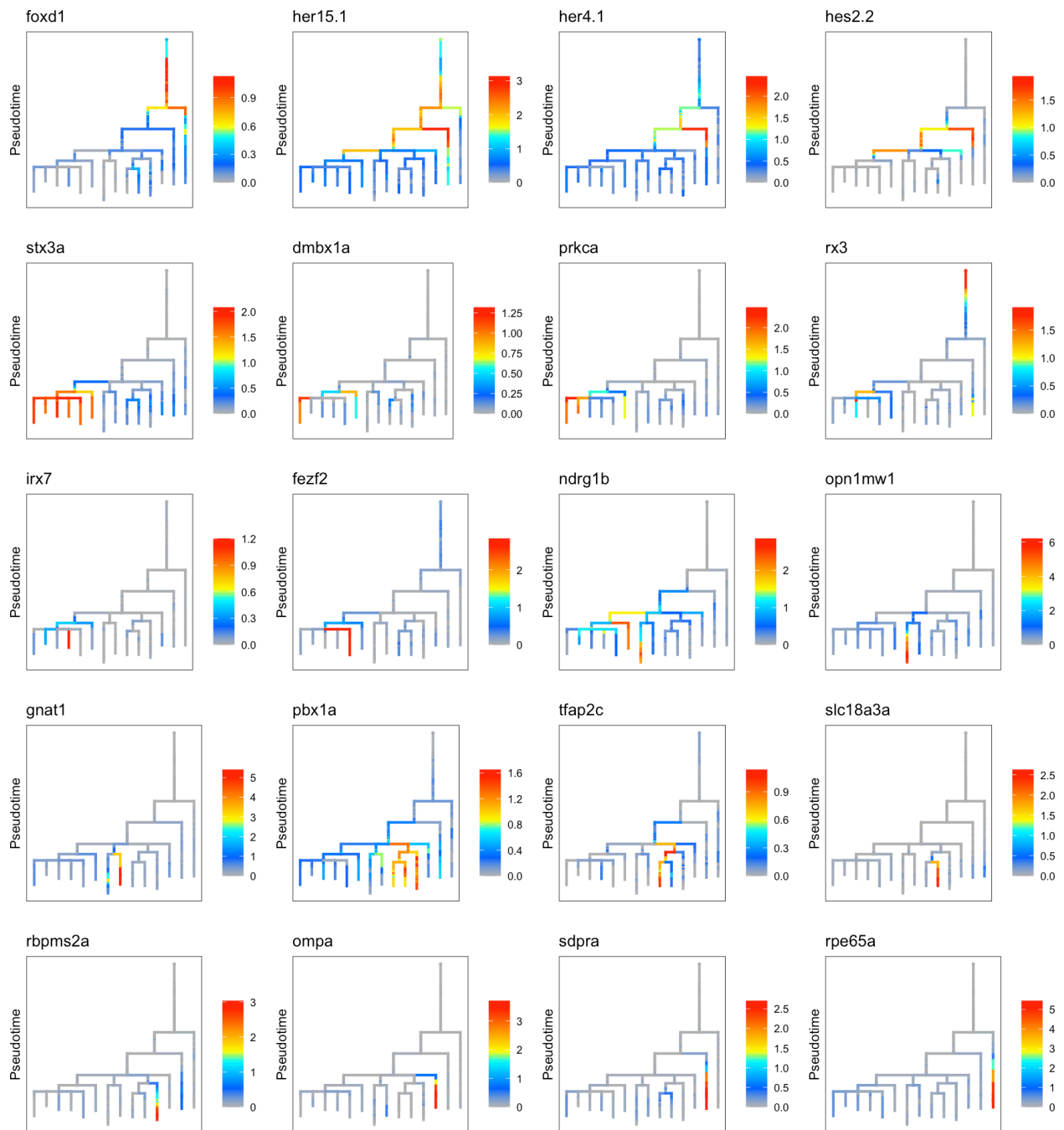
Plot tree with gene expression: main figures

```
gridExtra::grid.arrange(grobs = lapply(c("vsx1", "lmo4a", "pax6a", "rem1"), plotTree,
  object = obj, label.x = F, plot.cells = F), ncol = 2)
```



Plot tree with gene expression: supplemental figures

```
gridExtra::grid.arrange(grobs = lapply(c("foxd1", "her15.1", "her4.1", "hes2.2",
    "stx3a", "dmbx1a", "prkca", "rx3", "irx7", "fezf2", "ndrg1b", "opn1mw1", "gnat1",
    "pbx1a", "tfap2c", "slc18a3a", "rbpms2a", "ompa", "sdpra", "rpe65a"), plotTree,
    object = obj, label.x = F, plot.cells = F), ncol = 4)
```



Determine genes enriched in trajectories to particular cell types

Comparison between major cell types

We took each major group (“clade”) of branches from the end of the tree as a single entity (i.e. cone bipolar cells, photoreceptors, amacrine cells, retinal ganglion cells, horizontal cells) and compared them against each other pairwise to look for differentially expressed genes.

```
# Get the parent segment of each clade to consider as a group
combined.tips <- c("24", "25", "19", "8", "15")
```

```

# Get the cells in that segment and all child segments
cells.combined.tips <- lapply(combined.tips, function(t) whichCells(obj, label = "segment",
  value = segChildrenAll(obj, t, include.self = T)))
names(cells.combined.tips) <- combined.tips

# Loop through each of these clades and look for differentially expressed genes
combined.markers <- lapply(combined.tips, function(tip) {
  # Find all of the other clades
  opposing.tips <- setdiff(combined.tips, tip)
  # Perform pairwise comparisons to each other clade
  m.o <- lapply(opposing.tips, function(tip.opposing) {
    # message(paste0(Sys.time(), ': Comparing tip ', tip, ' to ', tip.opposing, '.'))
    # Find differentially expressed genes between the pair of clades
    ma <- markersAUCPR(object = obj, cells.1 = cells.combined.tips[[tip]], cells.2 = cells.combined.tips[[tip.opposing]],
      effect.size = 0.4, auc.factor = 1.1)
    # In order to facilitate combining all of the results later, add columns about
    # which two clades were compared and also a duplicate entry of the name of each
    # gene that's recovered.
    ma$gene <- rownames(ma)
    ma$tip1 <- tip
    ma$tip2 <- tip.opposing
    return(ma)
  })
  names(m.o) <- opposing.tips
  return(m.o)
})
names(combined.markers) <- combined.tips

# Require that genes are markers against at least 3 other clades
combined.markers.beatmult <- lapply(combined.markers, function(m) {
  names(which(table(unlist(lapply(m, rownames))) >= 3))
})

# Since genes might be a marker in a comparison to several other clades, combine
# the results into a single table, where each gene is listed only once with the
# info from the pairwise comparison where it had the strongest differential
# expression.
combined.markers.best <- lapply(1:length(combined.markers.beatmult), function(i) {
  cm <- do.call("rbind", combined.markers[[i]])
  cm <- cm[cm$gene %in% combined.markers.beatmult[[i]], ]
  cmb <- do.call("rbind", lapply(combined.markers.beatmult[[i]], function(g) {
    cmr <- cm[cm$gene == g, ]
    return(cmr[which.max(cmr$AUCPR.ratio), ])
  }))
  rownames(cmb) <- cmb$gene
  cmb <- cmb[order(cmb$AUCPR.ratio, decreasing = T), ]
  cmb$exp.global <- apply(obj@logupx.data[rownames(cmb), unlist(obj@tree$cells.in.segment)],
    1, mean.of.logs)
  cmb$exp.global.fc <- cmb$nTrans_1 - cmb$exp.global
  return(cmb)
})
names(combined.markers.best) <- combined.tips

```

AUCPR along tree

We also used the `AUCPRTTestAlongTree` function to ask for genes that are differential markers of a lineage using URD's tree structure. This makes a comparison at each branchpoint from a particular cell type up to the root.

```
# Get all of the tips from the tree
tips.in.tree <- as.character(obj@tree$tips)

# Tree segments to use as root for each particular cell population.
roots <- rep("29", length(tips.in.tree))
names(roots) <- tips.in.tree
roots["11"] <- "31"
roots["6"] <- "30"
roots[c("4", "17", "8")] <- "26"

# Perform a loop of tests with each tip.
markers <- lapply(tips.in.tree, function(t) {
  this.root <- roots[t]
  # message(paste0(Sys.time(), ': Starting tip ', t, ' and root ', this.root))
  these.markers <- aucprTestAlongTree(obj, pseudotime = "pseudotime", tips = as.character(t),
    genes.use = NULL, must.beat.sibs = 0.6, report.debug = F, root = this.root,
    auc.factor = 1.1, log.effect.size = 0.4)
  these.markers$gene <- rownames(these.markers)
  these.markers$tip <- t
  return(these.markers)
})
names(markers) <- tips.in.tree
```

Functions for curating differential expression results

We further curated those differentially expressed genes using the following functions:

`threshold.tree.markers`

Function to threshold markers from a `markersAUCPRAAlongTree` test with additional criteria

- **markers**: list of results from `markersAUCPRAAlongTree` tests
- **tip**: which tip (or element of the list to pursue)
- **global.fc**: fold.change that gene must have along the trajectory pursued vs. rest of the data
- **aucpr.ratio.all**: classifier score that gene must exhibit along trajectory test vs. rest of the data
- **branch.fc**: fold.change that gene must have (in best case) vs. the opposing branch at any branchpoint along the trajectory.
- Returns markers with only a subset of rows retained.

```
threshold.tree.markers <- function(markers, tip, global.fc = 0.1, branch.fc = 0.4,
  aucpr.ratio.all = 1.03) {
  m <- markers[[tip]]
  # First off -- lose global FC < x
  bye.globalfc <- rownames(m)[m$expfc.all < global.fc]
  # Second -- get rid of branch FC < x
  bye.branchfc <- rownames(m)[m$expfc.maxBranch < branch.fc]
  # Third -- get rid of stuff essentially worse than random classification on
  # global level
```

```

bye.badglobalaucpr <- rownames(m)[m$AUCPR.ratio.all < aucpr.ratio.all]
bye.all <- unique(c(bye.globalfc, bye.branchfc, bye.badglobalaucpr))
m.return <- m[setdiff(rownames(m), bye.all), ]
return(m.return)
}

```

threshold.clade.markers

Function to threshold markers of particular clades (see “Combined major branch families”) using additional criteria

- **markers:** result of markersAUCPR
- **global.fc:** fold.change that gene must have along the trajectory pursued vs. rest of the data (during testing, branches were compared pairwise. This compares one branch to all others together.)
- Returns markers with a subset of rows retained

```

threshold.clade.markers <- function(markers, global.fc = 0.1) {
  m <- markers
  # First off -- lose global FC < x
  bye.globalfc <- rownames(m)[m$exp.global.fc < global.fc]
  m.return <- m[setdiff(rownames(m), bye.globalfc), ]
  return(m.return)
}

```

divide.branches

Function to compare genes between two branches. Use this on a compiled list of markers to do a final selection of genes that are specific to one branch or another or markers of both (i.e. when making photoreceptor heatmap, use to divide into photoreceptor, cone, and rod markers)

- **object:** An URD object
- **genes:** (Character vector) Genes to test
- **clust.1:** (Character) Cluster 1
- **clust.2:** (Character) Cluster 2
- **clustering:** (Character) Clustering to pull from
- **exp.fc:** (Numeric) Minimum expression fold-change between branches to consider different
- **exp.thresh:** (Numeric) Minimum fraction of cells in order to consider gene expressed in a branch
- **exp.diff:** (Numeric) Minimum difference in fraction of cells expressing to consider gene differential
- Returns list of gene names (“specific.1” = specific to clust.1, “specific.2” = specific to clust.2, “markers” = all genes tested)

```

divide.branches <- function(object, genes, clust.1, clust.2, clustering = "segment",
  exp.fc = 0.4, exp.thresh = 0.1, exp.diff = 0.1) {
  # Double check which markers are unique to one or the other population
  mcomp <- markersAUCPR(object, clust.1 = clust.1, clust.2 = clust.2, clustering = clustering,
    effect.size = -Inf, auc.factor = 0, genes.use = genes, frac.min.diff = 0,
    frac.must.express = 0)
  specific.b <- rownames(mcomp)[abs(mcomp$exp.fc) > exp.fc & mcomp[, 4] < exp.thresh &
    mcomp[, 5] > pmin(mcomp[, 4] + exp.diff), 1)]
  specific.a <- rownames(mcomp)[abs(mcomp$exp.fc) > exp.fc & mcomp[, 5] < exp.thresh &
    mcomp[, 4] > pmin(mcomp[, 5] + exp.diff), 1)]
}

```



```

r <- list(specific.a, specific.b, mcomp)
names(r) <- c("specific.1", "specific.2", "markers")
return(r)
}

```

Functions for heatmap generation

These functions were used in the production of heatmaps:

Color scale

Generate color scale to use with heatmaps.

```

cols <- (scales::gradient_n_pal(RColorBrewer::brewer.pal(9, "YlOrRd")))(seq(0, 1,
length.out = 50))

```

determine.timing

Determines order to plot genes in heatmap. “Expression” is defined as 20% higher expression than the minimum observed value. “Peak” expression is defined as 50% higher expression than minimum observed value. The two longest stretches of “peak” expression are found, and then the later one is used. The onset time of the stretch of expression that contains that peak is also determined. Genes are then ordered by the pseudotime at which they enter “peak” expression, leave “peak” expression, start “expression”, and leave “expression”.

- **s**: result from `geneSmoothFit`
- **genes**: genes to order; default is all genes that were fit.
- Returns `s` but with an additional list entry (`$timing`) of the order to plot genes

```

determine.timing <- function(s, genes = rownames(s$mean.expression)) {
  s$timing <- as.data.frame(do.call("rbind", lapply(genes, function(g) {
    sv <- as.numeric(s$scaled.smooth[g, ])
    pt <- as.numeric(colnames(s$scaled.smooth))
    # Figure out baseline expression & threshold for finding peaks
    min.val <- max(min(sv), 0)
    peak.val <- ((1 - min.val)/2) + min.val
    exp.val <- ((1 - min.val)/5) + min.val
    # Run-length encoding of above/below the peak-threshold
    peak.rle <- rle(sv >= peak.val)
    peak.rle <- data.frame(lengths = peak.rle$lengths, values = peak.rle$values)
    peak.rle$end <- cumsum(peak.rle$lengths)
    peak.rle$start <- head(c(0, peak.rle$end) + 1, -1)
    # Run-length encoding of above/below the expressed-threshold
    exp.rle <- rle(sv >= exp.val)
    exp.rle <- data.frame(lengths = exp.rle$lengths, values = exp.rle$values)
    exp.rle$end <- cumsum(exp.rle$lengths)
    exp.rle$start <- head(c(0, exp.rle$end) + 1, -1)
    # Take top-two longest peak RLE & select later one. Find stretches that are
    # above peak value
    peak <- which(peak.rle$values)
    # Order by length and take 1 or 2 longest ones
    peak <- peak[order(peak.rle[peak, "lengths"], decreasing = T)][1:min(2, length(peak))]
  })]
}

```

```

# Order by start and take latest one.
peak <- peak[order(peak.rle[peak, "start"], decreasing = T)][1]
# Identify the actual peak value within that stretch
peak <- which.max(sv[peak.rle[peak, "start"]:peak.rle[peak, "end"]]) + peak.rle[peak,
  "start"] - 1
# Identify the start and stop of the expressed stretch that contains the peak
exp.start <- exp.rle[which(exp.rle$end >= peak & exp.rle$start <= peak),
  "start"]
exp.end <- exp.rle[which(exp.rle$end >= peak & exp.rle$start <= peak), "end"]
# Identify values of expression at start and stop
smooth.start <- sv[exp.start]
smooth.end <- sv[exp.end]
# Convert to pseudotime?
exp.start <- pt[exp.start]
exp.end <- pt[exp.end]
peak <- pt[peak]
# Return a vector
v <- c(exp.start, peak, exp.end, smooth.start, smooth.end)
names(v) <- c("pt.start", "pt.peak", "pt.end", "exp.start", "exp.end")
return(v)
})))
rownames(s$timing) <- genes

# Decide on ordering of genes
s$gene.order <- rownames(s$timing)[order(s$timing$pt.peak, s$timing$pt.start,
  s$timing$pt.end, s$timing$exp.end, decreasing = c(F, F, F, T), method = "radix")]

return(s)
}

```

filter.heatmap.genes

Removes undesired (mitochondrial, ribosomal, tandem duplicated genes) from heatmaps for presentation purposes.

- **genes:** (Character vector) genes to check
- Returns genes with undesired genes removed.

```

filter.heatmap.genes <- function(genes) {
  mt.genes <- grep("^mt-", ignore.case = T, genes, value = T)
  many.genes <- grep("\\(1 of many\\)", ignore.case = T, genes, value = T)
  ribo.genes <- grep("^rpl|^rps", ignore.case = T, genes, value = T)
  cox.genes <- grep("^cox", ignore.case = T, genes, value = T)
  return(setdiff(genes, c(mt.genes, many.genes, ribo.genes, cox.genes)))
}

```

Heatmaps of gene cascades

Using the genes that were determined as differentially expressed along the way to particular cell types, we generated expression cascades and plotted them as heatmaps.

Photoreceptors

Prepare cascade

```
## PHOTORECEPTORS: Seg 25 -> Cones (Seg 2) + Rods (Seg 12)

# Get markers from the two approaches:

# Lineage markers from above the combined clades
t25 <- threshold.clade.markers(combined.markers.best[["25"]], global.fc = 0.05)
# Cone markers from aucprTestAlongTree
m2 <- threshold.tree.markers(markers, "2", global.fc = 0.6)
# Rod markers from aucprTestAlongTree
m12 <- threshold.tree.markers(markers, "12", global.fc = 0.6)
pr.markers <- unique(c(rownames(t25), rownames(m2), rownames(m12)))

## Pseudotime for rods and cones is very different; for heatmaps, would like to
## normalize these, so that spline curves that consider both of them are not out
## of sync. Need to stretch pseudotime of cells in segment 12 / rods.

# Make a duplicate of the pseudotime measurement (pseudotime.212)
obj@pseudotime$pseudotime.212 <- obj@pseudotime$pseudotime
# Grab pseudotime of branchpoint
pt.start.212 <- as.numeric(obj@tree$segment.pseudotime.limits["2", "start"])
# Figure out lengths (and ratio) of the two branches in pseudotime
pt.end.212 <- as.numeric(obj@tree$segment.pseudotime.limits[c("2", "12"), "end"]) -
  pt.start.212
pt.ratio.212 <- pt.end.212[1]/pt.end.212[2]
# For cells in the shorter branch (12), subtract the starting pseudotime,
# multiply by the ratio of branch lengths, then add the starting pseudotime back
# in order to stretch the branch.
obj@pseudotime[cellsInCluster(obj, "segment", "12"), "pseudotime.212"] <- (obj@pseudotime[cellsInCluster(obj,
  "segment", "12"), "pseudotime.212"] - pt.start.212) * pt.ratio.212 + pt.start.212

# Calculate spline curves Using segments 29, 25, 2, and 12. Calculating a curve
# using only 29/25/2 for cone-specific genes, 29/25/12 for rod-specific genes,
# and 29/25/2+12 for genes that mark both. Should work now that pseudotimes are
# aligned.
spline.2 <- geneSmoothFit(obj, pseudotime = "pseudotime.212", cells = cellsInCluster(obj,
  "segment", c("29", "25", "2")), genes = pr.markers, method = "spline", moving.window = 5,
  cells.per.window = 25, pseudotime.per.window = 0.005, spar = 0.5, verbose = F)
spline.12 <- geneSmoothFit(obj, pseudotime = "pseudotime.212", cells = cellsInCluster(obj,
  "segment", c("29", "25", "12")), genes = pr.markers, method = "spline", moving.window = 5,
  cells.per.window = 25, pseudotime.per.window = 0.005, spar = 0.5, verbose = F)
spline.212 <- geneSmoothFit(obj, pseudotime = "pseudotime.212", cells = cellsInCluster(obj,
  "segment", c("29", "25", "2", "12")), genes = pr.markers, method = "spline",
  moving.window = 5, cells.per.window = 25, pseudotime.per.window = 0.005, spar = 0.5,
  verbose = F)

# Want to plot a heatmap that shows expression in photoreceptor progenitors and
# then each branch (i.e. rods, cones) as separate columns. Going to crop each
# spline fit to the correct pseudotime range and then combine them into a single
# one that can be plotted as a three-column heatmap.
```

```

pt.2v12 <- obj@tree$segment.pseudotime.limits["2", "start"] # pseudotime where the crop should happen
splines.pr <- list(cropSmoothFit(spline.212, pt.min = -Inf, pt.max = pt.2v12), cropSmoothFit(spline.2,
  pt.min = pt.2v12, pt.max = Inf), cropSmoothFit(spline.12, pt.min = pt.2v12, pt.max = Inf))
names(splines.pr) <- c("Photoreceptor Progenitors", "Rods", "Cones")
splines.pr.hm <- combineSmoothFit(splines.pr) # Combine into a single one

# Calculate gene expression timing for ordering rows
spline.212 <- determine.timing(s = spline.212)
spline.2 <- determine.timing(s = spline.2)
spline.12 <- determine.timing(s = spline.12)

# Decide which markers are specific to one cell type or both
d2v12 <- divide.branches(obj, pr.markers, clust.1 = "2", clust.2 = "12", exp.fc = 0.4,
  exp.thresh = 0.2, exp.diff = 0.1)

# Generate gene ordering based on timing & specificity
order.212 <- filter.heatmap.genes(setdiff(spline.212$gene.order, c(d2v12$specific.1,
  d2v12$specific.2)))
order.2 <- filter.heatmap.genes(intersect(spline.2$gene.order, d2v12$specific.1))
order.12 <- filter.heatmap.genes(intersect(spline.12$gene.order, d2v12$specific.2))
gene.order <- c(order.212, order.2, order.12)

# Output gene table
table.save <- data.frame(gene = gene.order, marks = c(rep("both", length(order.212)),
  rep("cone", length(order.2)), rep("rod", length(order.12))), stringsAsFactors = F)
table.save$clade.AUCPR.ratio <- t25[table.save$gene, "AUCPR.ratio"]
table.save$clade.exp.fc <- t25[table.save$gene, "exp.fc"]
table.save$clade.exp.fc.global <- t25[table.save$gene, "exp.global.fc"]
table.save$cone.AUCPR.ratio.all <- m2[table.save$gene, "AUCPR.ratio.all"]
table.save$cone.AUCPR.ratio.maxBranch <- m2[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$cone.exp.fc.all <- m2[table.save$gene, "expfc.all"]
table.save$cone.exp.fc.best <- m2[table.save$gene, "expfc.maxBranch"]
table.save$rod.AUCPR.ratio.all <- m12[table.save$gene, "AUCPR.ratio.all"]
table.save$rod.AUCPR.ratio.maxBranch <- m12[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$rod.exp.fc.all <- m12[table.save$gene, "expfc.all"]
table.save$rod.exp.fc.best <- m12[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/retina-photoreceptor.csv"))

```

Generate heatmap: all genes

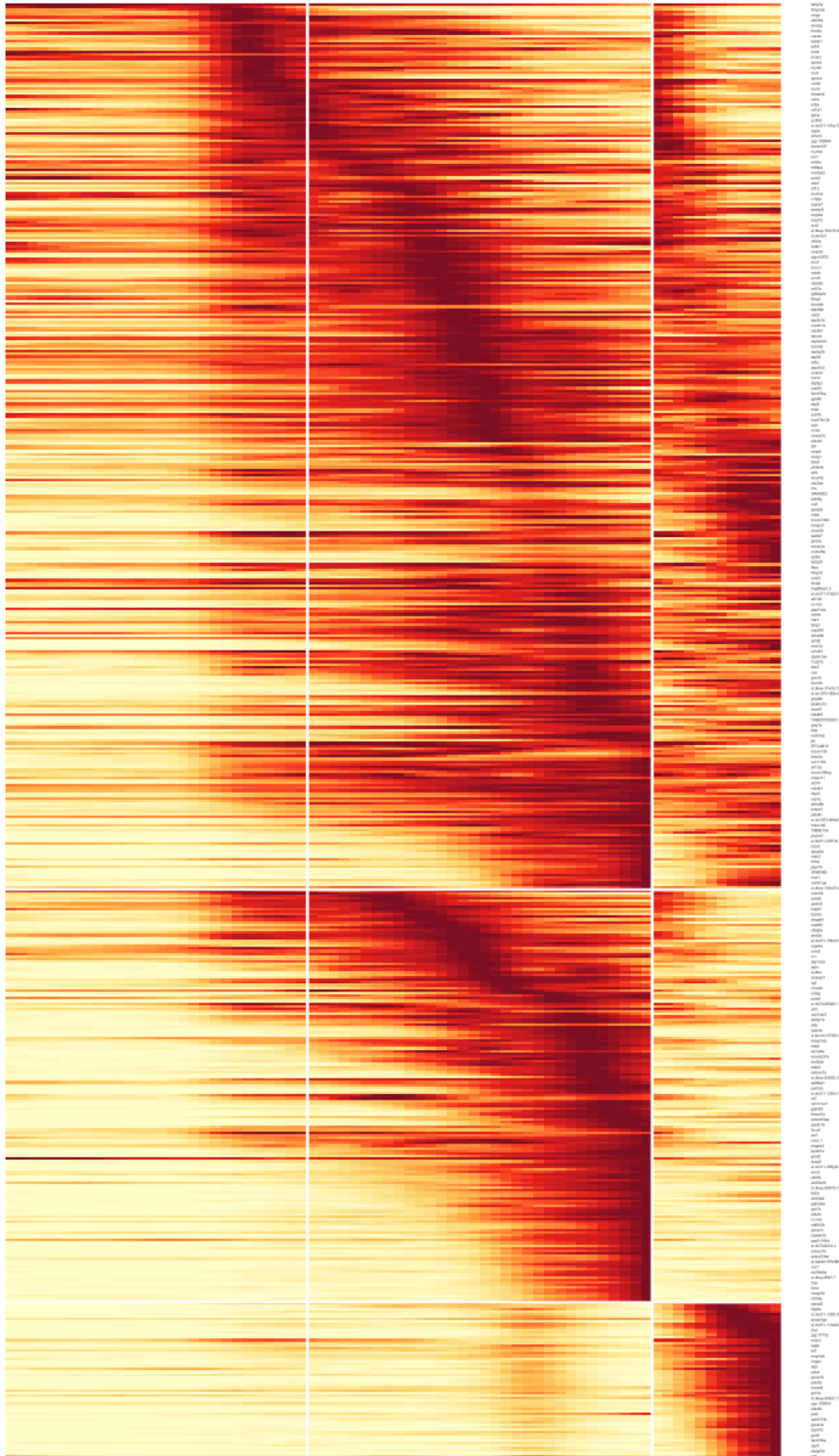
```

# Make sure any values <0 in the spline curves get set to 0 so that the heatmap
# scale doesn't get messed up.
splines.pr.hm$scaled.smooth[splines.pr.hm$scaled.smooth < 0] <- 0
# Determine where to place column separators (i.e. how many columns will each
# cell type occupy in the heatmap )
colsep <- cumsum(as.numeric(head(unlist(lapply(splines.pr, function(x) ncol(x$scaled.smooth))),
  -1)))
# Determine where to place row separators (i.e. how many common markers, and
# markers are specific to each cell type)
rowsep <- cumsum(c(length(order.212), length(order.2)))
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/retina-photoreceptor.pdf'), width=6, height=10)

```

```
gplots::heatmap.2(x = as.matrix(splines.pr.hm$scaled.smooth[gene.order, ]), Rowv = F,  
  Colv = F, dendrogram = "none", col = cols, trace = "none", density.info = "none",  
  key = F, cexCol = 0.8, cexRow = 0.15, margins = c(8, 8), lwid = c(0.3, 4), lhei = c(0.3,  
  4), labCol = NA, colsep = colsep, rowsep = rowsep, sepwidth = c(0.1, 0.2))  
title(main = "Photoreceptors")  
title(main = "Precursors", line = -41, adj = 0)  
title(main = "Cones", line = -41, adj = 0.45)  
title(main = "Rods", line = -41, adj = 0.76)
```

Photoreceptors



Precursors

Cones

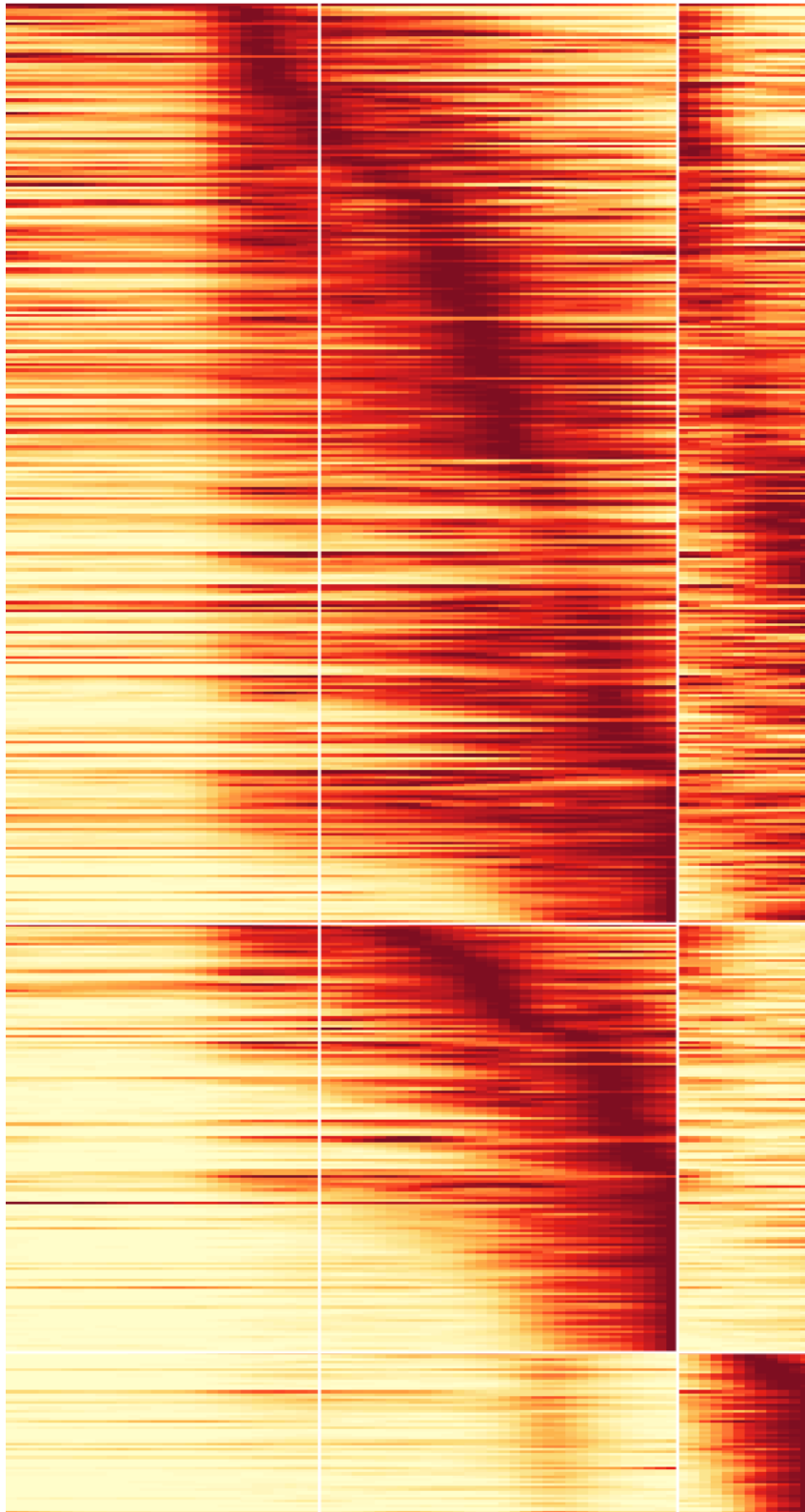
Rods

```
# dev.off()
```

Generate heatmap: main figure

```
## Generate heatmap with only particular genes labeled for main figure
genes.to.plot <- c("isl2a", "prdm1a", "otx5", "crx", "six7", "nr2f1b", "nr2e3", "aplnrb",
  "aplnra", "apln")
rownames.to.plot <- gene.order
rtp <- rownames.to.plot %in% genes.to.plot
rownames.to.plot[!rtp] <- ""
rownames.to.plot[rtp] <- paste0("- ", rownames.to.plot[rtp])
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/retina-photoreceptor-mainfig.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(splines.pr.hm$scaled.smooth[gene.order, ]), Rowv = F,
  Colv = F, dendrogram = "none", col = cols, trace = "none", density.info = "none",
  key = F, cexCol = 0.8, cexRow = 1.8, margins = c(8, 8), lwid = c(0.3, 4), lhei = c(0.3,
  4), labCol = NA, colsep = colsep, rowsep = rowsep, sepwidth = c(0.1, 0.2),
  labRow = rownames.to.plot)
title(main = "Photoreceptors")
title(main = "Precursors", line = -41, adj = 0)
title(main = "Cones", line = -41, adj = 0.45)
title(main = "Rods", line = -41, adj = 0.76)
```

Photoreceptors



Precursors

Cones

Rods


```
# dev.off()
```

Amacrine cells

Prepare cascade

```
## AMACRINE CELLS: Seg 19 -> Amacrine (Seg 4) + Starburst Amacrine (Seg 17)

# Get markers from the two approaches:

# Lineage markers from above the combined clades
t19 <- threshold.clade.markers(combined.markers.best[["19"]], global.fc = 0.05)
# Amacrine markers from aucprTestAlongTree
m4 <- threshold.tree.markers(markers, "4", global.fc = 0.6)
# Starburst amacrine markers from aucprTestAlongTree
m17 <- threshold.tree.markers(markers, "17", global.fc = 0.6)
am.markers <- unique(c(rownames(t19), rownames(m4), rownames(m17)))

## These have pretty equivalent pseudotimes, so don't need to worry about
## stretching them to match or anything.

# Calculate spline curves Using segments 29, 26, 19, and 4/17. Calculating a
# curve using only 29/26/19/4 for amacrine-specific genes, 29/26/19/4 for
# starburst-specific genes, and 29/26/19/4+17 for genes that mark both amacrine
# populations.
spline.4 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cellsInCluster(obj,
  "segment", c("29", "26", "19", "4")), genes = am.markers, method = "spline",
  moving.window = 5, cells.per.window = 25, pseudotime.per.window = 0.005, spar = 0.5,
  verbose = F)
spline.17 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cellsInCluster(obj,
  "segment", c("29", "26", "19", "17")), genes = am.markers, method = "spline",
  moving.window = 5, cells.per.window = 25, pseudotime.per.window = 0.005, spar = 0.5,
  verbose = F)
spline.417 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cellsInCluster(obj,
  "segment", c("29", "26", "19", "4", "17")), genes = am.markers, method = "spline",
  moving.window = 5, cells.per.window = 25, pseudotime.per.window = 0.005, spar = 0.5,
  verbose = F)

# Want to plot a heatmap that shows expression in amacrine progenitors and then
# each branch (i.e. amacrine_gaba, starburst_amacrine) as separate columns. Going
# to crop each spline fit to the correct pseudotime range and then combine them
# into a single one that can be plotted as a three-column heatmap.

pt.4v17 <- obj@tree$segment.pseudotime.limits["4", "start"] # pseudotime where the crop should happen
splines.am <- list(cropSmoothFit(spline.417, pt.min = -Inf, pt.max = pt.4v17), cropSmoothFit(spline.4,
  pt.min = pt.4v17, pt.max = Inf), cropSmoothFit(spline.17, pt.min = pt.4v17, pt.max = Inf))
names(splines.am) <- c("Amacrine Precursors", "Amacrine", "Starburst Amacrine")
splines.am.hm <- combineSmoothFit(splines.am) # Combine into a single one

# Calculate gene expression timing for ordering rows
spline.417 <- determine.timing(s = spline.417)
spline.4 <- determine.timing(s = spline.4)
spline.17 <- determine.timing(s = spline.17)
```

```

# Decide which markers are specific to one cell type or both
d4v17 <- divide.branches(obj, am.markers, clust.1 = "4", clust.2 = "17", exp.fc = 0.4,
  exp.thresh = 0.2, exp.diff = 0.1)

# Generate gene ordering based on timing & specificity
order.417 <- filter.heatmap.genes(setdiff(spline.417$gene.order, c(d4v17$specific.1,
  d4v17$specific.2)))
order.4 <- filter.heatmap.genes(intersect(spline.4$gene.order, d4v17$specific.1))
order.17 <- filter.heatmap.genes(intersect(spline.17$gene.order, d4v17$specific.2))
gene.order <- c(order.417, order.4, order.17)

# Output gene table
table.save <- data.frame(gene = gene.order, marks = c(rep("both", length(order.417)),
  rep("amacrine", length(order.4)), rep("starburst", length(order.17))), stringsAsFactors = F)
table.save$clade.AUCPR.ratio <- t19[table.save$gene, "AUCPR.ratio"]
table.save$clade.exp.fc <- t19[table.save$gene, "exp.fc"]
table.save$clade.exp.fc.global <- t19[table.save$gene, "exp.global.fc"]
table.save$am.AUCPR.ratio.all <- m4[table.save$gene, "AUCPR.ratio.all"]
table.save$am.AUCPR.ratio.maxBranch <- m4[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$am.exp.fc.all <- m4[table.save$gene, "expfc.all"]
table.save$am.exp.fc.best <- m4[table.save$gene, "expfc.maxBranch"]
table.save$star.AUCPR.ratio.all <- m17[table.save$gene, "AUCPR.ratio.all"]
table.save$star.AUCPR.ratio.maxBranch <- m17[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$star.exp.fc.all <- m17[table.save$gene, "expfc.all"]
table.save$star.exp.fc.best <- m17[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/retina-amacrine.csv"))

