

Mammary cell gene expression atlas links epithelial cell remodeling events to breast carcinogenesis

Kohei Saeki<sup>1</sup>, Gregory Chang<sup>1</sup>, Noriko Kanaya<sup>1</sup>, Xiwei Wu<sup>2</sup>, Jinhui Wang<sup>2</sup>, Lauren Bernal<sup>1</sup>, Desiree Ha<sup>1</sup>, Susan L Neuhausen<sup>3</sup>, Shiuan Chen<sup>\*1</sup>

<sup>1</sup>Department of Cancer Biology, Beckman Research Institute of City of Hope, Duarte, CA, USA,

<sup>2</sup>Integrative Genomics Core, Beckman Research Institute of City of Hope, Duarte, CA, USA,

<sup>3</sup>Department of Population Sciences, Beckman Research Institute of City of Hope, Duarte, CA, USA

Kohei Saeki *et al.*

**Correspondence:** Shiuan Chen, Ph.D.

Department of Cancer Biology

Beckman Research Institute of the City of Hope

1500 East Duarte Road, Duarte, CA, 91010

Tel: (626) 301-4673, Fax: (626) 301-8972

E-mail: [schen@coh.org](mailto:schen@coh.org)

**Supplementary Table 1 Software and packages.**

	Category	Version	Description	Availability
<i>Anaconda</i> <sup>1</sup>	Distribution	4.7.10	Package management	<a href="https://www.anaconda.com/">https://www.anaconda.com/</a>
<i>Bioconductor</i> <sup>2</sup>	Distribution	3.9	Package management	<a href="https://www.bioconductor.org/">https://www.bioconductor.org/</a>
<i>biomaRt</i> <sup>3</sup>	Package	2.40.4	Data retrieval	<a href="https://bioconductor.org/packages/biomaRt/">https://bioconductor.org/packages/biomaRt/</a>
<i>clusterProfile</i> <sup>4</sup>	Package	3.16.1	Enrichment analysis	<a href="https://bioconductor.org/packages/release/bioc/html/clusterProfiler.html">https://bioconductor.org/packages/release/bioc/html/clusterProfiler.html</a>
<i>cowplot</i> <sup>5</sup>	Package	1.0.0	Visualization	<a href="https://CRAN.R-project.org/package=cowplot">https://CRAN.R-project.org/package=cowplot</a>
<i>CytoTRACE</i> <sup>6</sup>	Package	0.1.0	Trajectory reconstruction	<a href="https://cytotrace.stanford.edu/">https://cytotrace.stanford.edu/</a>
<i>DoubletFinder</i> <sup>7</sup>	Package	2.0.2	Multiplet inference	<a href="https://github.com/chris-mcginnis-ucsf/DoubletFinder">https://github.com/chris-mcginnis-ucsf/DoubletFinder</a>
<i>dplyr</i> <sup>8</sup>	Package	0.8.3	Data manipulation	<a href="https://CRAN.R-project.org/package=dplyr">https://CRAN.R-project.org/package=dplyr</a>
<i>FactoMineR</i> <sup>9</sup>	Package	1.42	PCA	<a href="https://CRAN.R-project.org/package=FactoMineR">https://CRAN.R-project.org/package=FactoMineR</a>

