

1   **Differentiation dynamics of mammary epithelial cells revealed by single-cell**  
2   **RNA-sequencing**

3

4   Karsten Bach<sup>1,2,5</sup>, Sara Pensa<sup>1,5</sup>, Marta Grzelak<sup>2,5</sup>, James Hadfield<sup>2,5</sup>, David J.  
5   Adams<sup>3,5</sup>, John C. Marioni<sup>2,3,4</sup>, Walid T. Khaled<sup>1,5</sup>

6

7   1. Department of Pharmacology, University of Cambridge, Cambridge, UK  
8   2. Cancer Research UK Cambridge Institute, University of Cambridge, Cambridge, UK  
9   3. Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK  
10   4. European Bioinformatics Institute, European Molecular Biology Laboratory, Hinxton, UK  
11   5. Cancer Research UK Cambridge Cancer Centre, Cambridge, UK

12

13   Correspondence to: Walid T. Khaled ([wtk22@cam.ac.uk](mailto:wtk22@cam.ac.uk)) or John C. Marioni  
14   ([marioni@ebi.ac.uk](mailto:marioni@ebi.ac.uk))

15

16   **Keywords: Single-cell RNAseq, Development, Mammary gland**

17

18

19   **Abstract:**

20   The mammary gland is a unique organ as it undergoes most of its development

21   during puberty and adulthood. Characterising the hierarchy of the various

22   mammary epithelial cells (MECs) and how they are regulated in response to

23   gestation, lactation and involution is important for understanding breast

24   cancer. Recent studies have used only a handful of markers to isolate and

25   characterise the different MEC compartments within the adult gland. There is

26   however, a need for a comprehensive description of MECs at different

27   developmental stages. To this end we used single cell RNA-sequencing

28   (scRNAseq) to determine the gene expression profiles of individual MECs

29   across four adult developmental stages; nulliparous, mid gestation, lactation

30   and post weaning (full natural involution). Our data from 3,471 individual cells

31   identifies 9 distinct MEC clusters and allows their hierarchical structure across

32   development to be charted. We show that only few clusters could be fully

33   characterized by a single marker gene. We argue instead that the epithelial

34   cells – especially in the luminal compartment – should rather be

35   conceptualized as being part of a continuous spectrum of differentiation. This

36   view highlights the plasticity of the tissue and might help to explain some of

37   the conflicting results from lineage tracing studies.

38

39 **Main**

40 The purpose of the mammary gland is to provide nourishment and passive immunity  
41 to the young until they are capable of feeding themselves. From a developmental  
42 biology perspective, the mammary gland is a unique organ as it undergoes most of  
43 its development during puberty and adulthood<sup>1-4</sup>. In the pre-pubertal mouse the  
44 mammary gland consists of a rudimentary epithelial ductal structure embedded within  
45 a mammary fatpad, which is connected to the nipple<sup>5, 6</sup>. In response to the hormonal  
46 changes during puberty, the rudimentary ductal structure will proliferate and migrate  
47 to fill the entire mammary fatpad leaving a developed network of ductal structures  
48 that later serve as channels for milk transport during lactation. At the onset of  
49 pregnancy a highly proliferative stage is initiated characterised by further ductal side-  
50 branching and widespread lobuloalveolar development<sup>1</sup>. Differentiation of the  
51 epithelial cells within alveoli prepares the gland for milk production and secretion.  
52 Towards the end of pregnancy the gland is extremely dense and primarily occupied  
53 by epithelial cells and only few adipocytes. This morphology is largely maintained  
54 throughout lactation. However, in response to cessation of suckling the gland  
55 undergoes involution, which is characterised by extensive cell death and tissue  
56 remodelling<sup>4, 7</sup>. Towards the end of involution the gland reaches a morphology  
57 resembling that prior to pregnancy and subsequent pregnancies will trigger the same  
58 chain of events.

59

60 Recent efforts have focused on the identification and characterisation of the various  
61 MEC lineages within the gland that contribute to this developmental homeostasis.  
62 Pioneering fat-pad transplantation studies nearly 70 years ago were the first to  
63 demonstrate the regenerative and differentiation capacity of small numbers of cells<sup>8-</sup>

64 <sup>10</sup>. More recently the use of cell surface markers coupled with flow cytometry has  
65 been used to enrich for various progenitor and stem cell compartments<sup>11-14</sup> and  
66 shown that the imbalance of such cell populations results in cellular transformation  
67 and subsequently breast cancer<sup>15, 16</sup>. Other studies, inspired by breast cancer  
68 transcriptomic profiling, have identified transcriptional regulators of MEC types such  
69 as *GATA3* in the luminal cells<sup>14, 17</sup>. More recently, elegant lineage-tracing studies  
70 used key markers to address the contribution of each lineage to adult MEC  
71 homeostasis<sup>4, 18-20</sup>. However, in all of these studies only a handful of markers and  
72 genes were used to define the cellular hierarchy of the MECs, with a principal focus  
73 on the nulliparous developmental stage. Therefore, to properly understand its  
74 changing role throughout life, there is a need for an unbiased and comprehensive  
75 characterisation of MEC compartments at different developmental stages.

76

77 To this end we herein use single-cell RNA sequencing (scRNAseq) to map the  
78 cellular dynamics of MECs across four adult developmental stages; nulliparous, mid  
79 gestation, lactation and post weaning (full natural involution). Our data from 3,471  
80 individual cells identifies 9 distinct cell populations within the gland and allows their  
81 hierarchical structure across developmental time points to be charted.

82

### 83 **Results**

#### 84 **Single-cell RNA sequencing identifies 9 clusters of MECs**

85 First, we isolated MECs from four developmental time points; nulliparous (NP), day  
86 14.5 gestation (G), day 6 lactation (L) and two weeks post natural involution (PI). For  
87 each time point we sorted MECs based on the Epcam cell surface marker from two  
88 independent mice (Fig. 1a). All samples were then prepared for single-cell RNA-

89 sequencing using the 10x Chromium platform<sup>21</sup>. Following quality control (materials  
90 and methods), this yielded an average of 11,175 unique transcripts and 3,113 genes  
91 detected from 3,471 cells (1,681 from NP, 609 from G, 604 from L and 605 from PI)  
92 (Supplementary Fig. 1a).

93

94 Visual inspection of the data using t-distributed stochastic neighbour embedding (t-  
95 SNE) suggested that although there is some grouping of cells by time point, there are  
96 also other factors that underlie structure within the dataset (Fig. 1b and  
97 Supplementary Fig. 2a). We used unsupervised hierarchical clustering based on the  
98 expression of 704 highly variable genes to identify 9 clusters of epithelial cells (C1-  
99 C9). Importantly, these clusters were consistent with the groups observed in the t-  
100 SNE (Fig. 1c). Based on the expression of *Krt18*, *Krt8*, *Krt5*, *Krt14* and *Acta2* we  
101 grouped the clusters into 6 luminal (C1, C2, C3, C4, C5 & C8) and 3 basal clusters  
102 (C6, C7 & C9) (Fig. 1D, Supplementary Fig. 2b and 2c). Differential gene expression  
103 analysis allowed us to further characterise the clusters and infer putative identities  
104 (Fig. 1d, Table 1 and Supplementary Table 1). For example, C2 expresses high  
105 levels of beta-casein, *Csn2*, and is composed exclusively from cells from the  
106 gestation and lactation time points, suggesting that these cells are secretory alveolar  
107 cells (Fig. 1d). C1 and C3 on the other hand represent the clusters with the highest  
108 expression of hormone receptors. C5 expresses high levels of the luminal progenitor  
109 marker *Aldh1a3*<sup>22</sup> suggesting progenitor function, while C7 shows characteristics of  
110 myoepithelial cells such as high levels of *Acta2*, *Oxtr* and *Krt15* (Fig. 1d).

111

112 **Reconstruction of the luminal differentiation hierarchy**

113 We focused on cells from the NP and G time points to investigate mammary  
114 epithelial differentiation states of the gland. These differentiation states and the  
115 transitions between them can be computationally reconstructed using diffusion maps.  
116 Briefly, the method embeds the data in a low-dimensional space, where distances  
117 between cells represent the progression through a gradual but stochastic process  
118 such as differentiation. In diffusion maps constructed from all epithelial cells, we  
119 observed a clear segregation between the luminal and basal clusters, with virtually  
120 no transition states between the two (Fig. 2a). This supports the hypothesis that  
121 during normal tissue homoeostasis the two lineages are largely self-maintained,  
122 which is in agreement with the majority of the lineage tracing studies<sup>18, 23, 24</sup>. In  
123 contrast, the luminal compartment showed a distinct structure, with gradual  
124 transitions between different clusters and cells originating from a common origin (Fig.  
125 2b). We confirmed the robustness of this bifurcation by verifying that it was present  
126 when different methods of feature selection, algorithms for trajectory inference and  
127 down-sampling were employed (Supplementary Fig. 3). The expression of the  
128 progenitor marker *Aldh1a3* gradually decreased as cells progressed away from their  
129 common origin (Fig. 2c), which was largely composed of C5. We further noted that  
130 the left arm of the differentiation trajectory terminates at C2 and shows increasing  
131 expression of *Csn2* and *Glycam1* (Fig. 2c), consistent with a secretory phenotype.  
132 Between C5 and C2 we found cells from C8 (Fig. 2b). C8 is the second most rare cell  
133 population after C9 and is found in all developmental stages in approximately equal  
134 numbers (6 cells in NP, 5 cells in G, 8 cells in L and 7 cells in PI). C8 cells express  
135 *Aldh3a1*, *Eif5*, *Cd14* and *Kit* similar to C5 but in addition express high levels of  
136 prolactin induced protein (*Pip*) and apolipoprotein D (*Apod*) (Fig. 1d). On the right  
137 arm of the differentiation trajectory, cells in C5 transitioned to C3 and C1, during

138 which the expression of *Esr1* and *Pgr* increase suggesting that this branch  
139 represents differentiation towards hormone-sensing luminal cells (Fig. 2c).

140

141 Being confident that the diffusion map recapitulates the luminal differentiation  
142 process, we then computationally inferred the two branches and ordered the cells  
143 according to their progression through “pseudotime”<sup>25</sup> (see materials and methods for  
144 further explanation) (Fig. 3a). We found 828 genes that showed a pseudotime-  
145 dependent expression with the same directionality along both differentiation  
146 trajectories (Fig. 3b, top panel of heatmap, Supplementary Table 2, Supplementary  
147 Fig. 4a). These included genes associated with known progenitor characteristics  
148 such as *Aldh1a3*, *Kit*, *Cd14* and *Tspan8* as well as transcription factors that have not  
149 previously been associated with luminal differentiation such as *Creb5*, *Hey1* and  
150 *Fosl1* (Fig. 3c). In addition, we identified 1459 genes with branch-specific expression  
151 patterns (Fig. 3b, bottom panel of heatmap, Supplementary Table 3, Supplementary  
152 Fig. 4b), including milk related genes (*Wap*, *Csn1s1*, *Csn3*) as well as genes involved  
153 in hormone-receptor signalling (*Pgr*, *Prlr*, *S100a6*, *Cited1*, *Areg*). Amongst the  
154 branch-specific transcription factors we found for example *Runx1*, *Tox2* and *Bhlhe41*  
155 (also known as *Sharp1*) to be transcribed during differentiation towards the hormone-  
156 sensing lineage (Fig. 3d) and *Elf5*, *Foxs1* and *Ehf* in the secretory lineage (Fig. 3e).  
157 Interestingly, *Runx1* is a known repressor of *Elf5* and its deletion has been shown to  
158 be deleterious for ductal morphogenesis<sup>26</sup>. The expression of the known progenitor  
159 master regulator *Elf5* is maintained and further increased during secretory  
160 differentiation suggesting that its transcriptional level is fine-tuned in luminal  
161 progenitors (Fig. 3e).

162

163 **Cluster 4 represents a post involution specific progenitor cell**

164 From the initial clustering analysis we noticed that C4 appeared only after the  
165 animals had undergone a full pregnancy and was predominantly made up of cells  
166 from PI (0 cells from NP, 0 cells from G, 83 cells from L and 339 cells from PI) (Fig.  
167 1b). Moreover, multiple progenitor marker genes from C5 were also present in C4  
168 (Fig. 1c, Supplementary Fig. 5). The similarity between the two clusters also held true  
169 on a transcriptome-wide level (Fig. 4a, Supplementary Fig. 10c). We thus  
170 hypothesised that C4 represents the luminal progenitor population in the post-parous  
171 gland, which is supported by the observation that C5 virtually disappears after parity  
172 (Table 1). To further support this hypothesis, we identified genes that distinguish C5  
173 from the rest of the luminal compartment to see if they are also characteristic for the  
174 proposed post-parity progenitor population C4. Indeed, we find genes that are  
175 differentially expressed between C5 and the rest of the luminal compartment to show  
176 the same trend between C4 and the luminal compartment (Fig. 4b). In a similar  
177 manner we can distinguish C4 from the rest of the PI gland in a principal component  
178 analysis (PCA) using the identified progenitor genes (Supplementary Fig. 6a and b).  
179 Given this high degree of similarity we were not surprised that only 120 genes were  
180 differentially expressed between C4 and C5 (Supplementary Table 4). Interestingly,  
181 genes that were up-regulated in C4 were significantly enriched for pathways that are  
182 involved in the immune response and lactation (Fig. 4c, Supplementary Fig. 6c).  
183 Furthermore, the expression of milk genes was surprisingly maintained at high levels  
184 in C4 even after lactation (Fig. 4d). Of note, the genes of the casein locus (*Csn2*,  
185 *Csn1s1*, *Csn1s2a*, *Csn3*) have previously been reported to be up-regulated in the  
186 parous gland, most likely due to changes in chromatin accessibility<sup>27, 28</sup>. However, it  
187 has not been shown before that this effect is confined to the progenitor population of

188 the luminal compartment (Supplementary Fig. 7). Together, the data suggest that  
189 luminal progenitor cells maintain memory of having undergone gestation and  
190 involution. This memory could potentially prime progenitor cells towards the alveolar  
191 fate to facilitate alveogenesis in subsequent pregnancies.

192

### 193 **Discussion**

194 We have reported here the use of single-cell RNA sequencing to comprehensively  
195 map the transcriptomes of thousands of MECs across four developmental time  
196 points. Our analysis identified 9 clusters of epithelial cells, some of which are only  
197 present during specific developmental stage (e.g. C2 and C4). This study provides a  
198 rich dataset that can be mined online (see link in materials and methods) to identify  
199 marker genes and lineage specific factors that can be used to trace populations of  
200 cells *in vivo*. We note, however, that only some of the clusters can be fully  
201 characterized by a single marker gene. Instead we argue that the epithelial cells –  
202 especially in the luminal compartment – should rather be conceptualized as being  
203 part of a continuous spectrum of differentiation as visualised in Fig. 2. This view  
204 highlights the plasticity of the tissue and might help to explain some of the conflicting  
205 results from lineage tracing studies<sup>4</sup>. In this study we could not provide any evidence  
206 for contribution of a putative multipotent stem cell to the day-to-day homoeostasis of  
207 the gland. We do note, however, that the cluster C9 is placed between the luminal  
208 and basal cells in Fig. 2a. C9 represents the cluster with the fewest numbers of cells  
209 captured and expresses low levels of both luminal and basal markers, high levels of  
210 the stem cell marker *Procr*<sup>19</sup> and high levels of the luminal progenitor marker  
211 *Notch3*<sup>29</sup> (Fig. 1d). Yet, we did not capture cells in transition between C9 and any  
212 other clusters suggesting that if these cells are indeed, multipotent progenitors their

213 involvement in tissue homoeostasis is negligible. Alternatively, we cannot exclude the  
214 possibility that C9 is enriched for doublets.