```

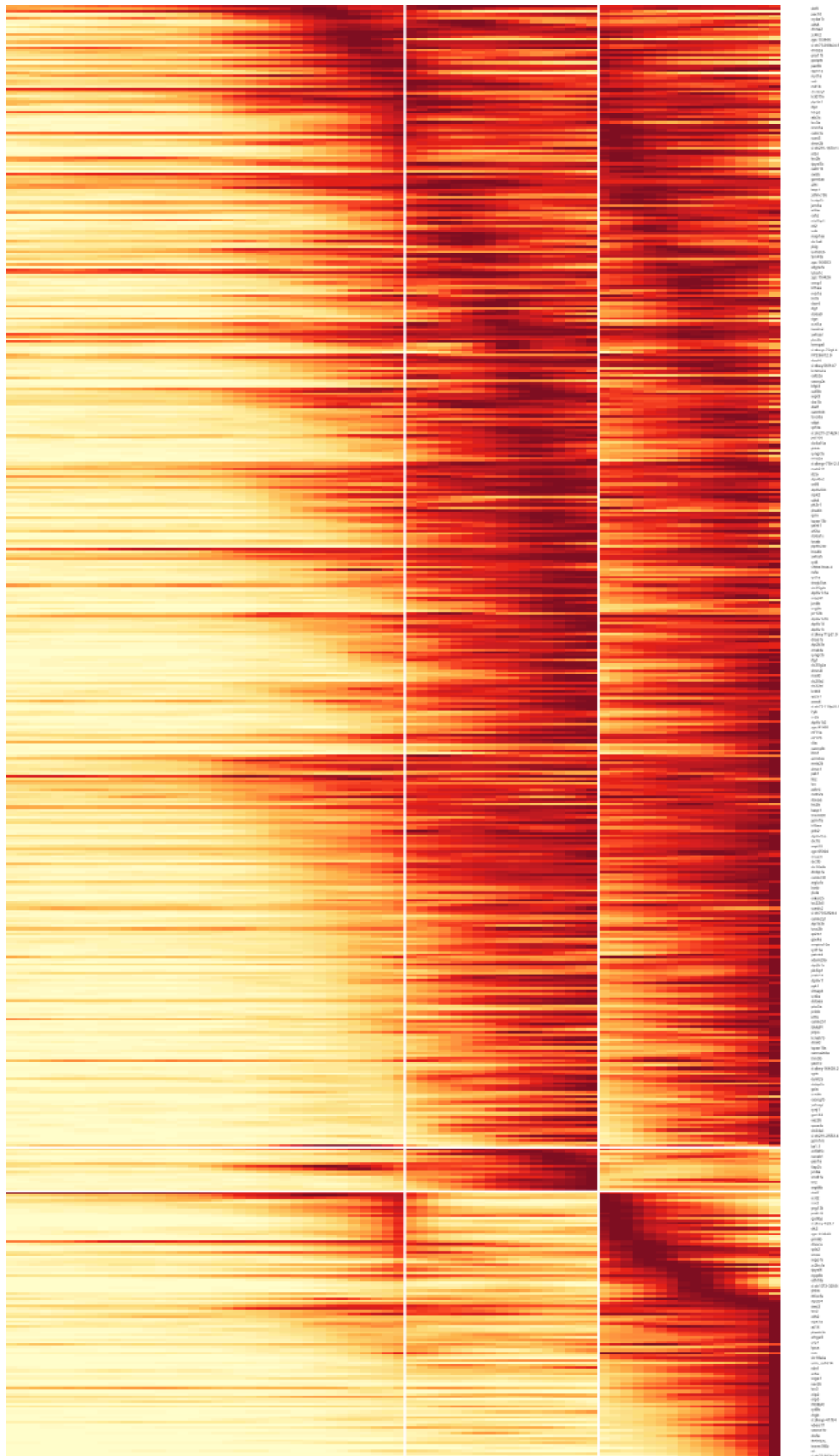
Generate heatmap: all genes

```

# Make sure any values <0 in the spline curves get set to 0 so that the heatmap
# scale doesn't get messed up.
splines.am.hm$scaled.smooth[splines.am.hm$scaled.smooth < 0] <- 0
# Determine where to place column separators (i.e. how many columns will each
# cell type occupy in the heatmap )
colsep <- cumsum(as.numeric(head(unlist(lapply(splines.am, function(x) ncol(x$scaled.smooth))),
  -1)))
# Determine where to place row separators (i.e. how many common markers, and
# markers are specific to each cell type)
rowsep <- cumsum(c(length(order.417), length(order.4)))
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/retina-amacrine.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(splines.am.hm$scaled.smooth[gene.order, ]), Rowv = F,
  Colv = F, dendrogram = "none", col = cols, trace = "none", density.info = "none",
  key = F, cexCol = 0.8, cexRow = 0.15, margins = c(8, 8), lwid = c(0.3, 4), lhei = c(0.3,
  4), labCol = NA, colsep = colsep, rowsep = rowsep, sepwidth = c(0.05, 0.2))
title(main = "Amacrine Cells")
title(main = "Precursors", line = -41, adj = 0.05)
title(main = "Amacrine", line = -41, adj = 0.475)
title(main = "Starburst", line = -41, adj = 0.75)

```

Amacrine Cells



Precursors

Amacrine

Starburst

```
# dev.off()
```

Retinal ganglion cells

Prepare cascade

```
## RGCs: Seg 8

# Get markers from the two approaches:

# Lineage markers from above the combined clades
t8 <- threshold.clade.markers(combined.markers.best[["8"]], global.fc = 0.05)
# RGC markers from aucprTestAlongTree
m8 <- threshold.tree.markers(markers, "8", global.fc = 0.6)
rgc.markers <- unique(c(rownames(t8), rownames(m8)))

# Calculate spline curves Using segments 29, 26, and 8.
spline.8 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cellsInCluster(obj,
  "segment", c("29", "26", "8")), genes = rgc.markers, method = "spline", moving.window = 5,
  cells.per.window = 25, pseudotime.per.window = 0.005, spar = 0.5, verbose = F)

# Calculate gene expression timing for ordering rows
spline.8 <- determine.timing(s = spline.8)
order.8 <- filter.heatmap.genes(spline.8$gene.order)

# Output gene table
table.save <- data.frame(gene = order.8, stringsAsFactors = F)
table.save$clade.AUCPR.ratio <- t8[table.save$gene, "AUCPR.ratio"]
table.save$clade.exp.fc <- t8[table.save$gene, "exp.fc"]
table.save$clade.exp.fc.global <- t8[table.save$gene, "exp.global.fc"]
table.save$rgc.AUCPR.ratio.all <- m8[table.save$gene, "AUCPR.ratio.all"]
table.save$rgc.AUCPR.ratio.maxBranch <- m8[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$rgc.exp.fc.all <- m8[table.save$gene, "expfc.all"]
table.save$rgc.exp.fc.best <- m8[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/retina-rgc.csv"))
```

Generate heatmap: all genes

```
# Make sure any values <0 in the spline curves get set to 0 so that the heatmap
# scale doesn't get messed up.
spline.8$scaled.smooth[spline.8$scaled.smooth < 0] <- 0
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/retina-rgc.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(spline.8$scaled.smooth[order.8, ]), Rowv = F, Colv = F,
  dendrogram = "none", col = cols, trace = "none", density.info = "none", key = F,
  cexCol = 0.8, cexRow = 0.15, margins = c(8, 8), lwid = c(0.3, 4), lhei = c(0.3,
  4), labCol = NA)
title(main = "Retinal Ganglion Cells")
```

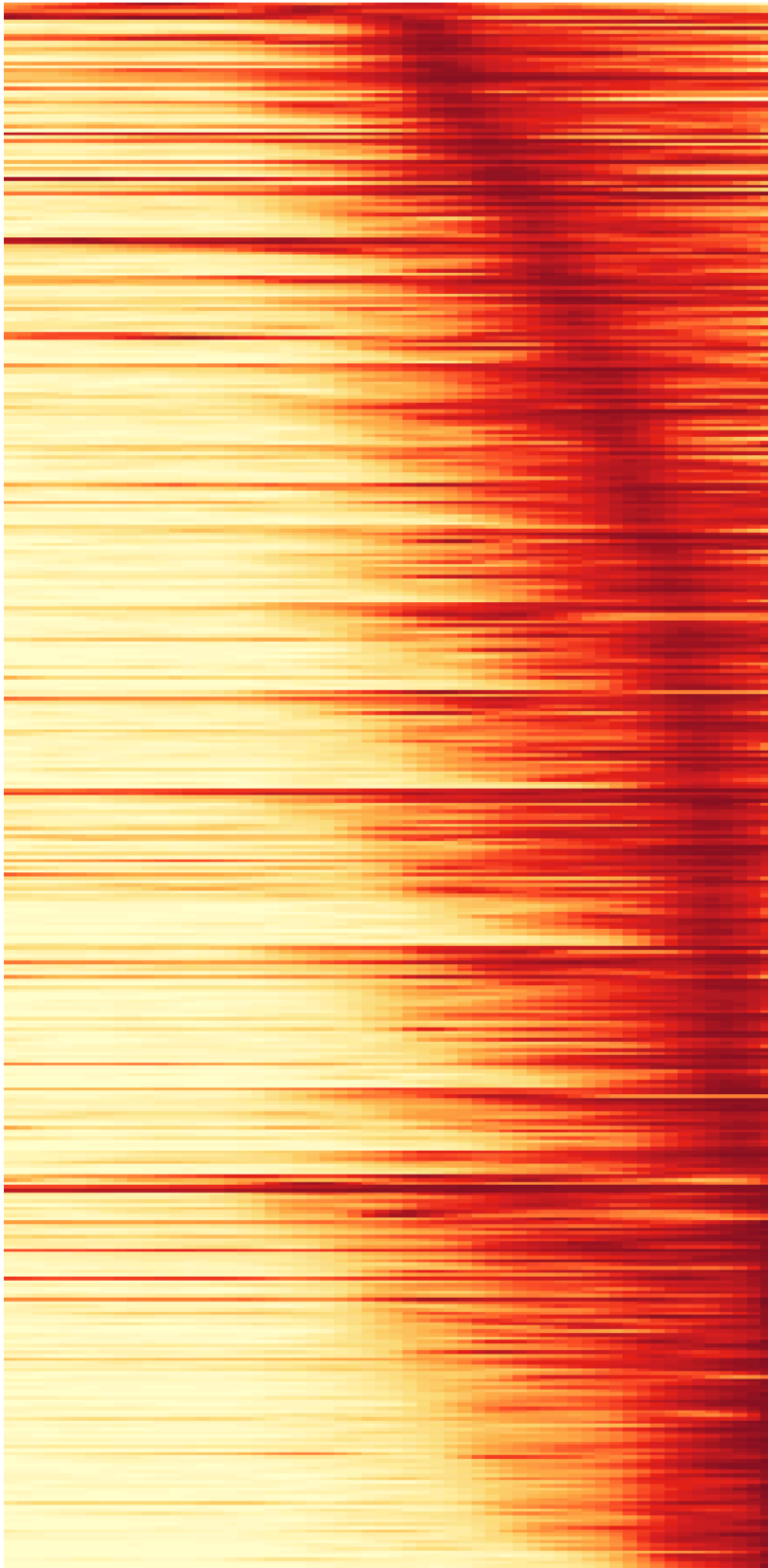


```
# dev.off()
```

Generate heatmap: main figure

```
genes.to.plot <- c("sox11a", "sox11b", "sox6", "irx4a", "pou4f2", "pou4f1", "rbpms2b",  
  "rbpms2a")  
rownames.to.plot <- order.8  
rtp <- rownames.to.plot %in% genes.to.plot  
rownames.to.plot[!rtp] <- ""  
rownames.to.plot[rtp] <- paste0("- ", rownames.to.plot[rtp])  
# Open a PDF and generate the heatmap pdf(paste0(base.path,  
# '/heatmaps/retina-rgc-mainfig.pdf'), width=6, height=10)  
gplots::heatmap.2(x = as.matrix(spline.8$scaled.smooth[order.8, ]), Rowv = F, Colv = F,  
  dendrogram = "none", col = cols, trace = "none", density.info = "none", key = F,  
  cexCol = 0.8, cexRow = 1.8, margins = c(8, 10), lwid = c(0.3, 4), lhei = c(0.3,  
  4), labCol = NA, labRow = rownames.to.plot)  
title(main = "Retinal Ganglion Cells")
```

Retinal Ganglion Cells



- sox11b

```
# dev.off()
```

Horizontal Cells

Prepare cascade

```
## Horizontal Cells: Seg 15

# Get markers from the two approaches:

# Lineage markers from above the combined clades
t15 <- threshold.clade.markers(combined.markers.best[["15"]], global.fc = 0.05)
# Horizontal Cell markers from aucprTestAlongTree
m15 <- threshold.tree.markers(markers, "15", global.fc = 0.6)
horiz.markers <- unique(c(rownames(t15), rownames(m15)))

# Calculate spline curves Using segments 29 and 15.
spline.15 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cellsInCluster(obj,
  "segment", c("29", "15")), genes = horiz.markers, method = "spline", moving.window = 5,
  cells.per.window = 25, pseudotime.per.window = 0.005, spar = 0.5, verbose = F)

# Calculate gene expression timing for ordering rows
spline.15 <- determine.timing(s = spline.15)
order.15 <- filter.heatmap.genes(spline.15$gene.order)

# Output gene table
table.save <- data.frame(gene = order.15, stringsAsFactors = F)
table.save$clade.AUCPR.ratio <- t15[table.save$gene, "AUCPR.ratio"]
table.save$clade.exp.fc <- t15[table.save$gene, "exp.fc"]
table.save$clade.exp.fc.global <- t15[table.save$gene, "exp.global.fc"]
table.save$horiz.AUCPR.ratio.all <- m15[table.save$gene, "AUCPR.ratio.all"]
table.save$horiz.AUCPR.ratio.maxBranch <- m15[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$horiz.exp.fc.all <- m15[table.save$gene, "expfc.all"]
table.save$horiz.exp.fc.best <- m15[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/retina-horiz.csv"))
```

Generate heatmap: all genes

```
# Make sure any values <0 in the spline curves get set to 0 so that the heatmap
# scale doesn't get messed up.
splines.am.hm$scaled.smooth[splines.am.hm$scaled.smooth < 0] <- 0
# Determine where to place column separators (i.e. how many columns will each
# cell type occupy in the heatmap )
colsep <- cumsum(as.numeric(head(unlist(lapply(splines.am, function(x) ncol(x$scaled.smooth))),
  -1)))
# Determine where to place row separators (i.e. how many common markers, and
# markers are specific to each cell type)
rowsep <- cumsum(c(length(order.417), length(order.4)))
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/retina-amacrine.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(splines.am.hm$scaled.smooth[order.15, ]), Rowv = F,
  Colv = F, dendrogram = "none", col = cols, trace = "none", density.info = "none",
```



```
key = F, cexCol = 0.8, cexRow = 0.15, margins = c(8, 8), lwid = c(0.3, 4), lhei = c(0.3,
4), labCol = NA, colsep = colsep, rowsep = rowsep, sepwidth = c(0.05, 0.2))
title(main = "Amacrine Cells")
title(main = "Precursors", line = -41, adj = 0.05)
title(main = "Amacrine", line = -41, adj = 0.475)
title(main = "Starburst", line = -41, adj = 0.75)
```



```
# dev.off()
```

Muller Glia

Prepare cascade

```
## Muller Glia: Seg 6

# Get markers from the two approaches
m6 <- threshold.tree.markers(markers, "6", global.fc = 0.6) # Muller Glia markers from aucprTestAlongTree
muller.markers <- rownames(m6)

# Calculate spline curves Using segments 29 and 15.
spline.6 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cellsInCluster(obj,
  "segment", c("30", "6")), genes = muller.markers, method = "spline", moving.window = 5,
  cells.per.window = 25, pseudotime.per.window = 0.005, spar = 0.5, verbose = F)

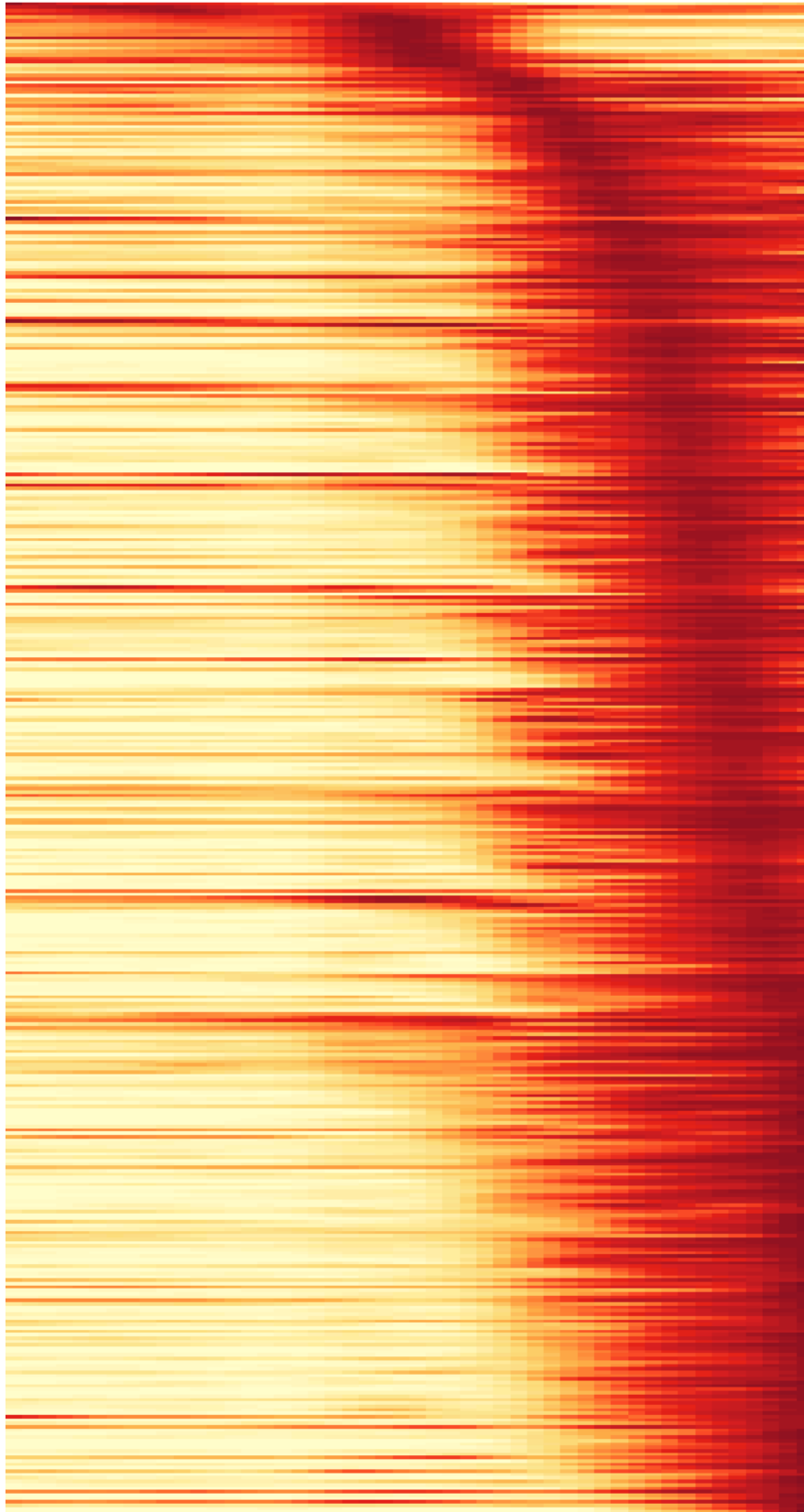
# Calculate gene expression timing for ordering rows
spline.6 <- determine.timing(s = spline.6)
order.6 <- filter.heatmap.genes(spline.6$gene.order)

# Output gene table
table.save <- data.frame(gene = order.6, stringsAsFactors = F)
table.save$muller.AUCPR.ratio.all <- m6[table.save$gene, "AUCPR.ratio.all"]
table.save$muller.AUCPR.ratio.maxBranch <- m6[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$muller.exp.fc.all <- m6[table.save$gene, "expfc.all"]
table.save$muller.exp.fc.best <- m6[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/retina-muller.csv"))
```

Generate heatmap: all genes

```
# Make sure any values <0 in the spline curves get set to 0 so that the heatmap
# scale doesn't get messed up.
spline.6$scaled.smooth[spline.6$scaled.smooth < 0] <- 0
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/retina-muller.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(spline.6$scaled.smooth[order.6, ]), Rowv = F, Colv = F,
  dendrogram = "none", col = cols, trace = "none", density.info = "none", key = F,
  cexCol = 0.8, cexRow = 0.15, margins = c(8, 8), lwid = c(0.3, 4), lhei = c(0.3,
  4), labCol = NA)
title(main = "Muller Glia")
```

Muller Glia



```
# dev.off()
```

Retinal Pigmented Epithelium

Prepare cascade

```
## RPE: Seg 11

# Get markers from the two approaches
m11 <- threshold.tree.markers(markers, "11", global.fc = 0.6) # RPE markers from aucprTestAlongTree
rpe.markers <- rownames(m11)

# Just want to plot part of cells from upstream segment 31, which is very long.
# Going to use cells from segment 11 and from segment 31 with pseudotime > 0.23
cells.rpe <- unique(c(whichCells(obj, "pseudotime", c(0.23, 0.30308134)), cellsInCluster(obj,
  "segment", "11")))

# Calculate spline curves
spline.11 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cells.rpe, genes = rpe.markers,
  method = "spline", moving.window = 5, cells.per.window = 25, pseudotime.per.window = 0.005,
  spar = 0.5, verbose = F)

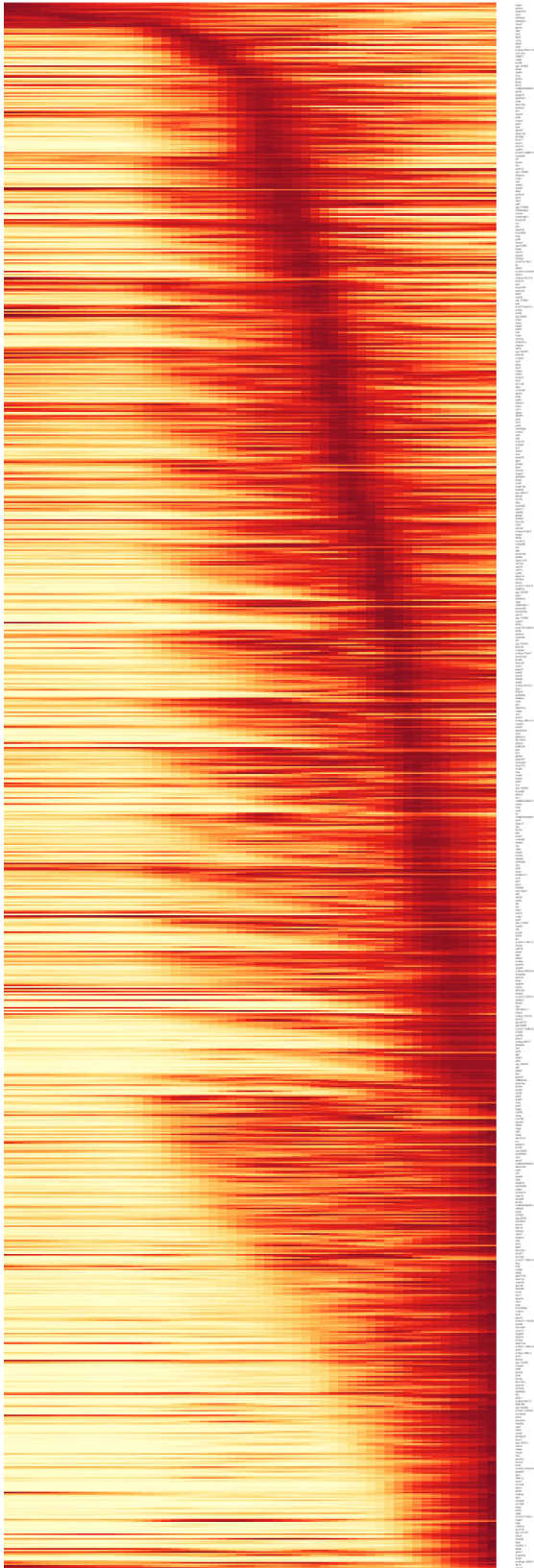
# Calculate gene expression timing for ordering rows
spline.11 <- determine.timing(s = spline.11)
order.11 <- filter.heatmap.genes(spline.11$gene.order)

# Output gene table
table.save <- data.frame(gene = order.11, stringsAsFactors = F)
table.save$rpe.AUCPR.ratio.all <- m11[table.save$gene, "AUCPR.ratio.all"]
table.save$rpe.AUCPR.ratio.maxBranch <- m11[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$rpe.exp.fc.all <- m11[table.save$gene, "expfc.all"]
table.save$rpe.exp.fc.best <- m11[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/retina-rpe.csv"))
```

Generate heatmap: all genes

```
# Make sure any values <0 in the spline curves get set to 0 so that the heatmap
# scale doesn't get messed up.
spline.11$scaled.smooth[spline.11$scaled.smooth < 0] <- 0
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/retina-rpe.pdf'), width=6, height=16)
gplots::heatmap.2(x = as.matrix(spline.11$scaled.smooth[order.11, ]), Rowv = F, Colv = F,
  dendrogram = "none", col = cols, trace = "none", density.info = "none", key = F,
  cexCol = 0.8, cexRow = 0.15, margins = c(8, 8), lwid = c(0.3, 4), lhei = c(0.3,
  4), labCol = NA)
title(main = "Retinal Pigmented Epithelium")
```

Retinal Pigmented Epithelium



```
# dev.off()
```

Continuous differentiation

Retinal cell types were often found with similar molecular states across many stages of development. This reflects that pseudotime accurately represents the asynchrony introduced by continuous differentiation.

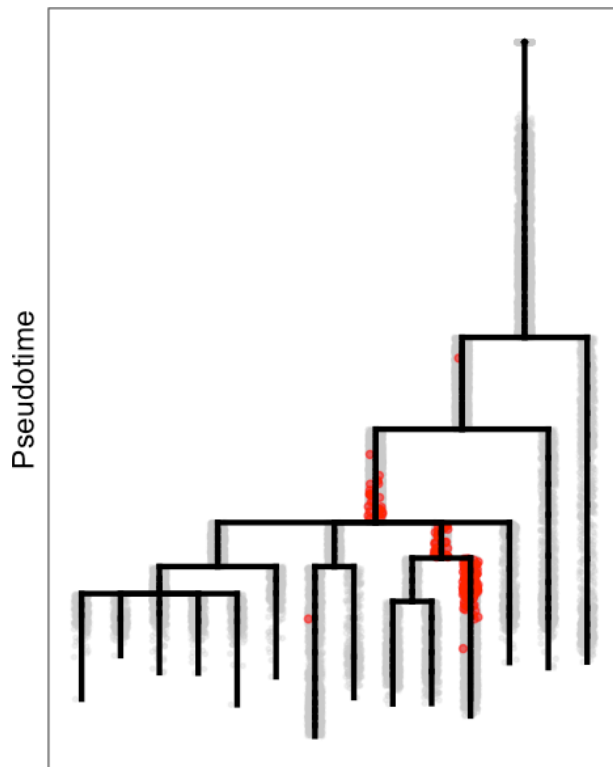
RGC cells

```
gridExtra::grid.arrange(grobs = list(plotTreeHighlight(obj, label.name = "clus.orig",  
  label.value = "7-36h_27", highlight.size = 1, title = "36 hpf Cluster 27: RGCs",  
  label.x = F), plotTreeHighlight(obj, label.name = "clus.orig", label.value = "12-15d_38",  
  highlight.size = 1, title = "15 dpf Cluster 38: RGCs", label.x = F)), ncol = 2)
```

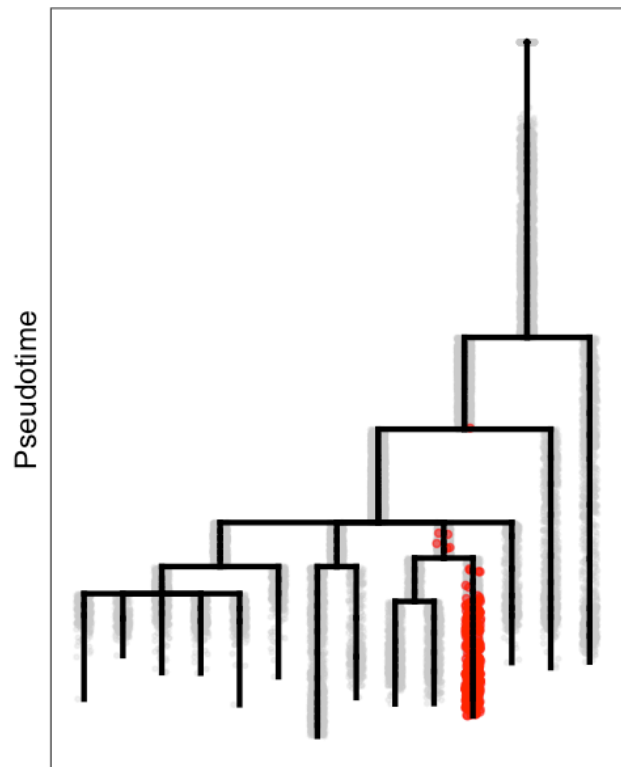
```
## Warning: Removed 5 rows containing missing values (geom_point).
```

```
## Warning: Removed 17 rows containing missing values (geom_point).
```

36 hpf Cluster 27: RGCs



15 dpf Cluster 38: RGCs



Progenitor cells

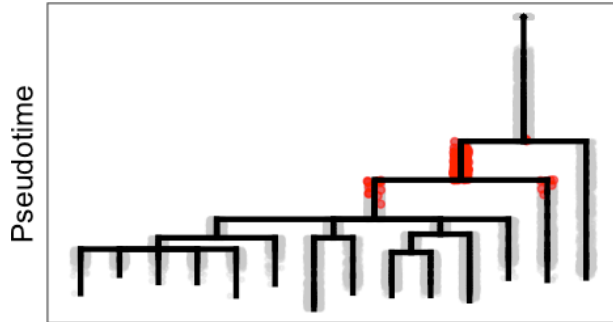
```
gridExtra::grid.arrange(grobs = list(plotTreeHighlight(obj, label.name = "clus.orig",  
  label.value = "6-24h_22", highlight.size = 1, title = "24 hpf Cluster 22: Progenitor Cells",  
  label.x = F), plotTreeHighlight(obj, label.name = "clus.orig", label.value = "7-36h_32",  
  highlight.size = 1, title = "36 hpf Cluster 32: Progenitor Cells", label.x = F)), ncol = 2)
```

```
plotTreeHighlight(obj, label.name = "clus.orig", label.value = "12-15d_39", highlight.size = 1,
  title = "15 dpf Cluster 39: Progenitor Cells", label.x = F)), ncol = 2)
```

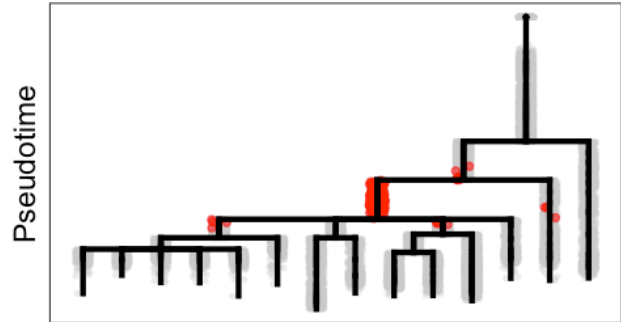
```
## Warning: Removed 9 rows containing missing values (geom_point).
```

```
## Warning: Removed 4 rows containing missing values (geom_point).
```

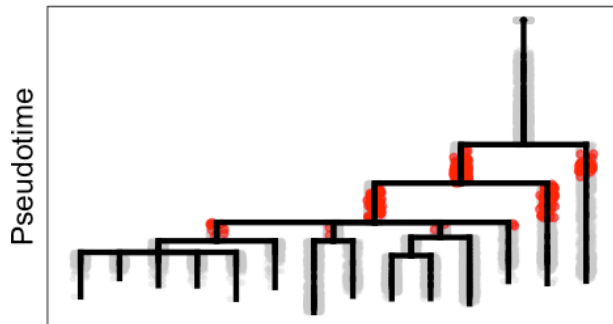
24 hpf Cluster 22: Progenitor Cells



36 hpf Cluster 32: Progenitor Cells



15 dpf Cluster 39: Progenitor Cells



Progenitors over time

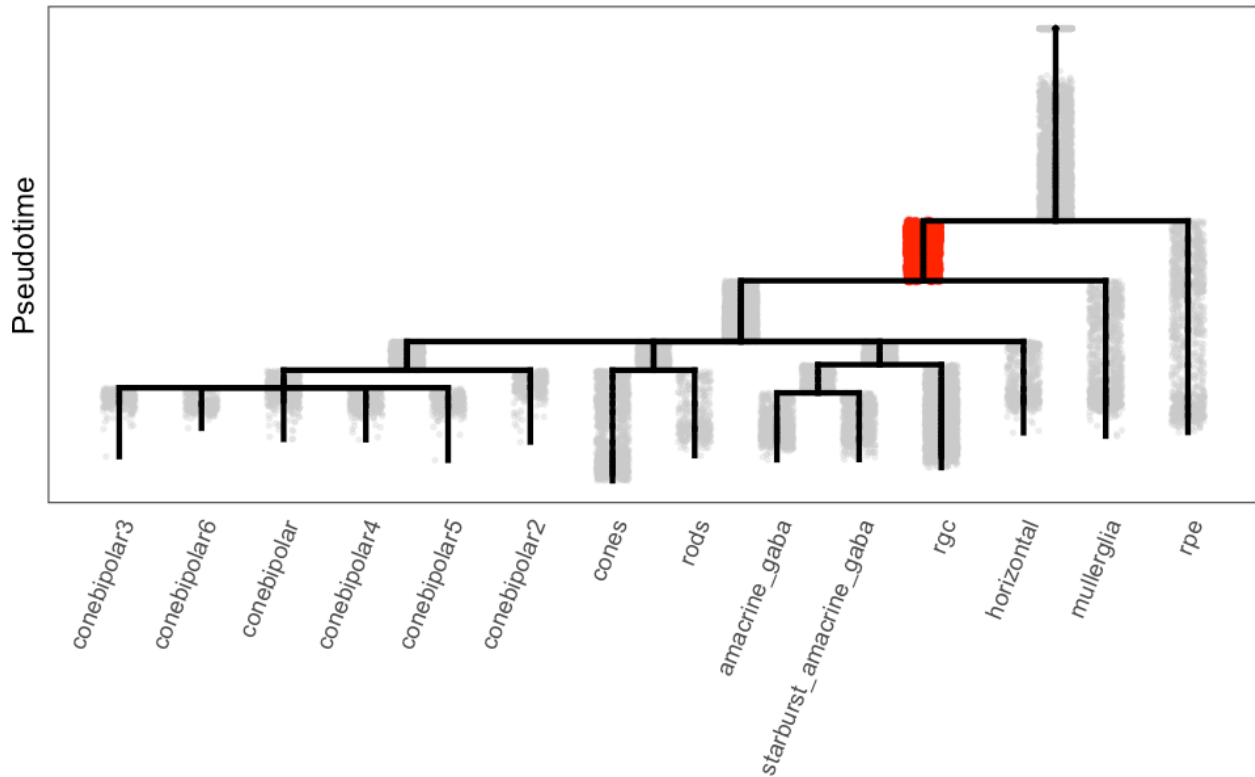
Retinal progenitors with similar transcriptional states are found across many different time points. We wanted to know whether there were significant transcriptional changes within those progenitors between early stages and late stages.