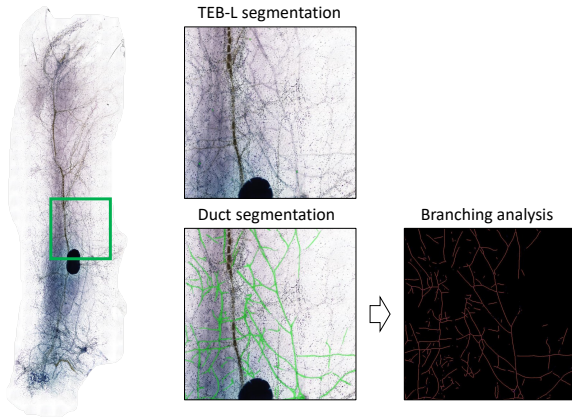
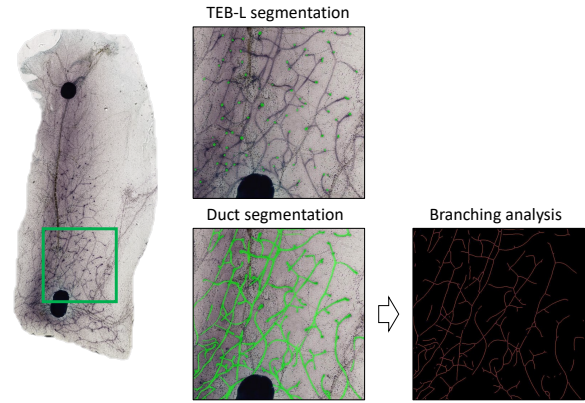
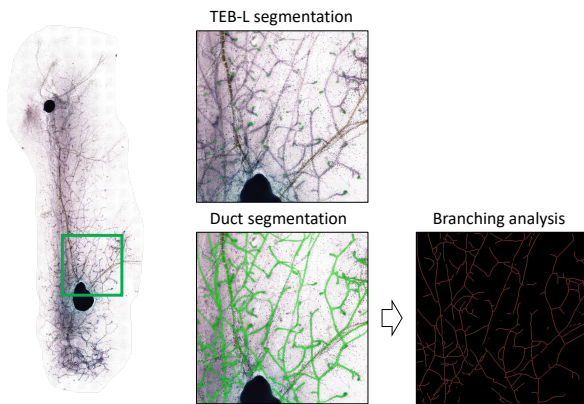
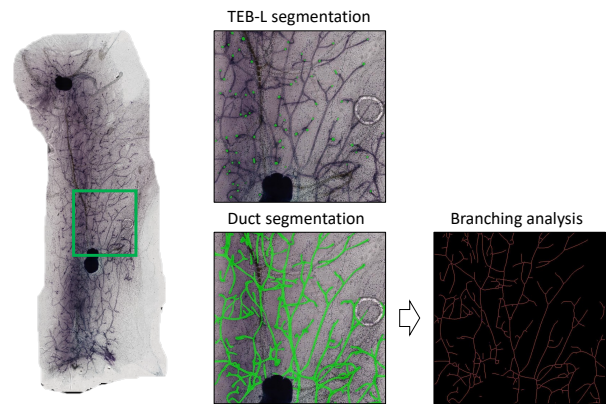
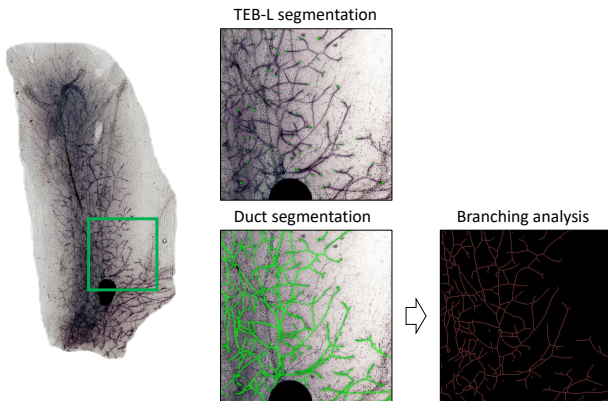
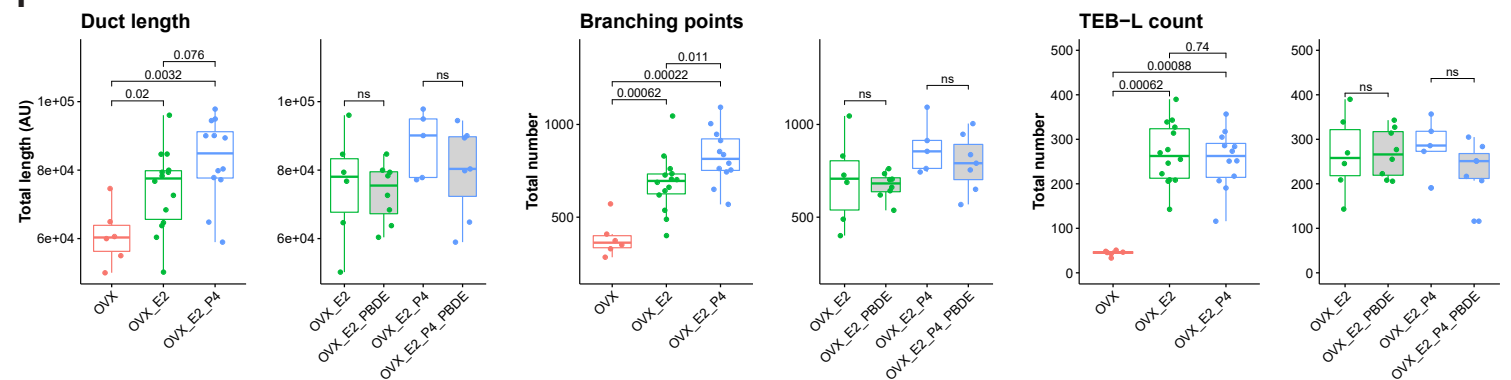
<i>factoextra</i> <sup>10</sup>	Package	1.0.5	Visualization	<a href="https://CRAN.R-project.org/package=factoextra">https://CRAN.R-project.org/package=factoextra</a>
<i>genefu</i> <sup>11</sup>	Package	2.16.0	Molecular subtyping	<a href="https://www.bioconductor.org/packages/release/bioc/html/genefu.html">https://www.bioconductor.org/packages/release/bioc/html/genefu.html</a>
<i>ggcorrplot</i> <sup>12</sup>	Package	0.1.3	Visualization	<a href="https://cran.r-project.org/package=ggcorrplot">https://cran.r-project.org/package=ggcorrplot</a>
<i>ggplot2</i> <sup>13</sup>	Package	3.2.1	Visualization	<a href="https://CRAN.R-project.org/package=ggplot2">https://CRAN.R-project.org/package=ggplot2</a>
<i>ggrepel</i> <sup>14</sup>	Package	0.8.1	Visualization	<a href="https://CRAN.R-project.org/package=ggrepel">https://CRAN.R-project.org/package=ggrepel</a>
<i>ggpubr</i> <sup>15</sup>	Package	0.2.3	Visualization, Statistics	<a href="https://CRAN.R-project.org/package=ggpubr">https://CRAN.R-project.org/package=ggpubr</a>
<i>ggtern</i> <sup>16</sup>	Package	3.1.0	Visualization	<a href="https://CRAN.R-project.org/package=ggtern">https://CRAN.R-project.org/package=ggtern</a>
<i>GSVA</i> <sup>17</sup>	Package	1.32.0	GSVA	<a href="https://bioconductor.org/packages/release/bioc/html/GSVA.html">https://bioconductor.org/packages/release/bioc/html/GSVA.html</a>
<i>JupyterLab</i> <sup>18</sup>	Environment	1.0.2	Environment	<a href="https://jupyter.org/index.html">https://jupyter.org/index.html</a>
<i>harmony</i> <sup>19</sup>	Package	1.0	scRNAseq integration	<a href="https://github.com/immunogenomics/harmony">https://github.com/immunogenomics/harmony</a>
<i>liger</i> <sup>20</sup>	Package	0.5.0	scRNAseq integration	<a href="https://github.com/welch-lab/liger">https://github.com/welch-lab/liger</a>