215

216 Based on the gene expression data presented here, the luminal compartment  
217 appears to have one common progenitor population (C5). We also identified a  
218 previously described progenitor population with expression of hormone-receptor  
219 related genes (C3) that appears to represent a meta-stable state on the  
220 differentiation path from luminal progenitors towards fully differentiated hormone-  
221 sensing cells (C1). It is worth noting that the expression of *Ly6a* (*Sca-1*) a marker  
222 recently described to differentiate between two sets of luminal progenitors<sup>22, 30</sup> is  
223 absent from C5 and is present in C3 (Fig. 1d). In the same study they describe that  
224 *Sca-1*<sup>+ve</sup> cells are also ER<sup>+ve</sup> and *Sca1*<sup>-ve</sup> progenitors are ER<sup>-ve</sup> this is in agreement  
225 with the expressions of both genes in C5 and C3. Furthermore, we characterised  
226 gene expression patterns along the differentiation hierarchy, thus enabling to  
227 disentangle the cellular hierarchy in the mammary gland. We were able to identify  
228 new transcription factors, which mark the progenitor populations and the two  
229 differentiated lineages (Fig. 3 and Supplementary Table 2 and 3). Investigating how  
230 these novel transcription factors are involved in cell fate decisions and cancer  
231 development should be investigated in the future.

232

233 By analysing the mammary gland at various stages of development we were also  
234 able to describe the molecular effects of parity at cellular resolution. We found that  
235 the luminal progenitor compartment undergoes lasting changes at the transcriptional  
236 level. This is especially interesting in light of the protective effect of pregnancies  
237 against breast cancer<sup>31, 32</sup> and the role of luminal progenitors as cell of origin<sup>15</sup>. The

238 majority of the changes were related to pathways involved in immunity and lactation,  
239 suggesting that, in particular, the luminal progenitors maintain a memory of gestation  
240 and involution. It is reasonable to assume that C4 overlaps with the previously  
241 described parity-induced MECs (PI-MECs)<sup>33</sup>.

242

243 In summary, this study provides a novel view of mammary gland development. Our  
244 unbiased approach helps support some previously formed hypotheses in the  
245 mammary gland field and describes differentiation processes at a high cellular  
246 resolution. The dataset will be a useful resource for future studies that aim to  
247 understand the relationship of the different cell types in the gland and how breast  
248 cancer develops and progresses.

249 **Figure legends**

250 **Figure 1. Single-cell RNA sequencing identifies 9 clusters of MECs**

251 **(a)** Schematic diagram highlighting the experimental setup for isolating and  
252 sequencing the RNA of single cells using the 10X Chromium system. **(b)** t-SNE plot  
253 of 3,471 cells visualizes general structure in the data. Cells are coloured by the four  
254 developmental time points Pink=NP, Dark Green=G, Light Green=L, Purple=PI. **(c)**  
255 Same as **(b)** but colouring cells by clusters. **(d)** Heatmap highlighting some key  
256 marker genes. Colour scale represents log-transformed and normalized UMI counts  
257 scaled to a maximum of 1 per row. The upper panel shows genes that were used to  
258 distinguish between luminal and basal cells. Upper bars represent the cluster  
259 assignment and stages for the individual cells. For visualization purposes only 100  
260 randomly selected cells were shown for each of the large clusters. The list of  
261 differentially expressed genes between clusters can be found in Supplementary  
262 Table 1.

263

264 **Figure 2. Computational reconstruction of differentiation processes in the**  
265 **mammary gland**

266 **(a)** Diffusion map of epithelial cells from the NP and G time points, showing the first  
267 three diffusion components. **(b)** Differentiation trajectory of the luminal compartment  
268 based on the first two diffusion components. **(c)** This panel depicts the log-  
269 transformed normalized expression levels of *Aldh1a3*, *Csn2*, *Glycam1*, *Pgr* and *Esr1*  
270 in all cells along the differentiation trajectory shown in **(b)**. Colour scale represents  
271 log-transformed and normalized UMI counts scaled to a maximum of 1.

272

273 **Figure 3. Pseudotime ordering identifies genes associated with luminal  
274 differentiation**

275 **(a)** Definition of the hormone-sensing and secretory differentiation branch. Cells are  
276 coloured by their progression through pseudotime, where low values represent  
277 undifferentiated cells. **(b)** Heatmap of the top 50 pseudotime-dependent genes (rows)  
278 with the same overall trend along both branches (top panel) or branch-specific trends  
279 (bottom panel, see materials and methods for definition). The black vertical line  
280 represents the common origin from the progenitor population. Arrows highlight  
281 differentiation direction. Pseudotime and the cluster assignment are annotated above  
282 the heatmap. The values in the heatmap represent z-scaled, spline-smoothed  
283 expression values. **(c-e)** Examples of transcription factors with pseudotime-  
284 dependent expression with the same overall trend on both branches **(c)** or branch  
285 specific trends **(d,e)**.

286

287 **Figure 4. The effect of parity on the transcriptomic landscape of the luminal  
288 progenitor compartment**

289 **(a)** Colour-coded similarity matrix between the log-transformed mean expression  
290 values of the 9 epithelial cell clusters. Similarities are based on Euclidean distances.  
291 **(b)** Comparison of fold changes from C5 versus luminal compartment and fold  
292 changes from C4 versus the luminal cells. The genes represent the top 500  
293 differentially expressed genes between C5 and luminal cells. **(c)** Volcano plot  
294 illustrates differential expression between C5 and C4, coloured dots represent  
295 significant genes with known function in lactation and immunity, and dashed lines  
296 highlight the P value threshold of 0.01 and a log fold change of 1. P values are  
297 adjusted for multiple testing using Benjamini-Hochberg. **(d)** Visualization of

298 expression difference for some lactation related genes from **(c)**, where C4 is split  
299 according to the developmental time point.

300

301 **Table 1. Summary of MEC clusters**

302 Overview of the different clusters including number of cells captured for each time  
303 point and key genes that were used to infer their identities.

304

305 **Supplementary Figure 1. Quality control of sequencing data**

306 **(a)** Table summarizing quality control criteria per sample. Number of unique  
307 molecules, genes detected and number of reads represent the median value for each  
308 sample. **(b-c)** Histograms for the four conditions showing the distributions of the  
309 number of genes detected **(b)** or total number of molecules **(c)** metric with indicated  
310 threshold (dashed line). **(d)** Scatterplot of number of genes detected versus  
311 percentage of mitochondrial RNA molecules, the threshold at 5% is indicated. For  
312 choice of thresholds see materials and methods.

313

314 **Supplementary Figure 2. Differences between replicates and expression of**  
315 **luminal and basal markers**

316 **(a)** t-SNE plot coloured by the eight different samples. Some minor structures (e.g.  
317 two groups PI1 and PI2 at the top) are mainly replicate dependent, but all of the  
318 clusters identified in the analysis are represented by multiple replicates. **(b-c)** t-SNE  
319 plots highlighting the expression of luminal **(b)** and basal **(c)** genes in all 3,471 cells.

320

321 **Supplementary Figure 3. Luminal bifurcation is robust to various alterations of**  
322 **the trajectory inference**

323 **(a-b)** Eigenvalues for the diffusion map used in Fig. 2a. **(a)** and Fig. 2b **(b)**. **(c)** The  
324 differentiation trajectory as determined by Monocle (see materials and methods)  
325 coloured by cluster assignment. **(d)** The diffusion map is robust towards down-  
326 sampling of cells (100, 50 or 25% of all cells were used in the left, middle or right  
327 panel, respectively) as well as the method of feature selection (all= all genes with  
328 mean expression level above 0.1, HVG= highly variable genes, PCA= first 50  
329 components from PCA, selected= a manually selected choice of genes that are  
330 known to be involved in luminal cell differentiation). The gene list included: *Csn2*,  
331 *Gata3*, *Prlr*, *Elf5*, *Esr1*, *Pgr*, *Aldh1a3*, *Wap*, *Tspan8*, *Krt18*, *Krt8*, *Areg*, *Fgfr1*, *Fgfr2*,  
332 *Notch1*, *Notch3*, *Foxc1* and *Zeb2*.

333

334 **Supplementary Figure 4. Full heatmaps of pseudotime dependent expression**  
335 Heatmaps as in Fig. 3, showing all genes that are pseudotime dependent in their  
336 expression along the trajectories with the same overall trend **(a)** or opposing trend  
337 **(b)**. Full gene list can be found in Supplementary Tables 2 and 3.

338

339 **Supplementary Figure 5. Expression of progenitors markers in C4**  
340 Panel of t-SNE plots (compare to Fig. 1) with coloured normalized log-expression  
341 values of progenitor markers.

342

343 **Supplementary Figure 6. C4 is a post-parity progenitor population**  
344 **(a)** Principal component analysis (PCA) on all luminal cells from the PI time point.  
345 The PCA was computed on the top 500 differentially expressed genes between C5  
346 and the rest of the luminal NP gland. PC1 separates the progenitors and the  
347 differentiated cells, with the progenitors showing negative PC1 values. **(b)** Genes that

348 are higher expressed in C5 compared to the rest of the NP gland also have negative  
349 PC1 loadings. **(c)** Top 20 GO-terms (biological processes) that are significantly  
350 enriched in upregulated genes from Fig. 4c. Dashed line indicates P value threshold  
351 at 0.001.

352

353 **Supplementary Figure 7. Expression of milk genes after parity is restricted to**  
354 **the progenitor compartment**

355 Panel of t-SNE plots (compare to Fig. 1) with coloured normalized log-expression  
356 values of milk genes from Fig. 4d.

357

358 **Supplementary Figure 8. Quality control of cells from the Lactation sample**

359 t-SNE plot computed only on the lactation sample. Cells that were flagged as low  
360 quality are highlighted in grey and cells that passed the QC are shown in black.  
361 Although a majority of the sample was removed, the cells that passed represent at  
362 least all major groups from this sample (see materials and methods).

363

364 **Supplementary Figure 9. Outline of the clustering strategy**

365 **(a)** Dendrogram of the hierarchical clustering using Euclidean distances and average  
366 linkage. Colour bars indicate the animal (binary, only comparable within one  
367 condition), the condition and cluster assignment. **(b)** Cluster stability as determined  
368 by the overlap between the original cluster and the most similar cluster in the  
369 bootstrap samples. The distributions of jaccard indeces are shown as a boxplot for  
370 100 bootstrap samples. Stability is shown for three values of the deepSplit parameter  
371 (0,1,2). Deep split 1 was chosen as this maintained overall a high cluster stability  
372 while identifying relevant cell types (see materials and methods).

373

374 **Supplementary Figure 10. Clustering Quality Control and removal of C10**

375 **(a)** Total molecule count grouped by different clusters. **(b)** Markers that were used to  
376 identify C10 as lymphocyte cluster<sup>34</sup> **(c)** Separation of C10 from all other clusters as  
377 shown by a PCA on the cluster average gene expression.

378

379 **Table S1.** Differentially expressed genes between clusters

380

381 **Table S2.** List of genes with pseudotime dependent gene expression with same  
382 overall trend in the two branches

383

384 **Table S3.** List of genes with pseudotime dependent gene expression with different  
385 trends in the two branches

386

387 **Table S4.** Genes differentially expressed between C4 and C5

388

389

390 **Materials and methods**

391 **Animals**

392 All experimental animal work was performed in accordance to the Animals (Scientific  
393 Procedures) Act 1986, UK and approved by the Ethics Committee at the Sanger  
394 Institute. C57BL/6N mice were housed in individually ventilated cages under a  
395 12:12 h light-dark cycle, with water and food available *ad libitum*. The experiment  
396 was set up to allow for all of the developmental time points to be collected and  
397 tissues to be processed at the same time. Mice were euthanized by terminal  
398 anaesthesia. Females were mated with studs and allowed to litter. Tissues were then  
399 harvested at gestation day 14.5 (G), lactation day 6 (L) and day 11 post natural  
400 weaning of the pups (PI). Tissue from NP females was harvested at 8 weeks of age.  
401 Two individual mice per developmental time point were used in the study.

402 **Mammary gland dissociation into single-cell suspension**

403 Lymph node divested mammary glands (excluding the cervical pair and the proximal  
404 region of one of the fourth glands) were dissected from the mice and mechanically  
405 dissociated. The finely minced tissue was transferred to a digestion mix consisting of  
406 DMEM/F12 (Gibco) + 10 mM HEPES (Gibco) + 1 mg ml<sup>-1</sup> collagenase (Roche) + 100  
407 U ml<sup>-1</sup> hyaluronidase (Sigma) + 50 µg ml<sup>-1</sup> gentamicin (Gibco) for 2.5 hours at 37 °C  
408 and vortexed every 30 minutes. After the lysis of red blood cells in NH<sub>4</sub>Cl, cells were  
409 briefly digested with warm 0.05% Trypsin-EDTA (Gibco), 5 mg ml<sup>-1</sup> dispase (Sigma)  
410 and 1 mg ml<sup>-1</sup> DNase (Sigma), and filtered through a 40 µm cell strainer (BD  
411 Biosciences).

412 **Cell labelling, flow cytometry and sorting**

413 Single cell suspensions were incubated in HF medium (Hank's balanced salt solution  
414 (Gibco) + 1% foetal bovine serum, Gibco) + 10% normal rat serum (Sigma) for 20  
415 min on ice to pre-block before antibody staining. All antibody incubations were  
416 performed for 10 min on ice in HF media. Mammary cells were stained with the  
417 following primary antibodies: 1  $\mu$ g ml<sup>-1</sup> CD31–biotin (eBioscience); 1  $\mu$ g ml<sup>-1</sup> CD45–  
418 biotin (eBioscience); 1  $\mu$ g ml<sup>-1</sup> Ter119–biotin (eBioscience) and 0.5  $\mu$ g ml<sup>-1</sup> EpCAM–  
419 PE (Biolegend). Cells were then stained with 0.4  $\mu$ g ml<sup>-1</sup> streptavidin–PE-CF594  
420 (BD-Biosciences). Propidium Iodide (PI, 1  $\mu$ g ml<sup>-1</sup>; Sigma) was used to detect dead  
421 cells. Cells were filtered through a 30  $\mu$ m cell strainer (Partec) before sorting. Sorting  
422 of cells was done using a SH800 sorter (SONY). Single-stained control cells were  
423 used to perform compensation manually and unstained cells were used to set gates.  
424 Chip alignment and sorting calibration was performed with automatic setup beads  
425 (SONY) immediately prior to sorting. Doublets, dead cells and contaminating  
426 haematopoietic, endothelial and stromal cells were gated out and EpCAM positive  
427 cells were sorted in FACS tubes (BD Biosciences) containing 1 mL of HF. After  
428 sorting, cells were transferred to a 1.5 ml tube through a 30  $\mu$ m cell strainer, spun  
429 down and resuspended in 15-50  $\mu$ l of HF according to the number of cells sorted per  
430 sample. Representative samples were manually counted using an improved  
431 Neubauer chamber to estimate sample loss from centrifugation and filtering. Equal  
432 numbers of cells per sample were processed for scRNA library preparation. Samples  
433 were processed for scRNA library preparation within 10 hours from tissue isolation.