Identify populations

First we grabbed early (24 / 36 hpf) and late (15 dpf) progenitors from two sections of the tree.

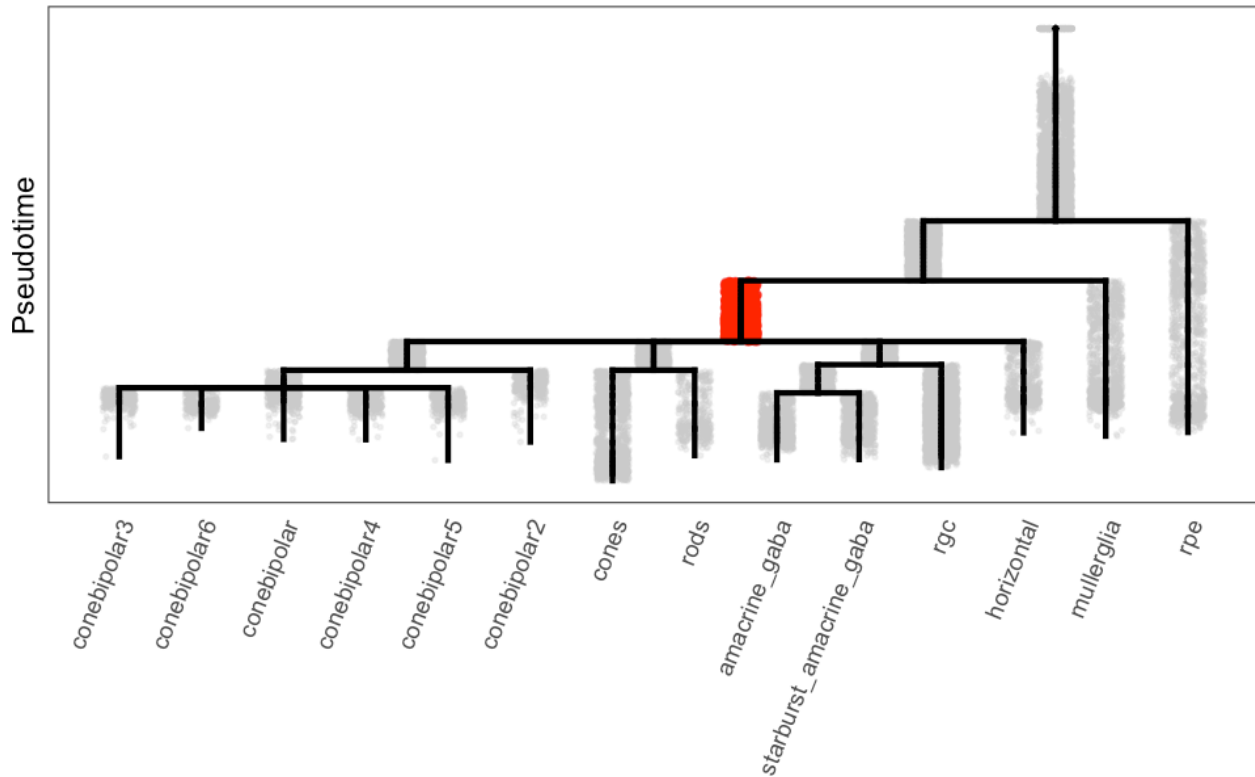
```
# Progenitors / 24-36 hpf / Segment 30
prog.early.s30 <- intersect(cellsInCluster(obj, "stage", c("06-24h", "07-36h")),
  cellsInCluster(obj, "segment", "30"))
obj <- groupFromCells(obj, group.id = "prog.early.s30", cells = prog.early.s30)
plotTreeHighlight(obj, "prog.early.s30", "TRUE", highlight.size = 1, highlight.alpha = 0.5,
  title = "Progenitors / 24-36 hpf / Segment 30")
```


Progenitors / 24-36 hpf / Segment 30



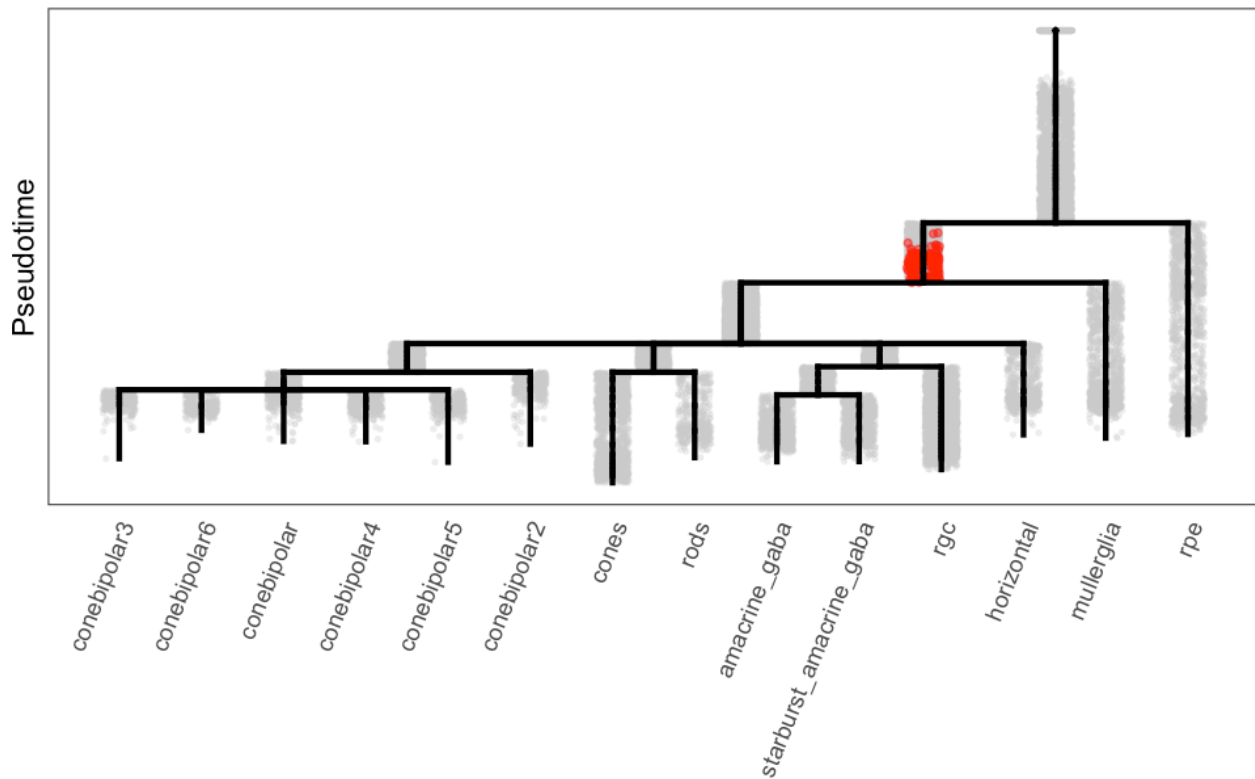
```
# Progenitors / 24-36 hpf / Segment 29
prog.early.s29 <- intersect(cellsInCluster(obj, "stage", c("06-24h", "07-36h")),
  cellsInCluster(obj, "segment", "29"))
obj <- groupFromCells(obj, group.id = "prog.early.s29", cells = prog.early.s29)
plotTreeHighlight(obj, "prog.early.s29", "TRUE", highlight.size = 1, highlight.alpha = 0.5,
  title = "Progenitors / 24-36 hpf / Segment 29")
```

Progenitors / 24-36 hpf / Segment 29



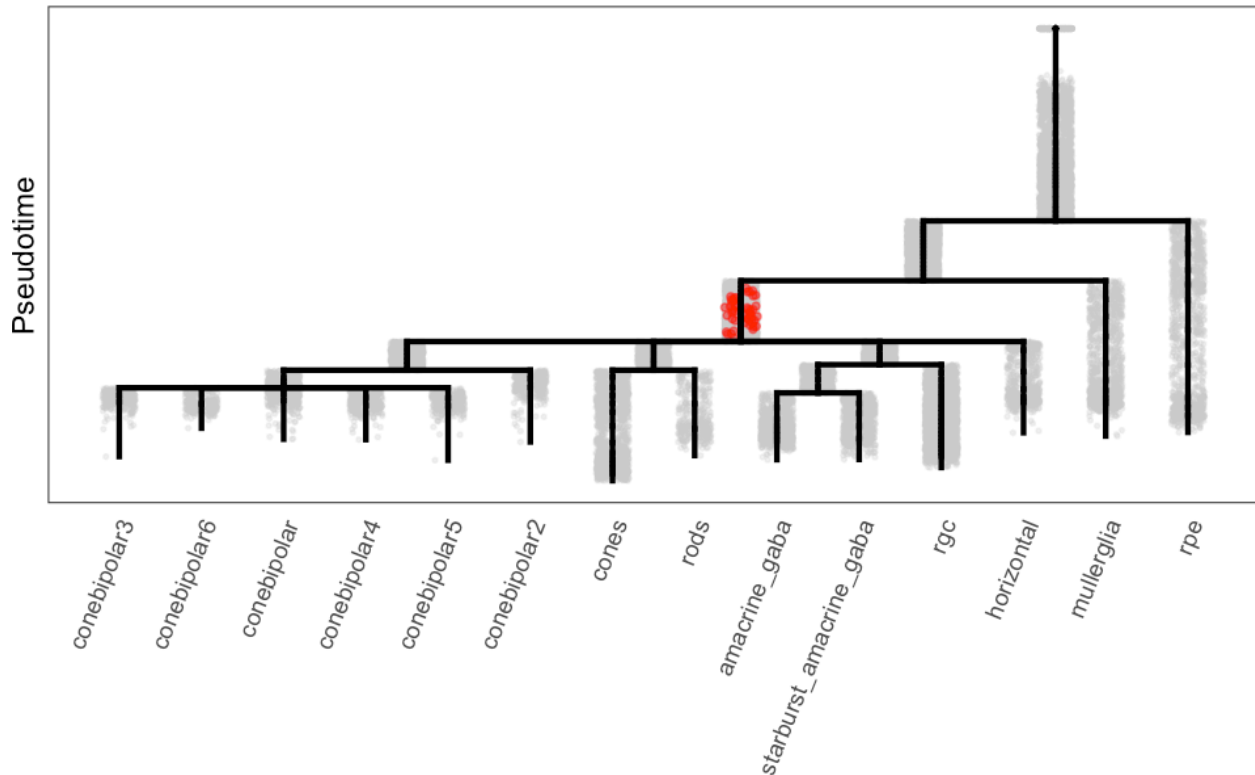
```
# Progenitors / 15 dpf / Cluster 39 / Segment 30
prog.15d.c39.s30 <- intersect(cellsInCluster(obj, "clus.orig", "12-15d_39"), cellsInCluster(obj,
"segment", "30"))
obj <- groupFromCells(obj, group.id = "prog.15d.c39.s30", cells = prog.15d.c39.s30)
plotTreeHighlight(obj, "prog.15d.c39.s30", "TRUE", highlight.size = 1, highlight.alpha = 0.5,
title = "Progenitors / 15 dpf / Cluster 39 / Segment 30")
```

Progenitors / 15 dpf / Cluster 39 / Segment 30



```
# Progenitors / 15 dpf / Cluster 39 / Segment 29
prog.15d.c39.s29 <- intersect(cellsInCluster(obj, "clus.orig", "12-15d_39"), cellsInCluster(obj,
"segment", "29"))
obj <- groupFromCells(obj, group.id = "prog.15d.c39.s29", cells = prog.15d.c39.s29)
plotTreeHighlight(obj, "prog.15d.c39.s29", "TRUE", highlight.size = 1, highlight.alpha = 0.5,
title = "Progenitors / 15 dpf / Cluster 39 / Segment 29")
```

Progenitors / 15 dpf / Cluster 39 / Segment 29



Differential expression between neural progenitor populations

Then we determined what was differentially expressed between them.

```
# Compare 15d ('late') vs. early - S29
markers.nb.lve.s29 <- markersAUCPR(obj, cells.1 = prog.15d.c39.s29, cells.2 = prog.early.s29,
  auc.factor = 1.1, effect.size = 0.4)

# Compare 15d ('late') vs. early - S30
markers.nb.lve.s30 <- markersAUCPR(obj, cells.1 = prog.15d.c39.s30, cells.2 = prog.early.s30,
  auc.factor = 1.1, effect.size = 0.4)
```

boot.fc

- **object**: An URD object
- **cells.1**: Cells from group 1 of the differential expression
- **cells.2**: Cells from group 2 of the differential expression
- **cells.segment**: All cells in the segment that can be pulled for bootstrapping
- **genes.test**: Genes to test in the bootstrapping
- **exp.fc**: Exp.fc from the original differential expression test to compare for bootstrap
- **exp.data**: Can pre-calculated un-logged expression data to pass to the function (`getUPXData`)
- **n**: (Numeric) Number of bootstrap simulations to run
- Returns list: p is the empirical p-value for each differential expression, `boot.fc` contains all of the test information.

```
# Function to bootstrap fold-change
boot.fc <- function(object, cells.1, cells.2, cells.segment, genes.test, exp.fc,
```

```

exp.data = NULL, n = 1000) {

# Pull random populations of equivalent sizes
l1 <- length(cells.1)
l2 <- length(cells.2)
random.pops <- lapply(1:1000, function(i) {
  y <- sample(x = cells.segment, size = l1 + l2, replace = F)
  return(list(a = y[1:l2], b = y[(l2 + 1):(l1 + l2)]))
})

# Get un-logged expression data, if not provided
if (is.null(exp.data))
  exp.data <- getUPXData(object)

# Calculate the expression fold-change for each random population
fc.boot <- as.data.frame(lapply(1:n, function(i) {
  exp.a <- exp.data[genes.test, random.pops[[i]][["a"]]
  exp.b <- exp.data[genes.test, random.pops[[i]][["b"]]
  exp.fc <- log2((rowMeans(exp.a)/rowMeans(exp.b)) + 1)
  return(exp.fc)
}))
names(fc.boot) <- paste0("rep", 1:n)

# Figure out p-value (proportion of these that beat provided exp.fc)
beat.boot <- sweep(fc.boot, 1, exp.fc, ">")
p.boot <- rowSums(beat.boot)/n

# Return information
return(list(p = p.boot, boot.fc = fc.boot))
}

```

Empirical p-value

Because these are relatively small populations, there's a decent chance that (due to the variability and noise inherent in scRNAseq data) that choosing any two similarly sized populations would find a number of differentially expressed genes also. Thus, we used an empirically-determined p-value to limit ourselves to differentially expressed genes that probably wouldn't arise by chance. We asked that our real comparison had a greater expression fold-change than two populations from a given segment of the same size chosen at random at least 99% of the time (i.e. $p < 0.01$).

```

# Try a bootstrapping approach to determine which markers are real, vs. which
# ones would arise just from small number of compared cells. Going to just do it
# on expression fc, so that the computation is reasonably fast.

# Isolate cells from each segment
cells.seg.29 <- cellsInCluster(obj, "segment", "29")
cells.seg.30 <- cellsInCluster(obj, "segment", "30")

# Get un-logged expression data to pass to the function
exp.data <- getUPXData(obj)

# Run the actual bootstrapping.
boot.s30.lve <- boot.fc(object, cells.1 = prog.15d.c39.s30, cells.2 = prog.early.s30,
  cells.segment = cells.seg.30, genes.test = rownames(markers.nb.lve.s30), exp.fc = markers.nb.lve.s30$exp.fc,

```

```

exp.data = exp.data, n = 1000)

boot.s29.lve <- boot.fc(object, cells.1 = prog.15d.c39.s29, cells.2 = prog.early.s29,
  cells.segment = cells.seg.29, genes.test = rownames(markers.nb.lve.s29), exp.fc = markers.nb.lve.s29$exp.fc,
  exp.data = exp.data, n = 1000)

# Limit markers to those that pass the bootstrap test
markers.nbb.lve.s30 <- markers.nb.lve.s30[which(boot.s30.lve$p <= 0.01), ]
markers.nbb.lve.s29 <- markers.nb.lve.s29[which(boot.s29.lve$p <= 0.01), ]

```

Tissue-specific changes

We also then divided genes based on whether they changed in all cells between 24/36 hpf and 15 dpf, or specifically in progenitors. Genes that change in all cells could represent either (a) global transcriptional changes in the tissue, or (b) changes in ambient RNA that is included with most cells based on highly expressed genes during different stages.

```

# Add global stage information to these – genes must change more in progenitors
# than just generally.

# Figure out early/late cells
cells.early <- cellsInCluster(obj, "stage", c("06–24h", "07–36h"))
cells.late <- cellsInCluster(obj, "stage", "12–15d")

# Calculate markers across stages generally with no restrictions
markers.nball.lve <- markersAUCPR(object = obj, cells.1 = cells.late, cells.2 = cells.early,
  effect.size = -Inf, frac.must.express = 0, auc.factor = 0, genes.use = unique(c(rownames(markers.nbb.lve.s30),
  rownames(markers.nbb.lve.s29))))

# Transfer information to NBB comparisons
markers.nbb.lve.s30$exp.fc.stage <- markers.nball.lve[rownames(markers.nbb.lve.s30),
  "exp.fc"]
markers.nbb.lve.s30$posFrac_stage1 <- markers.nball.lve[rownames(markers.nbb.lve.s30),
  "posFrac_1"]

markers.nbb.lve.s29$exp.fc.stage <- markers.nball.lve[rownames(markers.nbb.lve.s29),
  "exp.fc"]
markers.nbb.lve.s29$posFrac_stage1 <- markers.nball.lve[rownames(markers.nbb.lve.s29),
  "posFrac_1"]

# Calculate ratios (i.e. how much more does a gene change in progenitors than in
# the entire tissue)
markers.nbb.lve.s30$exp.fc.ratio <- pmin(markers.nbb.lve.s30$exp.fc, 1000) - pmin(markers.nbb.lve.s30$exp.fc.stage,
  1000)
markers.nbb.lve.s29$exp.fc.ratio <- pmin(markers.nbb.lve.s29$exp.fc, 1000) - pmin(markers.nbb.lve.s29$exp.fc.stage,
  1000)

markers.nbb.lve.s30$posFrac.ratio <- markers.nbb.lve.s30$posFrac_1/markers.nbb.lve.s30$posFrac_stage1
markers.nbb.lve.s29$posFrac.ratio <- markers.nbb.lve.s29$posFrac_1/markers.nbb.lve.s29$posFrac_stage1

```

Limit to well-expressed

We also limited ourselves to genes that had a decent level of expression. (In this case, they were detected in at least 20% of progenitor cells, and had a mean expression of at least 0.8.)

```
# All genes that change in segment 30
markers.nbbexp.lve.s30 <- markers.nbb.lve.s30[Reduce(intersect, list(which(markers.nbb.lve.s30$posFrac_1 >=
  0.2), which(markers.nbb.lve.s30$nTrans_1 >= 0.8))), ]

# All genes that change in segment 29
markers.nbbexp.lve.s29 <- markers.nbb.lve.s30[Reduce(intersect, list(which(markers.nbb.lve.s29$posFrac_1 >=
  0.2), which(markers.nbb.lve.s29$nTrans_1 >= 0.8))), ]

# Genes that change in segment 30 more than in the entire tissue
markers.nbbselect.lve.s30 <- markers.nbb.lve.s30[Reduce(intersect, list(which(markers.nbb.lve.s30$posFrac_1 >=
  0.2), which(markers.nbb.lve.s30$nTrans_1 >= 0.8), which(markers.nbb.lve.s30$exp.fc.ratio >=
  1.2), which(markers.nbb.lve.s30$posFrac.ratio >= 1.1))), ]

# Genes that change in segment 29 more than in the entire tissue
markers.nbbselect.lve.s29 <- markers.nbb.lve.s29[Reduce(intersect, list(which(markers.nbb.lve.s29$posFrac_1 >=
  0.2), which(markers.nbb.lve.s29$nTrans_1 >= 0.8), which(markers.nbb.lve.s29$exp.fc.ratio >=
  1.2), which(markers.nbb.lve.s29$posFrac.ratio >= 1.1))), ]
```

Result

That recovered a total of 71 genes that vary in progenitors between 24/36 hpf and 15 dpf, of which 16 change more in neural progenitors than the rest of the tissue.

```
# All genes that change in progenitors
unique(c(rownames(markers.nbbexp.lve.s29), rownames(markers.nbbexp.lve.s30)))
```

```
## [1] "hbbe2"          "hbz"           "ba1.1"
## [4] "rho"           "crabp1a"       "si:ch211-251b21.1"
## [7] "hbaa1"         "pde6h"        "tsc22d3"
## [10] "si:xx-by187g17.1" "ba1"          "rpe65a"
## [13] "zgc:153704"    "arr3a"        "lin7a"
## [16] "crygm1"        "gnat1"        "ptgdsb.1"
## [19] "gngt1"         "rgs16"        "cabp2a"
## [22] "junba"         "crygm2b"      "zgc:112320"
## [25] "si:dkey-183i3.5" "krt91"        "cabp5a"
## [28] "sagb"          "crygmx"       "scinla"
## [31] "rbp4l"         "gngt2b"       "rs1a"
## [34] "mt2"           "fosab"        "cebpd"
## [37] "snap25b"      "CNDP1"        "crabp2a"
## [40] "cryba4"        "jdp2b"        "cst3"
## [43] "higd1a"       "mt-nd3"       "si:dkey-16p21.8"
## [46] "crybb1"       "crygn2"       "gadd45ba"
## [49] "gapdhs"       "eno1a"        "mif"
## [52] "ggctb"        "ckbb"         "glula"
## [55] "tsc22d1"      "sod2"         "btg2"
## [58] "sod1"         "stmn1b"       "si:dkey-238o13.4"
## [61] "fabp11a"      "mdkb"         "gstp1"
## [64] "slc3a2b"      "si:dkey-238c7.12" "CABZ01102240.1"
## [67] "atp5ia"       "atpif1b"      "cadm3"
## [70] "h1f0"
```

```
# Genes that change in progenitors more than the rest of the tissue
unique(c(rownames(markers.nbbselect.lve.s29), rownames(markers.nbbselect.lve.s30)))

## [1] "si:ch211-114n24.6" "rps29"          "rrm2.1"
## [4] "si:ch211-193l2.6"  "si:dkey-238o13.4" "crabp1a"
## [7] "si:ch211-251b21.1" "CNDP1"          "junba"
## [10] "crabp2a"          "cryba4"         "crybb1"
## [13] "crygn2"           "fabp11a"        "si:dkey-238c7.12"
## [16] "cadm3"
```

Preservation of embryonic molecular profiles in larval progenitors

We were looking to compare the molecular profiles of embryonic and post-embryonic progenitors. Here, in the retina, we find progenitors at larval stages whose molecular signatures are preserved from embryonic stages. For comparison, in the hypothalamus, we find that progenitors at larval stages are transcriptionally different from embryonic progenitors (see Hypothalamus 3).

Identify populations to compare

We defined progenitor / precursor / neuron populations based on their location in the tree, cross-referenced with the expression of markers of each of these types

```
obj@group.ids$precursor.group <- NA

cells.s31 <- intersect(cellsInCluster(obj, "segment", "31"), whichCells(obj, "pseudotime",
  c(0.05, 1)))
obj@group.ids[cells.s31, "precursor.group"] <- "1a_prog_transient"

cells.prog.late <- cellsInCluster(obj, "segment", c("30", "29"))
obj@group.ids[cells.prog.late, "precursor.group"] <- "1b_prog_longterm"

cells.precursor <- intersect(cellsInCluster(obj, "segment", c("24", "25", "26", "15")),
  whichCells(obj, "pseudotime", c(0, 0.535)))
obj@group.ids[cells.precursor, "precursor.group"] <- "2_precursor"

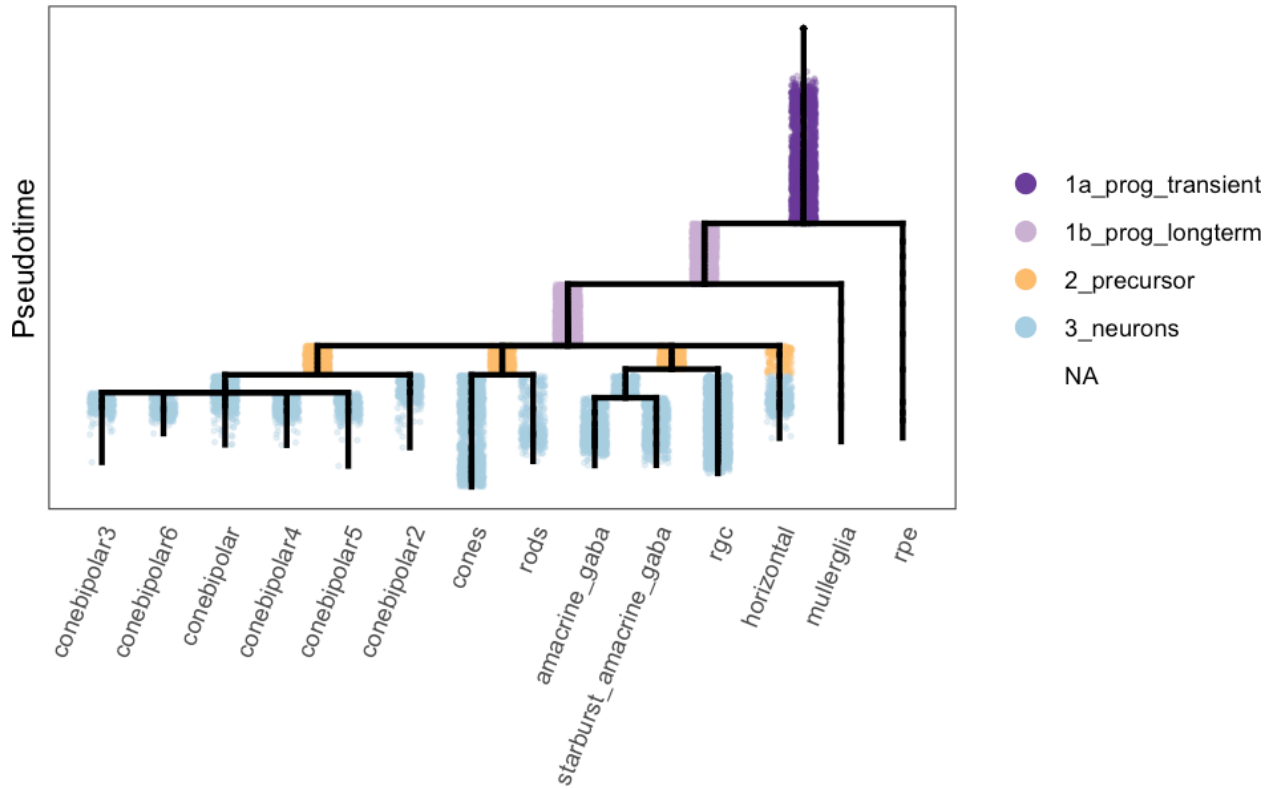
cells.neurons <- setdiff(whichCells(obj, "pseudotime", c(0.535, 1)), cellsInCluster(obj,
  "segment", c("6", "11")))
obj@group.ids[cells.neurons, "precursor.group"] <- "3_neurons"

# Colors for ggplot
stage.colors <- RColorBrewer::brewer.pal(12, "Paired")[c(10, 9, 7, 1)]

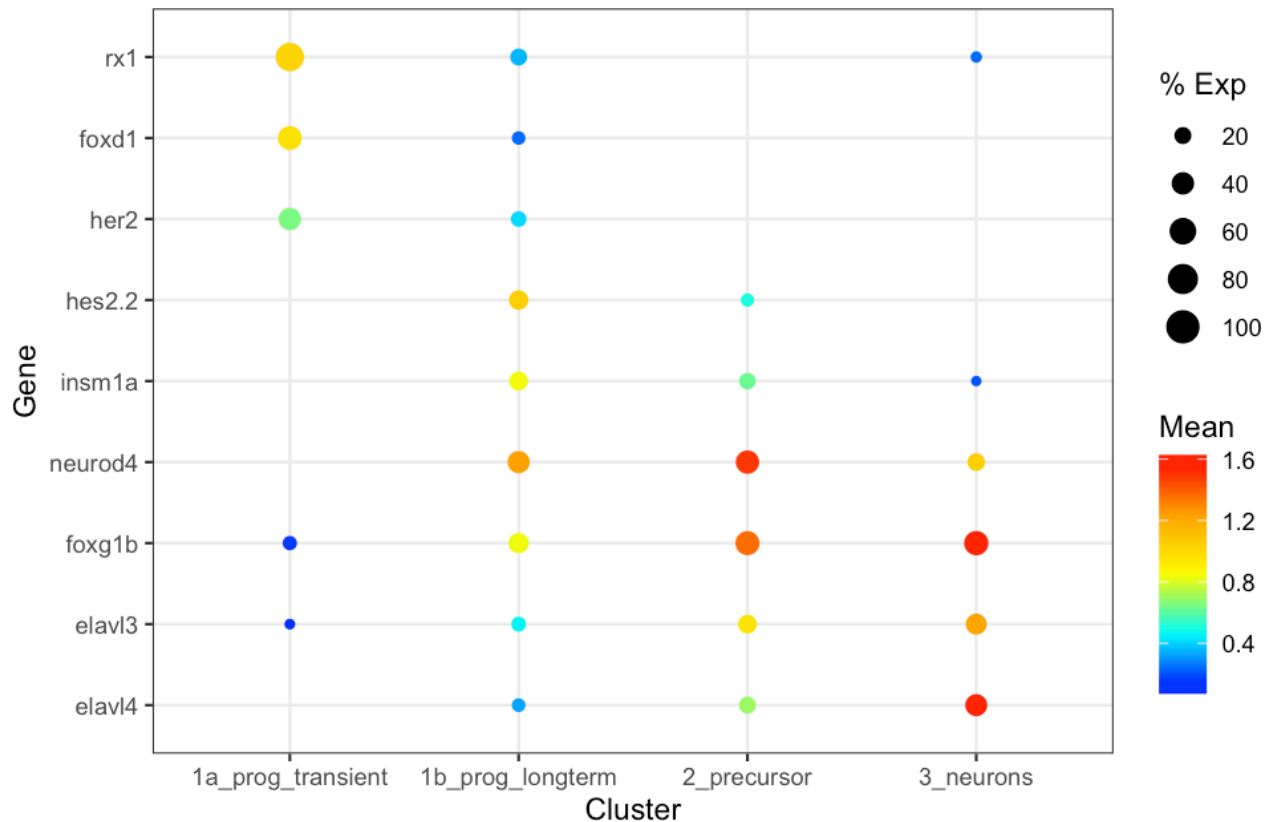
# Plot tree to show where the groups are
plotTree(obj, "precursor.group", discrete.colors = stage.colors)

## Warning: Removed 2530 rows containing missing values (geom_point).
```


precursor.group



```
# Plot genes in each group  
plotDot(obj, genes = c("rx1", "foxd1", "her2", "hes2.2", "insm1a", "neurod4", "foxg1b",  
  "elavl3", "elavl4"), clustering = "precursor.group", scale.by = "area") + theme_bw()
```



Determine proportion of cells in each state

We then determined the proportion of cells from different stages that fell into each of these transcriptional states.

```
# We combined stages to reduce number of plots
obj@group.ids$stage.collapsed <- plyr::mapvalues(x = obj@group.ids$stage, from = c("01-12h",
  "02-14h", "03-16h", "04-18h", "05-20h", "06-24h", "07-36h", "08-2d", "09-3d",
  "10-5d", "11-8d", "12-15d"), to = c(rep("01-12h-24h", 6), rep("02-36h-3d", 3),
  rep("03-5d-15d", 3)))
```

```
# Count number of cells from each stage group in each precursor group
stage.group.count <- plyr::count(obj@group.ids, vars = c("stage.collapsed", "precursor.group"))
```

```
# Remove cells that weren't part of a precursor group
stage.group.count <- stage.group.count[complete.cases(stage.group.count), ]
```

```
# Cast into a data frame and convert NA to 0 (no cells of that type observed)
stage.group.df <- reshape2::dcast(stage.group.count, formula = stage.collapsed ~
  precursor.group)
```

```
## Using freq as value column: use value.var to override.
```

```
stage.group.df[is.na(stage.group.df)] <- 0
```

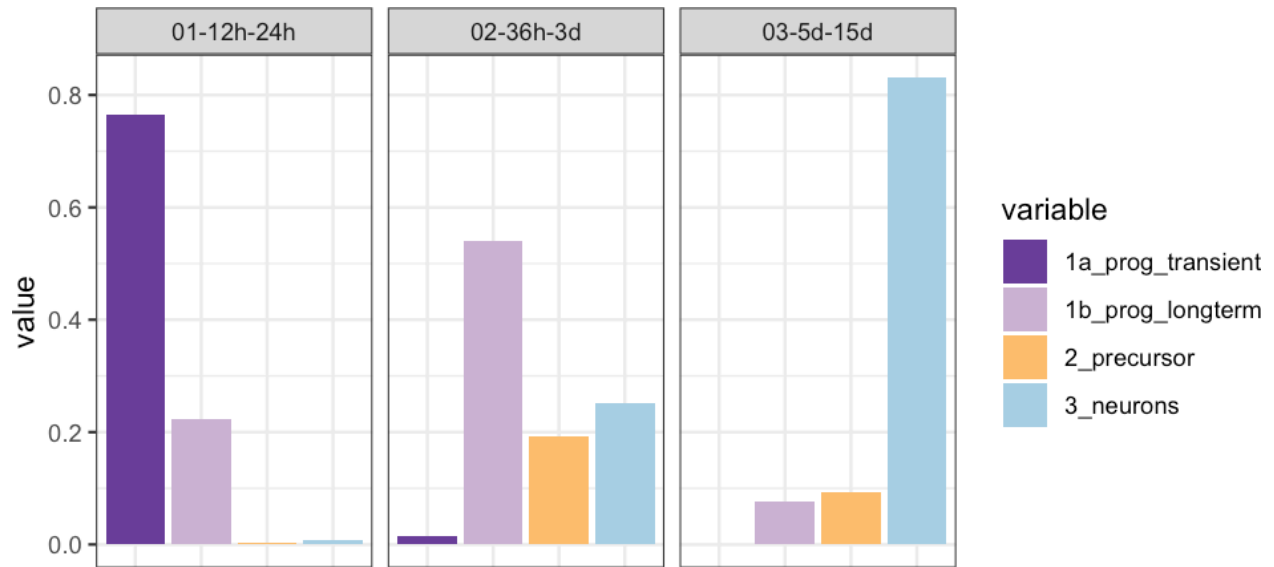
```
# Normalize by the number of precursors from each stage group
stage.group.df[, 2:5] <- sweep(stage.group.df[, 2:5], 1, rowSums(stage.group.df[,
  2:5]), "/")
```

```

# Melt for ggplot
stage.group.df.melt <- reshape2::melt(stage.group.df, id.vars = "stage.collapsed")

# Plot proportions
ggplot(stage.group.df.melt, aes(x = variable, y = value, group = stage.collapsed,
  fill = variable)) + geom_bar(stat = "identity") + facet_wrap(~stage.collapsed) +
  theme_bw() + scale_fill_manual(values = stage.colors) + theme(axis.title.x = element_blank(),
  axis.text.x = element_blank(), axis.ticks.x = element_blank())

```



Hypothalamus: 1 - URD object & doublet removal

Jeff Farrell

8/22/2019

Contents

Import data into URD	1
Convert Seurat object to URD	1
Combined individual stage clustering	1
Calculate highly variable genes	2
Calculate KNN graph and remove outliers	7
Remove cell type doublets	8
Add UMAP projection	8
Load NMF results and import into object	9
Select cell-type specific modules	9
Determine which module pairs to use for doublet removal	10

Import data into URD

Convert Seurat object to URD

We first loaded a Seurat object that contained just cells from the clusters that belonged to the hypothalamus from each stage.

```
suppressPackageStartupMessages(library(URD))
suppressPackageStartupMessages(library(Seurat))

base.path <- "~/urd-cluster-bushra/"

# Load Seurat object that has been cropped to hypothalamus cells
object.seurat <- readRDS(paste0(base.path, "obj/hypo_seurat.rds"))

# Convert to URD object
suburd <- seuratToURD(object.seurat)
```

Combined individual stage clustering

Bushra had performed individual clusterings across each stage with different resolutions. Here, it was better to create a single identifier that included stage + cluster information to combine all those clusterings (while preventing any overlap).

```

stages <- sort(unique(suburd@meta$stage))
clust.res.used <- paste0("res.", c("4.5", "4", "5", "5", "4.5", "5", "6",
  "6", "6", "5.5", "6", "5"))
names(clust.res.used) <- stages
suburd@group.ids$cluster <- NA
for (stage in stages) {
  suburd@group.ids[cellsInCluster(suburd, "stage", stage), "cluster"] <- paste0(stage,
    "-", suburd@group.ids[cellsInCluster(suburd, "stage", stage), clust.res.used[stage]])
}

```

Calculate highly variable genes

We calculated highly variable genes for each stage, used genes that were found as highly in at least two stages, but were not mitochondrial, ribosomal, heat-shock protein, or tandem duplicated genes.