<i>matplotlib</i> <sup>21</sup>	Package	3.0.2	Visualization	<a href="https://matplotlib.org/index.html">https://matplotlib.org/index.html</a>
<i>msigdb</i> <sup>22</sup>	Package	7.0.1	Data retrieval	<a href="https://CRAN.R-project.org/package=msigdb">https://CRAN.R-project.org/package=msigdb</a>
<i>NumPy</i> <sup>23</sup>	Package	1.17.2	Data manipulation	<a href="https://numpy.org/">https://numpy.org/</a>
<i>pandas</i> <sup>24</sup>	Package	0.25.1	Data manipulation	<a href="https://pandas.pydata.org/pandas-docs/stable/index.html">https://pandas.pydata.org/pandas-docs/stable/index.html</a>
<i>pheatmap</i> <sup>25</sup>	Package	1.0.12	Visualization	<a href="https://CRAN.R-project.org/package=pheatmap">https://CRAN.R-project.org/package=pheatmap</a>
<i>plotly</i> <sup>26</sup>	Package	4.9.0	Visualization	<a href="https://CRAN.R-project.org/package=plotly">https://CRAN.R-project.org/package=plotly</a>
<i>Python</i> <sup>27</sup>	Language	3.7.3	Language	<a href="https://www.python.org/">https://www.python.org/</a>
<i>R</i> <sup>28</sup>	Language	3.6.1	Language	<a href="https://www.r-project.org/">https://www.r-project.org/</a>
<i>reshape2</i> <sup>29</sup>	Package	1.4.3	Data manipulation	<a href="https://CRAN.R-project.org/package=reshape2">https://CRAN.R-project.org/package=reshape2</a>
<i>RStudio</i> <sup>30</sup>	Environment	1.2.5001	Environment	<a href="https://rstudio.com/">https://rstudio.com/</a>
<i>scAlign</i> <sup>31</sup>	Package	1.3.0	scRNAseq integration	<a href="https://github.com/quon-titative-biology/scAlign">https://github.com/quon-titative-biology/scAlign</a>
<i>Seurat</i> <sup>32</sup>	Package	3.1.0	scRNAseq processing	<a href="https://CRAN.R-project.org/package=Seurat">https://CRAN.R-project.org/package=Seurat</a>
<i>stream</i> <sup>33</sup>	Package	0.3.9	Trajectory reconstruction	<a href="https://github.com/pinellolab/STR-EAM">https://github.com/pinellolab/STR-EAM</a>