434 **Library preparation and sequencing**

435 Library preparation was performed according to instruction in the 10X chromium  
436 single cell kit.

437 **RNA-seq data processing**

438 Read processing was performed as previously reported<sup>21</sup>. Briefly, the Cell Ranger  
439 Single-Cell Software Suite was used for demultiplexing, barcode assignment and  
440 UMI quantification (<http://software.10xgenomics.com/single-cell/overview/welcome>).  
441 The reads were aligned to the mm10 reference genome using a pre-built annotation  
442 package obtained from the 10X Genomics website. All lanes per sample were  
443 processed using the “cellranger count” function. The outputs from different lanes  
444 were then aggregated using “cellranger aggr” with –normalize set to “none”.

445 **Quality control and preprocessing**

446 In total the Cell Ranger software identified 5,598 barcodes that contained enough  
447 unique molecules to be considered as cells (1,707 in NP, 623 in G, 2,648 in L and  
448 620 in PI). Libraries prepared from the NP, G and P time points all showed high  
449 quality that was reproducible between the two biological replicates (Supplementary  
450 Fig. 1a). We used the following metrics to flag poor quality cells: number of genes  
451 detected, total number of unique molecules (UMIs) and percentage of molecules  
452 mapped to mitochondrial genes. For these samples we then identified poor quality  
453 cells by setting a threshold on the number of genes and number of UMIs that was  
454 defined as four median absolute deviations (MAD) below the median for each sample  
455 (Supplementary Fig. 1b and c). The samples from the lactation sample, however,  
456 showed on average nearly one order of magnitude fewer genes detected and total  
457 number of UMIs. As this was specific to both lactation samples, we inferred that this  
458 was driven by the biological phenotype of the sample rather than technical artefacts.  
459 Nonetheless, we flagged a substantial amount of the lactation sample as poor quality  
460 cells by setting a fixed threshold of 500 detected genes and 1000 total molecules

461 detected (Supplementary Fig. 1b and c). This was necessary in order to allow for  
462 normalisation of all samples; otherwise, keeping more cells of lower quality would  
463 inflate the range of normalisation factors amplifying technical noise due to cell-  
464 specific biases and adversely impact downstream analyses. We ensured that we still  
465 sampled cells from all major structures present in the lactation sample  
466 (Supplementary Fig. 8). Finally, all cells with 5% or more of UMIs mapping to  
467 mitochondrial genes were defined as non-viable or apoptotic and removed form the  
468 analysis (Supplementary Fig. 1d). This left us with a total number of 3,499 cells  
469 (1,681 in NP, 609 in G, 604 in L and 605 in P). The cells were then normalised by  
470 size factors as previously described<sup>35</sup>. The log-transformed ( $\log_2(\text{counts}+1)$ ) counts  
471 of highly variable genes (HVGs) were used as features for dimensionality reduction  
472 and clustering. HVGs were defined using the method from Brennecke et al.<sup>36</sup>. Genes  
473 with a Benjamini-Hochberg adjusted P value smaller than 0.1 were considered to be  
474 significant. The t-SNE embedding in Fig. 1 was computed using the “Rtsne” package  
475 on scaled, log-transformed expression values with default settings and perplexity set  
476 to 25 (<https://github.com/jkrijthe/Rtsne>).

## 477 Clustering

478 As we expected the epithelial compartment to show a hierarchical organisation, we  
479 used agglomerative hierarchical clustering to identify cell types. We first computed  
480 the pairwise Euclidean distances between cells based on the log-transformed,  
481 normalized counts of HVGs. The dissimilarity matrix was used to perform hierarchical  
482 clustering (Supplementary Fig. 9a) with the 'hclust' function in R using average  
483 linkage, as this resulted in the highest cophenetic correlation coefficient between the  
484 tree and the dissimilarity matrix (average:0.91 single: 0.73, complete: 0.83, ward.D2:

485 0.73). Clusters were then defined with the 'cutreeDynamic' function from the  
486 'dynamicTreeCut' package (minClusterSize= 15, deepSplit=1, method="hybrid")<sup>37</sup>.  
487 The deepSplit parameter was set to 1 based on evaluating the robustness of the  
488 clustering. At a deepSplit of 0, all clusters were highly robust as determined by  
489 bootstrapping<sup>38</sup>. Increasing the parameter to 1, led to splitting C1 from C3 as well as  
490 C9 from C6. C1, C3 and C6 were highly robust in the bootstrap (see Supplementary  
491 Fig. 9b). C9 showed a below average robustness. However, due to the expression of  
492 marker genes for a previously reported cell type<sup>19</sup> we decided to keep it as a  
493 separate cluster. Interestingly, the putative cell type of C9 is known to have an  
494 expression profile that is in between luminal and basal cells (see main text), which  
495 might explain its assignment to other clusters during the bootstrap. An alternative  
496 explanation that cannot be excluded at this point is that C9 is enriched for doublets  
497 (see discussion). Further increasing the deepSplit parameter led to unstable clusters  
498 (Supplementary Fig. 9b). With this method, we identified ten clusters. In addition, the  
499 dynamic tree cut method allows for a “noise component” for cells that do not match to  
500 any cluster. 1 cell from the lactation sample was flagged as noise by the algorithm  
501 and removed from further analysis. C10 (27 cells) expressed lymphocyte related  
502 genes was also removed from all downstream analyses (as previously reported<sup>34</sup>,  
503 Supplementary Fig. 10).

#### 504 **Differential expression analysis**

505 Differential gene expression analysis was performed using “edgeR”<sup>39</sup>. Genes with a  
506 mean expression level below 0.1 counts were removed from the analysis. A negative  
507 binomial generalized log-linear model was fitted to the remaining genes with the

508 cluster assignments as covariate(s). The 'glmTreat' function was used to identify  
509 genes that have a significantly higher log fold change than 1 at an FDR of 0.01.

510 **Diffusion maps and pseudotime inference**

511 For inferring the differentiation trajectory we used diffusion maps. First, we selected  
512 all cells from the NP and G time point (Fig. 2a) and detected the HVGs as described  
513 above. The log-transformed ( $\log_2(\text{count}+1)$ ) gene counts were then used to compute  
514 the first twenty diffusion components using the 'DiffusionMap' function ( $k=50$ ,  
515 otherwise the default parameters as in 'destiny'<sup>40</sup>. After inspecting the fraction of  
516 explained variation for the first twenty components, we decided to retain the first 3. In  
517 Fig. 2b we then focused on the luminal compartment and recomputed the diffusion  
518 map based only on the luminal cells, using the aforementioned procedure. Here we  
519 retained only the first two components (Supplementary Fig. 3a). Notably, the  
520 structure inferred by the diffusion map algorithm was robust to the choice of features  
521 and down-sampling of cells (Supplementary Fig. 3d). The structure of a common  
522 origin and the two branches could also be inferred using Monocle with standard  
523 settings<sup>41</sup> (Supplementary Fig. 3c). For inferring the branches and pseudotime  
524 ordering, we defined the following three tips, the cell with the smallest value for the  
525 second eigenvector (which was set as root) and the cells with the largest and  
526 smallest values for the first eigenvector (compare Fig. 2b).

527 **Pseudotime-dependent expression**

528 To identify genes whose expression was significantly associated with the pseudotime  
529 we first fitted a natural cubic spline with three degrees of freedom to the log-  
530 transformed ( $\log_2(\text{counts}+1)$ ) expression data in each branch. A likelihood ratio test

531 was then used to assess statistical significance of the fit compared to a null  
532 (pseudotime-independent) model. Genes with a Benjamini-Hochberg corrected P  
533 value below 0.001 were considered to be significantly pseudotime-dependent. We  
534 then used a heuristic definition of branch-specific expression instead of modelling the  
535 branch assignment explicitly. This was motivated as follows. We were interested in  
536 general expression trends of genes, i.e. increase or decrease along the  
537 differentiation towards one of the two cell types, rather than comparing the exact  
538 timing of gene activation/inactivation between the branches. Any approach trying to  
539 do the latter would have been complicated by the different cell densities along the  
540 branches and by the difficulty of verifying any such hypotheses *in vivo*. Hence, we  
541 defined genes to be branch-specific when they were pseudotime-dependent in their  
542 expression in at least one of the two branches and when the gradient differed in  
543 signs between two branches. The gradient was determined as the coefficient of a  
544 linear model fit to the spline-smoothed expression values, which was set to 0 if the  
545 coefficient was not significantly different from 0 at alpha=0.01. Consequently, the  
546 gradient of a gene could either be -1 (decreasing), 0 (flat) or 1 (increasing).

547 **Gene set enrichment analysis**

548 A gene set enrichment analysis based on gene-ontology (GO) terms was conducted  
549 to characterize genes that were up-regulated in C4 compared to C5. Genes with  
550 positive log fold changes were compared to all genes that were tested for differential  
551 expression using topGO with default settings<sup>42</sup>.

552

553 **Code availability**

554 All computational analyses were performed in R (Version 3.3.3) using standard  
555 functions unless otherwise indicated. Code is available online at  
556 <https://github.com/MarioniLab/MammaryGland>. The processed data can be browsed  
557 and accessed at <https://karstenbach.shinyapps.io/webapp/>

558 **Author contribution**

559 K.B. performed most of the experiments and all the computational analysis. S.P. and  
560 K.B. setup and collected the MECs. M.G. and J.H performed the 10X library  
561 production and sequencing. J.C.M, D.A., K.B and W.T.K conceptualised the study  
562 and wrote the manuscript.

563

564 **Acknowledgements**

565 We would like to thank the staff at Sanger Institute, Research Service Facility (RSF)  
566 for their assistance. We would like to thank Dr. Aaron T. Lun (CRUK CI) for helpful  
567 discussions and comments on the manuscript. K.B. is funded by a Cambridge  
568 Cancer Centre studentship. D.A. is funded by CRUK and Wellcome Trust. S.P. is  
569 funded by CRUK. J.C.M is funded by CRUK and EMBL. W.T.K is funded by a CRUK  
570 career establishment award (C47525/A17348), University of Cambridge and  
571 Magdalene College, Cambridge.

572

573 **References**

- 574 1. Watson, C.J. & Khaled, W.T. Mammary development in the embryo and adult:  
575 a journey of morphogenesis and commitment. *Development (Cambridge, England)* **135**, 995-1003 (2008).  
576 2. Hennighausen, L. & Robinson, G.W. Information networks in the mammary  
577 gland. *Nat Rev Mol Cell Biol* **6**, 715-725 (2005).  
578 3. Hennighausen, L. & Robinson, G.W. Think globally, act locally: the making of  
579 a mouse mammary gland. *Genes & development* **12**, 449-455 (1998).

- 581 4. Inman, J.L., Robertson, C., Mott, J.D. & Bissell, M.J. Mammary gland  
582 development: cell fate specification, stem cells and the microenvironment.  
583 *Development (Cambridge, England)* **142**, 1028-1042 (2015).
- 584 5. Mikkola, M.L. & Millar, S.E. The mammary bud as a skin appendage: unique  
585 and shared aspects of development. *Journal of mammary gland biology and*  
586 *neoplasia* **11**, 187-203 (2006).
- 587 6. Hens, J.R. & Wysolmerski, J.J. Key stages of mammary gland development:  
588 molecular mechanisms involved in the formation of the embryonic mammary  
589 gland. *Breast Cancer Res* **7**, 220-224 (2005).
- 590 7. Watson, C.J. Post-lactational mammary gland regression: molecular basis and  
591 implications for breast cancer. *Expert Reviews in Molecular Medicine* **8**, 1  
592 (2006).
- 593 8. Kordon, E.C. & Smith, G.H. An entire functional mammary gland may  
594 comprise the progeny from a single cell. *Development (Cambridge, England)*  
595 **125**, 1921-1930 (1998).
- 596 9. Faulkin, L.J., Jr. & Deome, K.B. Regulation of growth and spacing of gland  
597 elements in the mammary fat pad of the C3H mouse. *J Natl Cancer Inst* **24**,  
598 953-969 (1960).
- 599 10. Daniel, C.W. Regulation of cell division in aging mouse mammary epithelium.  
600 *Adv Exp Med Biol* **61**, 1-19 (1975).
- 601 11. Smalley, M.J., Titley, J. & O'Hare, M.J. Clonal characterization of mouse  
602 mammary luminal epithelial and myoepithelial cells separated by  
603 fluorescence-activated cell sorting. *In Vitro Cell Dev Biol Anim* **34**, 711-721  
604 (1998).
- 605 12. Stingl, J. et al. Purification and unique properties of mammary epithelial stem  
606 cells. *Nature* **439**, 993-997 (2006).
- 607 13. Shackleton, M. et al. Generation of a functional mammary gland from a single  
608 stem cell. *Nature* **439**, 84-88 (2006).
- 609 14. Asselin-Labat, M.L. et al. Gata-3 is an essential regulator of mammary-gland  
610 morphogenesis and luminal-cell differentiation. *Nature cell biology* **9**, 201-209  
611 (2007).
- 612 15. Lim, E. et al. Aberrant luminal progenitors as the candidate target population  
613 for basal tumor development in BRCA1 mutation carriers. *Nature medicine* **15**,  
614 907-913 (2009).
- 615 16. Molyneux, G. et al. BRCA1 Basal-like Breast Cancers Originate from Luminal  
616 Epithelial Progenitors and Not from Basal Stem Cells. *Cell Stem Cell* **7**, 403-  
617 417 (2010).
- 618 17. Kourous-Mehr, H., Slorach, E.M., Sternlicht, M.D. & Werb, Z. GATA-3 maintains  
619 the differentiation of the luminal cell fate in the mammary gland. *Cell* **127**,  
620 1041-1055 (2006).
- 621 18. Van Keymeulen, A. et al. Distinct stem cells contribute to mammary gland  
622 development and maintenance. *Nature* **479**, 189-193 (2011).
- 623 19. Wang, D. et al. Identification of multipotent mammary stem cells by protein C  
624 receptor expression. *Nature* **517**, 81-84 (2015).
- 625 20. Rios, A.C., Fu, N.Y., Lindeman, G.J. & Visvader, J.E. In situ identification of  
626 bipotent stem cells in the mammary gland. *Nature* **506**, 322-327 (2014).
- 627 21. Zheng, G.X. et al. Massively parallel digital transcriptional profiling of single  
628 cells. *Nature communications* **8**, 14049 (2017).
- 629 22. Shehata, M. et al. Phenotypic and functional characterisation of the luminal  
630 cell hierarchy of the mammary gland. *Breast Cancer Res* **14**, R134 (2012).