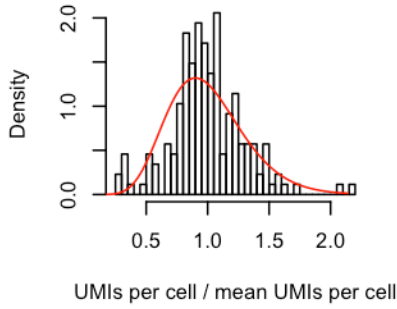
```

# Calculated on each stage separaely, final gene list was all genes
# that were 'variable' in at least two stages NB: For a couple of
# stages, the gamma fit was poor -- the library size distribution
# seemed bimodal. Have seen this before in 10X data, but not sure what
# it means.
var.genes.by.stage <- lapply(stages, function(stage) {
  findVariableGenes(suburd, cells.fit = cellsInCluster(suburd, "stage",
    stage), set.object.var.genes = F, diffCV.cutoff = 0.3, main.use = stage,
    do.plot = T)
})

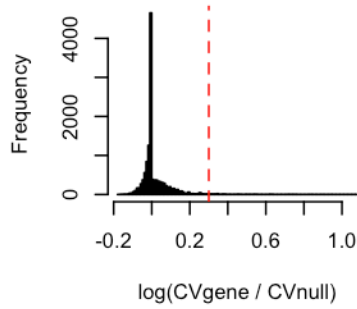
```

01-12h

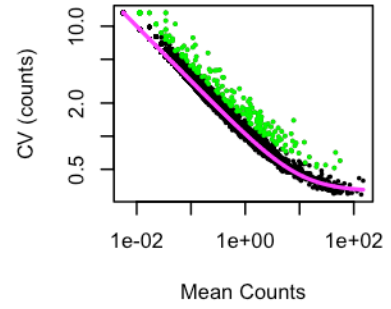
Size Factors & Gamma Fit ($\alpha=10.1$)



Diff CV

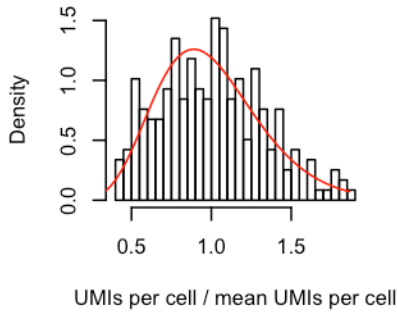


Selection of Variable Genes

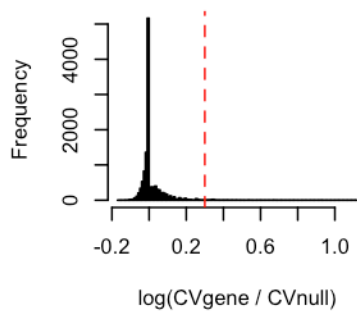


02-14h

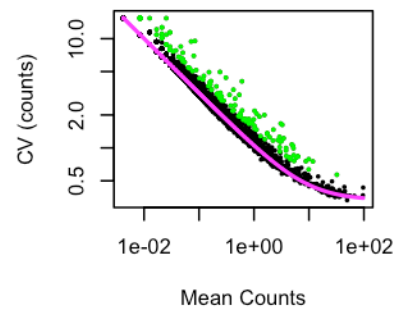
Size Factors & Gamma Fit ($\alpha=9.0$)



Diff CV

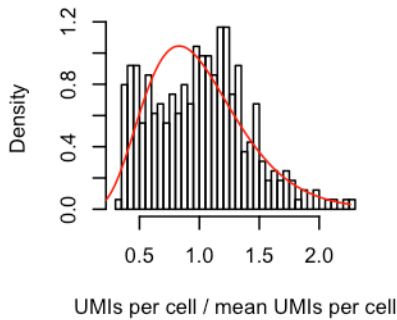


Selection of Variable Genes

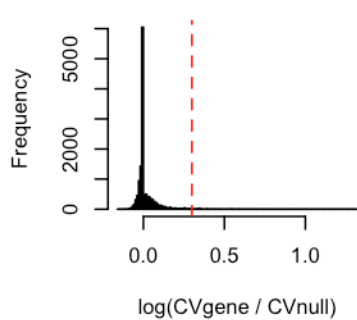


03-16h

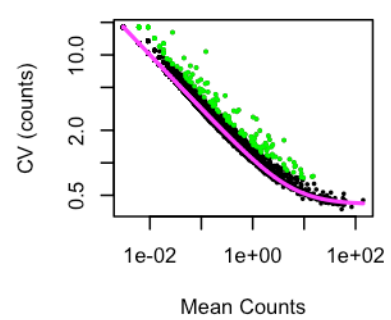
Size Factors & Gamma Fit ($\alpha=5.1$)



Diff CV

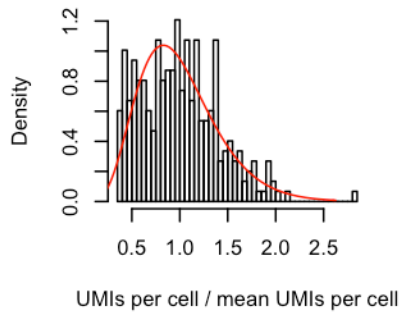


Selection of Variable Genes

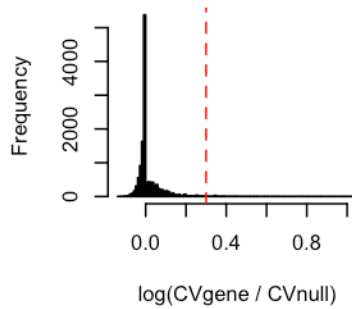


04-18h

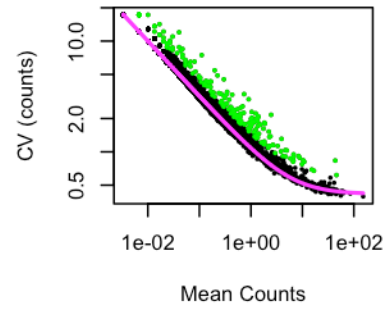
Size Factors & Gamma Fit ($\alpha=5.8$)



Diff CV

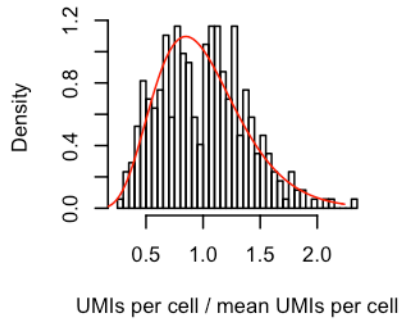


Selection of Variable Genes

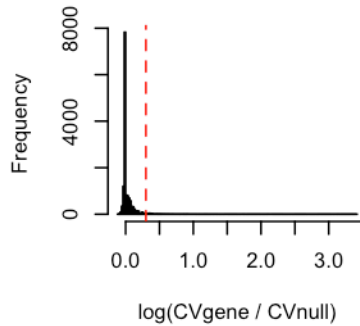


05-20h

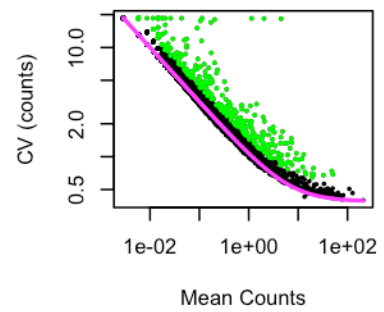
Size Factors & Gamma Fit ($\alpha=6.6$)



Diff CV

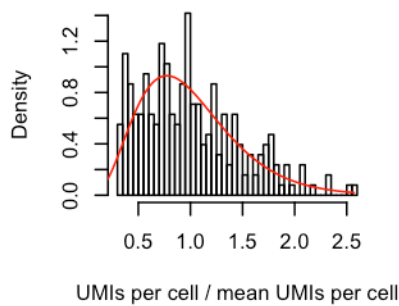


Selection of Variable Genes

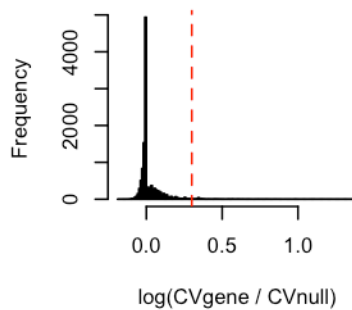


06-24h

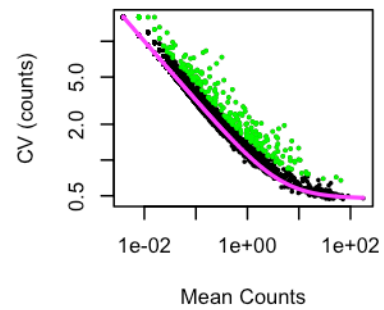
Size Factors & Gamma Fit ($\alpha=4.4$)



Diff CV

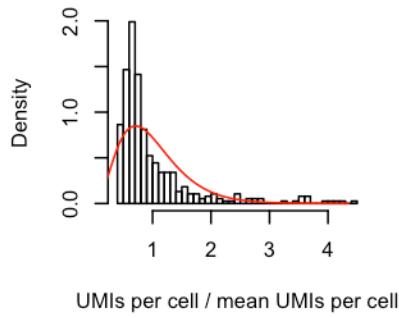


Selection of Variable Genes

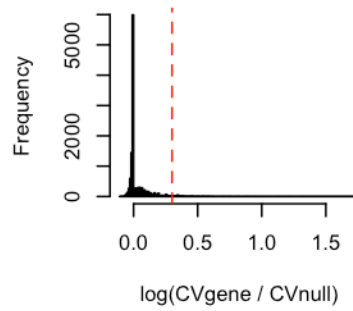


07-36h

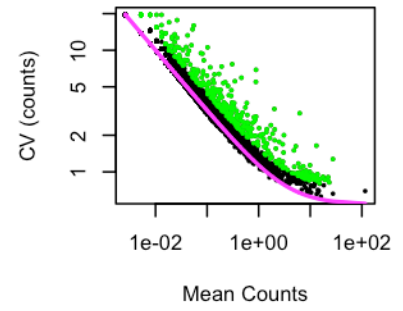
Size Factors & Gamma Fit ($\alpha=3.4$)



Diff CV

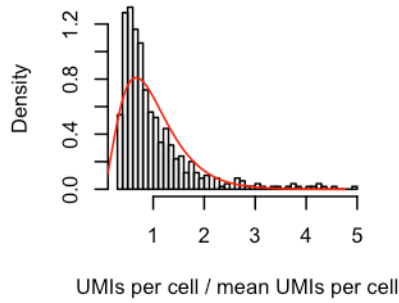


Selection of Variable Genes

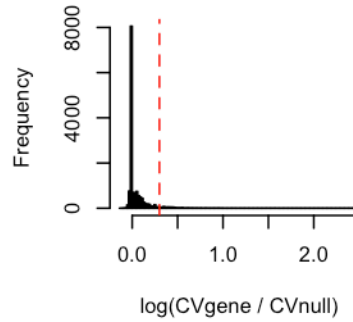


08-2d

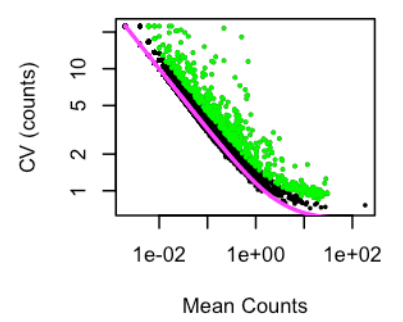
Size Factors & Gamma Fit ($\alpha=2.9$)



Diff CV

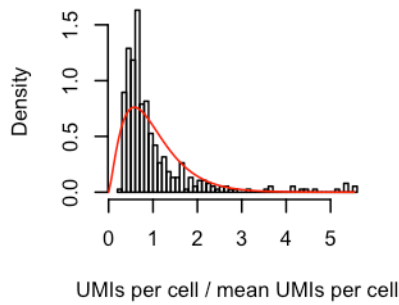


Selection of Variable Genes

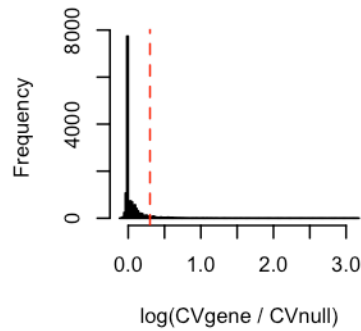


09-3d

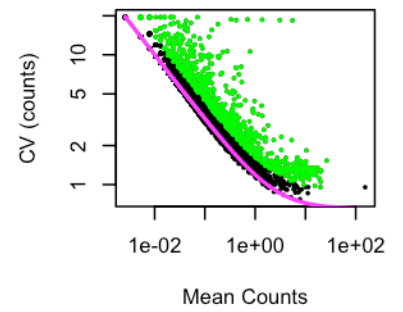
Size Factors & Gamma Fit ($\alpha=2.3$)



Diff CV

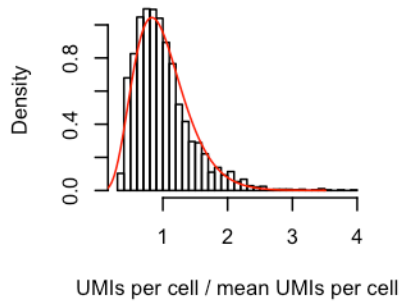


Selection of Variable Genes

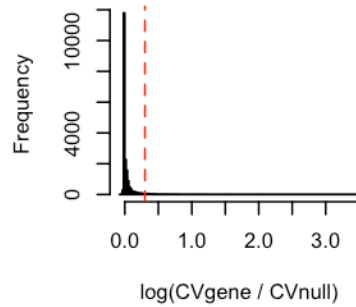


10-5d

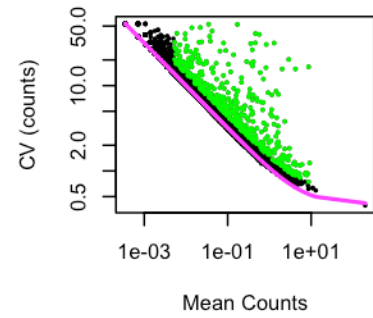
Size Factors & Gamma Fit (a=5.8)



Diff CV

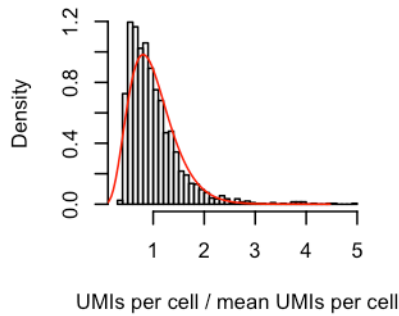


Selection of Variable Genes

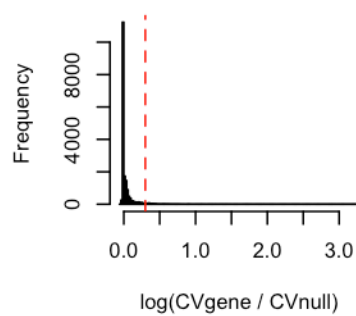


11-8d

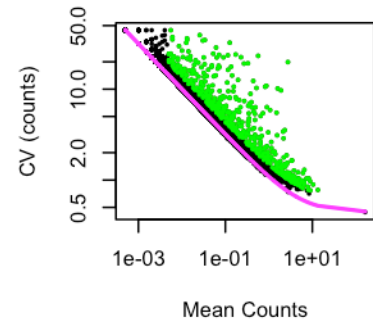
Size Factors & Gamma Fit (a=5.0)



Diff CV

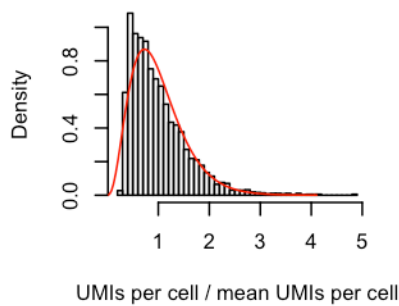


Selection of Variable Genes

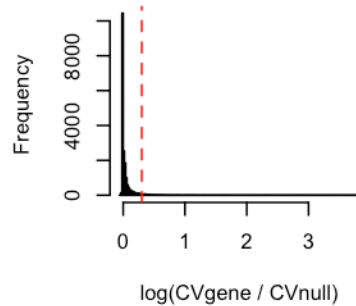


12-15d

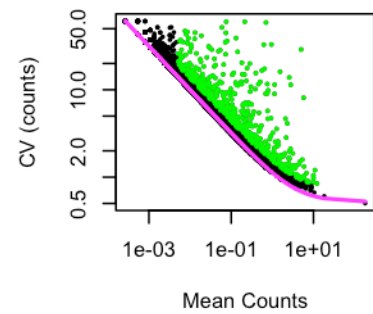
Size Factors & Gamma Fit (a=3.6)



Diff CV



Selection of Variable Genes



```
names(var.gen.es.by.stage) <- stages
var.gen.es <- sort(unique(unlist(var.gen.es.by.stage)))
print(paste0("Length of variable genes is ", length(var.gen.es)))
```

```
## [1] "Length of variable genes is 1783"
```

```
var.gen.es.twice <- names(which(table(unlist(var.gen.es.by.stage)) >= 2))
print(paste0("Length of variable genes shared across at least 2 stages is ",
length(var.gen.es.twice)))
```

```
## [1] "Length of variable genes shared across at least 2 stages is 957"
```

```
# Remove mitochondrial genes
var.mito <- grep("^mt-|^AC0", var.genes.twice, value = T)
# Remove ribosomal genes
var.ribo <- grep("^rps|^rpl", var.genes.twice, value = T)
# Remove hsp genes
var.hsp <- grep("^hsp", var.genes.twice, value = T)
# Remove genes with duplicates
var.dups <- grep("of many", var.genes.twice, value = T)
suburd@var.genes <- setdiff(var.genes.twice, c(var.mito, var.ribo, var.hsp,
var.dups))
print(paste0("Length of final variable genes list (after removing mito, ribo, hsp genes) is ",
length(suburd@var.genes)))
```

```
## [1] "Length of final variable genes list (after removing mito, ribo, hsp genes) is 856"
```

To prevent downstream problems, we also removed any cells from the data that had the exact same expression of the variable genes (i.e. cells with completely duplicated coordinates in the high-dimensional space we would use for analysis downstream).

```
# Check for duplicate data points - cells with exact same expression of
# variable genes
vg.dups <- duplicated(as.data.frame(as.matrix(t(suburd@logupx.data[suburd@var.genes,
])))
if (length(which(vg.dups)) > 0) {
  print(paste("Removing", length(which(vg.dups)), "cell(s) with duplicated variable gene expression."))
  not.dup.cells <- colnames(suburd@logupx.data)[!vg.dups]
  suburd <- urdSubset(suburd, not.dup.cells)
}
```

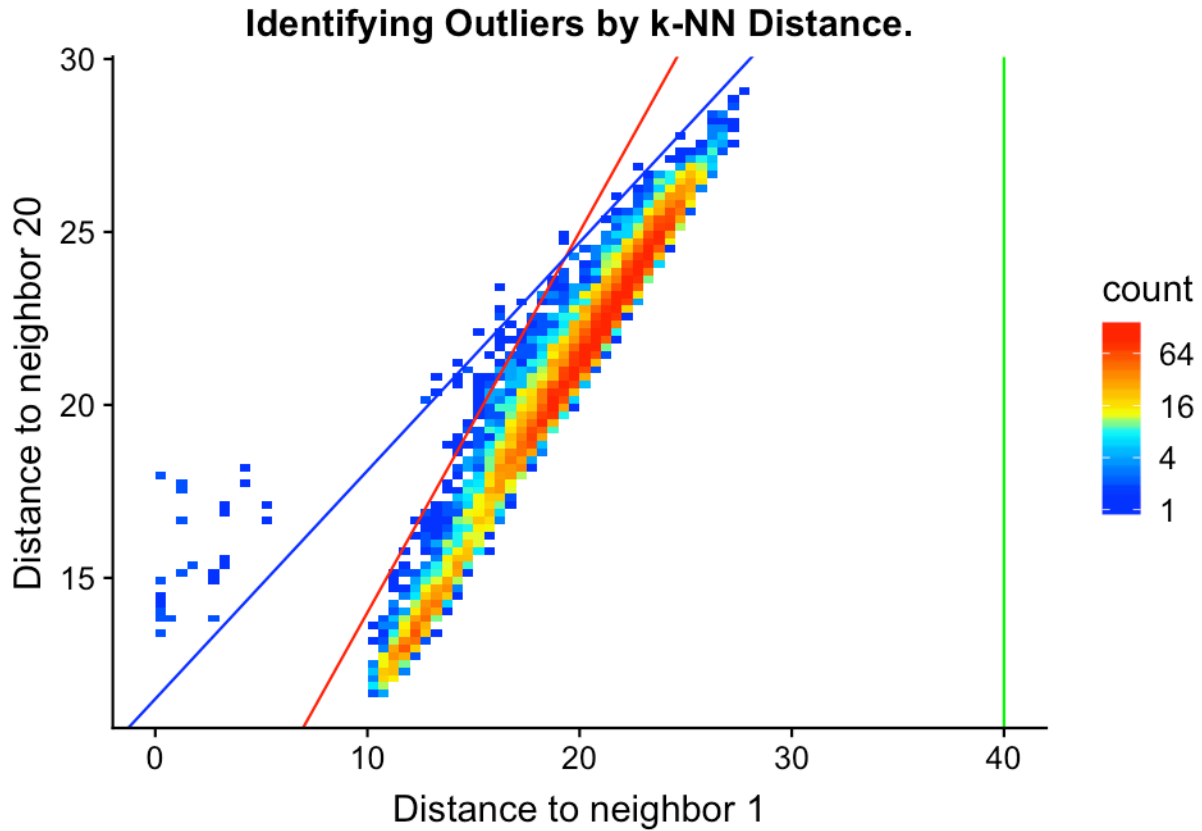
```
## [1] "Removing 1 cell(s) with duplicated variable gene expression."
```

Calculate KNN graph and remove outliers

We then calculated a k-nearest neighbor graph and removed cells that had unusual distance to their nearest neighbor, or unusual distance to their 20th nearest neighbor (given their distance to their nearest neighbor). These sorts of outliers often cause problems or skew diffusion maps (used downstream).

```
# Calculate k-nn
suburd <- calcKNN(suburd)

# Check what the outliers are
outliers <- knnOutliers(suburd, nn.1 = 1, nn.2 = 20, x.max = 40, slope.r = 1.1,
  int.r = 3, slope.b = 0.66, int.b = 11.5, title = "Identifying Outliers by k-NN Distance.")
```



```
length(outliers)
```

```
## [1] 87
```

```
suburd <- urdSubset(suburd, cells.keep = setdiff(colnames(suburd@logupx.data),
  outliers))
```

Remove cell type doublets

Add UMAP projection

While not strictly required, a UMAP projection can make it easier to assess the expression of NMF modules and whether thresholds for overlap are set correctly.

```
## add UMAP command
```

```
# Load pre-calculated UMAP
umap <- readRDS(paste0(base.path, "/umap/umap_hypo.rds"))
```

```
# Add projection to URD object
suburd@tsne.y <- umap
```

Load NMF results and import into object

NMF results were calculated by providing `suburd@logupx.data` to an external NMF pipeline written in Python. The output results are imported here, scaled, and added to the URD object.

```
# Load the NMF results
load(paste0(base.path, "/NMF/hypo/result_tbls.Robj"))

# The results object contains NMF runs for several K values. k=28 was
# chosen for this tissue, so this extracts the results for that
# particular parameter
k.use <- "28"
nmf.cells <- result_obj[[paste0("K=", k.use)]][[1]]$C
rownames(nmf.cells) <- paste0("nmf", 1:nrow(nmf.cells))
colnames(nmf.cells) <- gsub("\\.", "-", colnames(nmf.cells))
nmf.genes <- result_obj[[paste0("K=", k.use)]][[1]]$G
colnames(nmf.genes) <- paste0("nmf", 1:nrow(nmf.cells))

# Scale NMF results 0-1
nmf.cells.scaled <- sweep(nmf.cells, 1, apply(nmf.cells, 1, max), "/")

# Add scaled NMF results to the URD object
suburd@nmf.c1 <- as(t(as.matrix(nmf.cells.scaled)), "dgCMatrix")
```

Select cell-type specific modules

Several NMF modules will be poor markers of cell types — these are often modules driven mostly by the expression of 1-2 genes (where the gene loading of the first gene is much greater than that of the fourth gene, for instance), or modules that don't exhibit any restriction in a tSNE or UMAP projection.

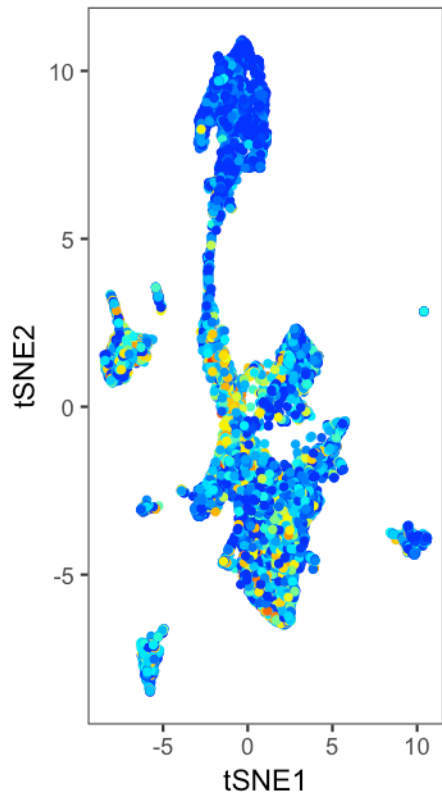
```
# Plot size parameters
plot.height = 6
plot.width = 16
dpi = 150

# Plot every module to determine which exhibit cell-type specificity
# This saves directly to the hard drive: two example plots are shown
# below.

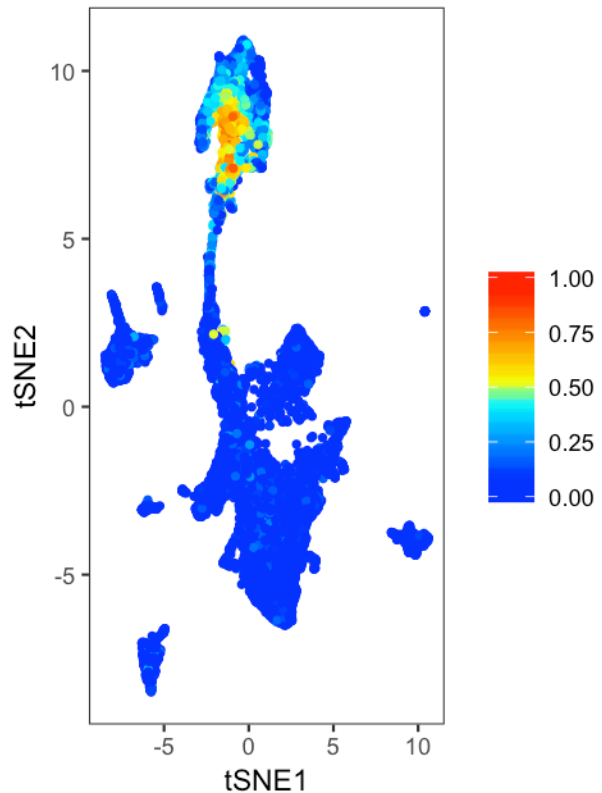
# for (n in colnames(suburd@nmf.c1)) { png(paste0(path, '/doublets/',
# subset, '-plots/', n, '.png'), width=dpi*plot.width,
# height=dpi*plot.height) plot(plotDim(suburd, n)) dev.off() }

gridExtra::grid.arrange(grobs = list(plotDim(suburd, "nmf2", plot.title = "nmf2: exhibits poor restriction"),
  plotDim(suburd, "nmf27", plot.title = "nmf27: exhibits good cell-type")),
  ncol = 2)
```

nmf2: exhibits poor restriction



nmf27: exhibits good cell-type



```
# Module Gene 1 : Gene 4 Ratios
top.genes <- result_obj[[paste0("K=", k.use)]][[1]]$top30genes
top.weights <- top.genes[, grep("Weights", colnames(top.genes), value = T)]
colnames(top.weights) <- paste0("nmf", 1:nrow(nmf.cells))
top.weights.ratio <- top.weights[1, ]/top.weights[4, ]

# Which modules exhibit cell-type restriction?
modules.bad.ratio <- names(top.weights.ratio)[which(top.weights.ratio >
3)]
unrestricted.modules <- paste0("nmf", c("12", "19", "24", "28"))
good.modules <- setdiff(colnames(suburd@nmf.c1), c(modules.bad.ratio, unrestricted.modules))
```

Determine which module pairs to use for doublet removal

We consider NMF modules pairwise and only use those pairs that don't are non-overlapping in the data. (In other words, NMF modules that are mutually exclusive in the majority of the data.) Here, we determine thresholds for selecting those module pairs.