<i>stringr</i> <sup>34</sup>	Package	1.4.0	Data manipulation	<a href="https://CRAN.R-project.org/package=stringr">https://CRAN.R-project.org/package=stringr</a>
<i>SummarizedExperiment</i> <sup>35</sup>	Package	1.14.1	Data manipulation	<a href="http://bioconductor.org/packages/release/bioc/html/SummarizedExperiment.html">http://bioconductor.org/packages/release/bioc/html/SummarizedExperiment.html</a>
<i>TCGAbiolinks</i> <sup>36</sup>	Package	2.1.2.6	Data retrieval	<a href="https://bioconductor.org/packages/release/bioc/html/TCGAbiolinks.html">https://bioconductor.org/packages/release/bioc/html/TCGAbiolinks.html</a>
<i>TFEA.ChIP</i> <sup>37</sup>	Package	1.4.2	Data retrieval	<a href="https://bioconductor.org/packages/release/bioc/html/TFEA.ChIP.html">https://bioconductor.org/packages/release/bioc/html/TFEA.ChIP.html</a>
<i>umap</i> <sup>38</sup>	Package	0.2.5.0	UMAP	<a href="https://CRAN.R-project.org/package=umap">https://CRAN.R-project.org/package=umap</a>
<i>viridis</i> <sup>39</sup>	Package	0.5.1	Color palette	<a href="https://CRAN.R-project.org/package=viridis">https://CRAN.R-project.org/package=viridis</a>
<i>wesanderson</i> <sup>40</sup>	Package	0.3.6	Color palette	<a href="https://CRAN.R-project.org/package=wesanderson">https://CRAN.R-project.org/package=wesanderson</a>

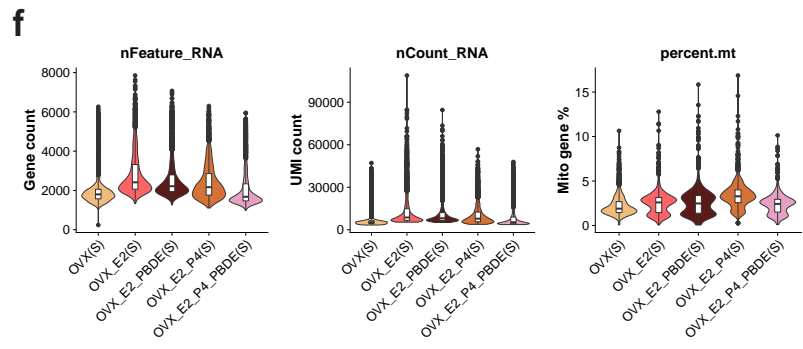
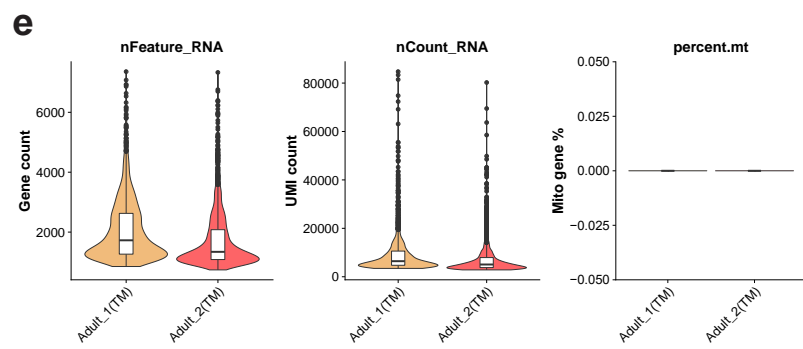
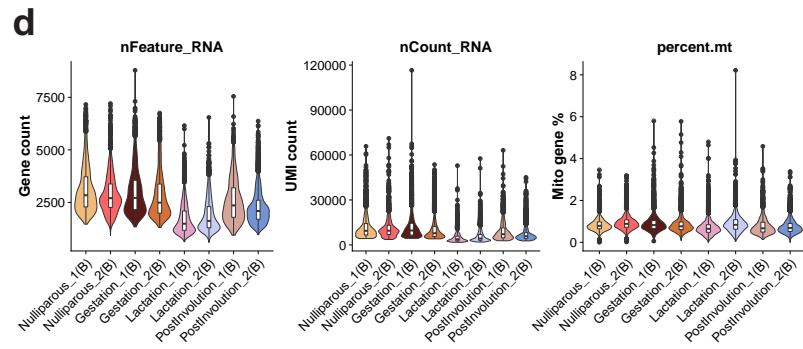
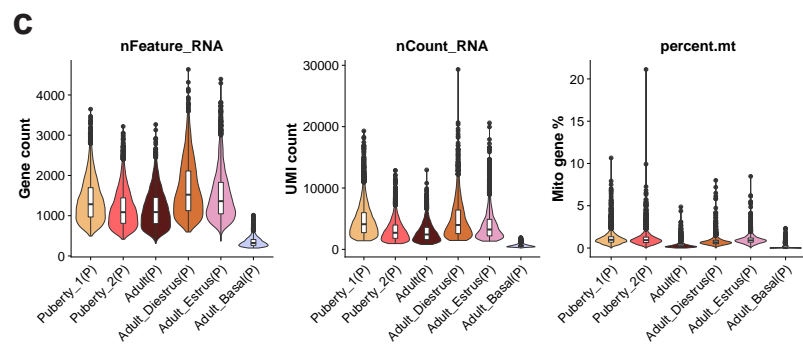
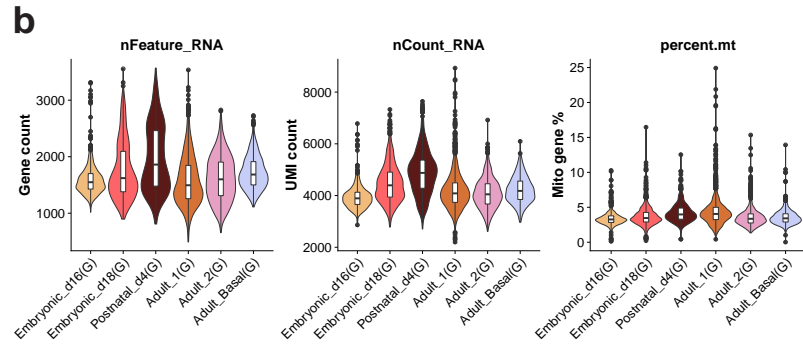
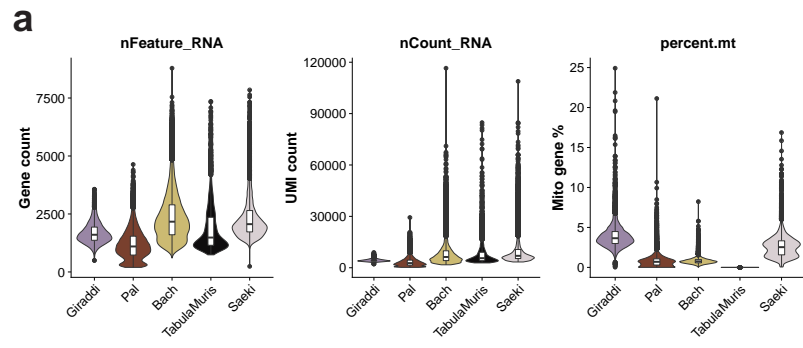
Links were checked on 4-13-2021 for validity.