- 631 23. van Amerongen, R., Bowman, A.N. & Nusse, R. Developmental stage and  
632 time dictate the fate of Wnt/beta-catenin-responsive stem cells in the  
633 mammary gland. *Cell Stem Cell* **11**, 387-400 (2012).
- 634 24. Davis, F.M. *et al.* Single-cell lineage tracing in the mammary gland reveals  
635 stochastic clonal dispersion of stem/progenitor cell progeny. *Nature  
636 communications* **7**, 13053 (2016).
- 637 25. Haghverdi, L., Buttner, M., Wolf, F.A., Buettner, F. & Theis, F.J. Diffusion  
638 pseudotime robustly reconstructs lineage branching. *Nature methods* **13**, 845-  
639 848 (2016).
- 640 26. van Bragt, M.P., Hu, X., Xie, Y. & Li, Z. RUNX1, a transcription factor mutated  
641 in breast cancer, controls the fate of ER-positive mammary luminal cells. *Elife*  
642 **3**, e03881 (2014).
- 643 27. Dos Santos, C.O., Dolzhenko, E., Hodges, E., Smith, A.D. & Hannon, G.J. An  
644 epigenetic memory of pregnancy in the mouse mammary gland. *Cell Rep* **11**,  
645 1102-1109 (2015).
- 646 28. Rijnkels, M. *et al.* Epigenetic modifications unlock the milk protein gene loci  
647 during mouse mammary gland development and differentiation. *PloS one* **8**,  
648 e53270 (2013).
- 649 29. Lafkas, D. *et al.* Notch3 marks clonogenic mammary luminal progenitor cells in  
650 vivo. *The Journal of cell biology* **203**, 47-56 (2013).
- 651 30. Giraddi, R.R. *et al.* Stem and progenitor cell division kinetics during postnatal  
652 mouse mammary gland development. *Nature communications* **6**, 8487 (2015).
- 653 31. Russo, I.H. & Russo, J. Pregnancy-induced changes in breast cancer risk.  
654 *Journal of mammary gland biology and neoplasia* **16**, 221-233 (2011).
- 655 32. Lyons, T.R., Schedin, P.J. & Borges, V.F. Pregnancy and breast cancer: when  
656 they collide. *Journal of mammary gland biology and neoplasia* **14**, 87-98  
657 (2009).
- 658 33. Wagner, K.U. *et al.* An adjunct mammary epithelial cell population in parous  
659 females: its role in functional adaptation and tissue renewal. *Development  
660 (Cambridge, England)* **129**, 1377-1386 (2002).
- 661 34. Scheele, C.L. *et al.* Identity and dynamics of mammary stem cells during  
662 branching morphogenesis. *Nature* **542**, 313-317 (2017).
- 663 35. Lun, A.T., Bach, K. & Marioni, J.C. Pooling across cells to normalize single-  
664 cell RNA sequencing data with many zero counts. *Genome biology* **17**, 75  
665 (2016).
- 666 36. Brennecke, P. *et al.* Accounting for technical noise in single-cell RNA-seq  
667 experiments. *Nature methods* **10**, 1093-1095 (2013).
- 668 37. Langfelder, P., Zhang, B. & Horvath, S. Defining clusters from a hierarchical  
669 cluster tree: the Dynamic Tree Cut package for R. *Bioinformatics (Oxford,  
670 England)* **24**, 719-720 (2008).
- 671 38. Hennig, C. fpc: Flexible Procedures for Clustering. R package version 2.1-10.  
672 <https://CRAN.R-project.org/package=fpc> ((2015)).
- 673 39. Robinson, M.D., McCarthy, D.J. & Smyth, G.K. edgeR: a Bioconductor  
674 package for differential expression analysis of digital gene expression data.  
675 *Bioinformatics (Oxford, England)* **26**, 139-140 (2010).
- 676 40. Angerer, P. *et al.* destiny: diffusion maps for large-scale single-cell data in R.  
677 *Bioinformatics (Oxford, England)* **32**, 1241-1243 (2016).
- 678 41. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are  
679 revealed by pseudotemporal ordering of single cells. *Nature biotechnology* **32**,  
680 381-386 (2014).

681 42. J, A.A.a.R. topGO: Enrichment Analysis for Gene Ontology. *R package*  
682 *version 2.26.0.* (2016).  
683

Figure 1

Bach et al.

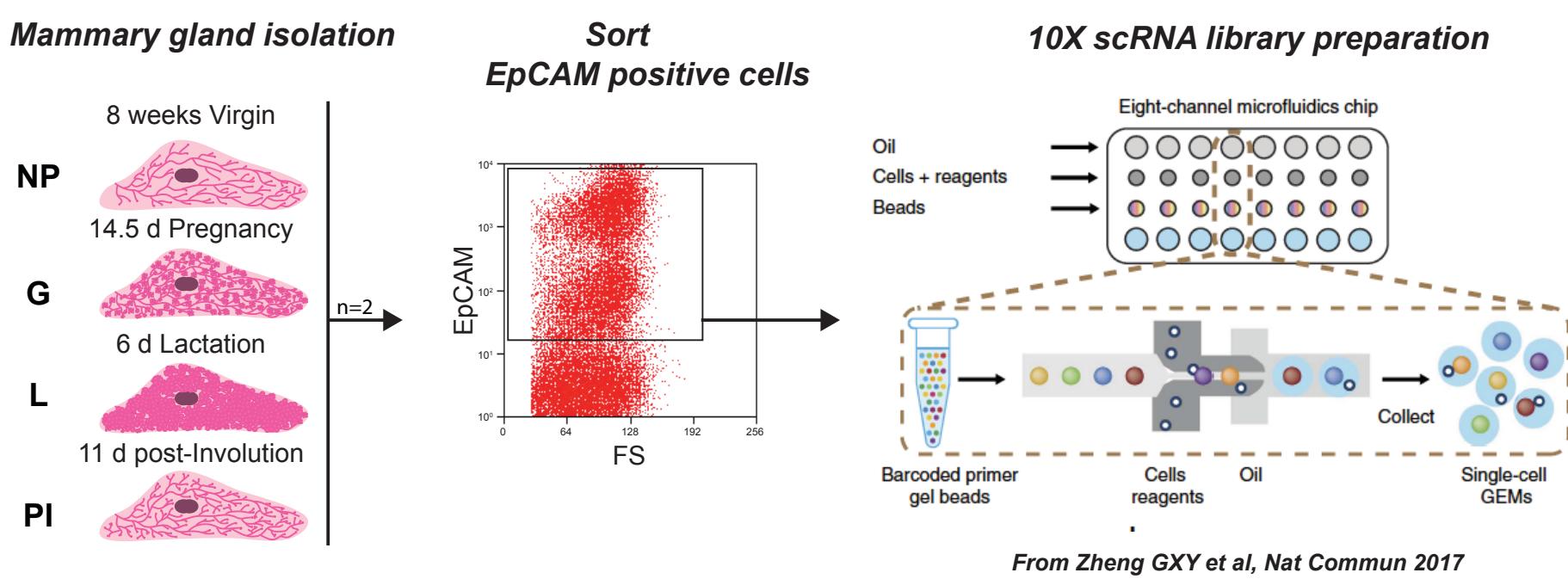
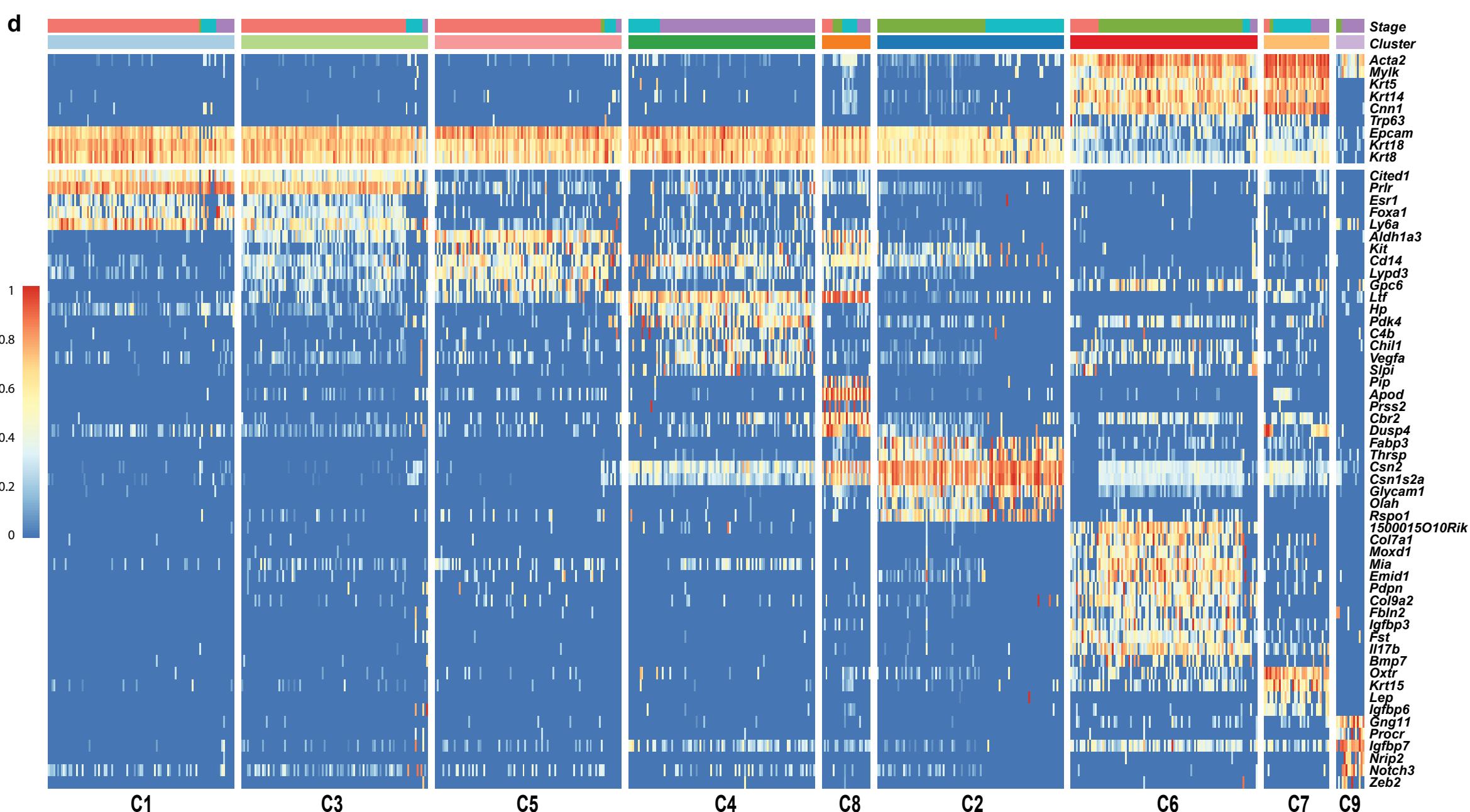
**a****b****t-SNE coloured by timepoint****c****t-SNE coloured by cluster****d**

Figure 2

Bach et al.

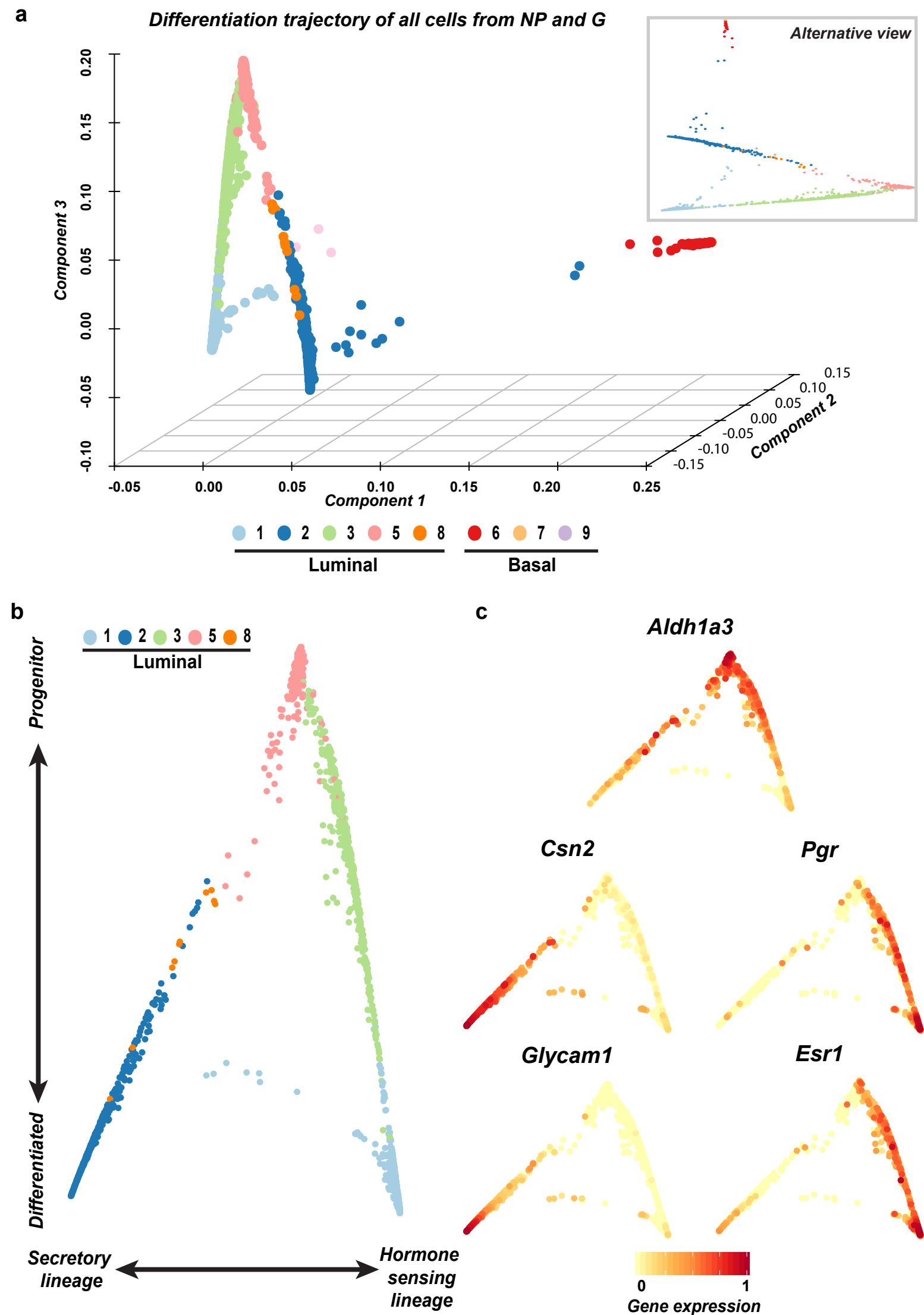


Figure 3

Bach et al.