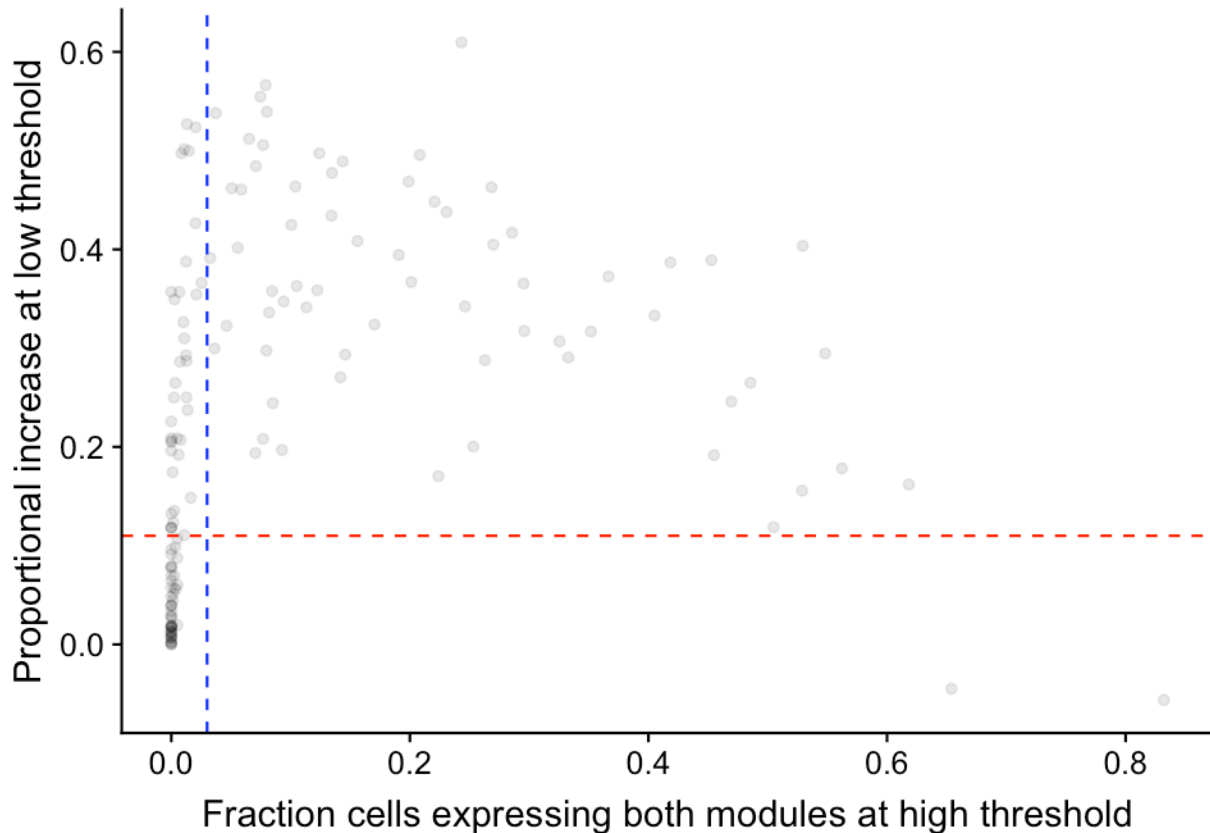
```
# Determine overlaps between module pairs
nmf.doublet.combos <- NMFDoubletsDefineModules(suburd, modules.use = good.modules,
module.thresh.high = 0.4, module.thresh.low = 0.15)

# Determine thresholds for NMF modules
frac.overlap.max = 0.03
frac.overlap.diff.max = 0.11
```

```
module.expressed.thresh = 0.33
```

```
# Determine which module pairs to use for doublets
```

```
NMFDoubletsPlotModuleThresholds(nmf.doublet.combos, frac.overlap.max = frac.overlap.max,  
  frac.overlap.diff.max = frac.overlap.diff.max)
```



```
# These commands save plots directly to the hard-drive.
```

```
# Make plots to see how your thresholds are
```

```
NMFDoubletsPlotModuleCombos(suburd, path = paste0(path, "/doublets/"), subset,  
  "-doublet-combos/"), module.combos = nmf.doublet.combos, module.expressed.thresh = module.expressed.thresh,  
  frac.overlap.max = frac.overlap.max, frac.overlap.diff.max = frac.overlap.diff.max,  
  boundary = "pass", sort = "near", n.plots = 25)
```

```
NMFDoubletsPlotModuleCombos(suburd, path = paste0(path, "/doublets/"), subset,  
  "-ok-combos/"), module.combos = nmf.doublet.combos, module.expressed.thresh = module.expressed.thresh,  
  frac.overlap.max = frac.overlap.max, frac.overlap.diff.max = frac.overlap.diff.max,  
  boundary = "discarded", sort = "near", n.plots = 25)
```

```
# Define doublet cells
```

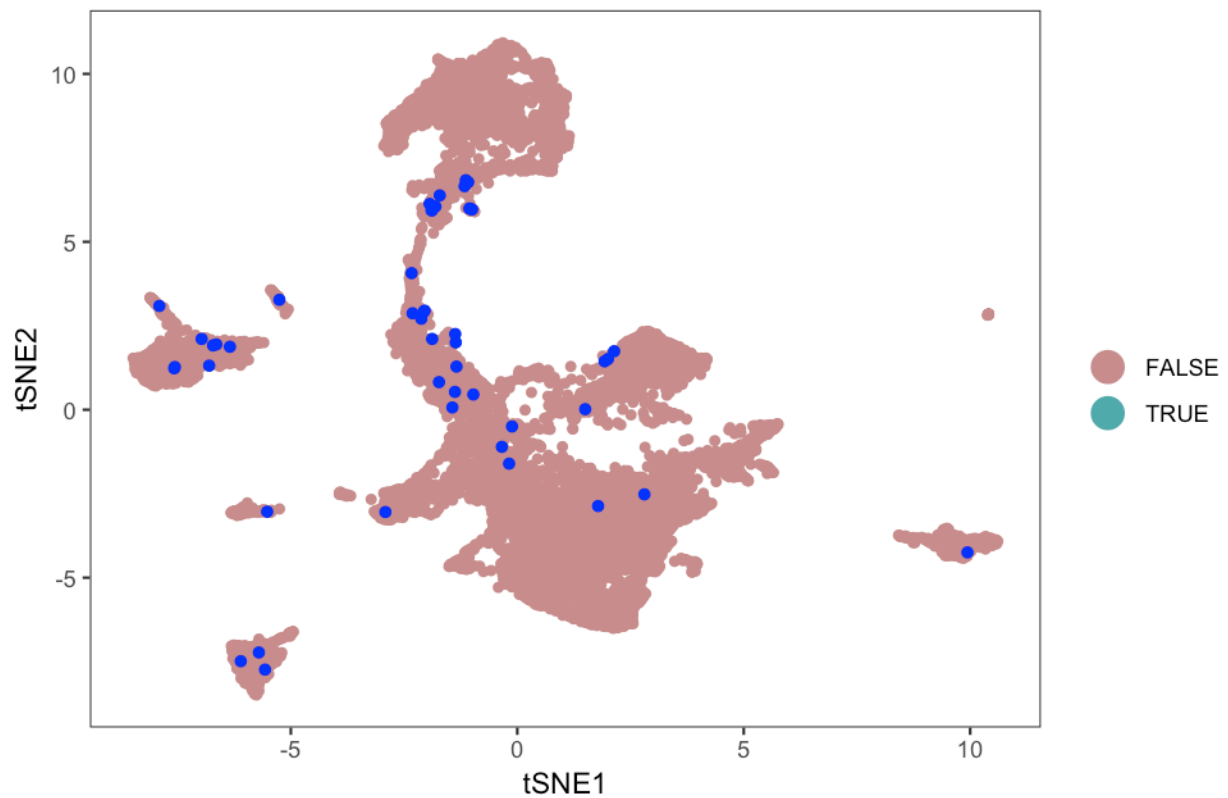
```
nmf.doublets <- NMFDoubletsDetermineCells(suburd, nmf.doublet.combos, module.expressed.thresh = module.expressed.thresh,  
  frac.overlap.max = frac.overlap.max, frac.overlap.diff.max = frac.overlap.diff.max) # 49 cells / 11307 cells
```

```
# Plot doublet cells on the UMAP
```

```
suburd <- groupFromCells(suburd, "nmf.doublets", cells = nmf.doublets)  
plot(plotDimHighlight(suburd, clustering = "nmf.doublets", cluster = "TRUE",  
  plot.title = paste0("NMF doublets: ", length(nmf.doublets), " cells"),
```

```
point.size = 2, highlight.color = "blue"))
```

NMF doublets: 49 cells (Highlight TRUE)



```
# Crop object to exclude doublets
```

```
suburd.cropped <- urdSubset(suburd, cells.keep = setdiff(colnames(suburd@logupx.data),  
nmf.doublets))
```

And then save the completed object for use downstream in building a tree using URD.

```
saveRDS(suburd.cropped, file = paste0(base.path, "/obj/URD_hypo_ND.rds"))
```

Hypothalamus: 2 - URD tree

Jeff Farrell

9/07/2019

Contents

Load data	1
Processed on the cluster	1
Calculate diffusion map and pseudotime	2
Calculate diffusion map	2
Calculate pseudotime	4
Calculate biased transition matrix	9
Perform biased random walks	9
Determine tips	9
Perform the biased random walks	11
Process the random walks	12
Build the URD tree	12
Save the URD tree	13

Load data

```
suppressPackageStartupMessages(library(URD))
suppressPackageStartupMessages(library(Seurat))

base.path <- "~/urd-cluster-bushra/"

# Load processed URD object
object <- readRDS(paste0(base.path, "obj/URD_hypo_ND.rds"))
```

Processed on the cluster

Most of the following steps were run on a computing cluster. These individual tissue subsets can be run on a modern, well-equipped laptop. The use of a computing cluster allows multiple parameter choices to be tried in parallel, and also allows further parallelization of the random walk procedure, speeding it up. Below, we should the commands that one would run on their laptop, and then generally load the pre-processed results from the cluster that were used in the paper. If you want to parallelize your own processing on a compute cluster, the scripts we used will be available at <http://github.com/farrellja/URD/cluster/>

Calculate diffusion map and pseudotime

These two steps are run in the cluster script URD-DM-PT.R.

Calculate diffusion map

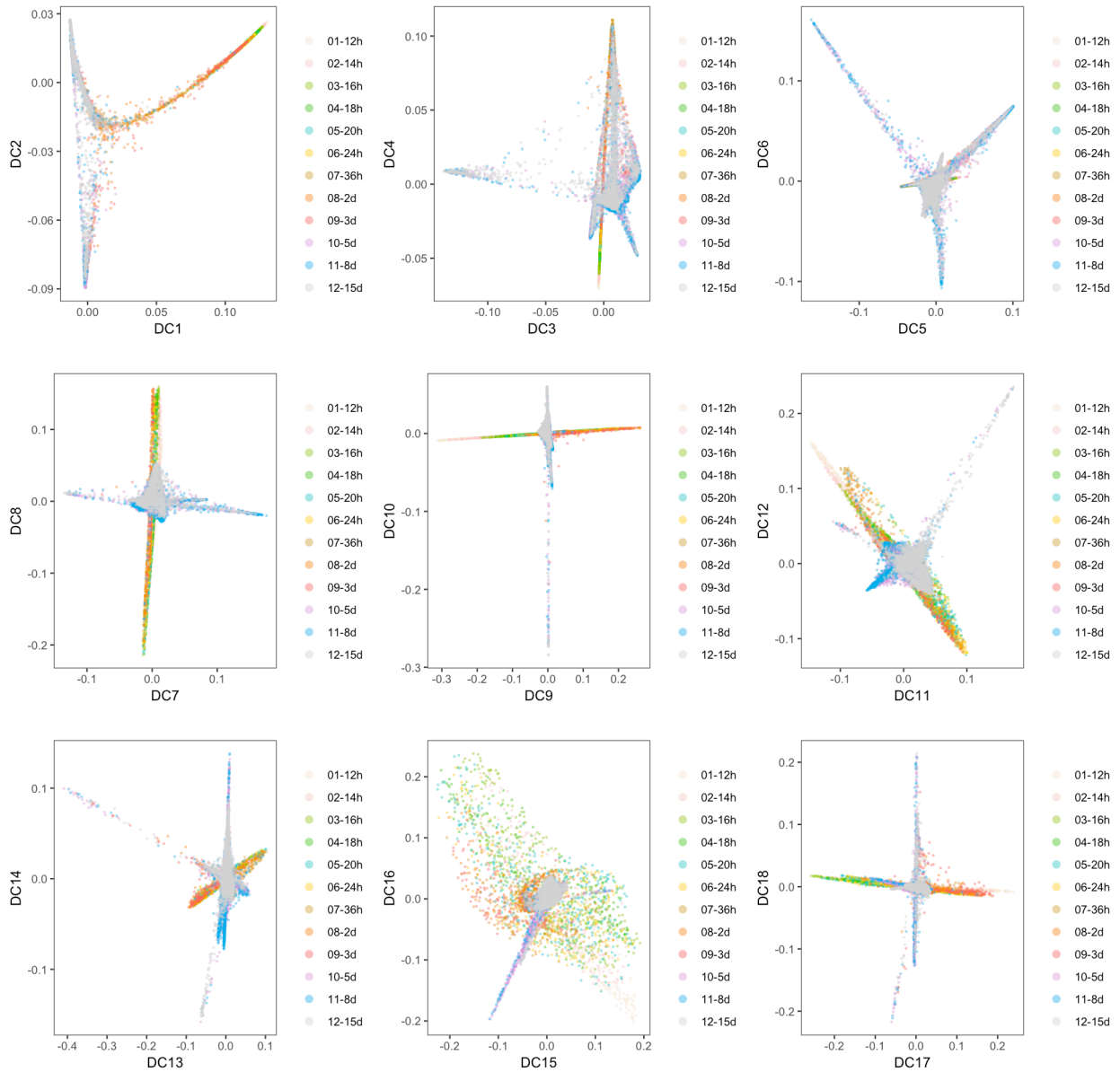
```
# To run locally: Calculate a diffusion map projection
object <- calcDM(object, knn = 100, sigma.use = 8)

# Or: Load a pre-computed diffusion map projection
dm <- readRDS(paste0(base.path, "dm/dm_hypoND_knn-100_sigma-8.rds"))
object <- importDM(object, dm)

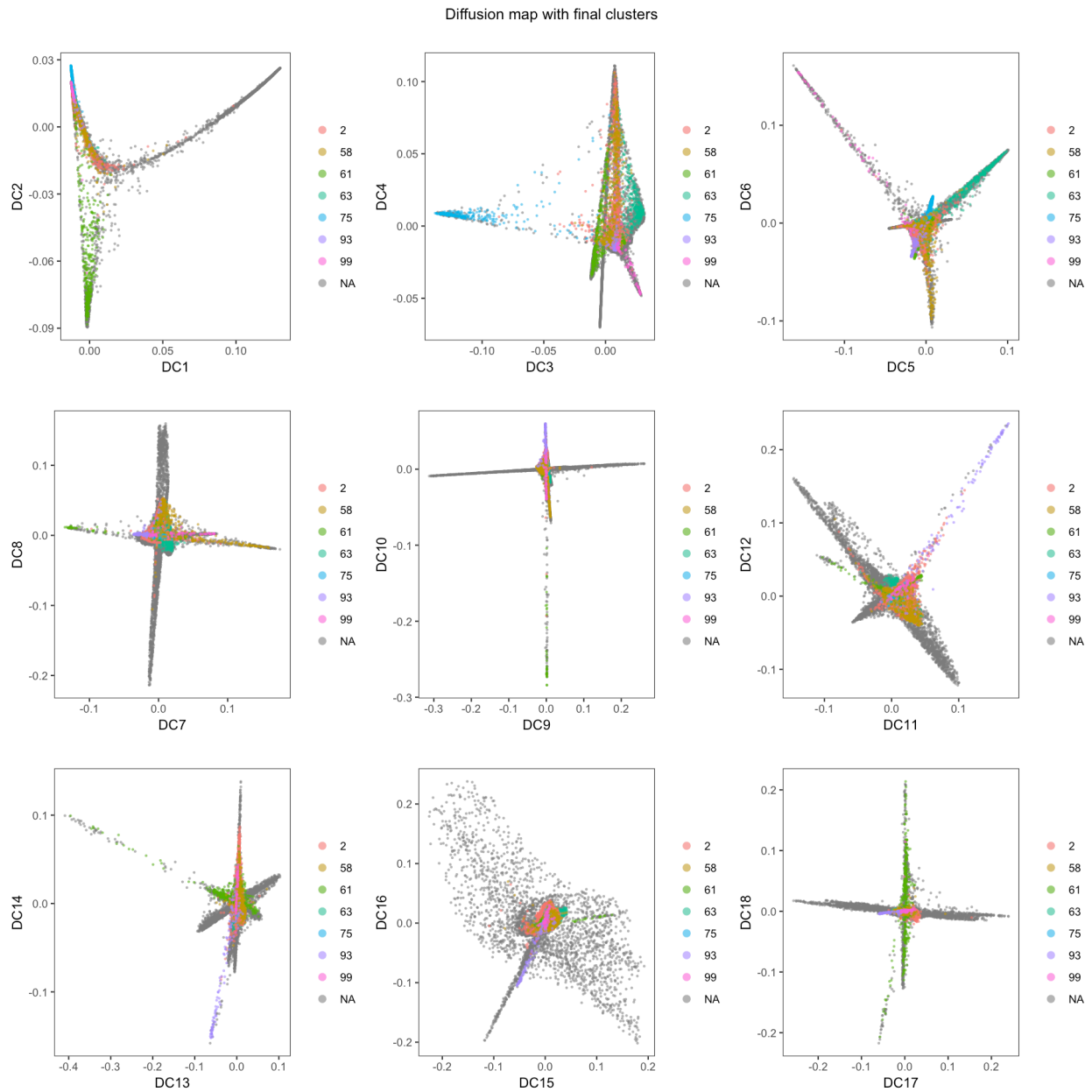
# Plot diffusion maps
stage.colors <- c("antiquewhite", "#FFCCCC", "#99CC00", "#33CC00", "cyan3",
  "gold", "goldenrod", "darkorange", "indianred1", "plum", "deepskyblue2",
  "lightgrey")

# Plot by stage
plotDimArray(object = object, reduction.use = "dm", dims.to.plot = 1:18,
  label = "stage", plot.title = "", outer.title = "Diffusion map labeled by Stage",
  legend = T, alpha = 0.45, discrete.colors = stage.colors)
```

Diffusion map labeled by Stage



```
# Plot with final cell types labeled
object@group.ids$final.cluster <- NA
object@group.ids[cellsInCluster(object, "stage", "12-15d"), "final.cluster"] <- object@group.ids[cellsInCluster(
  "stage", "12-15d"), "res.5"]
plotDimArray(object = object, reduction.use = "dm", dims.to.plot = 1:18,
  label = "final.cluster", plot.title = "", outer.title = "Diffusion map with final clusters",
  legend = T, alpha = 0.6)
```

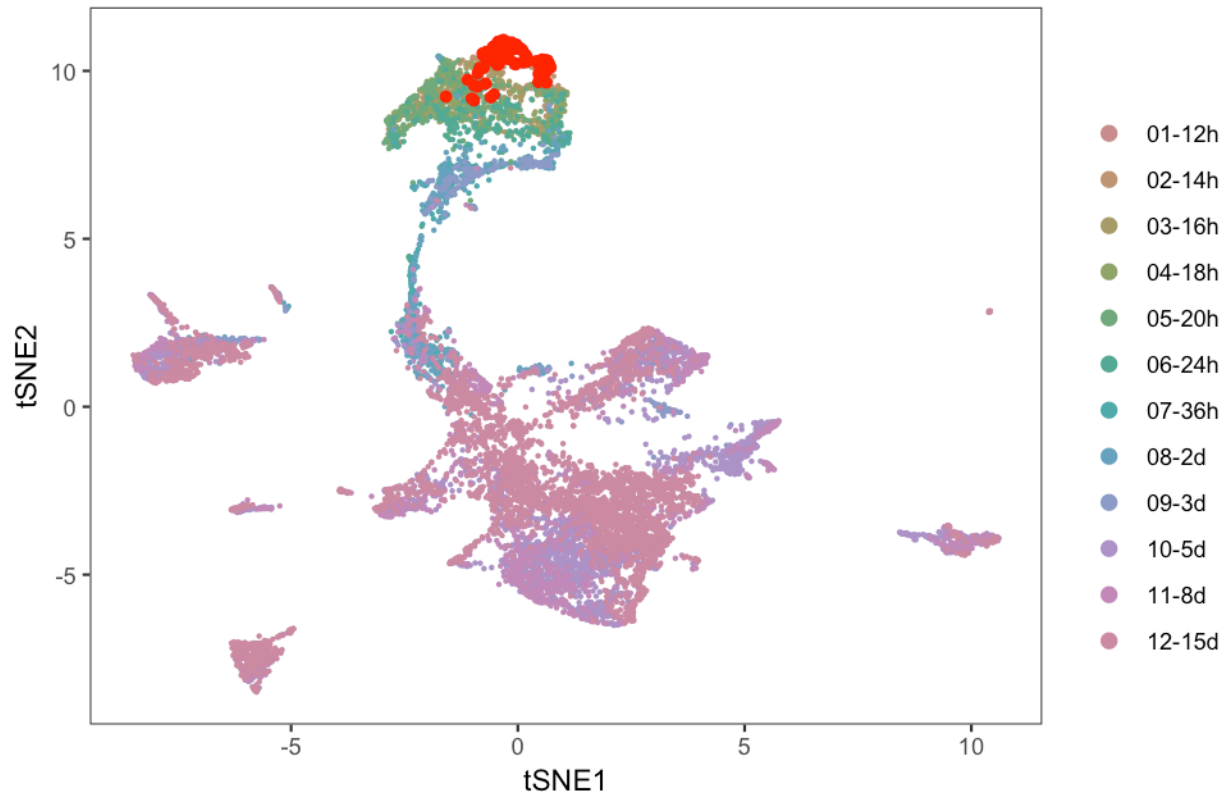


Calculate pseudotime

URD requires a starting point or 'root' for determining pseudotime. Here, we used all cells from the first timepoint (i.e. 12 hpf) as the root.

```
# Here, we used all cells from the first timepoint (i.e. 12 hours) as
# the root.
root.cells <- cellsInCluster(object, "stage", "01-12h")
plotDimHighlight(object, "stage", "01-12h", plot.title = "Root is 12 hpf cells")
```

Root is 12 hpf cells (Highlight 01-12h)



```
# To run locally: Run graph-search simulations to determine pseudotime
flood.result <- floodPseudotime(object, root.cells = root.cells, n = 100,
  minimum.cells.flooded = 2, verbose = T)
```

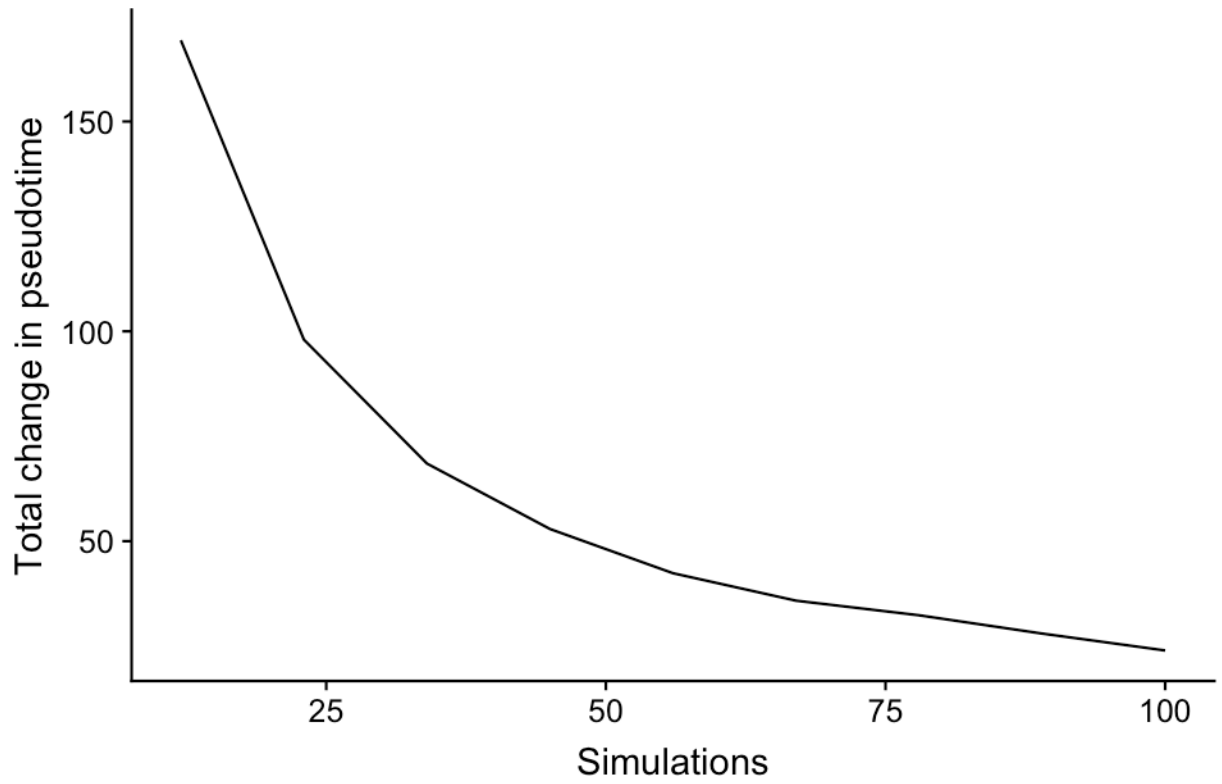
```
# Or load a pre-computed graph-search simulation result
flood.result <- readRDS(paste0(base.path, "flood/flood_hypoND_knn-100_sigma-8.rds"))
```

```
# Process the graph-search simulations to determine the pseudotime of
# each cell
object <- floodPseudotimeProcess(object, flood.result, floods.name = "pseudotime",
  max.frac.NA = 0.4, pseudotime.fun = mean, stability.div = 10)
```

```
# If enough simulations have been run, then as additional simulations
# are added, the overall change in pseudotime of cells should reach an
# asymptote. If it does not, then floodPseudotime should be run with a
# higher n.
```

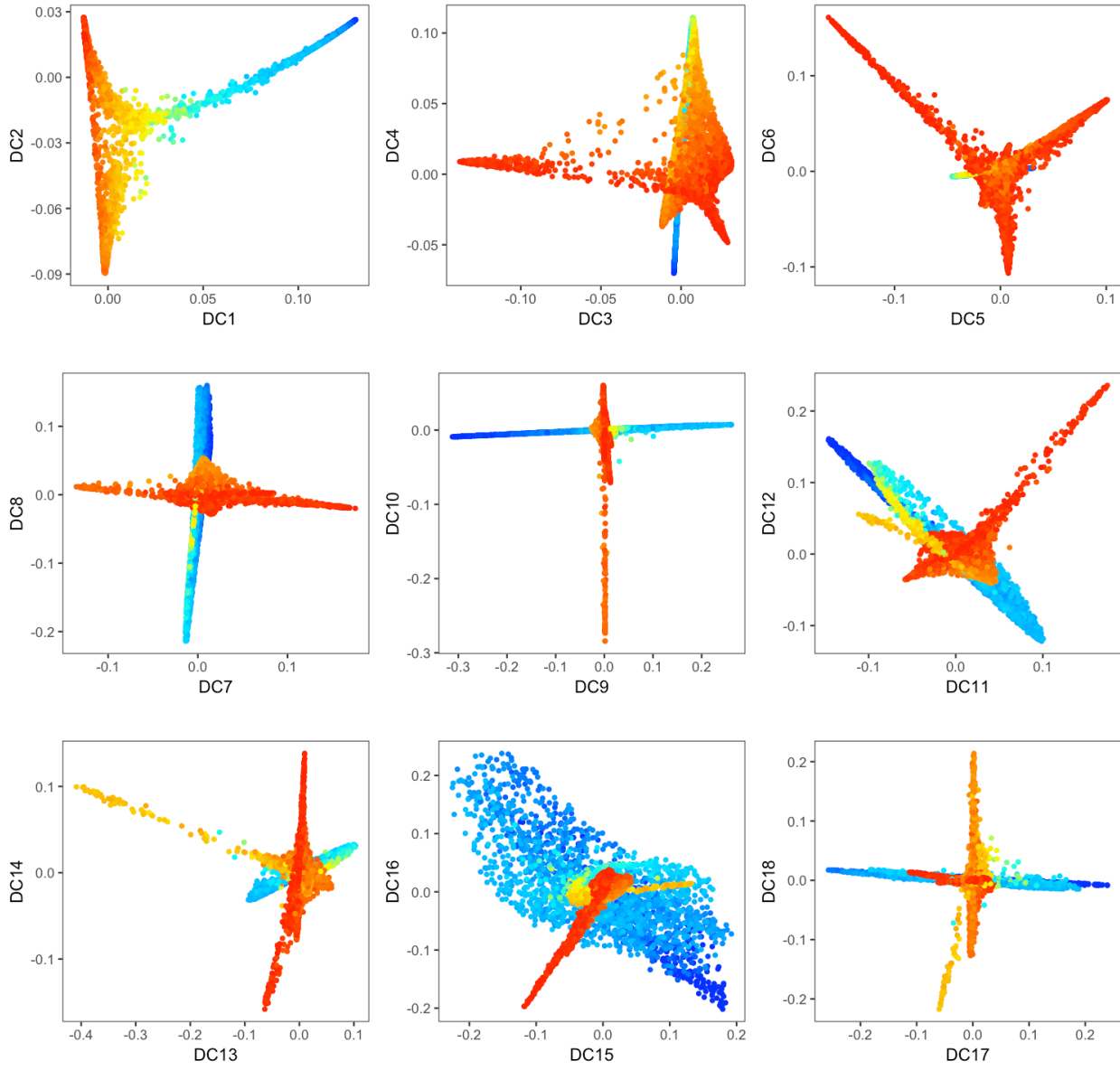
```
pseudotimePlotStabilityOverall(object)
```

Overall Pseudotime Stability



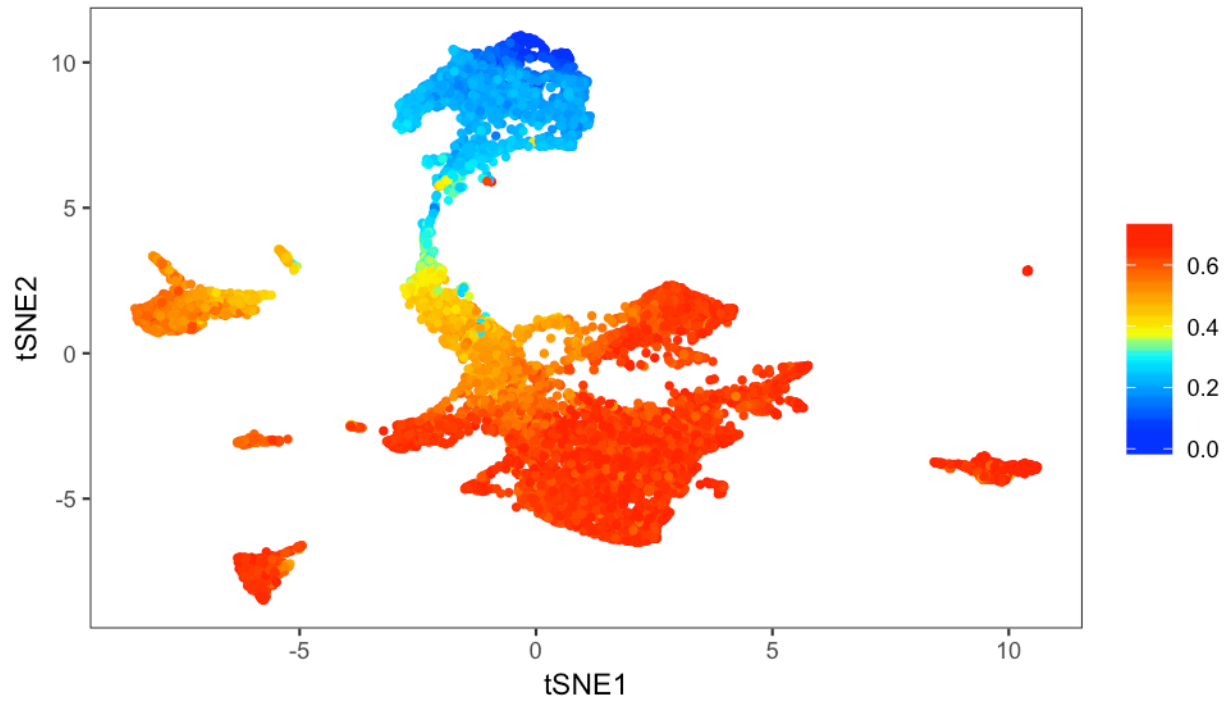
```
plotDimArray(object = object, reduction.use = "dm", dims.to.plot = 1:18,  
  label = "pseudotime", plot.title = "", outer.title = "Diffusion Map labeled by pseudotime",  
  legend = F, alpha = 0.4)
```

Diffusion Map labeled by pseudotime



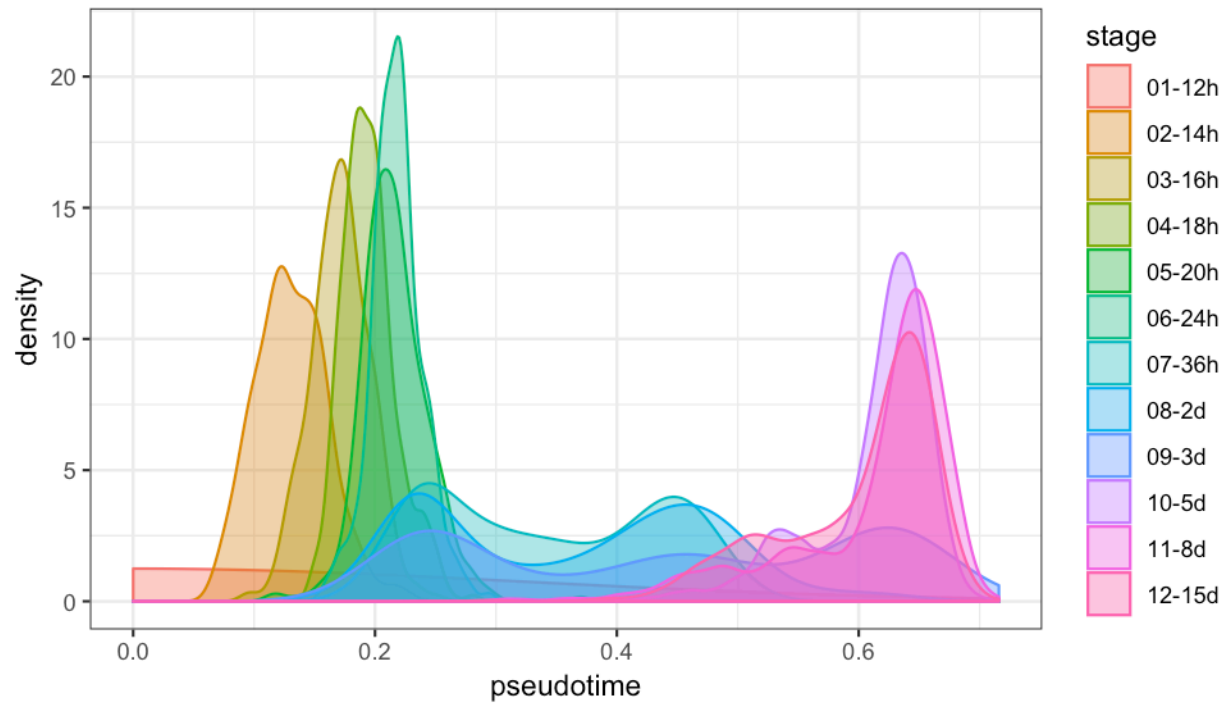
```
plotDim(object, "pseudotime", plot.title = "UMAP projection colored by pseudotime")
```

UMAP projection colored by pseudotime



```
plotDists(object, "pseudotime", "stage", plot.title = "Pseudotime by stage")
```

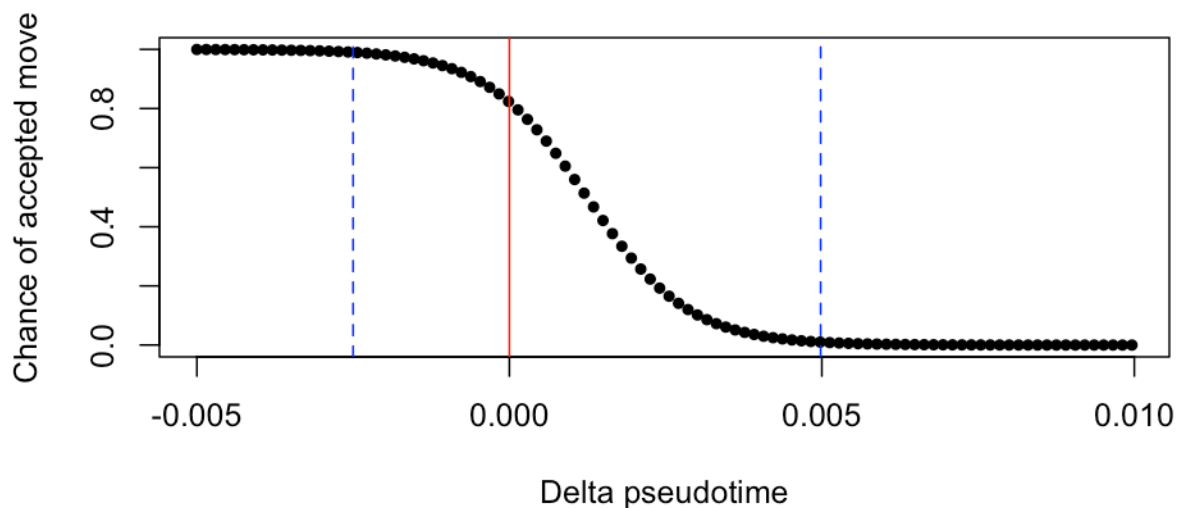
Pseudotime by stage



Calculate biased transition matrix

In order to perform biased random walks, we must first bias the transition matrix to ensure that walks proceed towards the root and do not turn into other differentiated cell types. This is performed in the cluster script `URD-TM.R`.

```
# Calculate parameters for biasing the transition matrix.
diffusion.logistic <- pseudotimeDetermineLogistic(object, "pseudotime",
  optimal.cells.forward = 40, max.cells.back = 80, pseudotime.direction = "<",
  do.plot = T, print.values = T)
```



```
## [1] "Mean pseudotime back (~80 cells) 0.00498088811519173"
## [1] "Chance of accepted move to equal pseudotime is 0.821561374686937"
## [1] "Mean pseudotime forward (~40 cells) -0.00250030667341253"
```

```
# Calculate the biased matrix.
biased.tm <- pseudotimeWeightTransitionMatrix(object, pseudotime = "pseudotime",
  logistic.params = diffusion.logistic, pseudotime.direction = "<")
```

Perform biased random walks

Then, we perform biased walks starting from each tip. Visited cells are inferred to lie along the trajectory that connects the root to each cell type. This is performed in the cluster script `URD-Walk.R`.