## Supplementary Figures

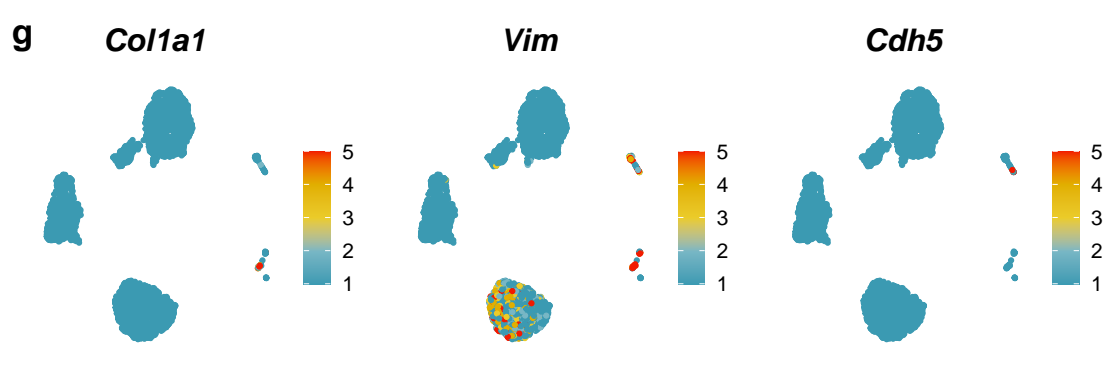
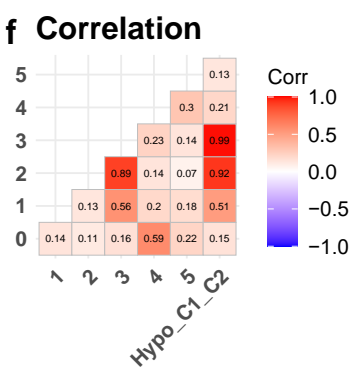
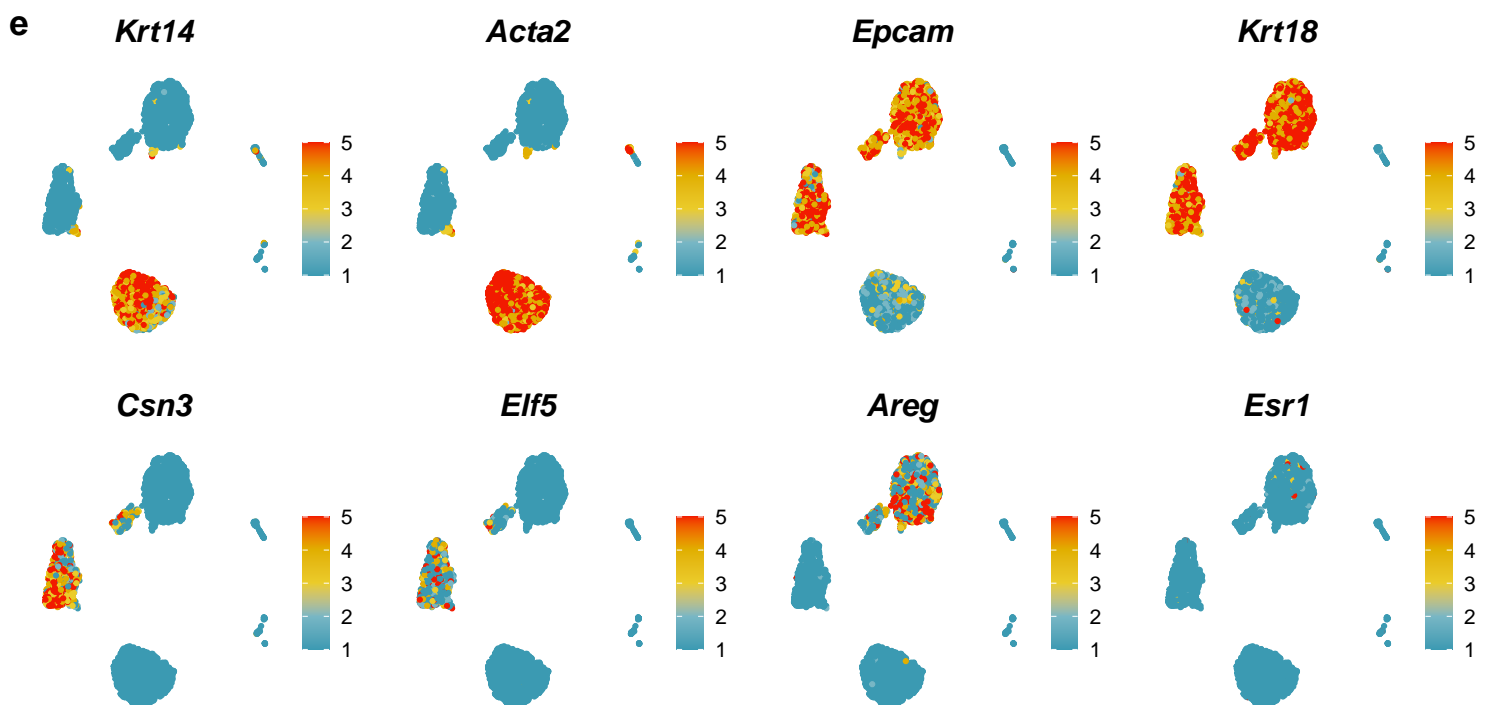
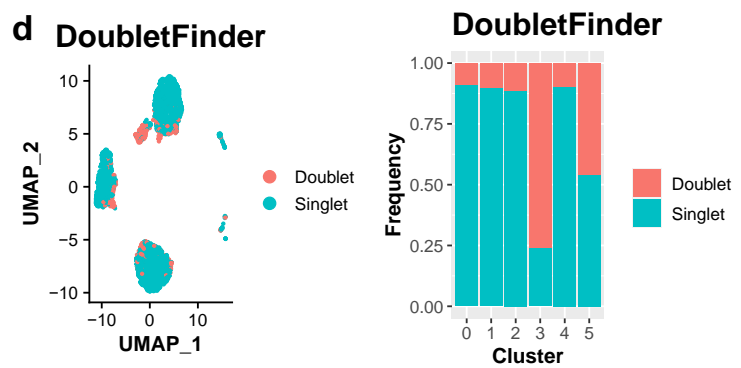
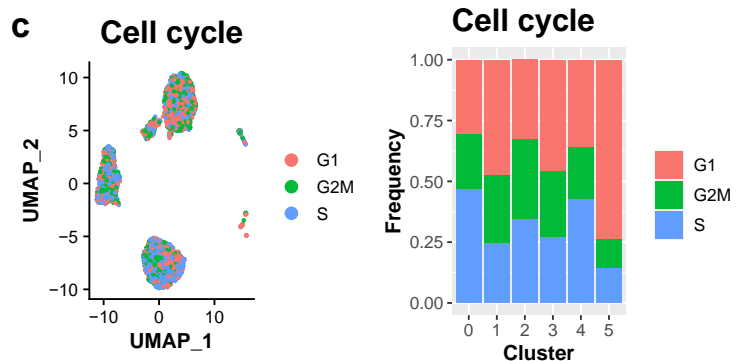