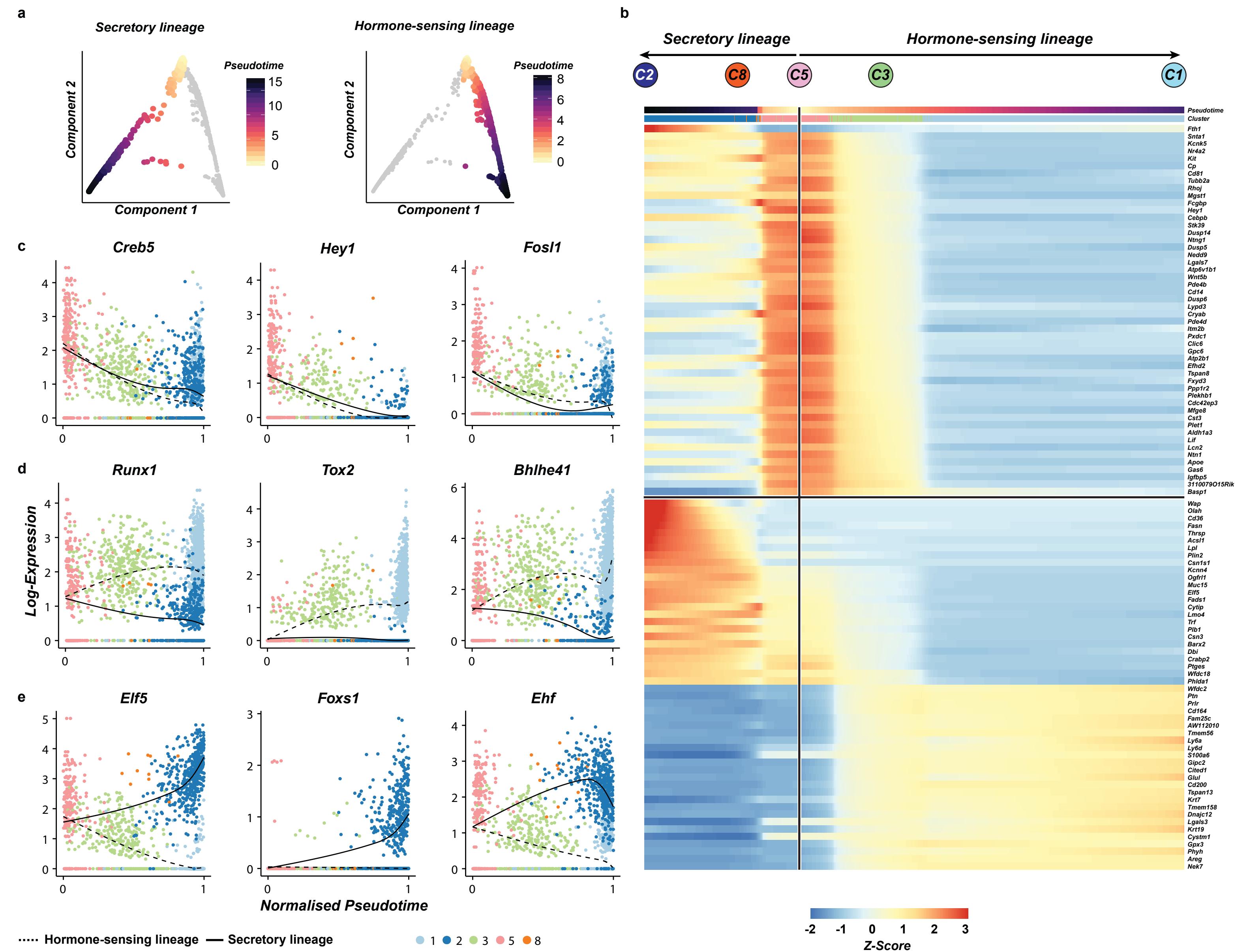
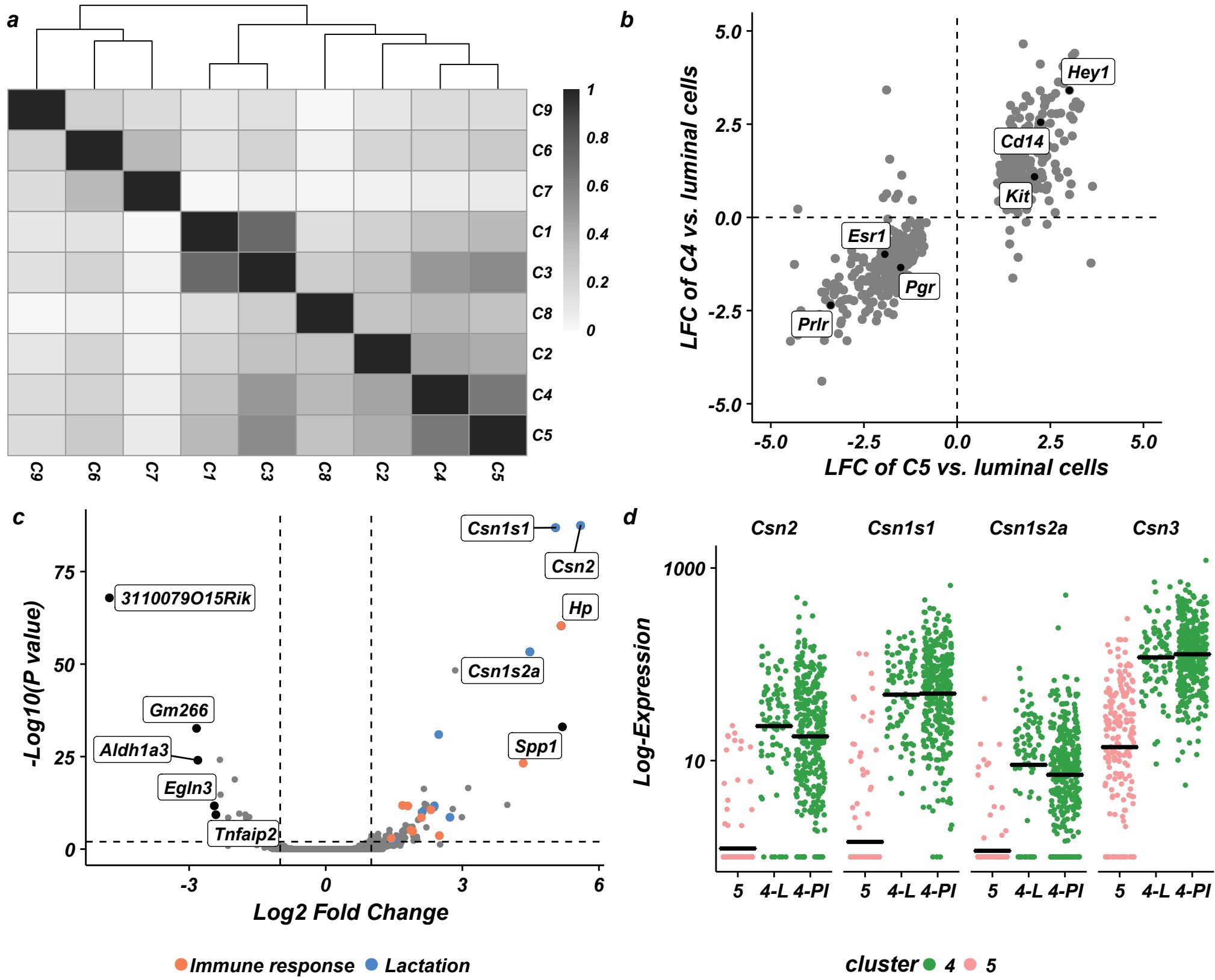


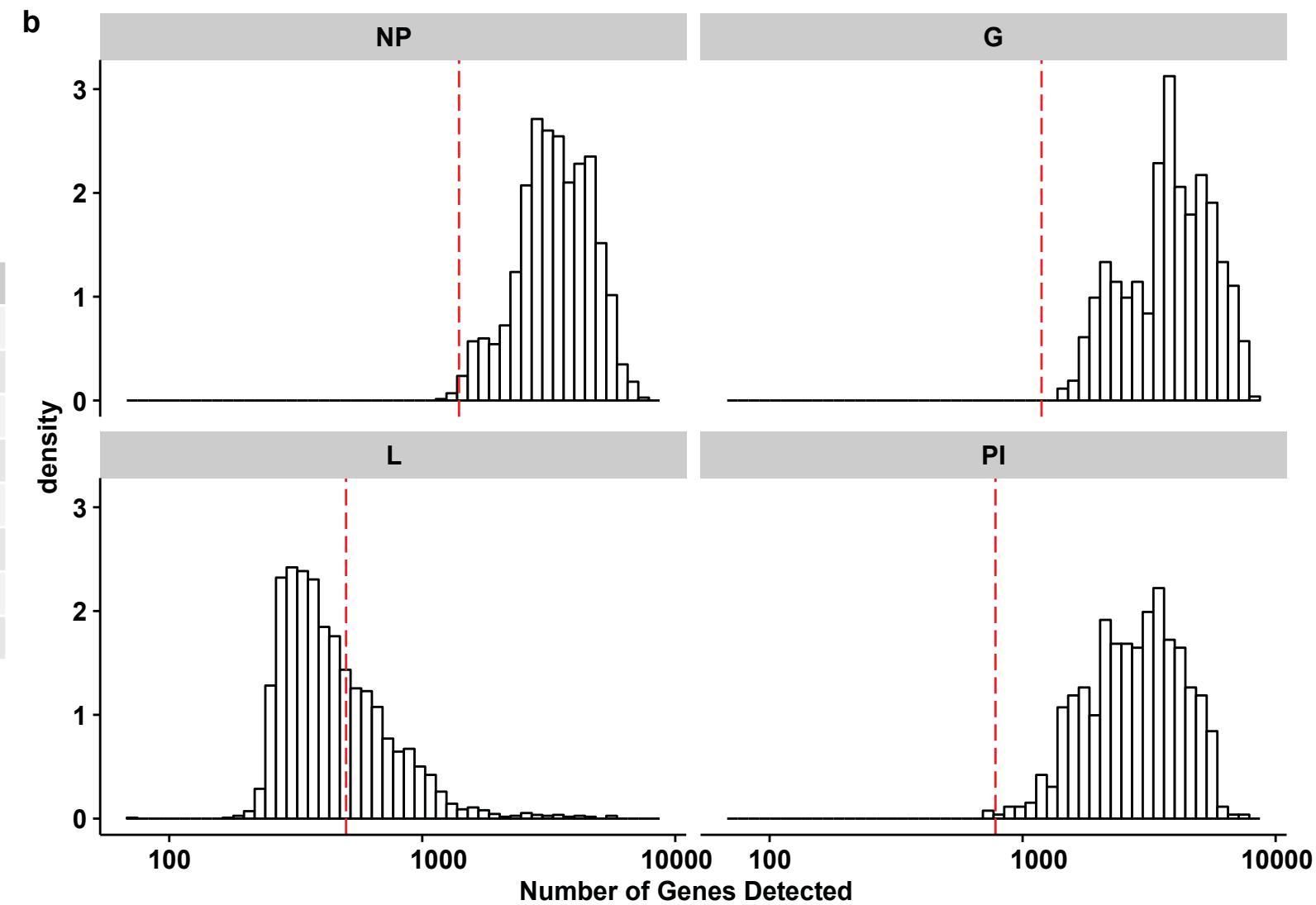
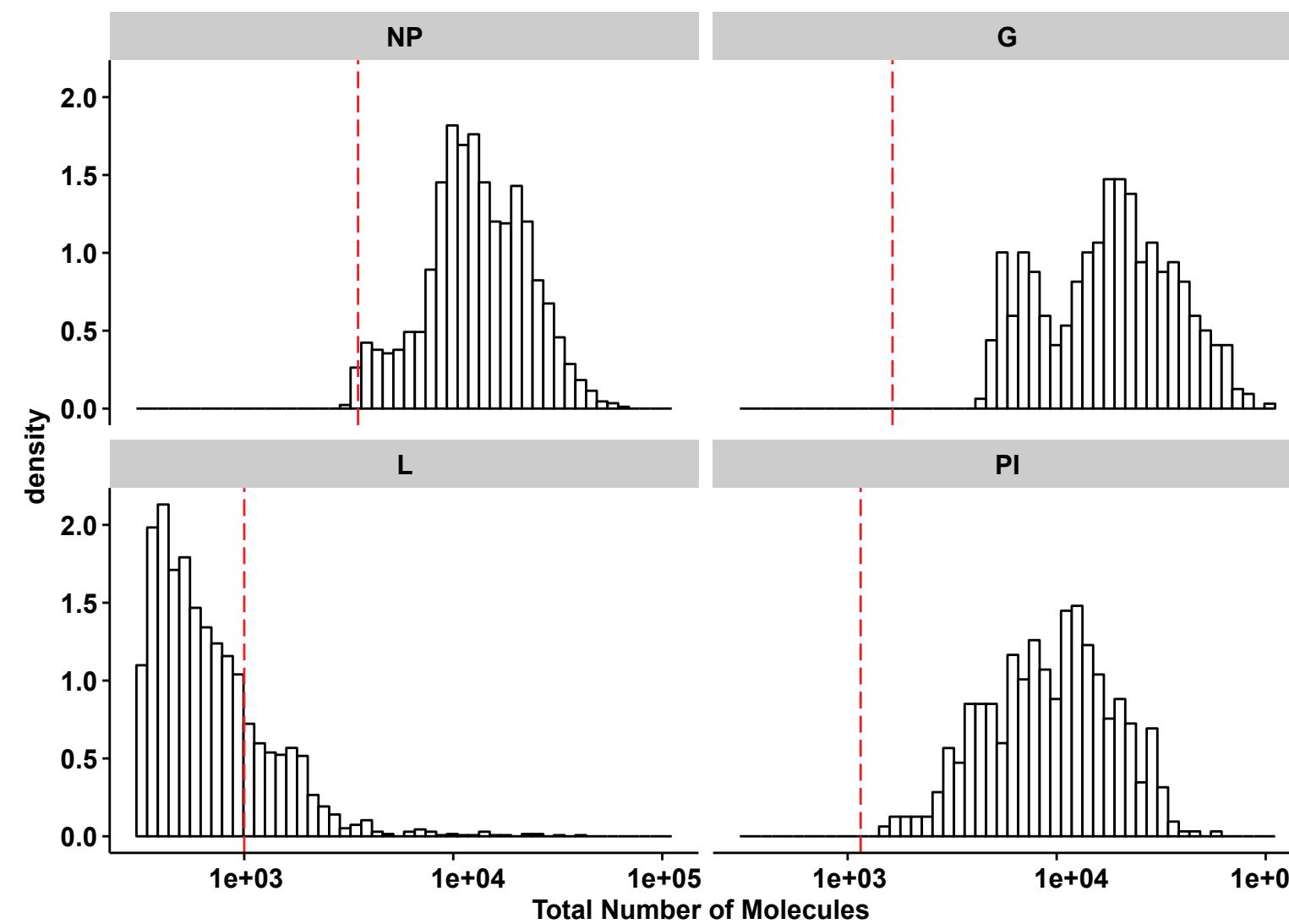
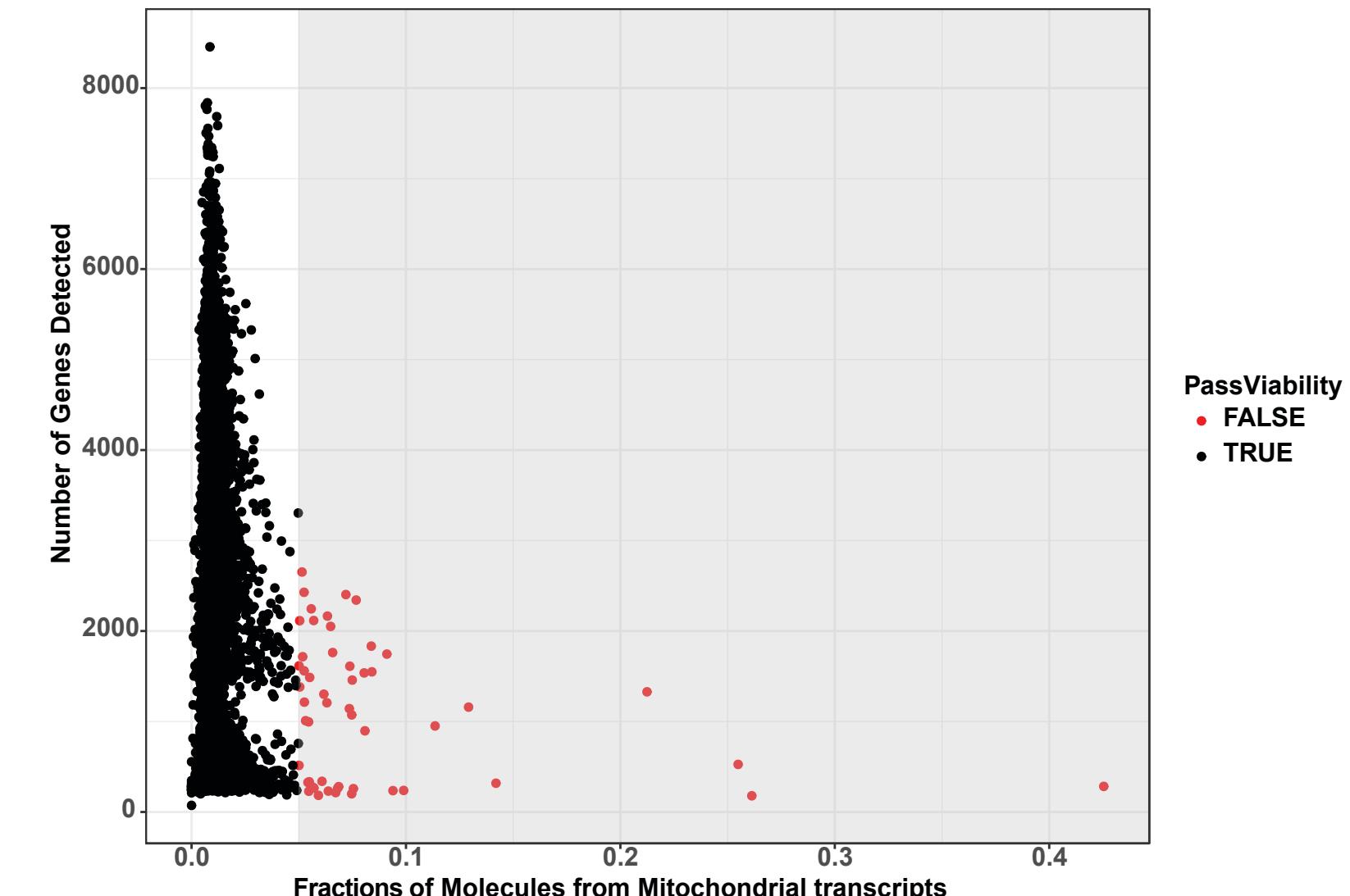
Figure 4