Determine tips

We used clusters from 15 dpf as the tips for performing biased random walks. Here we define the cells belonging to each of those clusters.


```

# All clusters at 15 days
clusters.15day <- unique(object@group.ids[grep("15d", object@group.ids$stage),
      "res.5"])
# All cells at 15 days
cells.15day <- rownames(object@group.ids)[grep("15d", object@group.ids$stage)]
# Cell lists of each cluster at 15dpf
cells.15dpf.clusters <- lapply(clusters.15day, function(clust) intersect(cells.15day,
      cellsInCluster(object, "res.5", clust)))
names(cells.15dpf.clusters) <- paste0("15d-", clusters.15day)

```

We also load a .csv file that contains information about the tips. It has four columns:

- id: Cluster ID for the tip
- use: Whether this cluster should be used when building the tree
- name: The name for this tip, which will be used on 2D plots
- short.name: The 'short' name for this tip, which would be used on 3D plots (though we did not use that feature in this study).

```

# Load CSV
tip.names <- read.csv(paste0(base.path, "tips/tip_names_hypoND.csv"), header = F,
      stringsAsFactors = F, colClasses = c("character", "logical", "character",
      "character"))

# Name columns and rows
names(tip.names) <- c("id", "use", "name", "short.name")
rownames(tip.names) <- gsub("_", "-", tip.names$id)

# Sort alphabetically
tip.names <- tip.names[order(rownames(tip.names)), ]

```

These are the tips that were considered during the construction of the retina URD tree (some were excluded during tree construction later in the `buildTree` command).

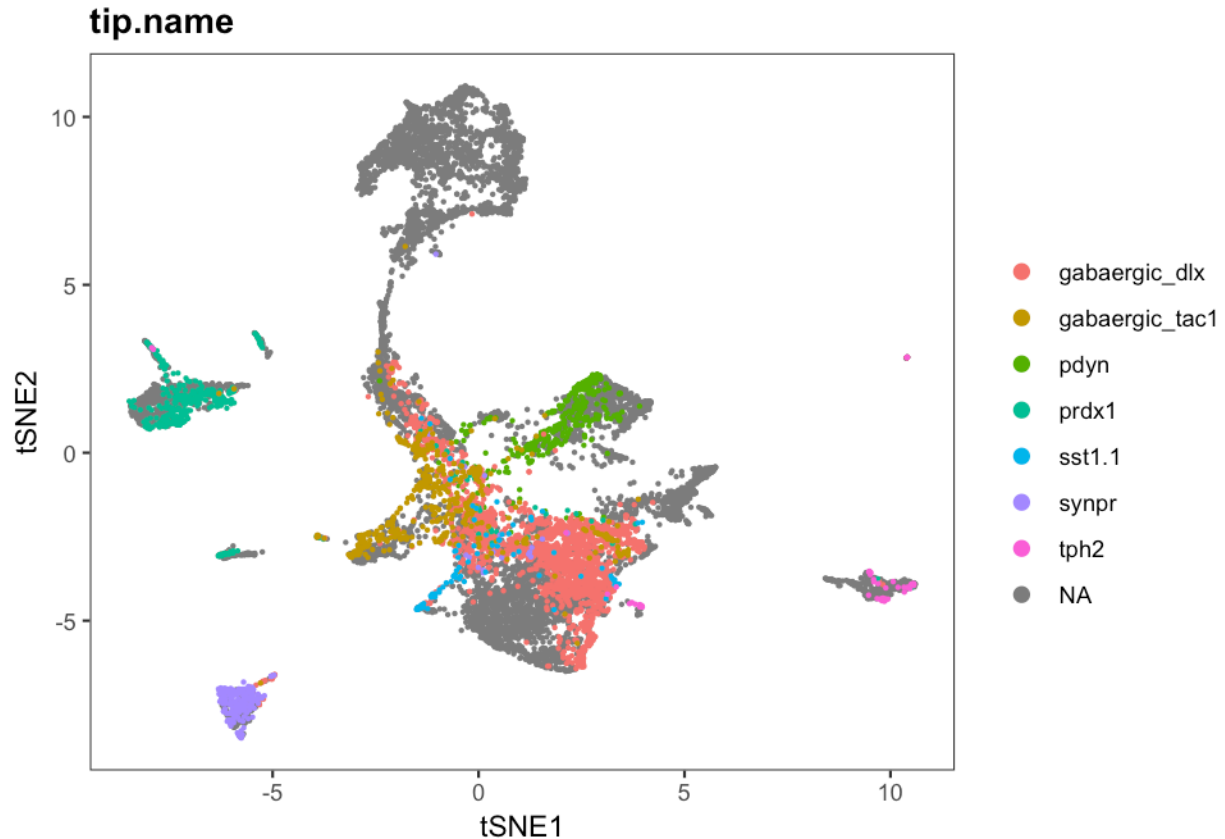
```

# Define a 'tips' clustering
object@group.ids$tip <- NA
object@group.ids$tip.id <- NA
object@group.ids$tip.name <- NA

# If the tip will be used in the tree, define its cells in the
# clustering
for (i in 1:nrow(tip.names)) {
  tip.cells <- cells.15dpf.clusters[[rownames(tip.names)[i]]]
  object@group.ids[tip.cells, "tip"] <- as.character(i)
  object@group.ids[tip.cells, "tip.id"] <- rownames(tip.names)[i]
  object@group.ids[tip.cells, "tip.name"] <- as.character(tip.names[i,
      "name"])
}

# Plot the tips
plotDim(object, "tip.name")

```



Perform the biased random walks

Biased random walks then need to be run starting from each tip. This can be performed on a laptop, but is an ideal candidate for parallelization on a cluster. (The walks from each tip can be run as a separate job.)

```
## IF RUNNING LOCALLY

# Loop through each cluster
walks <- lapply(rownames(tip.names), function(c) {
  # Exclude any tip cells that for whatever reason didn't end up in the
  # biased TM (e.g. maybe not assigned a pseudotime).
  tip.cells <- intersect(cells.15dpf.clusters[[c]], rownames(biased.tm))
  # Perform the random walk simulation
  this.walk <- simulateRandomWalk(start.cells = tip.cells, transition.matrix = biased.tm,
    end.cells = root.cells, n = 50000, end.visits = 1, verbose.freq = 1000,
    max.steps = 5000)
  return(this.walk)
})
names(walks) <- rownames(tip.names)

# Alternatively, this loop is automated by the function
# simulateRandomWalksFromTips
```

Alternatively, a set of pre-calculated walks can be loaded. Since the walks are a simulation (and

therefore not deterministic), this is particularly crucial for reproducing results.

```
## IF LOADING PRE-CALCULATED WALKS

# Get list of files in the walks directory
walks.files <- list.files(paste0(base.path, "/walks/hypoND/"), pattern = ".rds")

# Load the walks previously performed for each cluster
walks <- lapply(rownames(tip.names), function(c) {
  walk.file <- grep(pattern = paste0("_tip-", c, "_"), x = walks.files,
    value = T)[1]
  return(readRDS(paste0(base.path, "/walks/hypoND/", walk.file)))
})
names(walks) <- rownames(tip.names)
```

Process the random walks

The walks are then converted to visitation frequency by importing them into the URD object.

```
for (i in 1:nrow(tip.names)) {
  # Load the individual walk visitation frequencies into the object
  object <- processRandomWalks(object, walks = walks[[i]], walks.name = i,
    n.subsample = 1, verbose = F)
}
```

Build the URD tree

Then, a branching tree is constructed, by joining trajectories in an agglomerative fashion when cells are highly visited by walks from multiple tips. The following steps were performed in the cluster script URD-Tree.R.

```
# Tree building is destructive, so create a copy of the object
object.tree <- object
```

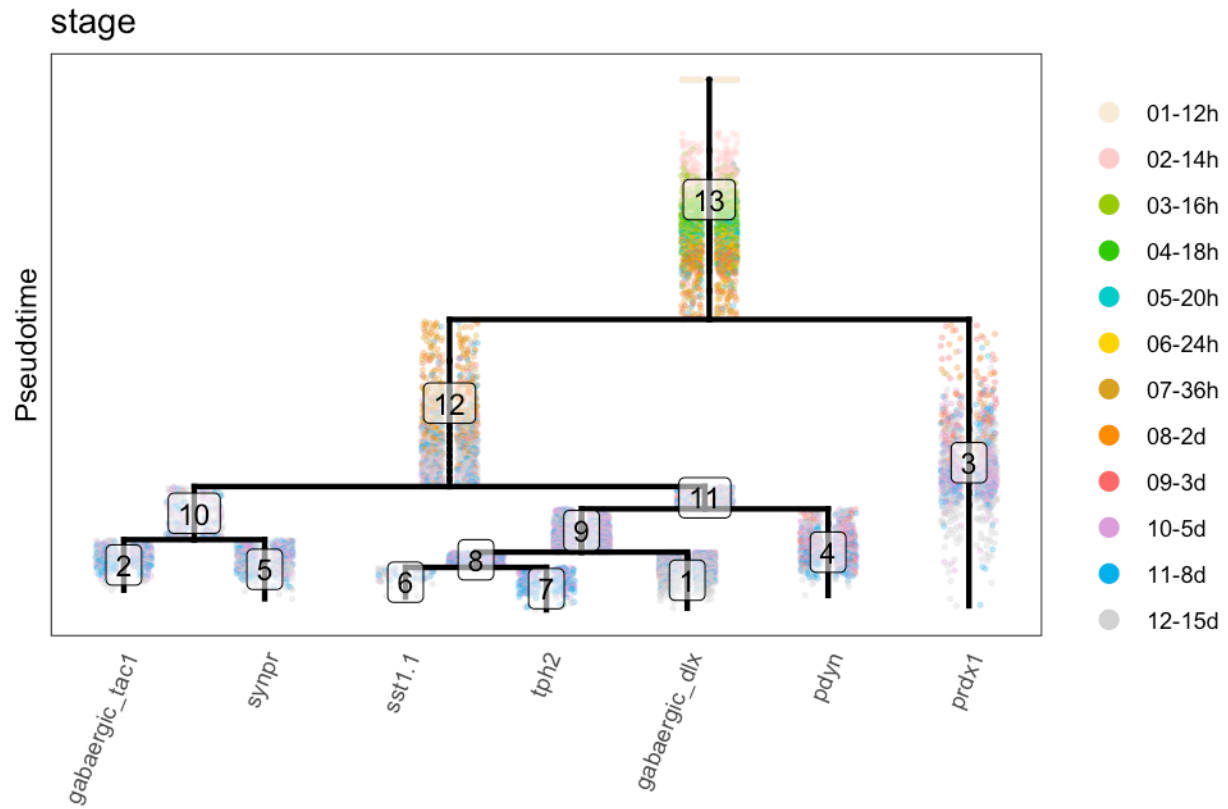
```
# Load tip cells
object.tree <- loadTipCells(object.tree, "tip")
```

```
# Determine tips to use
tips.to.use <- which(tip.names$use)
```

```
# Build the tree
object.tree <- buildTree(object.tree, pseudotime = "pseudotime", divergence.method = "ks",
  cells.per.pseudotime.bin = 40, bins.per.pseudotime.window = 5, save.all.breakpoint.info = T,
  p.thresh = 1e-04, verbose = F, tips.use = as.character(tips.to.use))
```

```
# Name the tips of the tree
object.tree <- nameSegments(object.tree, segments = tips.to.use, segment.names = as.character(tip.names[tips.to.use,
  "name"]), short.names = as.character(tip.names[tips.to.use, "short.name"]))
```

```
plotTree(object.tree, "stage", discrete.colors = stage.colors, label.segments = T)
```



Save the URD tree

The tree is then saved for use in downstream analysis, and can easily be loaded for further perusal.

```
saveRDS(object.tree, file = paste0(base.path, "tree/URD-Tree-Hypo.rds"))
```

Hypothalamus: 3 - URD Cascades and Figures

Jeff Farrell

10/08/2019, updated 07/30/2020

Contents

Load data	2
Plot gene expression on the tree	2
Plot tree by stage	2
Plot tree with gene expression: main figures	3
Plot tree with gene expression: supplemental figures	5
Determine genes enriched in trajectories to particular cell types	6
Comparison between major cell types	6
AUCPR along tree	8
Markers of the prdx1- neuron clade	8
Functions for curating differential expression results	9
threshold.tree.markers	9
divide.branches	9
divide.branches.triple	10
Functions for heatmap generation	11
Color scale	11
determine.timing	12
filter.heatmap.genes	13
Heatmaps of gene cascades	13
Pdyn+ neurons	13
Prepare cascade	13
Generate heatmap: all genes	14
Prdx1+ neurons vs. other neurons	16
Prepare cascade	16
Generate heatmap: all genes	17
nrgna+ neurons	19
Prepare cascade	19
Generate heatmap: all genes	20
Generate heatmap: main figure	22
GABAergic neurons	24
Prepare cascade	24
Generate heatmap: all genes	25
Embryonic molecular profiles are not found in larval progenitors	28
Identify populations to compare	28
Determine proportion of cells in each state	30

Load data

```
# Load URD
library(URD)

## Loading required package: ggplot2
## Loading required package: Matrix
## Registered S3 method overwritten by 'xts':
##   method      from
##   as.zoo.xts  zoo

# Basic location
base.path <- "~/Documents/R sessions/urd-cluster-bushra/"

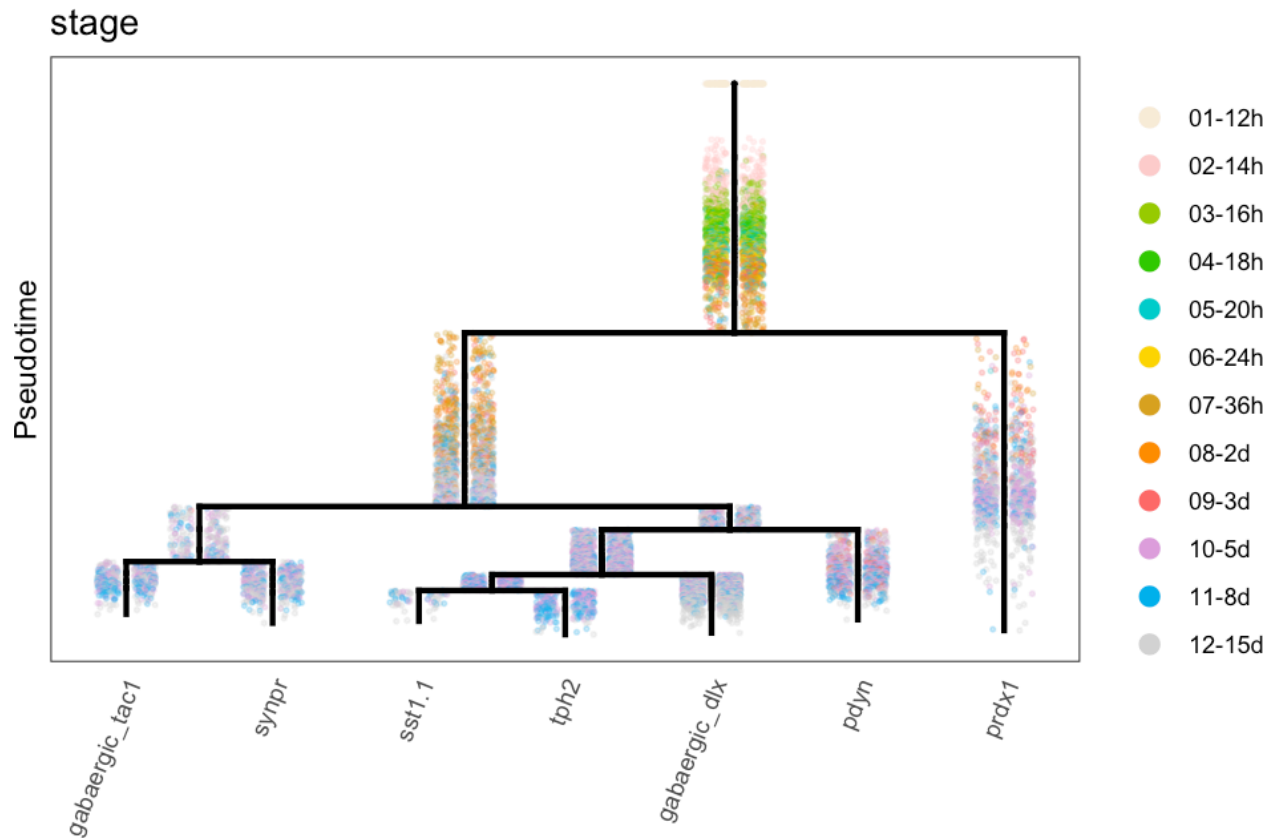
# Load completed hypothalamus tree object
obj.path <- paste0(base.path, "tree/hypoND/tree-hypoND_knn-100_sigma-8_40F-80B_N0-_ks_0001.rds")
obj <- readRDS(obj.path)
```

Plot gene expression on the tree

Plot tree by stage

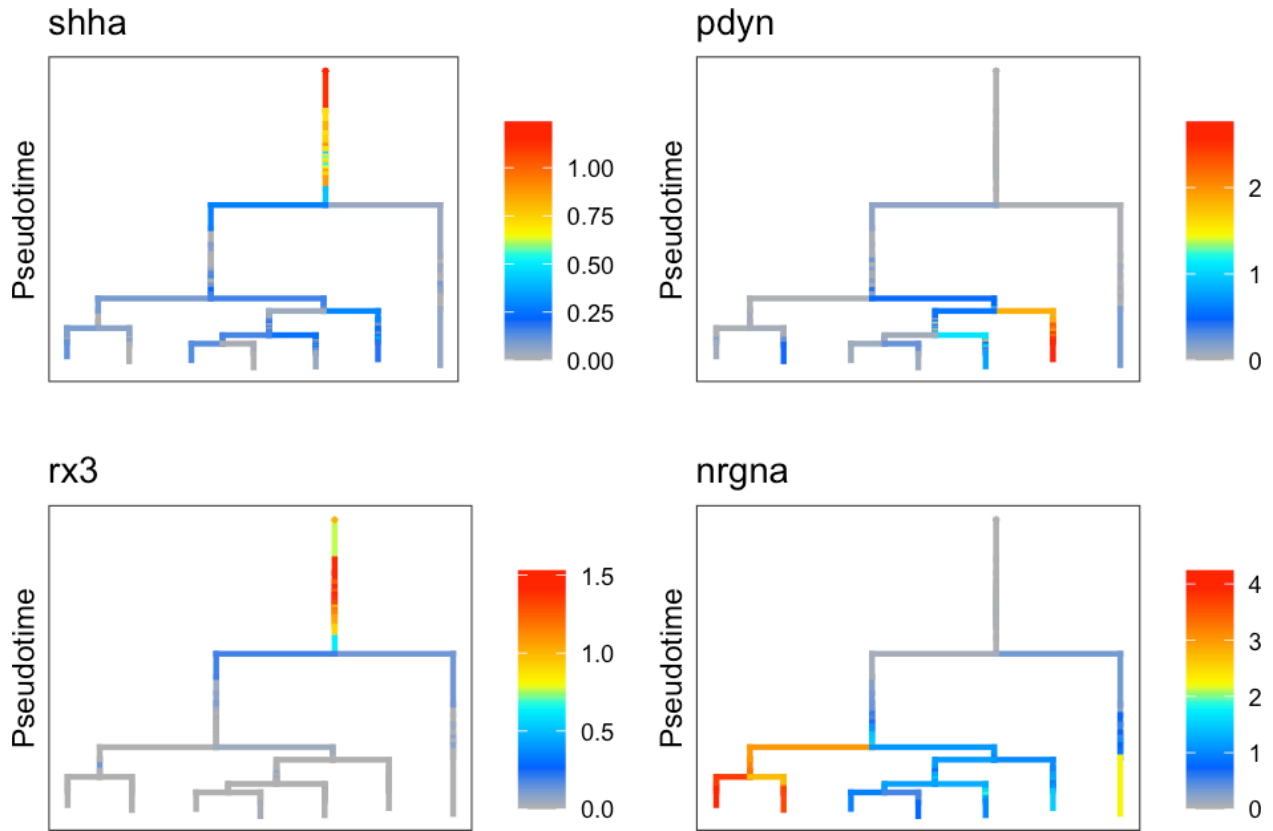
```
stage.colors <- c("antiquewhite", "#FFCCCC", "#99CC00", "#33CC00", "cyan3",
  "gold", "goldenrod", "darkorange", "indianred1", "plum", "deepskyblue2",
  "lightgrey")

plotTree(obj, "stage", label.type = "group", discrete.colors = stage.colors)
```

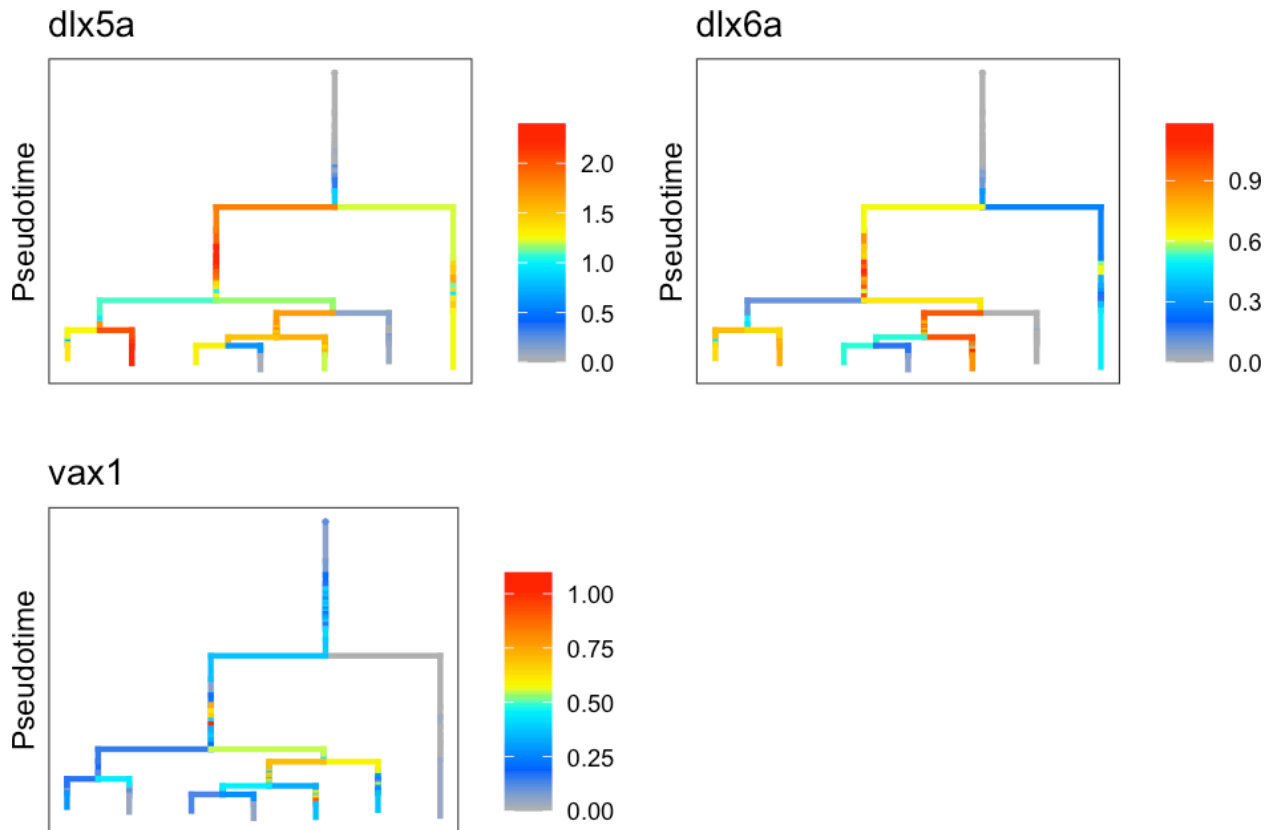


Plot tree with gene expression: main figures

```
gridExtra::grid.arrange(grobs = lapply(c("shha", "pdyn", "rx3", "nrgna"),
  plotTree, object = obj, label.x = F, plot.cells = F), ncol = 2)
```

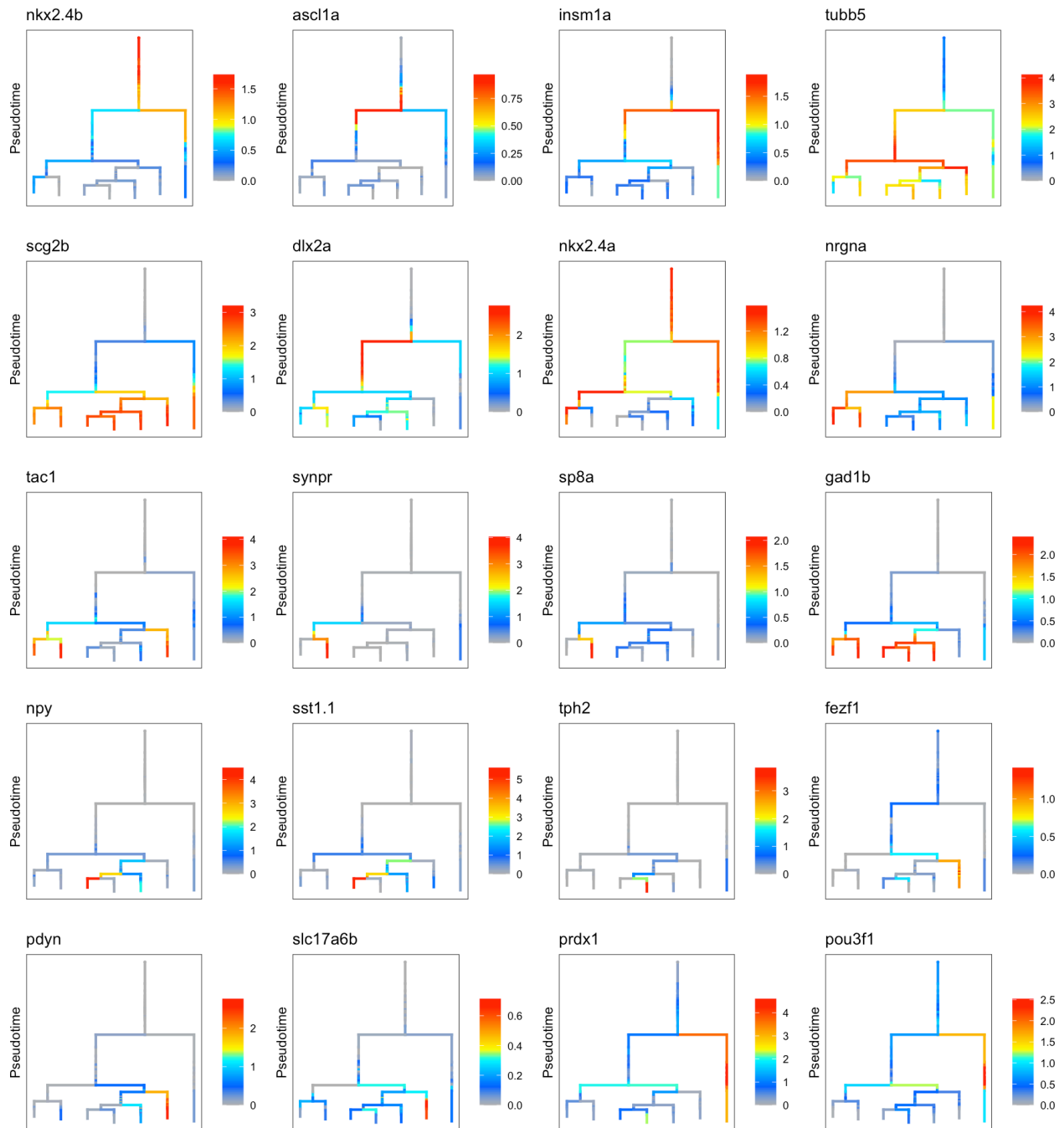


```
gridExtra::grid.arrange(grobs = lapply(c("dlx5a", "dlx6a", "vax1"), plotTree,
  object = obj, label.x = F, plot.cells = F), ncol = 2)
```

Plot tree with gene expression: supplemental figures

```
gridExtra::grid.arrange(grobs = lapply(c("nkx2.4b", "ascl1a", "insm1a",
  "tubb5", "scg2b", "dlx2a", "nkx2.4a", "nrgna", "tac1", "synpr", "sp8a",
  "gad1b", "npy", "sst1.1", "tph2", "fezf1", "pdyn", "slc17a6b", "prdx1",
  "pou3f1"), plotTree, object = obj, label.x = F, plot.cells = F), ncol = 4)
```



Determine genes enriched in trajectories to particular cell types

Comparison between major cell types

We took each major group (“clade”) of branches from the end of the tree as a single entity (i.e. prdx1+ neurons, pdyn+ neurons, GABAergic dlx+ neurons, nrgna+ neurons) and compared them against each other to look for differentially expressed genes.

```
# Get the parent segment of each clade to consider as a group
combined.tips <- c("3", "4", "9", "10")
```

```

# Get the cells in that segment and all child segments
cells.combined.tips <- lapply(combined.tips, function(t) whichCells(obj,
  label = "segment", value = segChildrenAll(obj, t, include.self = T)))
names(cells.combined.tips) <- combined.tips

# Loop through each of these clades and look for differentially
# expressed genes
combined.markers <- lapply(combined.tips, function(tip) {
  # Find all of the other clades
  opposing.tips <- setdiff(combined.tips, tip)
  # Perform pairwise comparisons to each other clade
  m.o <- lapply(opposing.tips, function(tip.opposing) {
    # message(paste0(Sys.time(), ': Comparing tip ', tip, ' to ',
    # tip.opposing, '.')) Find differentially expressed genes between the
    # pair of clades
    ma <- markersAUCPR(object = obj, cells.1 = cells.combined.tips[[tip]],
      cells.2 = cells.combined.tips[[tip.opposing]], effect.size = 0.5,
      auc.factor = 1.1)
    # In order to facilitate combining all of the results later, add
    # columns about which two clades were compared and also a duplicate
    # entry of the name of each gene that's recovered.
    if (nrow(ma) > 0) {
      ma$gene <- rownames(ma)
      ma$tip1 <- tip
      ma$tip2 <- tip.opposing
    }
    return(ma)
  })
  names(m.o) <- opposing.tips
  return(m.o)
})
names(combined.markers) <- combined.tips

# Require that genes are markers against at least 2 other clades
combined.markers.beatmult <- lapply(combined.markers, function(m) {
  names(which(table(unlist(lapply(m, rownames))) >= 2))
})

# Since genes might be a marker in a comparison to several other
# clades, combine the results into a single table, where each gene is
# listed only once with the info from the pairwise comparison where it
# had the strongest differential expression.
combined.markers.best <- lapply(1:length(combined.markers.beatmult), function(i) {
  cm <- do.call("rbind", combined.markers[[i]])
  cm <- cm[cm$gene %in% combined.markers.beatmult[[i]], ]
  cmb <- do.call("rbind", lapply(combined.markers.beatmult[[i]], function(g) {
    cmr <- cm[cm$gene == g, ]
    return(cmr[which.max(cmr$AUCPR.ratio), ])
  }))
  rownames(cmb) <- cmb$gene
  if (!is.null(cmb)) {
    cmb <- cmb[order(cmb$AUCPR.ratio, decreasing = T), ]
    cmb$exp.global <- apply(obj@logupx.data[rownames(cmb), unlist(obj@tree$cells.in.segment)],

```

```

    1, mean.of.logs)
    cmb$exp.global.fc <- cmb$nTrans_1 - cmb$exp.global
  }
  return(cmb)
})
names(combined.markers.best) <- combined.tips

```

AUCPR along tree

We also used the `AUCPRTTestAlongTree` function to ask for genes that are differential markers of a lineage using URD's tree structure. This makes a comparison at each branchpoint from a particular cell type up to the root.

```

# Get all of the tips from the tree
tips.in.tree <- as.character(obj@tree$tips)

# Tree segments to use as root
roots <- rep("12", length(tips.in.tree))
names(roots) <- tips.in.tree
roots["3"] <- "13"

# Define parameters to use for calculation Used more permissive values
# in the sst1.1+ / tph2+ / gabaergic dlx+ neuronal comparisons due to
# the small number of cells in these populations
auc.use <- rep(1.2, length(tips.in.tree))
names(auc.use) <- tips.in.tree
auc.use[c("1", "6", "7")] <- 1.15
log.effect.use <- rep(0.8, length(tips.in.tree))
names(log.effect.use) <- tips.in.tree
log.effect.use[c("1", "6", "7")] <- 0.6

# Perform a loop of tests with each tip.
markers <- lapply(tips.in.tree, function(t) {
  this.root <- roots[t]
  this.auc <- auc.use[t]
  this.log <- log.effect.use[t]
  # message(paste0(Sys.time(), ': Starting tip ', t, ' and root ',
  # this.root, ' with params ', this.auc, ' AUC and ', this.log, ' effect
  # size.'))
  these.markers <- aucprTestAlongTree(obj, pseudotime = "pseudotime",
    tips = as.character(t), genes.use = NULL, must.beat.sibs = 0.6,
    report.debug = F, root = this.root, auc.factor = this.auc, log.effect.size = this.log)
  these.markers$gene <- rownames(these.markers)
  these.markers$tip <- t
  return(these.markers)
})
names(markers) <- tips.in.tree

```

Markers of the prdx1- neuron clade

```

# Calculate from segment 12 against segment 3 specifically
nonprdx.markers <- markersAUCPR(obj, clust.1 = "12", clust.2 = "3", clustering = "segment",

```

```

effect.size = 0.8, auc.factor = 1.2)
# Also look at segment 12 vs. rest of the hypothalamus with lower
# thresholds
nonprdx.markers.global <- markersAUCPR(obj, clust.1 = "12", clust.2 = as.character(c(1:11,
13)), clustering = "segment", effect.size = 0.4, auc.factor = 1.1)

## Warning in names(genes.data)[4:7] <- paste(c("posFrac", "posFrac", "nTrans", :
## number of items to replace is not a multiple of replacement length

```

Functions for curating differential expression results

We further curated those differentially expressed genes using the following functions:

threshold.tree.markers

Function to threshold markers from a markersAUCPRAlongTree test with additional criteria

- **markers**: list of results from markersAUCPRAlongTree tests
- **tip**: which tip (or element of the list to pursue)
- **global.fc**: fold.change that gene must have along the trajectory pursued vs. rest of the data
- **aucpr.ratio.all**: classifier score that gene must exhibit along trajectory test vs. rest of the data
- **branch.fc**: fold.change that gene must have (in best case) vs. the opposing branch at any branchpoint along the trajectory.
- Returns markers with only a subset of rows retained.

```

threshold.tree.markers <- function(markers, tip, global.fc = 0.1, branch.fc = 0.4,
aucpr.ratio.all = 1.03) {
  m <- markers[[tip]]
  # First off -- lose global FC < x
  bye.globalfc <- rownames(m)[m$expfc.all < global.fc]
  # Second -- get rid of branch FC < x
  bye.branchfc <- rownames(m)[m$expfc.maxBranch < branch.fc]
  # Third -- get rid of stuff essentially worse than random
  # classification on global level
  bye.badglobalaucpr <- rownames(m)[m$AUCPR.ratio.all < aucpr.ratio.all]
  bye.all <- unique(c(bye.globalfc, bye.branchfc, bye.badglobalaucpr))
  m.return <- m[setdiff(rownames(m), bye.all), ]
  return(m.return)
}

```

divide.branches

Function to compare genes between two branches. Use this on a compiled list of markers to do a final selection of genes that are specific to one branch or another or markers of both (i.e. when making photoreceptor heatmap, use to divide into photoreceptor, cone, and rod markers)

- **object**: An URD object
- **genes**: (Character vector) Genes to test
- **clust.1**: (Character) Cluster 1
- **clust.2**: (Character) Cluster 2
- **clustering**: (Character) Clustering to pull from
- **exp.fc**: (Numeric) Minimum expression fold-change between branches to consider different

- **exp.thresh**: (Numeric) Minimum fraction of cells in order to consider gene expressed in a branch
- **exp.diff**: (Numeric) Minimum difference in fraction of cells expressing to consider gene differential
- Returns list of gene names (“specific.1” = specific to clust.1, “specific.2” = specific to clust.2, “markers” = all genes tested)

```
divide.branches <- function(object, genes, clust.1, clust.2, clustering = "segment",
  exp.fc = 0.4, exp.thresh = 0.1, exp.diff = 0.1) {
  # Double check which markers are unique to one or the other population
  mcomp <- markersAUCPR(object, clust.1 = clust.1, clust.2 = clust.2,
    clustering = clustering, effect.size = -Inf, auc.factor = 0, genes.use = genes,
    frac.min.diff = 0, frac.must.express = 0)
  specific.b <- rownames(mcomp)[abs(mcomp$exp.fc) > exp.fc & mcomp[,
    4] < exp.thresh & mcomp[, 5] > pmin((mcomp[, 4] + exp.diff), 1)]
  specific.a <- rownames(mcomp)[abs(mcomp$exp.fc) > exp.fc & mcomp[,
    5] < exp.thresh & mcomp[, 4] > pmin((mcomp[, 5] + exp.diff), 1)]
  r <- list(specific.a, specific.b, mcomp)
  names(r) <- c("specific.1", "specific.2", "markers")
  return(r)
}
```

divide.branches.triple

Function to compare genes between three branches. Use this on a compiled list of markers to do a final selection of genes that are specific to one branch or another or markers of both (i.e. when making gad2+ heatmap, use to divide into general, dlx+, sst+, and tph2+ markers).