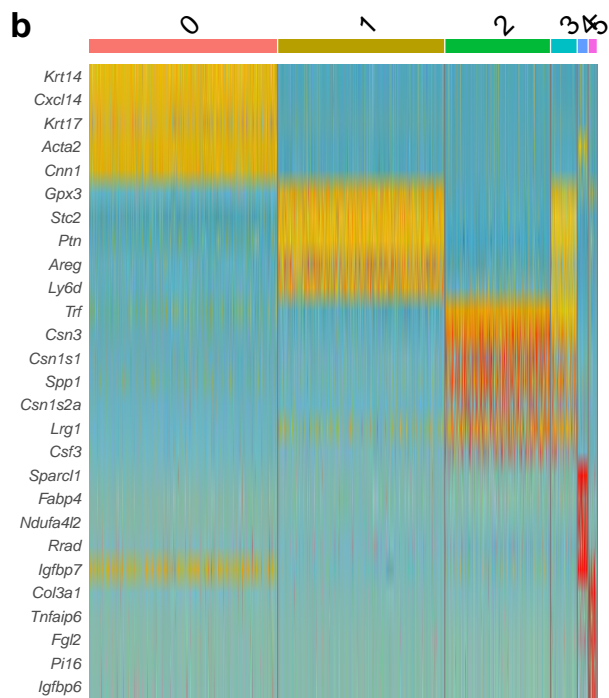
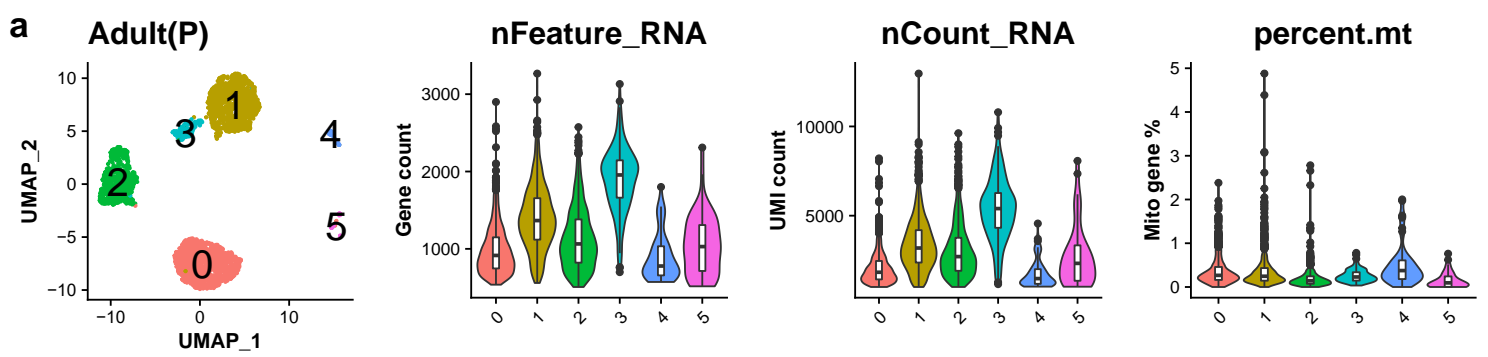
**a****OVX (Animal # 252)****b****OVX\_E2 (Animal # 233)****c****OVX\_E2\_PBDE (Animal # 141)****d****OVX\_E2\_P4 (Animal # 142)****e****OVX\_E2\_P4\_PBDE (Animal # 264)****f**