Bach et al.



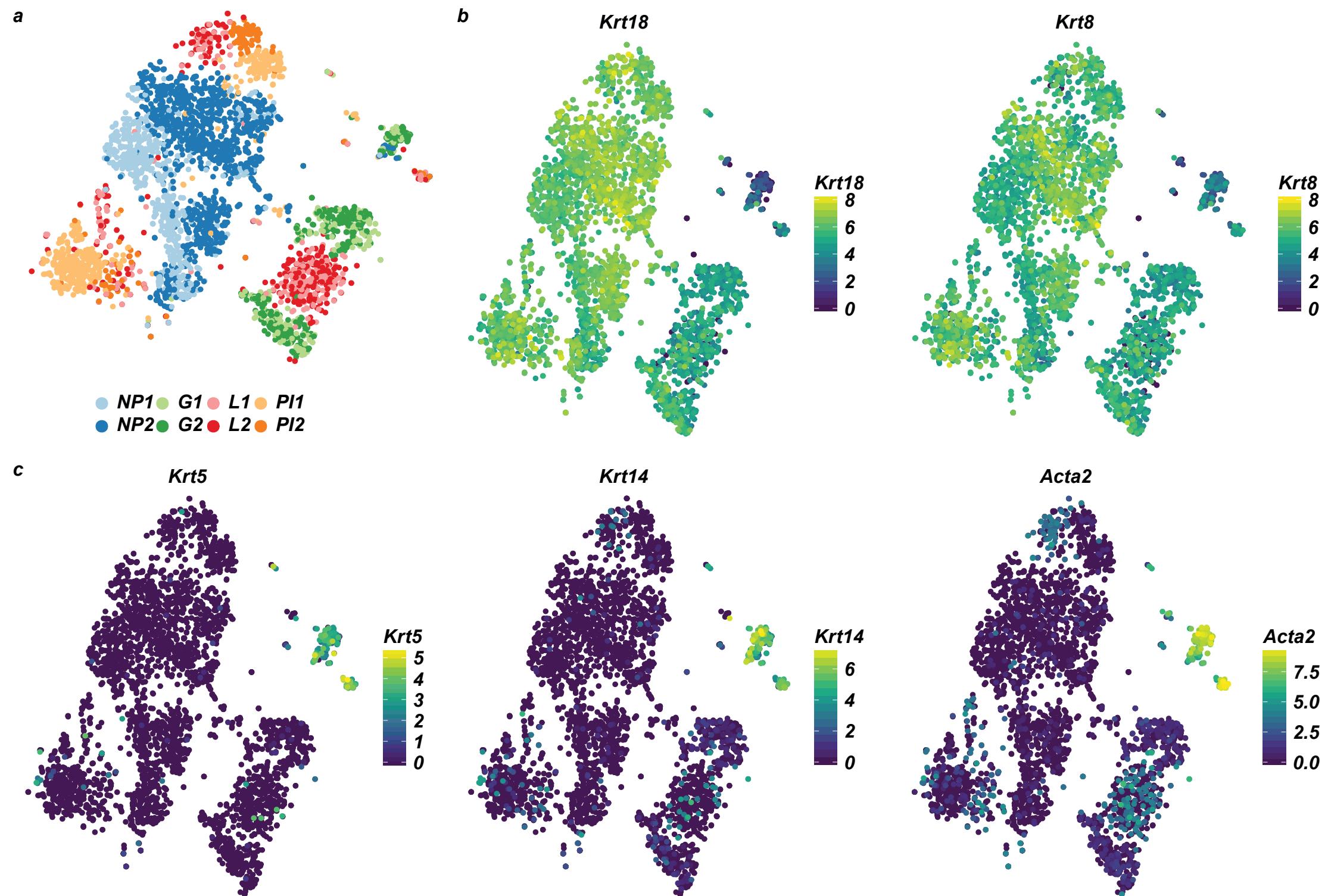
**a**

Sample	Number of cells	Unique molecules	Genes Detected	Number of reads	Sequencing Saturation
NP1	596	12778.5	3326.0	257639	86.3
NP2	1111	12640.0	3358.0	124686	75.3
G1	248	20043.5	4106.5	574902	91.4
G2	375	18068.0	3754.0	372905	88.3
L1	1312	620.5	414.0	106789	97.7
L2	1336	571.5	389.5	111926	98.0
PI1	469	9329.0	2773.0	264971	90.2
PI2	151	10473.0	3061.0	914353	97.2

**b****c****d**

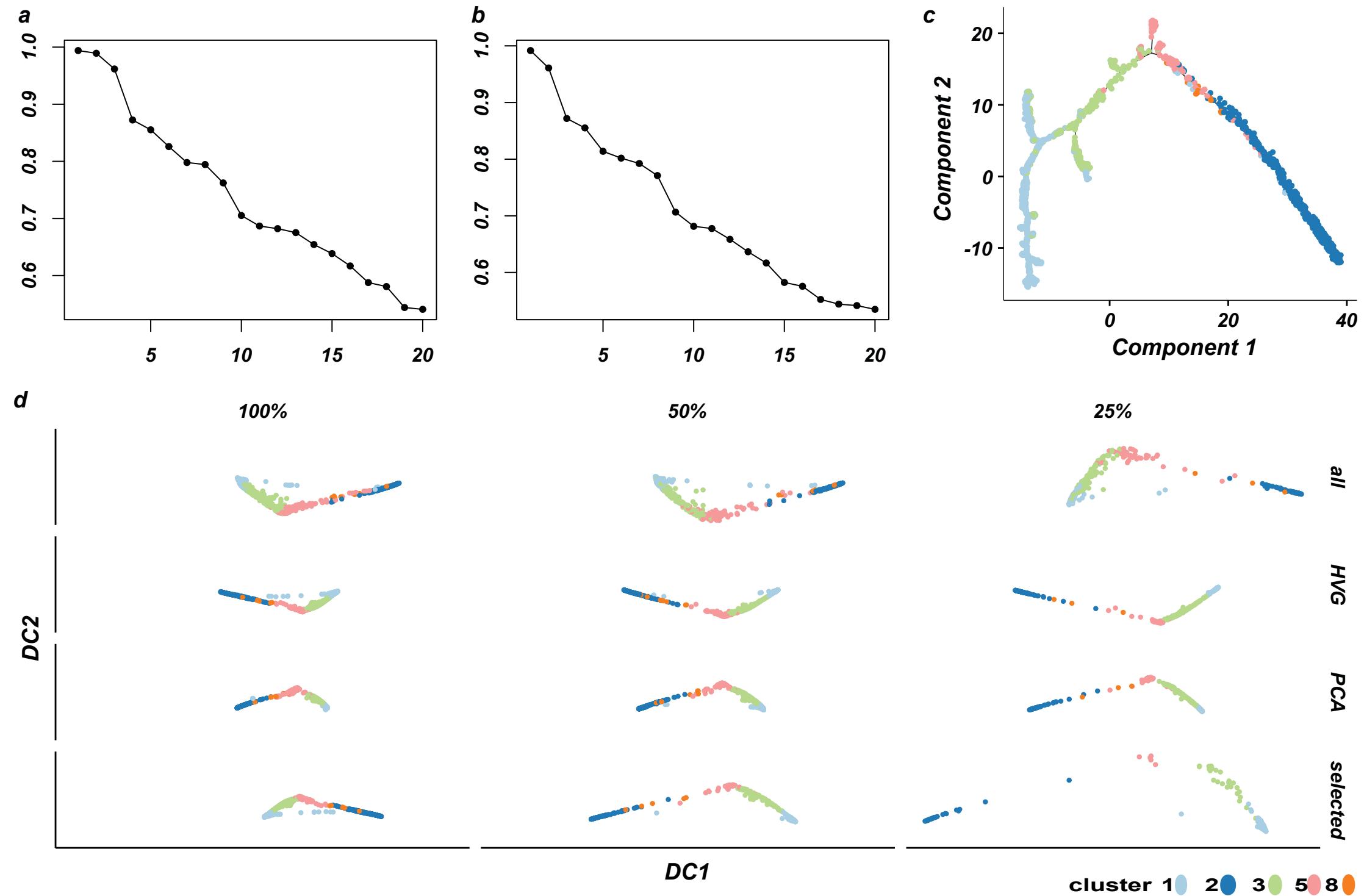
Supplementary Figure 2

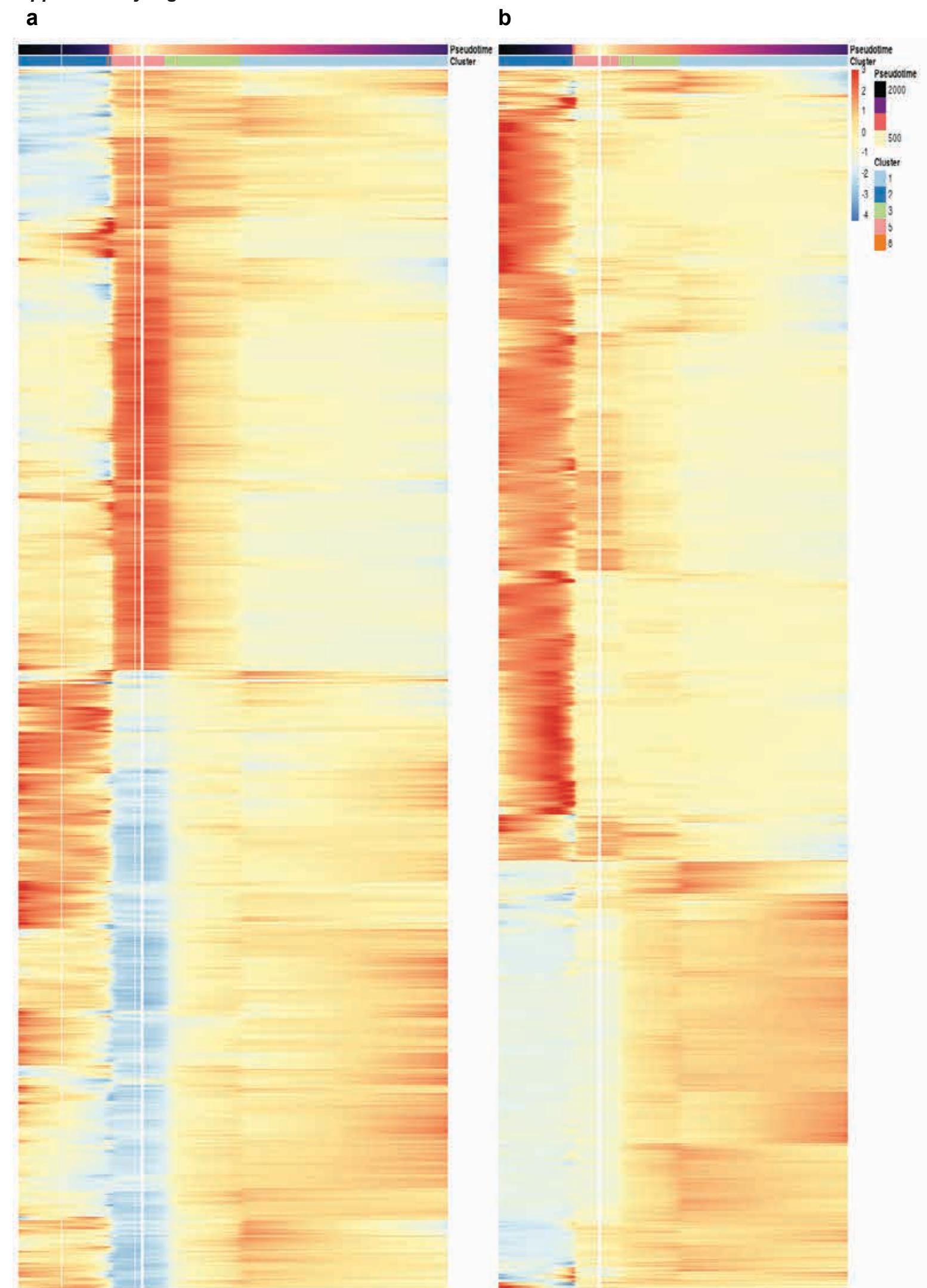
Bach et al.