- **object**: An URD object
- **genes**: (Character vector) Genes to test
- **clust.1**: (Character) Cluster 1
- **clust.2**: (Character) Cluster 2
- **clust.3**: (Character) Cluster 3
- **clustering**: (Character) Clustering to pull from
- **exp.fc**: (Numeric) Minimum expression fold-change between branches to consider different
- **exp.thresh**: (Numeric) Minimum fraction of cells in order to consider gene expressed in a branch
- **exp.diff**: (Numeric) Minimum difference in fraction of cells expressing to consider gene differential
- Returns list of gene names (“nonspecific” = genes not in specific.1/2/3, “specific.1” = specific to clust.1, “specific.2” = specific to clust.2, “specific.3” = specific to clust.3, each pairwise comparison, and “markers” = all genes tested)

```
divide.branches.triple <- function(object, genes, clust.1, clust.2, clust.3,
  clustering = "segment", exp.fc = 0.4, exp.thresh = 0.1, exp.diff = 0.1) {
  # Double check which markers are unique to one or the other population
  mcomp12 <- markersAUCPR(object, clust.1 = clust.1, clust.2 = clust.2,
    clustering = clustering, effect.size = -Inf, auc.factor = 0, genes.use = genes,
    frac.min.diff = 0, frac.must.express = 0)
  mcomp23 <- markersAUCPR(object, clust.1 = clust.2, clust.2 = clust.3,
    clustering = clustering, effect.size = -Inf, auc.factor = 0, genes.use = genes,
    frac.min.diff = 0, frac.must.express = 0)
  mcomp13 <- markersAUCPR(object, clust.1 = clust.1, clust.2 = clust.3,
    clustering = clustering, effect.size = -Inf, auc.factor = 0, genes.use = genes,
    frac.min.diff = 0, frac.must.express = 0)
```

```

specific.1v2 <- rownames(mcomp12)[abs(mcomp12$exp.fc) > exp.fc & mcomp12[,
  5] < exp.thresh & mcomp12[, 4] > pmin((mcomp12[, 5] + exp.diff),
  1)]
specific.2v1 <- rownames(mcomp12)[abs(mcomp12$exp.fc) > exp.fc & mcomp12[,
  4] < exp.thresh & mcomp12[, 5] > pmin((mcomp12[, 4] + exp.diff),
  1)]

specific.2v3 <- rownames(mcomp23)[abs(mcomp23$exp.fc) > exp.fc & mcomp23[,
  5] < exp.thresh & mcomp23[, 4] > pmin((mcomp23[, 5] + exp.diff),
  1)]
specific.3v2 <- rownames(mcomp23)[abs(mcomp23$exp.fc) > exp.fc & mcomp23[,
  4] < exp.thresh & mcomp23[, 5] > pmin((mcomp23[, 4] + exp.diff),
  1)]

specific.1v3 <- rownames(mcomp13)[abs(mcomp13$exp.fc) > exp.fc & mcomp13[,
  5] < exp.thresh & mcomp13[, 4] > pmin((mcomp13[, 5] + exp.diff),
  1)]
specific.3v1 <- rownames(mcomp13)[abs(mcomp13$exp.fc) > exp.fc & mcomp13[,
  4] < exp.thresh & mcomp13[, 5] > pmin((mcomp13[, 4] + exp.diff),
  1)]

specific.1 <- unique(setdiff(c(specific.1v2, specific.1v3), c(specific.2v3,
  specific.3v2)))
specific.2 <- unique(setdiff(c(specific.2v1, specific.2v3), c(specific.1v3,
  specific.3v1)))
specific.3 <- unique(setdiff(c(specific.3v2, specific.3v1), c(specific.2v1,
  specific.1v2)))
nonspecific <- setdiff(genes, c(specific.1, specific.2, specific.3))

markers.comp <- list(mcomp12, mcomp13, mcomp23)
names(markers.comp) <- c("1v2", "1v3", "2v3")

r <- list(nonspecific, specific.1, specific.2, specific.3, specific.1v2,
  specific.1v3, specific.2v1, specific.2v3, specific.3v1, specific.3v2,
  markers.comp)
names(r) <- c("nonspecific", "specific.1", "specific.2", "specific.3",
  "specific.1v2", "specific.1v3", "specific.2v1", "specific.2v3",
  "specific.3v1", "specific.3v2", "markers")
return(r)
}

```

Functions for heatmap generation

These functions were used in the production of heatmaps:

Color scale

Generate color scale to use with heatmaps.

```

cols <- (scales::gradient_n_pal(RColorBrewer::brewer.pal(9, "YlOrRd")))(seq(0,
  1, length.out = 50))

```

determine.timing

Determines order to plot genes in heatmap. "Expression" is defined as 20% higher expression than the minimum observed value. "Peak" expression is defined as 50% higher expression than minimum observed value. The two longest stretches of "peak" expression are found, and then the later one is used. The onset time of the stretch of expression that contains that peak is also determined. Genes are then ordered by the pseudotime at which they enter "peak" expression, leave "peak" expression, start "expression", and leave "expression".

- **s**: result from `geneSmoothFit`
- **genes**: genes to order; default is all genes that were fit.
- Returns `s` but with an additional list entry (`s$timing`) of the order to plot genes

```
determine.timing <- function(s, genes = rownames(s$mean.expression)) {
  s$timing <- as.data.frame(do.call("rbind", lapply(genes, function(g) {
    sv <- as.numeric(s$scaled.smooth[g, ])
    pt <- as.numeric(colnames(s$scaled.smooth))
    # Figure out baseline expression & threshold for finding peaks
    min.val <- max(min(sv), 0)
    peak.val <- ((1 - min.val)/2) + min.val
    exp.val <- ((1 - min.val)/5) + min.val
    # Run-length encoding of above/below the peak-threshold
    peak.rle <- rle(sv >= peak.val)
    peak.rle <- data.frame(lengths = peak.rle$lengths, values = peak.rle$values)
    peak.rle$end <- cumsum(peak.rle$lengths)
    peak.rle$start <- head(c(0, peak.rle$end) + 1, -1)
    # Run-length encoding of above/below the expressed-threshold
    exp.rle <- rle(sv >= exp.val)
    exp.rle <- data.frame(lengths = exp.rle$lengths, values = exp.rle$values)
    exp.rle$end <- cumsum(exp.rle$lengths)
    exp.rle$start <- head(c(0, exp.rle$end) + 1, -1)
    # Take top-two longest peak RLE & select later one. Find stretches
    # that are above peak value
    peak <- which(peak.rle$values)
    # Order by length and take 1 or 2 longest ones
    peak <- peak[order(peak.rle[peak, "lengths"], decreasing = T)][1:min(2,
      length(peak))]
    # Order by start and take latest one.
    peak <- peak[order(peak.rle[peak, "start"], decreasing = T)][1]
    # Identify the actual peak value within that stretch
    peak <- which.max(sv[peak.rle[peak, "start"]:peak.rle[peak, "end"]]) +
      peak.rle[peak, "start"] - 1
    # Identify the start and stop of the expressed stretch that contains
    # the peak
    exp.start <- exp.rle[which(exp.rle$end >= peak & exp.rle$start <=
      peak), "start"]
    exp.end <- exp.rle[which(exp.rle$end >= peak & exp.rle$start <=
      peak), "end"]
    # Identify values of expression at start and stop
    smooth.start <- sv[exp.start]
    smooth.end <- sv[exp.end]
    # Convert to pseudotime?
    exp.start <- pt[exp.start]
    exp.end <- pt[exp.end]
    peak <- pt[peak]
  })})
```



```

# Return a vector
v <- c(exp.start, peak, exp.end, smooth.start, smooth.end)
names(v) <- c("pt.start", "pt.peak", "pt.end", "exp.start", "exp.end")
return(v)
}))
rownames(s$timing) <- genes

# Decide on ordering of genes
s$gene.order <- rownames(s$timing)[order(s$timing$pt.peak, s$timing$pt.start,
s$timing$pt.end, s$timing$exp.end, decreasing = c(F, F, F, T),
method = "radix")]

return(s)
}

```

filter.heatmap.genes

Removes undesired (mitochondrial, ribosomal, tandem duplicated genes) from heatmaps for presentation purposes.

- **genes:** (Character vector) genes to check
- Returns genes with undesired genes removed.

```

filter.heatmap.genes <- function(genes) {
  mt.genes <- grep("^mt-", ignore.case = T, genes, value = T)
  many.genes <- grep("\\(1 of many\\)", ignore.case = T, genes, value = T)
  ribo.genes <- grep("^rpl|^rps", ignore.case = T, genes, value = T)
  cox.genes <- grep("^cox", ignore.case = T, genes, value = T)
  hsp.genes <- grep("^hsp", ignore.case = T, genes, value = T)
  return(setdiff(genes, c(mt.genes, many.genes, ribo.genes, cox.genes,
hsp.genes)))
}

```

Heatmaps of gene cascades

Using the genes that were determined as differentially expressed along the way to particular cell types, we generated expression cascades and plotted them as heatmaps.

Pdyn+ neurons

Prepare cascade

```

## Pdyn+ neurons: Seg 4

# Get markers from the two approaches
t <- combined.markers.best[["4"]] # pdyn+ markers from above the combined clades
t <- t[t$exp.global.fc >= 0.8, ] # limited to those with good global parameters
m <- threshold.tree.markers(markers, "4", global.fc = 0.6) # Pdyn+ Cell markers from aucprTestAlongTree
pdyn.markers <- unique(c(rownames(t), rownames(m)))

# Just want to plot part of cells from upstream segment 12, which is
# very long. Going to use cells from segments 4, 11, and from segment

```

```

# 12 with pseudotime > 0.23
cells.plot <- unique(c(intersect(cellsInCluster(obj, "segment", "12"),
  whichCells(obj, "pseudotime", c(0.45, Inf))), cellsInCluster(obj, "segment",
  c("11", "4"))))

# Calculate spline curve
spline.plot <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cells.plot,
  genes = pdyn.markers, method = "spline", moving.window = 5, cells.per.window = 25,
  pseudotime.per.window = 0.005, spar = 0.5, verbose = F)

# Calculate gene expression timing for ordering rows
spline.plot <- determine.timing(s = spline.plot)
order.plot <- filter.heatmap.genes(spline.plot$gene.order)

# Output gene table
table.save <- data.frame(gene = order.plot, stringsAsFactors = F)
table.save$clade.AUCPR.ratio <- t[table.save$gene, "AUCPR.ratio"]
table.save$clade.exp.fc <- t[table.save$gene, "exp.fc"]
table.save$clade.exp.fc.global <- t[table.save$gene, "exp.global.fc"]
table.save$pdyn.AUCPR.ratio.all <- m[table.save$gene, "AUCPR.ratio.all"]
table.save$pdyn.AUCPR.ratio.maxBranch <- m[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$pdyn.exp.fc.all <- m[table.save$gene, "expfc.all"]
table.save$pdyn.exp.fc.best <- m[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/hypo-pdyn.csv"))

```

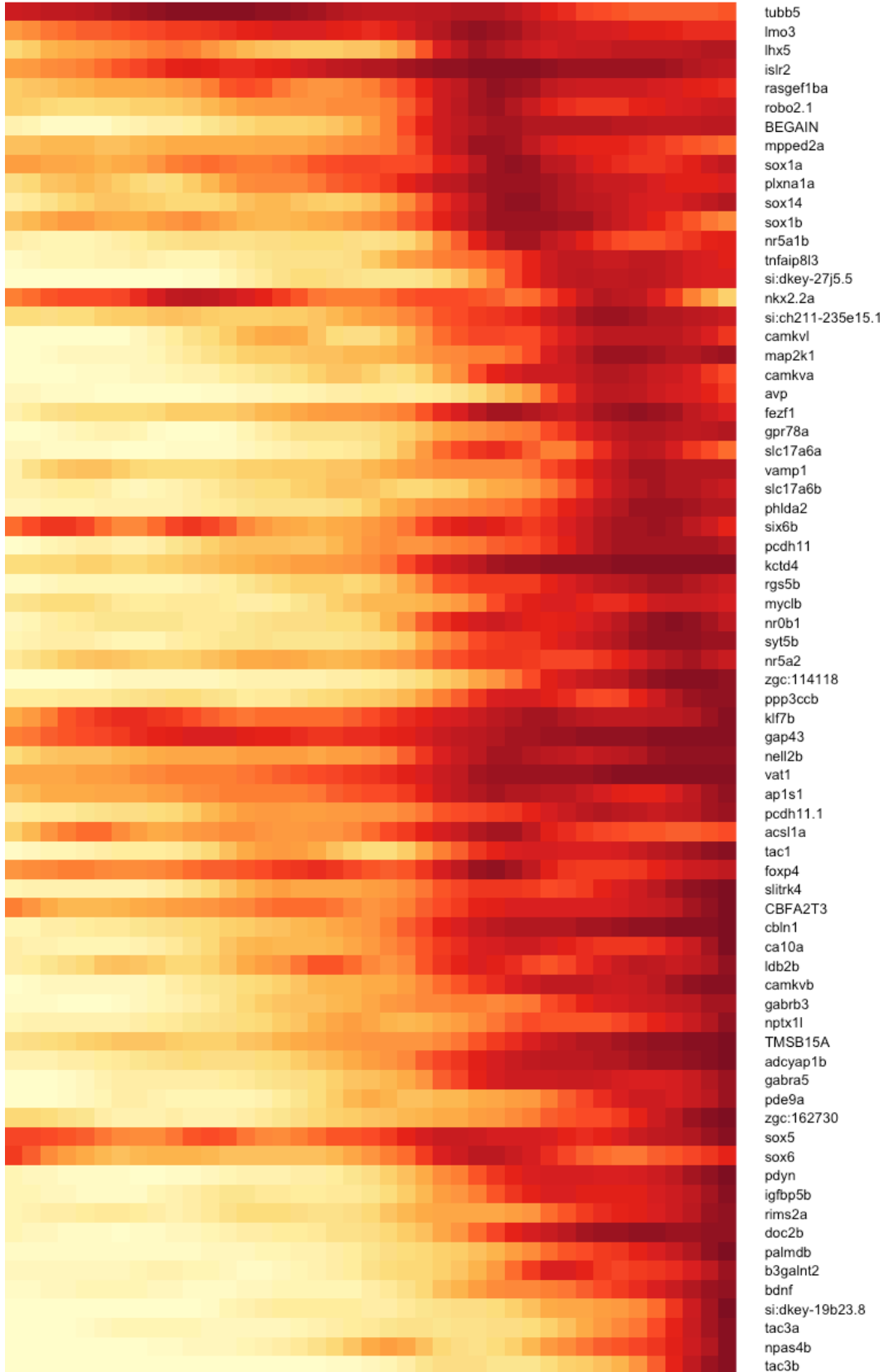
Generate heatmap: all genes

```

## Generate heatmap Make sure any values <0 in the spline curves get set
## to 0 so that the heatmap scale doesn't get messed up.
spline.plot$scaled.smooth[spline.plot$scaled.smooth < 0] <- 0
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/hypo-pdyn.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(spline.plot$scaled.smooth[order.plot, ]),
  Rowv = F, Colv = F, dendrogram = "none", col = cols, trace = "none",
  density.info = "none", key = F, cexCol = 0.8, cexRow = 0.6, margins = c(8,
  8), lwid = c(0.3, 4), lhei = c(0.3, 4), labCol = NA)
title(main = "pdyn+ Neurons")

```

pdyn+ Neurons



```
# dev.off()
```

Prdx1+ neurons vs. other neurons

Prepare cascade

```
## Prdx1+ neurons: Seg 3

# Get markers from the two approaches
t <- combined.markers.best[["3"]] # prdx1+ markers from above the combined clades
t <- t[t$exp.global.fc >= 0.8, ] # limited to those with good global parameters
m <- threshold.tree.markers(markers, "3", global.fc = 0.6) # prdx1+ Cell markers from aucprTestAlongTree
prdx.markers <- unique(c(rownames(t), rownames(m)))

# Get markers for the opposing segment
opposing.prdx.markers <- intersect(rownames(nonprdx.markers), rownames(nonprdx.markers.global))

prdx.hm.markers <- unique(c(prdx.markers, opposing.prdx.markers))

# Calculate spline curves Using segments 13 and 12 or 3.
spline.12 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cellsInCluster(obj,
  "segment", c("13", "12")), genes = prdx.hm.markers, method = "spline",
  moving.window = 5, cells.per.window = 25, pseudotime.per.window = 0.005,
  spar = 0.5, verbose = F)
spline.3 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cellsInCluster(obj,
  "segment", c("13", "3")), genes = prdx.hm.markers, method = "spline",
  moving.window = 5, cells.per.window = 25, pseudotime.per.window = 0.005,
  spar = 0.5, verbose = F)
spline.123 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cellsInCluster(obj,
  "segment", c("13", "12", "3")), genes = prdx.hm.markers, method = "spline",
  moving.window = 5, cells.per.window = 25, pseudotime.per.window = 0.005,
  spar = 0.5, verbose = F)

# Want to plot a heatmap that shows expression in most upstream
# progenitors and then each branch (i.e. prdx1- vs. prdx1+ neurons) as
# separate columns. Going to crop each spline fit to the correct
# pseudotime range and then combine them into a single one that can be
# plotted as a three-column heatmap.

pt.12v3 <- obj@tree$segment.pseudotime.limits["3", "start"] # pseudotime where the crop should happen
splines.prdx <- list(cropSmoothFit(spline.123, pt.min = -Inf, pt.max = pt.12v3),
  cropSmoothFit(spline.12, pt.min = pt.12v3, pt.max = Inf), cropSmoothFit(spline.3,
  pt.min = pt.12v3, pt.max = Inf))
names(splines.prdx) <- c("Hypo Precursors", "Prdx1-", "Prdx1+")
splines.prdx.hm <- combineSmoothFit(splines.prdx) # Combine into a single one

# Calculate gene expression timing for ordering rows
spline.12 <- determine.timing(s = spline.12)
spline.3 <- determine.timing(s = spline.3)
spline.123 <- determine.timing(s = spline.123)

# Decide which markers are specific to one cell type or both
d12v3 <- divide.branches(obj, prdx.hm.markers, clust.1 = "12", clust.2 = "3",
```

```

exp.fc = 0.4, exp.thresh = 0.2, exp.diff = 0.1)

# Generate gene ordering based on timing & specificity
order.123 <- filter.heatmap.genes(setdiff(spline.123$gene.order, c(d12v3$specific.1,
  d12v3$specific.2)))
order.12 <- filter.heatmap.genes(intersect(spline.12$gene.order, d12v3$specific.1))
order.3 <- filter.heatmap.genes(intersect(spline.3$gene.order, d12v3$specific.2))
gene.order <- c(order.123, order.12, order.3)

# Output gene table
table.save <- data.frame(gene = gene.order, marks = c(rep("both", length(order.123)),
  rep("prdx1-", length(order.12)), rep("prdx1+", length(order.3))), stringsAsFactors = F)
table.save$clade.AUCPR.ratio <- t[table.save$gene, "AUCPR.ratio"]
table.save$clade.exp.fc <- t[table.save$gene, "exp.fc"]
table.save$clade.exp.fc.global <- t[table.save$gene, "exp.global.fc"]
table.save$non.AUCPR.ratio <- nonprdx.markers[table.save$gene, "AUCPR.ratio"]
table.save$non.exp.fc <- nonprdx.markers[table.save$gene, "exp.fc"]
table.save$non.exp.fc.global <- nonprdx.markers.global[table.save$gene,
  "exp.fc"]
table.save$prdx.AUCPR.ratio.all <- m[table.save$gene, "AUCPR.ratio.all"]
table.save$prdx.AUCPR.ratio.maxBranch <- m[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$prdx.exp.fc.all <- m[table.save$gene, "expfc.all"]
table.save$prdx.exp.fc.best <- m[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/hypo-prdx1-vs-non.csv"))

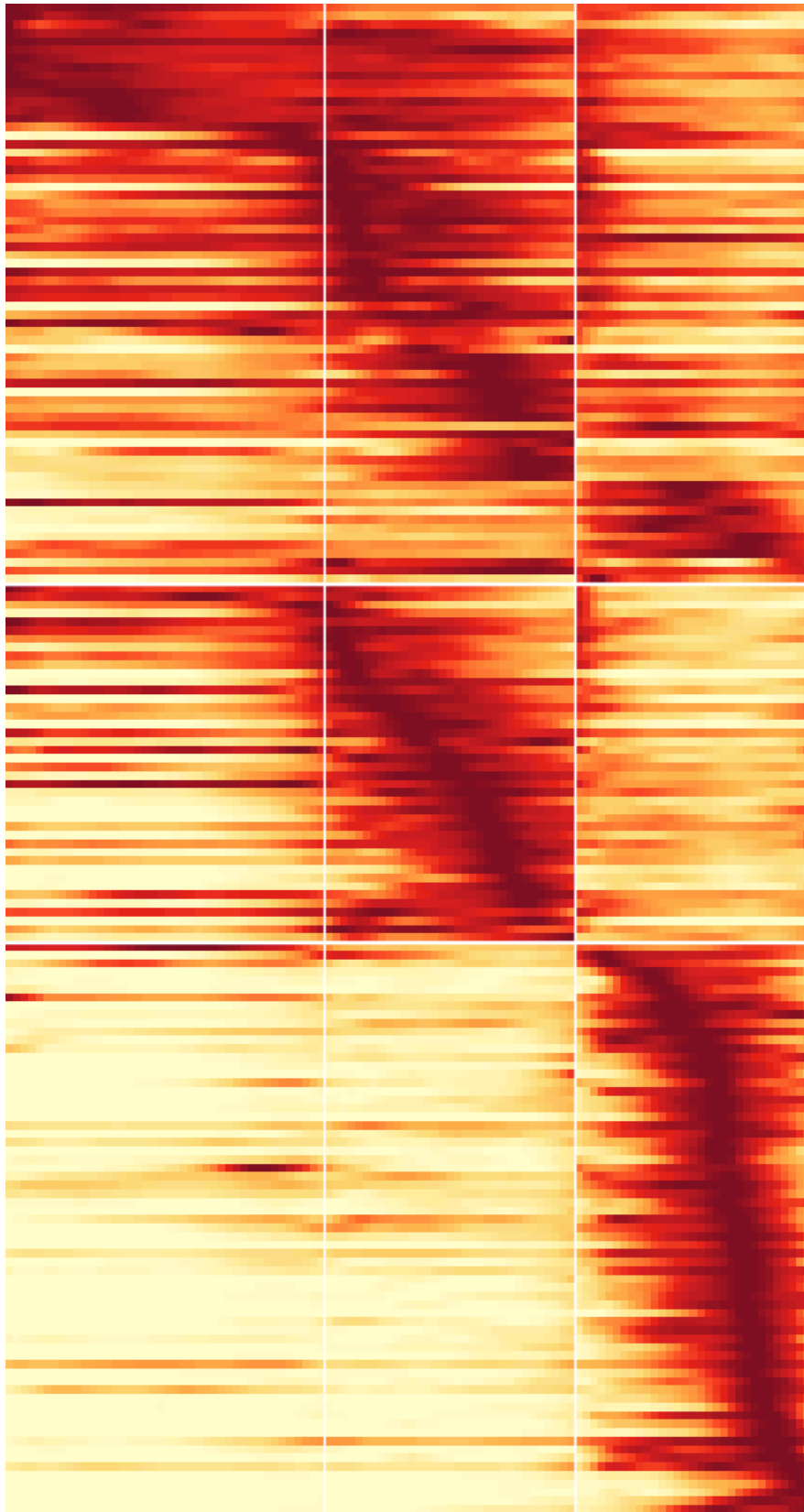
```

Generate heatmap: all genes

```

# Make sure any values <0 in the spline curves get set to 0 so that the
# heatmap scale doesn't get messed up.
splines.prdx.hm$scaled.smooth[splines.prdx.hm$scaled.smooth < 0] <- 0
# Determine where to place column separators (i.e. how many columns
# will each cell type occupy in the heatmap )
colsep <- cumsum(as.numeric(head(unlist(lapply(splines.prdx, function(x) ncol(x$scaled.smooth))),
  -1)))
# Determine where to place row separators (i.e. how many common
# markers, and markers are specific to each cell type)
rowsep <- cumsum(c(length(order.123), length(order.12)))
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/hypo-prdx-vs-non.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(splines.prdx.hm$scaled.smooth[gene.order,
  ]), Rowv = F, Colv = F, dendrogram = "none", col = cols, trace = "none",
  density.info = "none", key = F, cexCol = 0.8, cexRow = 0.35, margins = c(8,
    8), lwid = c(0.3, 4), lhei = c(0.3, 4), labCol = NA, colsep = colsep,
  rowsep = rowsep, sepwidth = c(0.05, 0.2))
title(main = "Precursors", line = -41, adj = 0)
title(main = "prdx1 -", line = -41, adj = 0.425)
title(main = "prdx1 +", line = -41, adj = 0.725)

```



scch73-281n10.2
 nkx2-4b
 rcor2
 ifi2
 ybx1
 marksab
 ube2b
 sf3b5
 hnrnpa0l.1
 ddx39ab
 sirpf
 cct5
 cct3
 hsp1
 dnajh1b
 scch211-132d3.4
 ppiia
 six5b
 gfh10
 nero
 mbx3b
 dbp
 fatp3
 esbp3b
 scdkay-56m19.5
 ncor2
 brd7
 rarg1
 CABZ01075268.2
 chd4a
 myl1b
 hnrngaba
 com2
 cct4
 smarca1
 hmx3a
 fam168a
 rbb4
 hsp11a
 casia
 foxg1a
 csnrbp1
 zc4h2
 bcam
 tubb2b
 tmsb
 tuba1a
 tmeff1b
 scch73-46j18.5
 zgc:35493
 tuba1c
 myl1a
 tuba1b
 tuba5
 ccf9l2
 mfsd2ab
 sp5-2a
 sox1a
 nkx2-4a
 sp5a
 scch211-260e23.9
 zc9-2b
 tub1b
 hsp11b
 csnrb1
 nr2f2
 tmsb-4x
 left
 hcf2
 fam60a1.1
 scch211-193j2.7
 myl1a
 ddx38aa
 rcc2
 sox11a
 pih2a
 smarcb1b
 sox11b
 dti2b
 elav3
 actf6a
 dlx2a
 vma2l1
 zfx4
 dlx1a
 wnt9
 nfix3
 meox3b
 cxxc4
 hnrnpa1a
 sox4a.1
 zgc:153867
 nfix-5
 zfx3
 atf1
 dti6a
 hnrnp
 abv1a
 nkain1
 dck1b
 cnp
 ppp1r14ba
 ataf1
 kyo9l
 dpy9l3
 cct5
 gsa1
 ddc34b
 hnc3
 ccng1
 insm1a
 tead1b
 rrm5
 pccng
 enkur
 rfx2
 plp5
 msc2b
 zgc:162707
 prdx1
 ton
 scdkayp-122a12.1
 hspan18a
 ak5l
 cksa
 slp1a1b
 pifad3
 kcnj3a
 rgs8
 rgsn
 hysm2
 b3gal2
 tvrt1
 syf6b
 pira
 scdkayp-110a12.4
 pou3f5a
 pou3f1
 rap1b
 zgc:73340
 klf12.1
 pou3f3b
 insm1b
 sog3
 ndufa42a
 plpnr2
 gpn1
 glpc2
 dac
 fash2a
 csh2
 zgc:162995
 tgh1a
 hsa
 scch211-56a11.2
 rap3db
 mt2
 scch211-270n8.4
 hsp2
 nkx1
 rgs3a
 junba
 zgc:152698
 npr2b
 sic354
 th2
 npas4b
 zgc:52140
 hectd2
 scch211-131k2.2
 hcf3
 htr1ab
 scdkay-85n7.8
 aoc1
 chg9
 scdkayp-72h1.1

Precursors

prdx1 -

prdx1 +

```
# dev.off()
```

nrgna+ neurons

Prepare cascade

```
# Get markers from the two approaches
t <- combined.markers.best[["10"]] # Markers from above the combined clades
t <- t[t$exp.global.fc >= 0.8, ] # limited to those with good global parameters
m2 <- threshold.tree.markers(markers, "2", global.fc = 0.6) # Markers from aucprTestAlongTree
m5 <- threshold.tree.markers(markers, "5", global.fc = 0.6) # Markers from aucprTestAlongTree
tacsyn.markers <- unique(c(rownames(t), rownames(m2), rownames(m5)))

# Just want to plot part of cells from upstream segment 12, which is
# very long. Going to use cells from segments 2 or 5, 10, and from
# segment 12 with pseudotime > 0.23
cells.plot.2 <- unique(c(intersect(cellsInCluster(obj, "segment", "12"),
  whichCells(obj, "pseudotime", c(0.45, Inf))), cellsInCluster(obj, "segment",
  c("10", "2"))))
cells.plot.5 <- unique(c(intersect(cellsInCluster(obj, "segment", "12"),
  whichCells(obj, "pseudotime", c(0.45, Inf))), cellsInCluster(obj, "segment",
  c("10", "5"))))
cells.plot.25 <- unique(c(intersect(cellsInCluster(obj, "segment", "12"),
  whichCells(obj, "pseudotime", c(0.45, Inf))), cellsInCluster(obj, "segment",
  c("10", "2", "5"))))

# Calculate spline curves
spline.2 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cells.plot.2,
  genes = tacsyn.markers, method = "spline", moving.window = 5, cells.per.window = 25,
  pseudotime.per.window = 0.005, spar = 0.5, verbose = F)
spline.5 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cells.plot.5,
  genes = tacsyn.markers, method = "spline", moving.window = 5, cells.per.window = 25,
  pseudotime.per.window = 0.005, spar = 0.5, verbose = F)
spline.25 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cells.plot.25,
  genes = tacsyn.markers, method = "spline", moving.window = 5, cells.per.window = 25,
  pseudotime.per.window = 0.005, spar = 0.5, verbose = F)

# Want to plot a heatmap that shows expression in upstream progenitors
# and then each branch (i.e. gabaergic_tacl vs. synpr+) as separate
# columns. Going to crop each spline fit to the correct pseudotime
# range and then combine them into a single one that can be plotted as
# a three-column heatmap.

pt.2v5 <- obj@tree$segment.pseudotime.limits["2", "start"] # pseudotime where the crop should happen
splines.tacsyn <- list(cropSmoothFit(spline.25, pt.min = -Inf, pt.max = pt.2v5),
  cropSmoothFit(spline.2, pt.min = pt.2v5, pt.max = Inf), cropSmoothFit(spline.5,
  pt.min = pt.2v5, pt.max = Inf))
names(splines.tacsyn) <- c("Precursors", "Synpr-", "Synpr+")
splines.tacsyn.hm <- combineSmoothFit(splines.tacsyn) # Combine into a single one

# Calculate gene expression timing for ordering rows
spline.2 <- determine.timing(s = spline.2)
spline.5 <- determine.timing(s = spline.5)
```

```

spline.25 <- determine.timing(s = spline.25)

# Decide which markers are specific to one cell type or both
d2v5 <- divide.branches(obj, tacsyn.markers, clust.1 = "2", clust.2 = "5",
  exp.fc = 0.4, exp.thresh = 0.2, exp.diff = 0.1)

# Generate gene ordering based on timing & specificity
order.25 <- filter.heatmap.genes(setdiff(spline.25$gene.order, c(d2v5$specific.1,
  d2v5$specific.2)))
order.2 <- filter.heatmap.genes(intersect(spline.2$gene.order, d2v5$specific.1))
order.5 <- filter.heatmap.genes(intersect(spline.5$gene.order, d2v5$specific.2))
gene.order <- c(order.25, order.2, order.5)

# Output gene table
table.save <- data.frame(gene = gene.order, marks = c(rep("both", length(order.25)),
  rep("synpr-", length(order.2)), rep("synpr+", length(order.5))), stringsAsFactors = F)
table.save$clade.AUCPR.ratio <- t[table.save$gene, "AUCPR.ratio"]
table.save$clade.exp.fc <- t[table.save$gene, "exp.fc"]
table.save$clade.exp.fc.global <- t[table.save$gene, "exp.global.fc"]
table.save$nonynpr.AUCPR.ratio.all <- m2[table.save$gene, "AUCPR.ratio.all"]
table.save$nonynpr.AUCPR.ratio.maxBranch <- m2[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$nonynpr.exp.fc.all <- m2[table.save$gene, "expfc.all"]
table.save$nonynpr.exp.fc.best <- m2[table.save$gene, "expfc.maxBranch"]
table.save$synpr.AUCPR.ratio.all <- m5[table.save$gene, "AUCPR.ratio.all"]
table.save$synpr.AUCPR.ratio.maxBranch <- m5[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$synpr.exp.fc.all <- m5[table.save$gene, "expfc.all"]
table.save$synpr.exp.fc.best <- m5[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/hypo-nrgna.csv"))

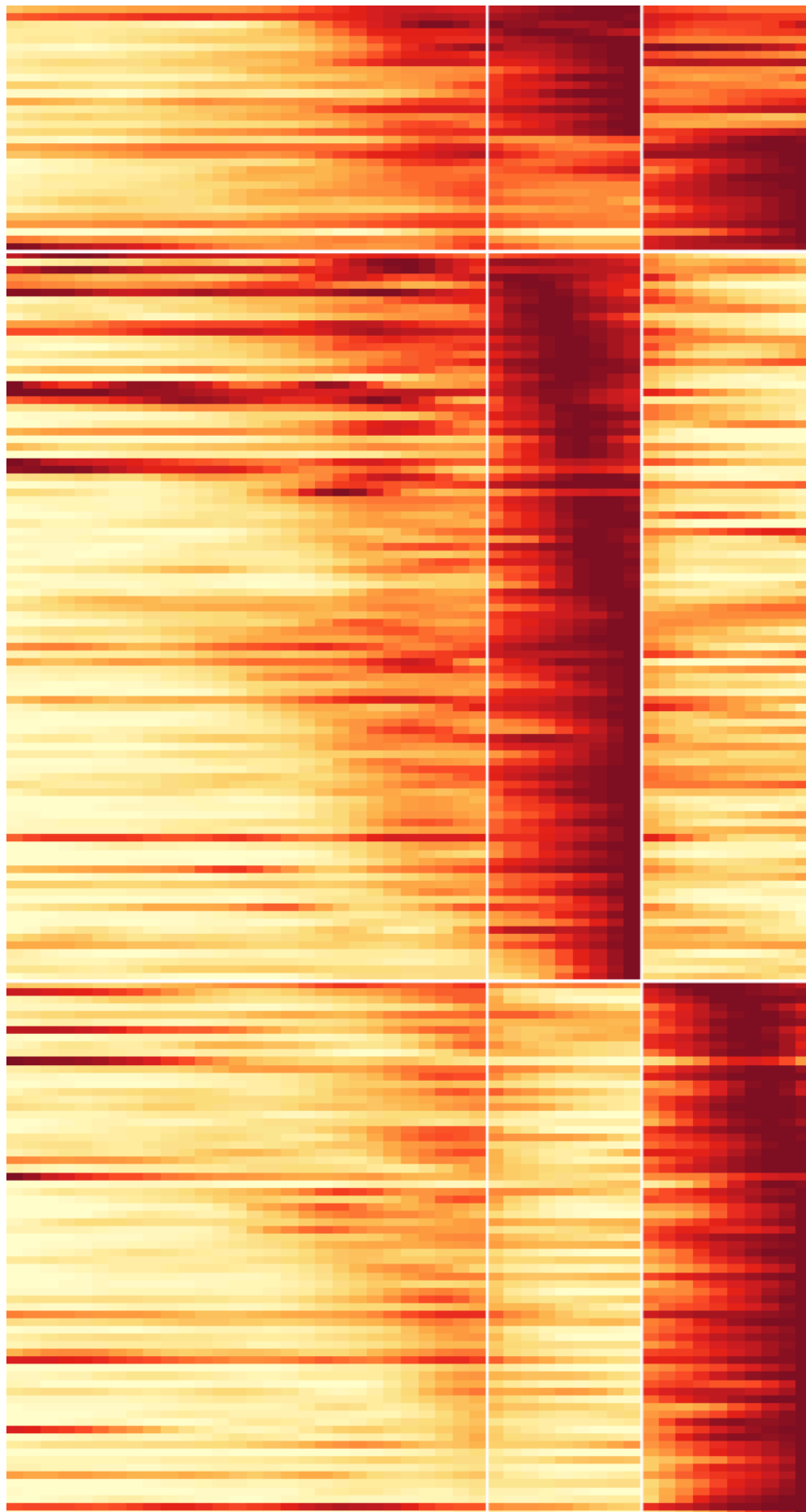
```