**Supplementary Figure 1 Reorganization of the surgically menopausal glands after the HRTs and PBDEs exposure.** (a-e) The representative whole gland staining images from OVX (n=6) (a), OVX\_E2 (n=6) (b), OVX\_E2\_PBDE (n=8) (c), OVX\_E2\_P4 (n=5) (d), and OVX\_E2\_P4\_PBDE (n=7) groups (e). In each image, the entire gland, the results of the segmentation of the terminal end bud-like structures and the ductal structure, and the skeletonized image for the branching analysis are presented. (f) The changes by different HRTs and the addition of PBDEs on the total ductal length, total number of the branching points, and TEB-L counts were summarized and compared between the groups. The n number is described above (biologically independent samples). Cliff's delta values for the comparisons of duct length between OVX and OVX\_E2, OVX and OVX\_E2\_P4, and OVX\_E2 and OVX\_E2\_P4 were 0.67 (CI: 0.17-0.89), 0.83 (CI: 0.39-0.96), and 0.42 (0.07-0.74), respectively. Cliff's delta values for the comparisons of branching points between OVX and OVX\_E2, OVX and OVX\_E2\_P4, and OVX\_E2 and OVX\_E2\_P4 were 0.90 (CI: 0.59-0.98), 0.97 (CI: 0.83-0.99), and 0.58 (CI: 0.11-0.84), respectively. Cliff's delta values for the comparisons of TEB-L counts between OVX and OVX\_E2, OVX and OVX\_E2\_P4 were 1.00 (CI: 0.96-1.00) and 1.00 (CI: 0.96-1.00), respectively. ns; not significant. CI; confidence interval. The box-plot elements were defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers.