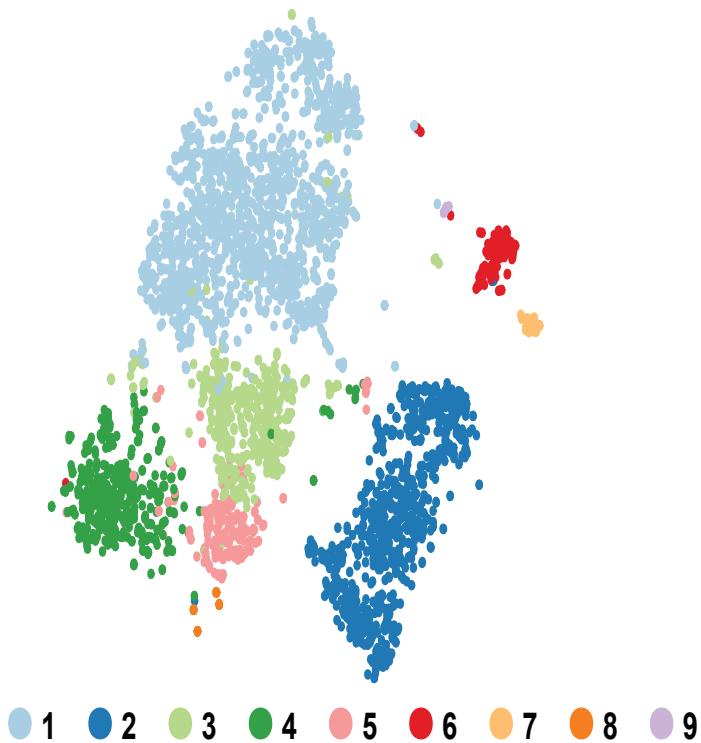
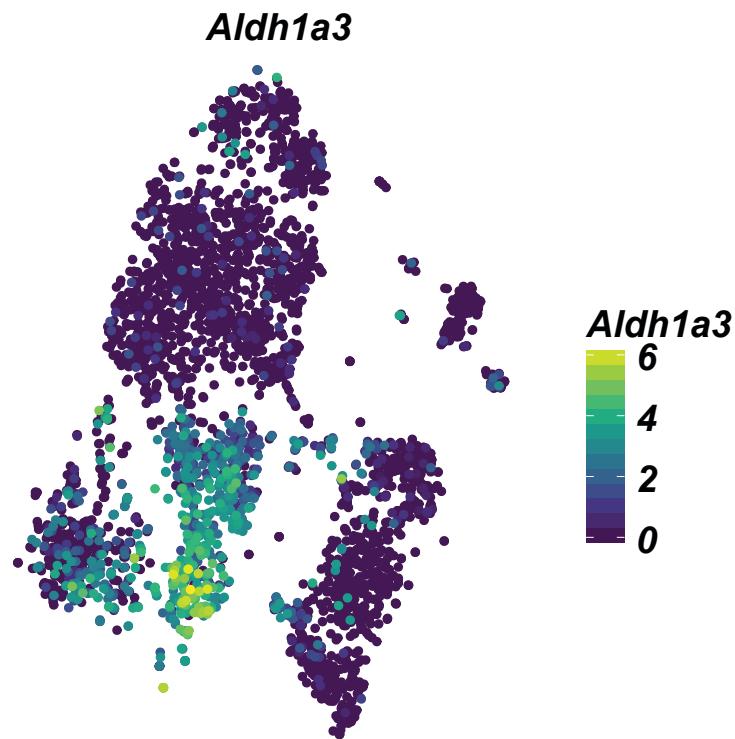
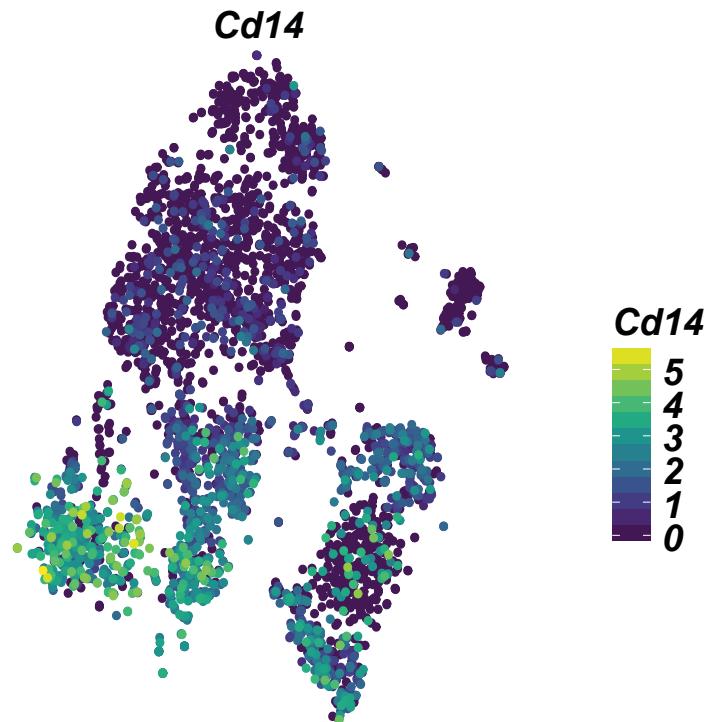
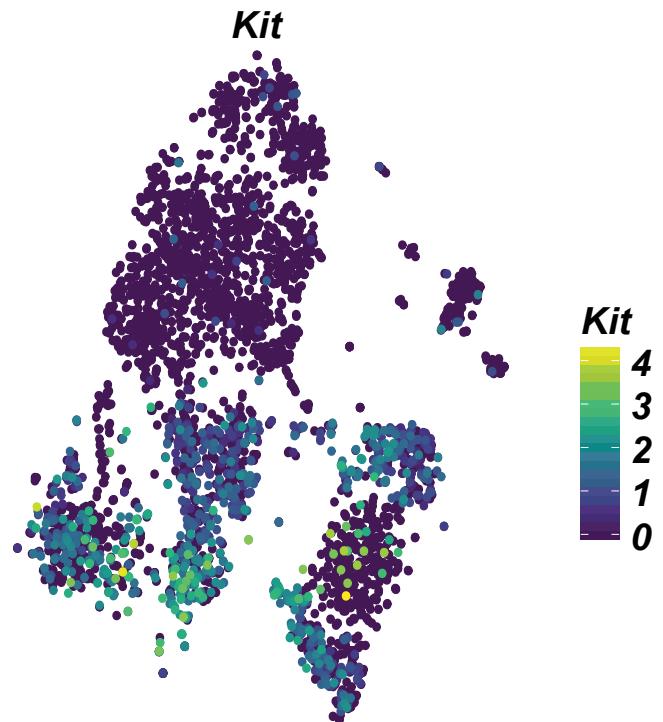


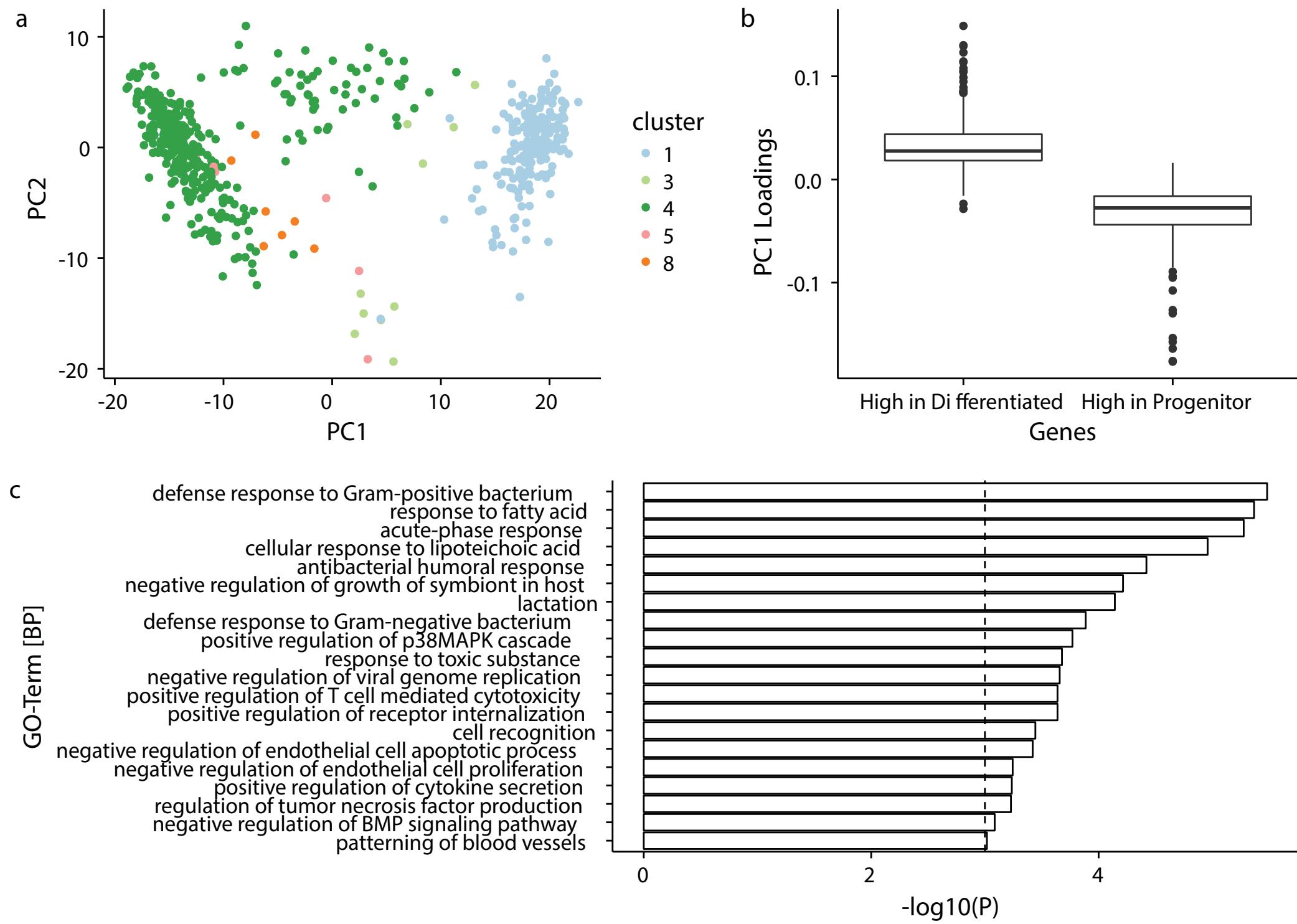
Supplementary Figure 3

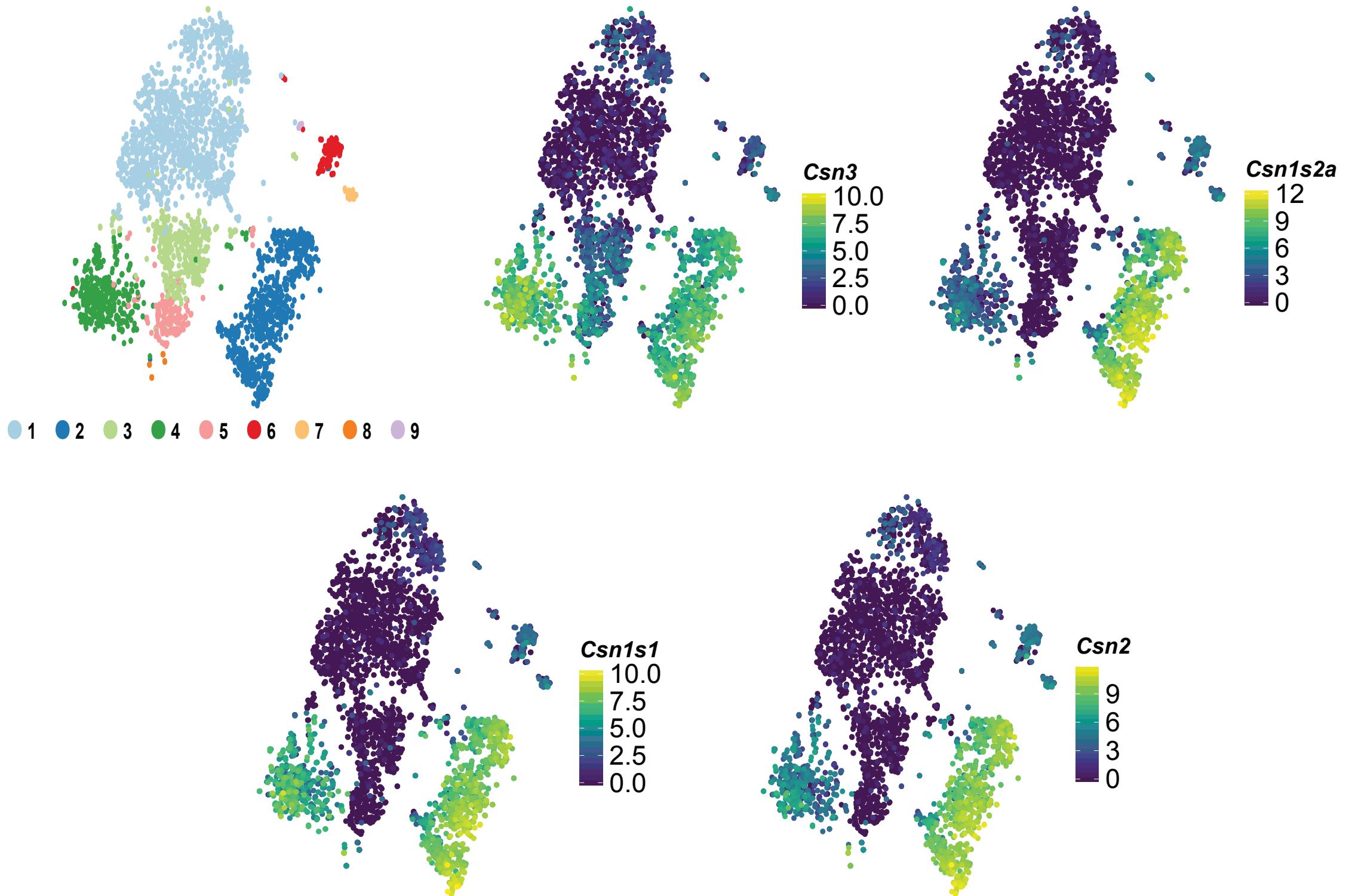
Bach et al.