Generate heatmap: all genes

```

# Make sure any values <0 in the spline curves get set to 0 so that the
# heatmap scale doesn't get messed up.
splines.tacsyn.hm$scaled.smooth[splines.tacsyn.hm$scaled.smooth < 0] <- 0
# Determine where to place column separators (i.e. how many columns
# will each cell type occupy in the heatmap )
colsep <- cumsum(as.numeric(head(unlist(lapply(splines.tacsyn, function(x) ncol(x$scaled.smooth))),
  -1)))
# Determine where to place row separators (i.e. how many common
# markers, and markers are specific to each cell type)
rowsep <- cumsum(c(length(order.25), length(order.2)))
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/hypo-nrgna.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(splines.tacsyn.hm$scaled.smooth[gene.order,
  ]), Rowv = F, Colv = F, dendrogram = "none", col = cols, trace = "none",
  density.info = "none", key = F, cexCol = 0.8, cexRow = 0.35, margins = c(8,
    8), lwid = c(0.3, 4), lhei = c(0.3, 4), labCol = NA, colsep = colsep,
  rowsep = rowsep, sepwidth = c(0.05, 0.2))
title(main = "synpr-", line = -41, adj = 0.535)
title(main = "synpr+", line = -41, adj = 0.75)

```

hsbp1a
 an4d
 nptna
 sidkey-35f13.1
 reep2
 smd11b
 cpw2l
 ndufc4
 rfoxf1
 atp1b3a
 map1ab
 eno2
 egr4
 kznip1b
 mdkb
 mcl1a
 isl1
 pfdn6
 foad2
 CR047844.2
 prmi8b
 nr2f1a
 cyth3a
 camkvb
 sidkey-27f5.5
 tax
 rgma
 kong3b
 vip
 zgc:114118
 hmx3a
 map2k1
 scn3b
 penkb
 add3b
 PRIM1
 cyth1b
 lmp
 npr3
 spna2
 gabrd
 gng13b
 gfm1a
 tmem9
 gabrb3
 ppp3r1a
 oprd1b
 calcrla
 sicch73-103b11.2
 sicch73-305a9.3
 atp2b4
 tmed8
 prkcb
 RAMP1
 nr2f2
 sicch211-168d1.3
 rasgef1ba
 penka
 ntn1b
 camk1db
 tbbp1
 fam107b
 sidkey-112a7.4
 zgc:64022
 foxg1a
 mmel1
 flbp1b
 ncs1b
 dhrax
 CU929259.1
 syt6a
 sic38a3a
 chodi
 synpr
 met2cb
 dli2b
 mxd3a
 mt2
 ddc
 pkig
 egr1
 sidkey-253a1.2
 cacna1ea
 gdpd5b
 hhz1
 plex3b
 onh3
 fam126a
 sicch211-132d3.4
 VSTM2B
 zgc:91860
 gulp1a
 konf1b
 gm1a
 ppyb
 sail3a
 cart2
 syt10
 sidkey-423.7

synpr- synpr+

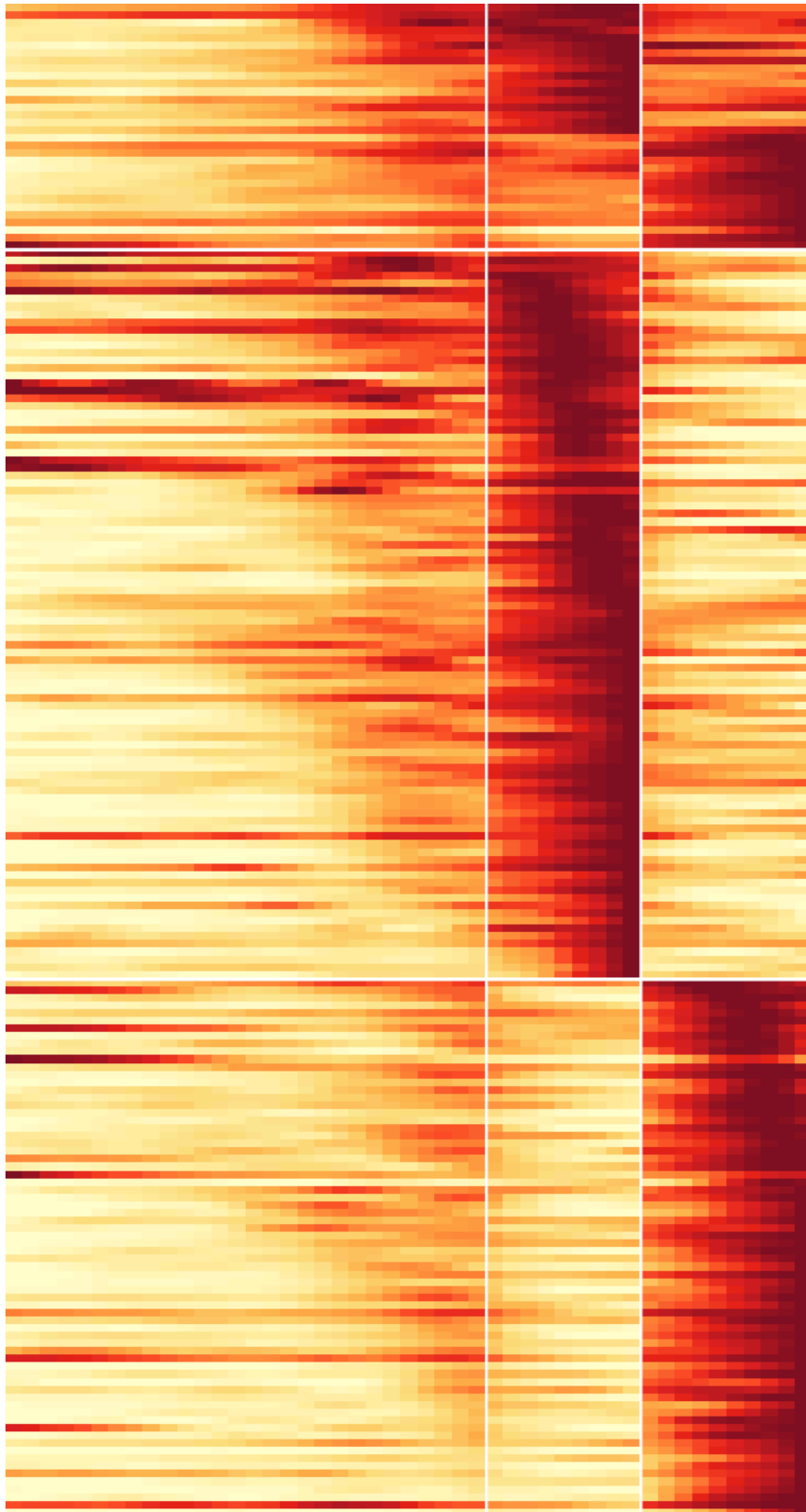
```
# dev.off()
```

Generate heatmap: main figure

```
# Plot heatmap with only certain genes labeled for main figure

genes.to.plot <- c("rgs5b", "dlx5a", "isl1", "nkx2.4a", "nkx2.2a", "hmx3a",
  "dlx1a", "dlx2b", "sp8a", "pbx3b")
rownames.to.plot <- gene.order
rtp <- rownames.to.plot %in% genes.to.plot
rownames.to.plot[!rtp] <- ""
rownames.to.plot[rtp] <- paste0("- ", rownames.to.plot[rtp])

# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/hypo-nrgna-mainfig.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(splines.tacsyn.hm$scaled.smooth[gene.order,
  ]), Rowv = F, Colv = F, dendrogram = "none", col = cols, trace = "none",
  density.info = "none", key = F, cexCol = 0.8, cexRow = 1.8, margins = c(8,
  8), lwid = c(0.3, 4), lhei = c(0.3, 4), labCol = NA, colsep = colsep,
  rowsep = rowsep, sepwidth = c(0.05, 0.2), labRow = rownames.to.plot)
title(main = "synpr-", line = -41, adj = 0.535)
title(main = "synpr+", line = -41, adj = 0.75)
```



synpr- synpr+

```
# dev.off()
```

GABAergic neurons

Prepare cascade

```
# Get markers from the two approaches
t <- combined.markers.best[["9"]] # Markers from above the combined clades
m1 <- threshold.tree.markers(markers, "1", global.fc = 0.6) # Markers from aucprTestAlongTree
m6 <- threshold.tree.markers(markers, "6", global.fc = 0.6) # Markers from aucprTestAlongTree
m7 <- threshold.tree.markers(markers, "7", global.fc = 0.6) # Markers from aucprTestAlongTree
gaba.markers <- unique(c(rownames(t), rownames(m1), rownames(m6), rownames(m7)))

# Defining cell populations to use in the heatmap Just want to use a
# tiny bit of segment 12 and then the rest of the gabaergic clade
cells.plot.1 <- unique(c(intersect(cellsInCluster(obj, "segment", "12"),
  whichCells(obj, "pseudotime", c(0.5, Inf))), cellsInCluster(obj, "segment",
  c("11", "9", "1"))))
cells.plot.6 <- unique(c(intersect(cellsInCluster(obj, "segment", "12"),
  whichCells(obj, "pseudotime", c(0.5, Inf))), cellsInCluster(obj, "segment",
  c("11", "9", "8", "6"))))
cells.plot.7 <- unique(c(intersect(cellsInCluster(obj, "segment", "12"),
  whichCells(obj, "pseudotime", c(0.5, Inf))), cellsInCluster(obj, "segment",
  c("11", "9", "8", "7"))))
cells.plot.167 <- unique(c(intersect(cellsInCluster(obj, "segment", "12"),
  whichCells(obj, "pseudotime", c(0.5, Inf))), cellsInCluster(obj, "segment",
  c("11", "9", "8", "1", "6", "7"))))

# Calculate spline curves
spline.1 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cells.plot.1,
  genes = gaba.markers, method = "spline", moving.window = 5, cells.per.window = 25,
  pseudotime.per.window = 0.005, spar = 0.5, verbose = F)
spline.6 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cells.plot.6,
  genes = gaba.markers, method = "spline", moving.window = 5, cells.per.window = 25,
  pseudotime.per.window = 0.005, spar = 0.5, verbose = F)
spline.7 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cells.plot.7,
  genes = gaba.markers, method = "spline", moving.window = 5, cells.per.window = 25,
  pseudotime.per.window = 0.005, spar = 0.5, verbose = F)
spline.167 <- geneSmoothFit(obj, pseudotime = "pseudotime", cells = cells.plot.167,
  genes = gaba.markers, method = "spline", moving.window = 5, cells.per.window = 25,
  pseudotime.per.window = 0.005, spar = 0.5, verbose = F)

# Want to plot a heatmap that shows expression in upstream progenitors
# and then each branch as separate columns. Going to crop each spline
# fit to the correct pseudotime range and then combine them into a
# single one that can be plotted as a three-column heatmap.

pt.1 <- obj@tree$segment.pseudotime.limits["1", "start"] # pseudotime where the crop should happen
splines.gaba <- list(cropSmoothFit(spline.167, pt.min = -Inf, pt.max = pt.1),
  cropSmoothFit(spline.1, pt.min = pt.1, pt.max = Inf), cropSmoothFit(spline.6,
  pt.min = pt.1, pt.max = Inf), cropSmoothFit(spline.7, pt.min = pt.1,
  pt.max = Inf))
names(splines.gaba) <- c("Precursors", "dlx+", "sst1.1+", "tph2+")
```

```

splines.gaba.hm <- combineSmoothFit(splines.gaba) # Combine into a single one

# Calculate gene expression timing for ordering rows
spline.1 <- determine.timing(s = spline.1)
spline.6 <- determine.timing(s = spline.6)
spline.7 <- determine.timing(s = spline.7)
spline.167 <- determine.timing(s = spline.167, genes = setdiff(gaba.markers,
  "sst1.2"))

# Decide which markers are specific to one cell type or not specific
d <- divide.branches.triple(obj, gaba.markers, clust.1 = "1", clust.2 = "6",
  clust.3 = "7", exp.fc = 0.4, exp.thresh = 0.2, exp.diff = 0.1)

# Generate gene ordering based on timing & specificity
order.167 <- filter.heatmap.genes(intersect(spline.167$gene.order, d$nonpecific))
order.1 <- filter.heatmap.genes(intersect(spline.1$gene.order, d$specific.1))
order.6 <- filter.heatmap.genes(intersect(spline.6$gene.order, d$specific.2))
order.7 <- filter.heatmap.genes(intersect(spline.7$gene.order, d$specific.3))
gene.order <- c(order.167, order.1, order.6, order.7)

# Output gene table
table.save <- data.frame(gene = gene.order, marks = c(rep("multiple", length(order.167)),
  rep("dlx+", length(order.1)), rep("sst+", length(order.6)), rep("tph2+",
  length(order.7))), stringsAsFactors = F)
table.save$clade.AUCPR.ratio <- t[table.save$gene, "AUCPR.ratio"]
table.save$clade.exp.fc <- t[table.save$gene, "exp.fc"]
table.save$clade.exp.fc.global <- t[table.save$gene, "exp.global.fc"]
table.save$dlx.AUCPR.ratio.all <- m1[table.save$gene, "AUCPR.ratio.all"]
table.save$dlx.AUCPR.ratio.maxBranch <- m1[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$dlx.exp.fc.all <- m1[table.save$gene, "expfc.all"]
table.save$dlx.exp.fc.best <- m1[table.save$gene, "expfc.maxBranch"]
table.save$sst.AUCPR.ratio.all <- m6[table.save$gene, "AUCPR.ratio.all"]
table.save$sst.AUCPR.ratio.maxBranch <- m6[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$sst.exp.fc.all <- m6[table.save$gene, "expfc.all"]
table.save$sst.exp.fc.best <- m6[table.save$gene, "expfc.maxBranch"]
table.save$tph.AUCPR.ratio.all <- m7[table.save$gene, "AUCPR.ratio.all"]
table.save$tph.AUCPR.ratio.maxBranch <- m7[table.save$gene, "AUCPR.ratio.maxBranch"]
table.save$tph.exp.fc.all <- m7[table.save$gene, "expfc.all"]
table.save$tph.exp.fc.best <- m7[table.save$gene, "expfc.maxBranch"]
write.csv(table.save, quote = F, file = paste0(base.path, "/heatmaps/hypo-gaba.csv"))

```

Generate heatmap: all genes

```

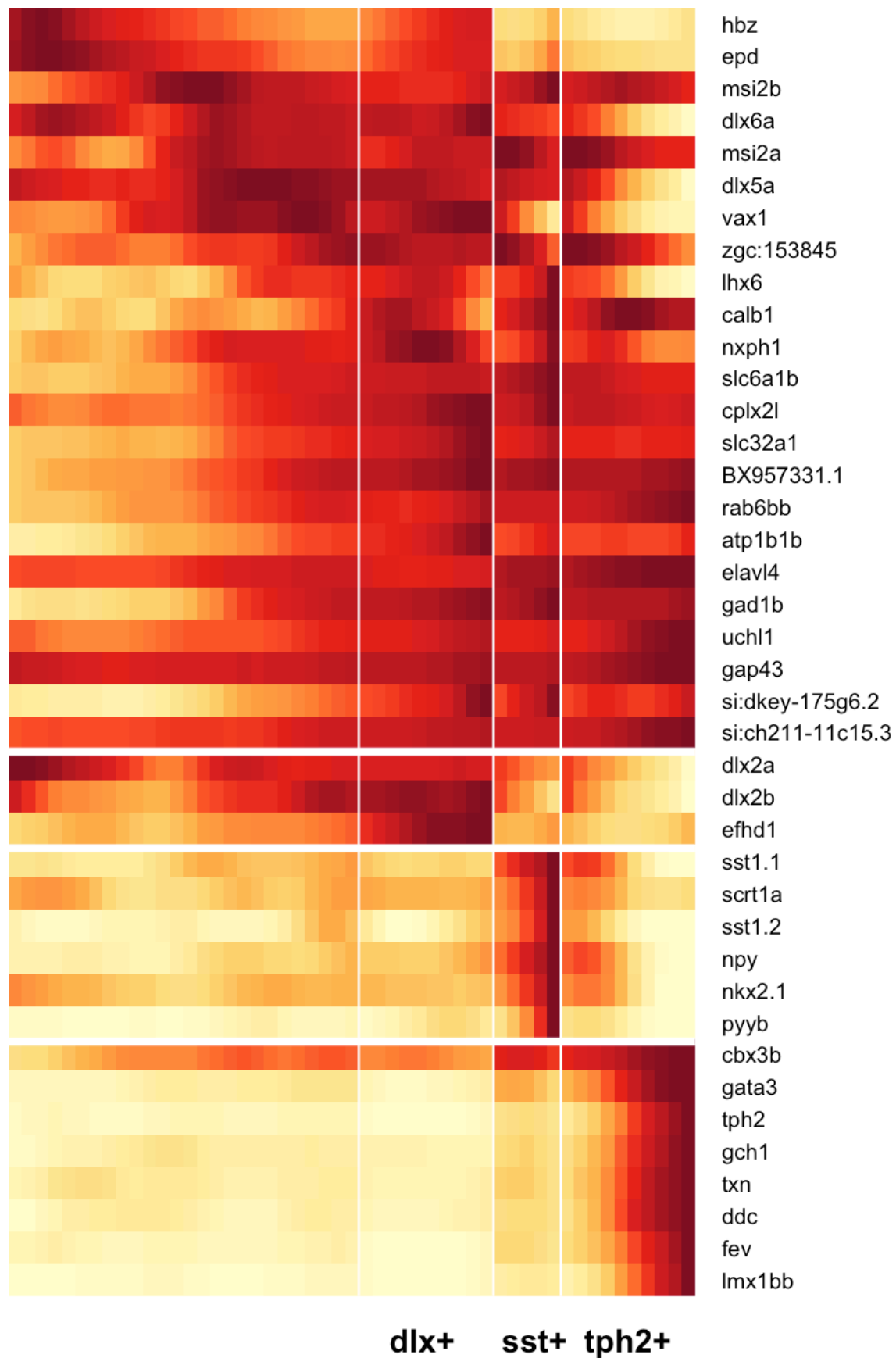
# Make sure any values < 0 in the spline curves get set to 0 so that the
# heatmap scale doesn't get messed up.
splines.gaba.hm$scaled.smooth[splines.gaba.hm$scaled.smooth < 0] <- 0
# Determine where to place column separators (i.e. how many columns
# will each cell type occupy in the heatmap )
colsep <- cumsum(as.numeric(head(unlist(lapply(splines.gaba, function(x) ncol(x$scaled.smooth))),
  -1)))
# Determine where to place row separators (i.e. how many common
# markers, and markers are specific to each cell type)

```

```

rowsep <- cumsum(c(length(order.167), length(order.1), length(order.6)))
# Open a PDF and generate the heatmap pdf(paste0(base.path,
# '/heatmaps/hypo-gaba.pdf'), width=6, height=10)
gplots::heatmap.2(x = as.matrix(splines.gaba.hm$scaled.smooth[gene.order,
]), Rowv = F, Colv = F, dendrogram = "none", col = cols, trace = "none",
density.info = "none", key = F, cexCol = 0.8, cexRow = 1, margins = c(8,
8), lwid = c(0.3, 4), lhei = c(0.3, 4), labCol = NA, colsep = colsep,
rowsep = rowsep, sepwidth = c(0.05, 0.2))
title(main = "dlx+", line = -41, adj = 0.45)
title(main = "sst+", line = -41, adj = 0.61)
title(main = "tph2+", line = -41, adj = 0.75)

```



```
# dev.off()
```

Embryonic molecular profiles are not found in larval progenitors

We were looking to compare the molecular profiles of embryonic and post-embryonic progenitors. In the retina, we found progenitors at larval stages whose molecular signatures were preserved from embryonic stages (see Retina 3). Here, we find that, in the hypothalamus, progenitors at larval stages are transcriptionally different from embryonic progenitors.

Identify populations to compare

We defined progenitor / precursor / neuron populations based on their location in the tree, cross-referenced with the expression of markers of each of these types.

```
obj@group.ids$precursor.group <- NA

cells.s13 <- intersect(cellsInCluster(obj, "segment", "13"), whichCells(obj,
  "pseudotime", c(0.05, 1)))

cells.s3neurons <- intersect(cellsInCluster(obj, "segment", "3"), whichCells(obj,
  "pseudotime", c(0.5, 1)))

cells.s3precursors <- intersect(cellsInCluster(obj, "segment", "3"), whichCells(obj,
  "pseudotime", c(0.3, 0.5)))

obj@group.ids[cells.s13, "precursor.group"] <- "1_progenitors"

obj@group.ids[c(cells.s3precursors, cellsInCluster(obj, "segment", c("12",
  "11", "10"))), "precursor.group"] <- "2_precursors"

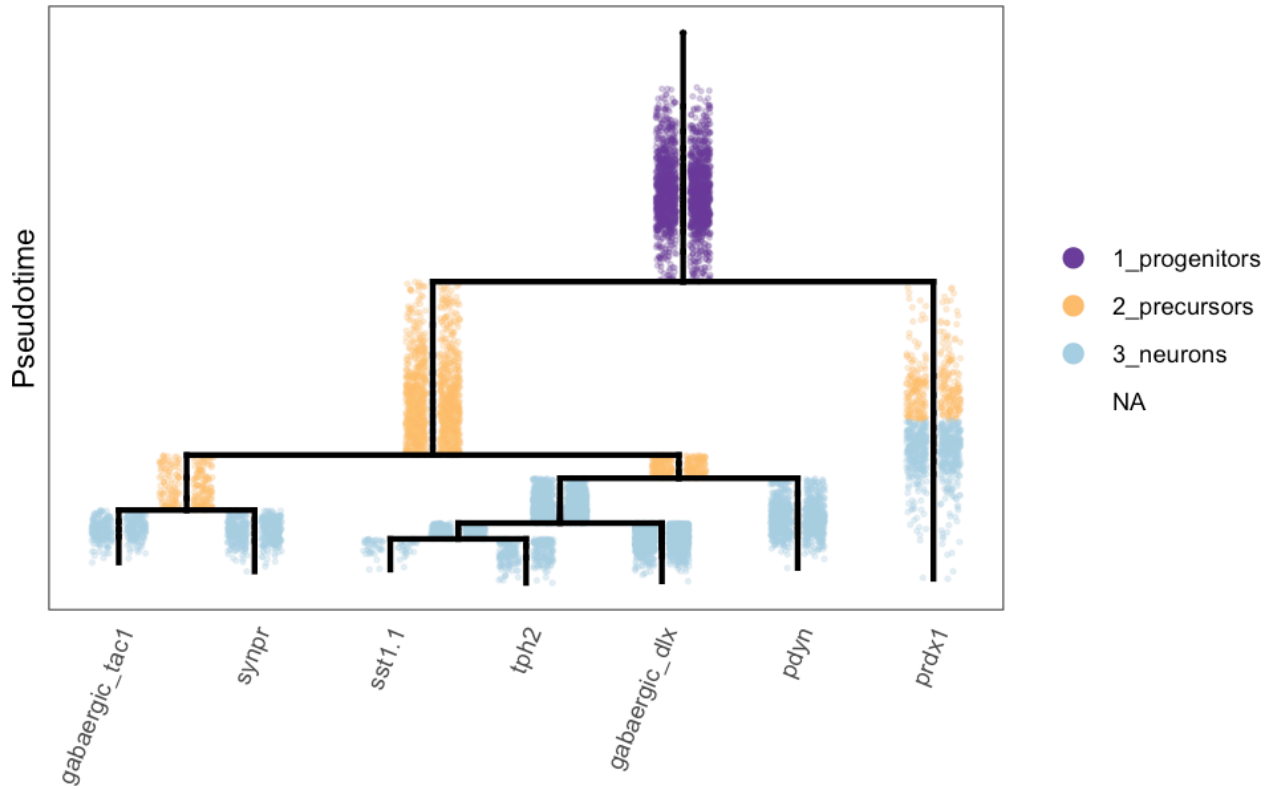
obj@group.ids[c(cells.s3neurons, cellsInCluster(obj, "segment", c("9",
  "1", "4", "8", "6", "7", "2", "5"))), "precursor.group"] <- "3_neurons"

# Colors for ggplot
stage.colors <- RColorBrewer::brewer.pal(12, "Paired")[c(10, 7, 1)]

# Plot tree to show where the groups are
plotTree(obj, "precursor.group", discrete.colors = stage.colors)

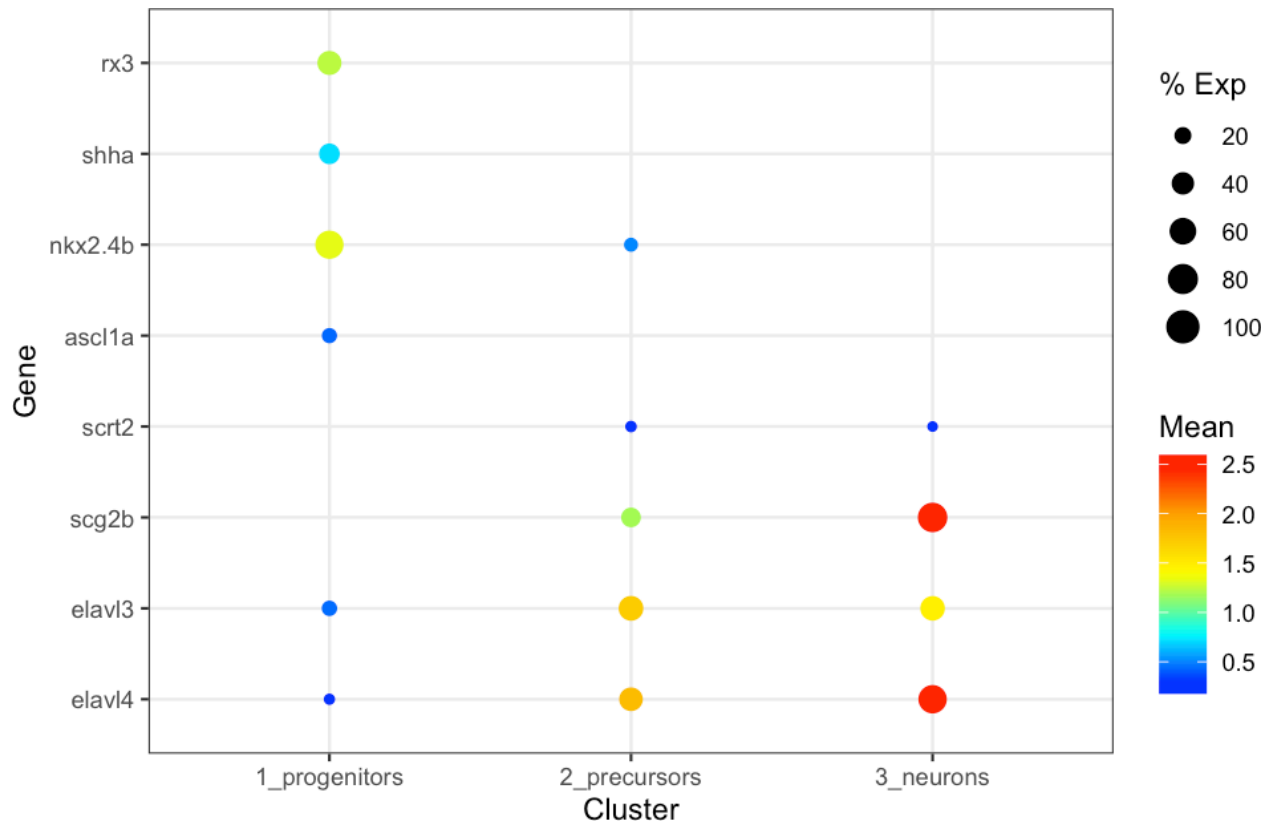
## Warning: Removed 175 rows containing missing values (geom_point).
```


precursor.group



Plot genes in each group

```
plotDot(obj, genes = c("rx3", "shha", "nkx2.4b", "ascl1a", "scrt2", "scg2b",  
"elavl3", "elavl4"), clustering = "precursor.group", scale.by = "area") +  
theme_bw()
```



Determine proportion of cells in each state

We then determined the proportion of cells from different stages that fell into each of these transcriptional states.

```
# Combine stages to reduce number of plots
obj@group.ids$stage.collapsed <- plyr::mapvalues(x = obj@group.ids$stage,
  from = c("01-12h", "02-14h", "03-16h", "04-18h", "05-20h", "06-24h",
    "07-36h", "08-2d", "09-3d", "10-5d", "11-8d", "12-15d"), to = c(rep("01-12h-24h",
    6), rep("02-36h-3d", 3), rep("03-5d-15d", 3)))

# Count number of cells from each stage group in each precursor group
stage.group.count <- plyr::count(obj@group.ids, vars = c("stage.collapsed",
  "precursor.group"))

# Remove cells that weren't part of a precursor group
stage.group.count <- stage.group.count[complete.cases(stage.group.count),
  1]

# Cast into a data frame and convert NA to 0 (no cells of that type
# observed)
stage.group.df <- reshape2::dcast(stage.group.count, formula = stage.collapsed ~
  precursor.group)

## Using freq as value column: use value.var to override.
stage.group.df[is.na(stage.group.df)] <- 0
```

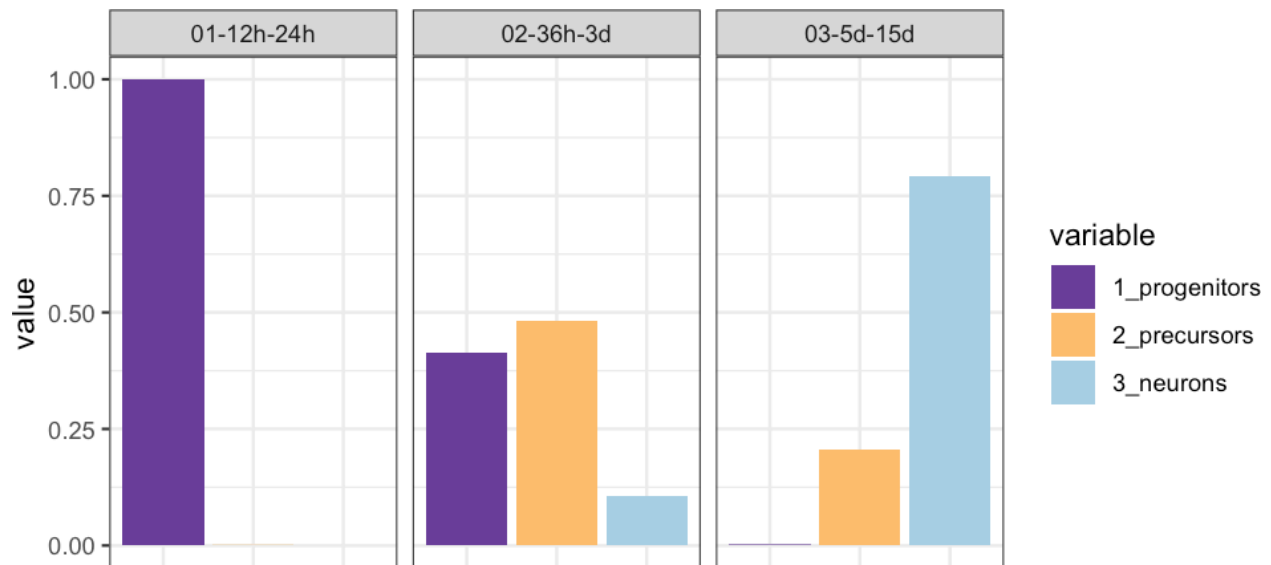
```

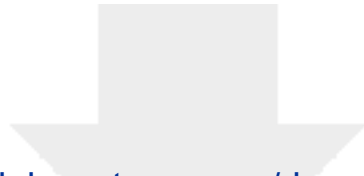
# Normalize by the number of precursors from each stage group
stage.group.df[, 2:4] <- sweep(stage.group.df[, 2:4], 1, rowSums(stage.group.df[,
  2:4]), "/")

# Melt for ggplot
stage.group.df.melt <- reshape2::melt(stage.group.df, id.vars = "stage.collapsed")

# Plot proportions
ggplot(stage.group.df.melt, aes(x = variable, y = value, group = stage.collapsed,
  fill = variable)) + geom_bar(stat = "identity") + facet_wrap(~stage.collapsed) +
  theme_bw() + scale_fill_manual(values = stage.colors) + theme(axis.title.x = element_blank(),
  axis.text.x = element_blank(), axis.ticks.x = element_blank())

```





[Click here to access/download](#)

Supplemental Videos and Spreadsheets
Raj2020_SupTable.xlsx