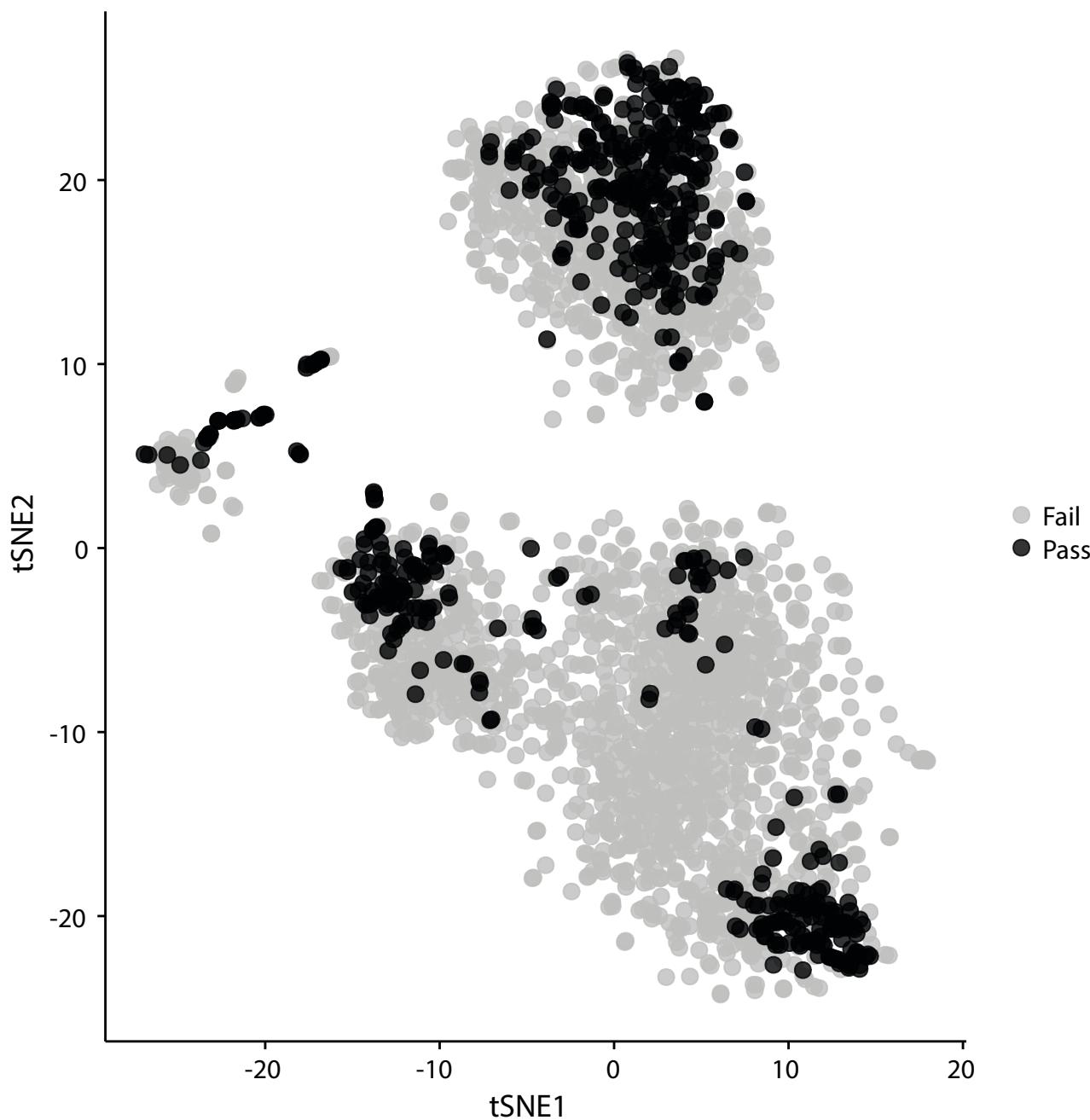




*t-SNE coloured by cluster**Aldh1a3**Cd14**Kit*

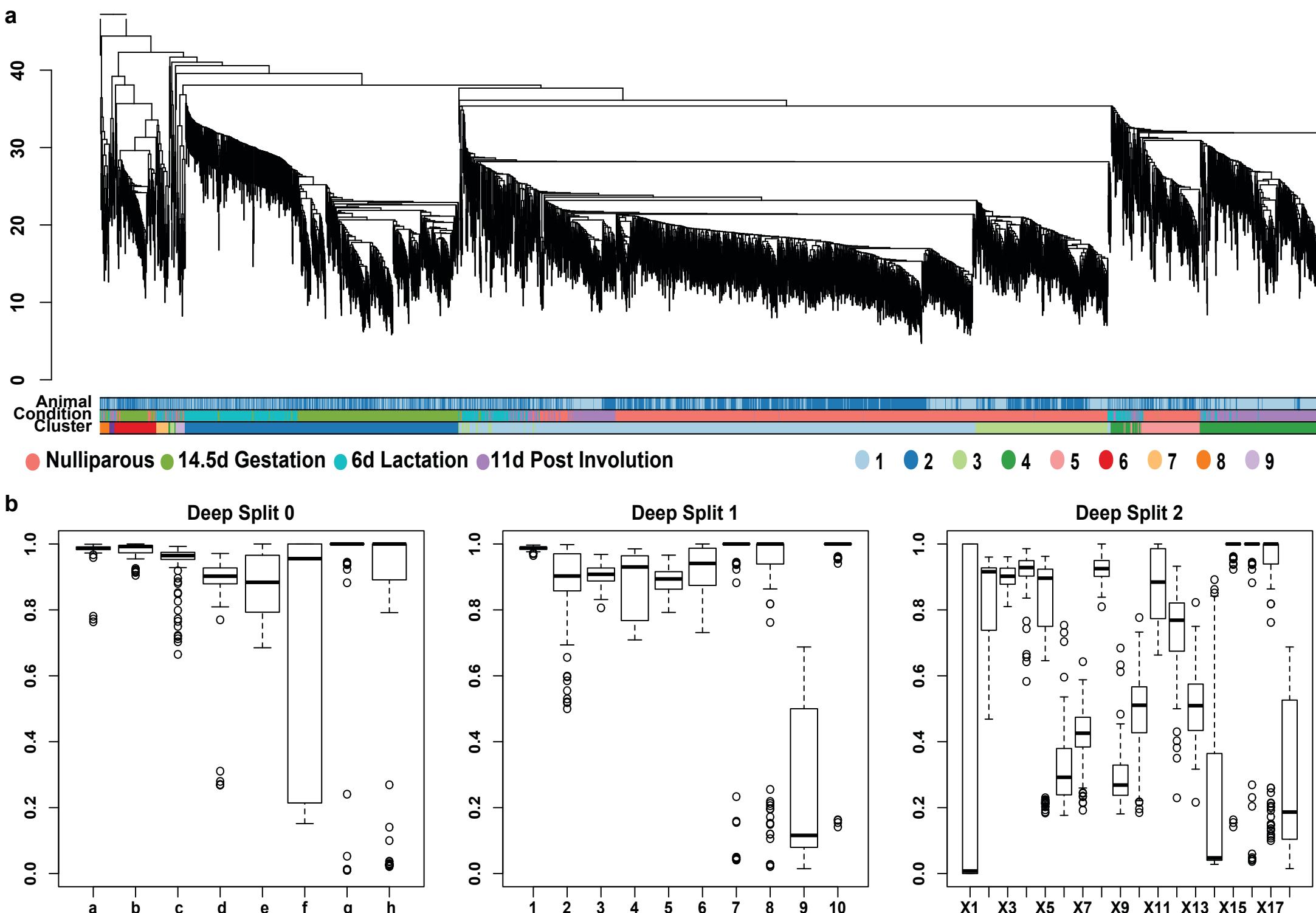


*t*-SNE coloured by cluster



Supplementary Figure 9

Bach et al.



Supplementary Figure 10

Bach et al.

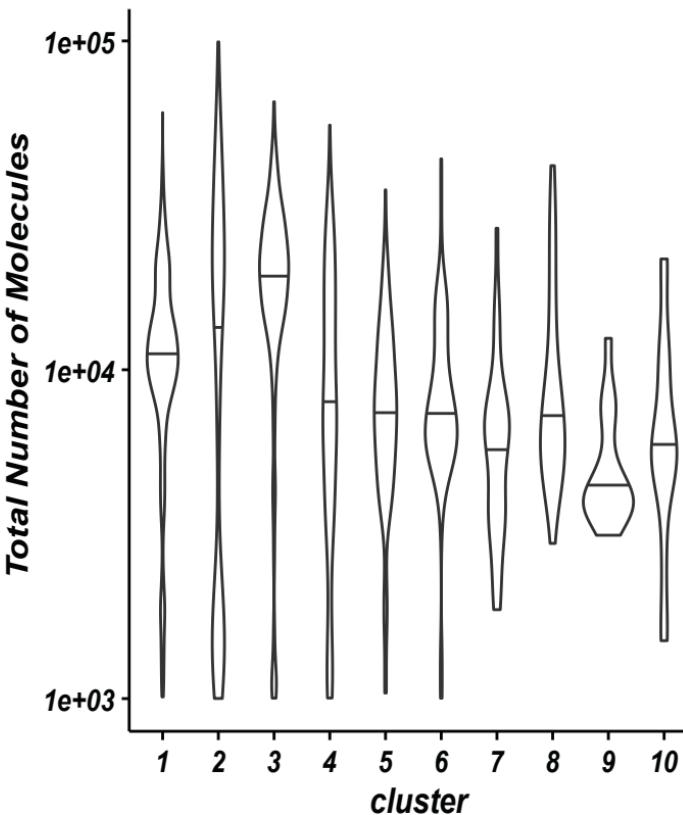
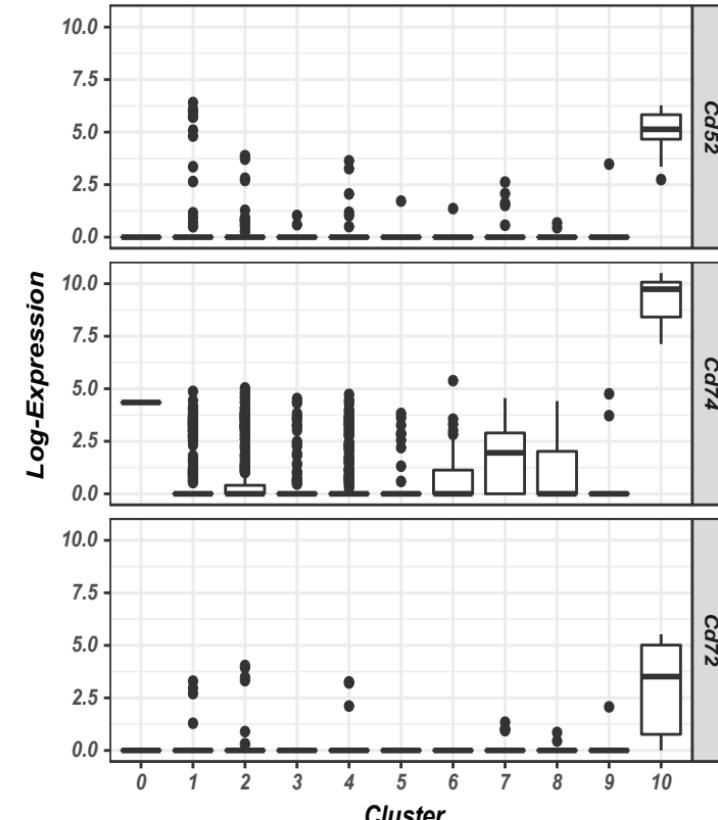
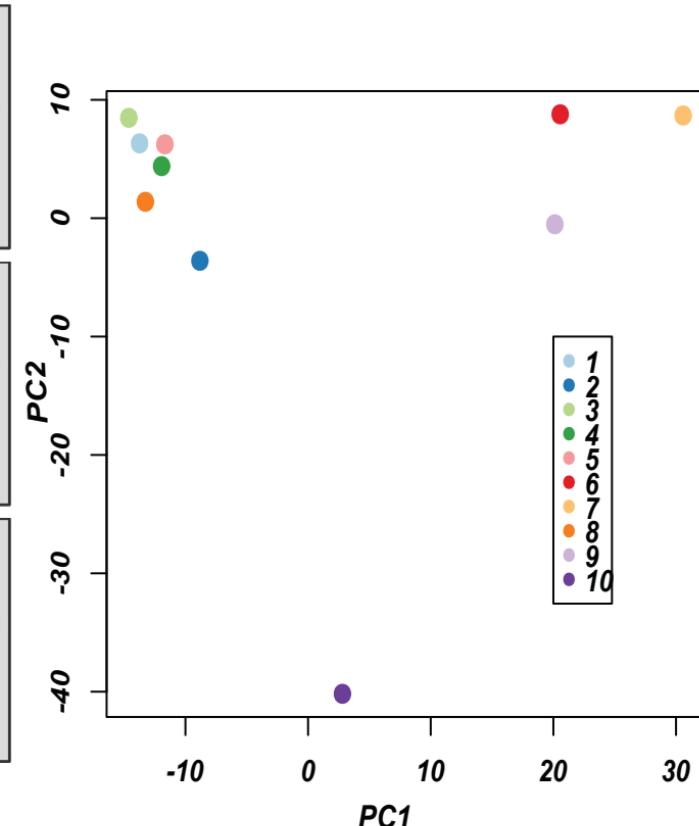
**a****b****c**

Table 1

Bach et al.

Cluster	Key genes expressed	Number of cells captured	Putative identity
1	<i>Esr1, Prlr, Pgr, S100a6, Cited1</i>	1103(NP), 12(G), 126(L), 211(PI)	Differentiated hormone sensing luminal cells
2	<i>Wap, Csn2, Glycam1</i>	480(G), 305(L)	Secretory luminal cells
3	<i>Esr1, Prlr, Pgr, Aldh1a3, Cd14, Kit, Ly6a (Sca-1)</i>	388(NP), 36(L), 12(PI)	Hormone sensing luminal progenitor cells
4	<i>Aldh1a3, Cd14, Kit</i>	83(L), 339(PI)	Post-parity luminal progenitor cell
5	<i>Aldh1a3, Cd14, Kit</i>	161(NP), 2(G), 13(L), 6(PI)	Luminal progenitor cells
6	<i>Krt4, Krt14, Pdpn, Etv5, Acta2</i>	18(NP), 94(G), 4(L), 4(PI)	Basal progenitor cells
7	<i>Oxtr, Acta2, Krt4, Krt14</i>	3(NP), 2(G), 20(L), 10(PI)	Myoepithelial cells
8	<i>Elf5, Aldh1a3, Kit</i>	6(NP), 5(G), 8 (L), 7 (PI)	Luminal progenitor cell with bias towards the alveolar fate
9	<i>Procr, Igfbp4, Gng11, and Zeb2</i>	3(G), 12(PI)	Procr+ basal cells