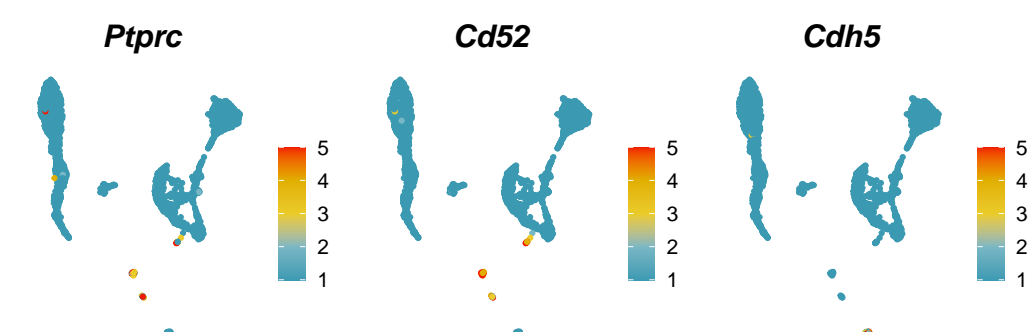
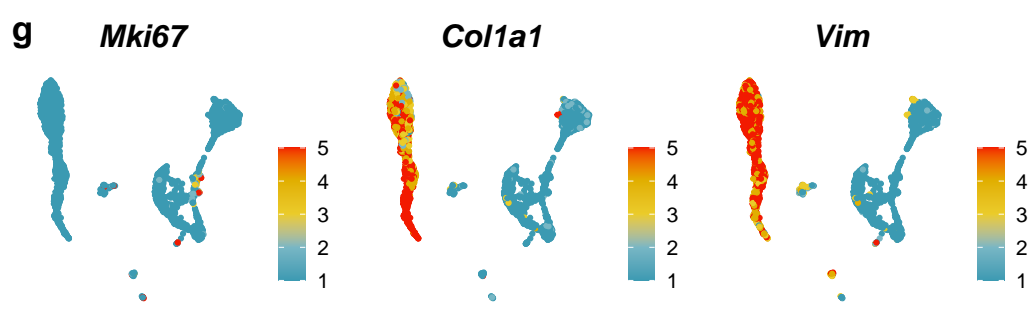
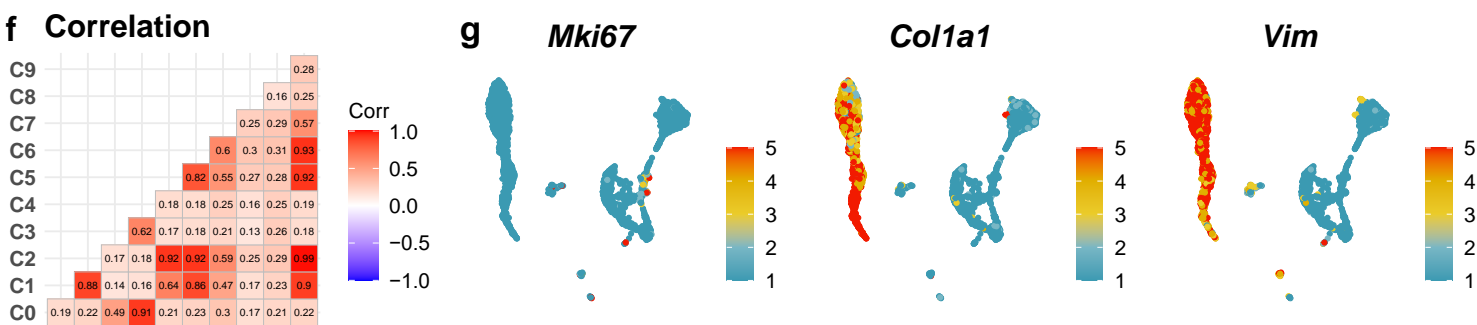
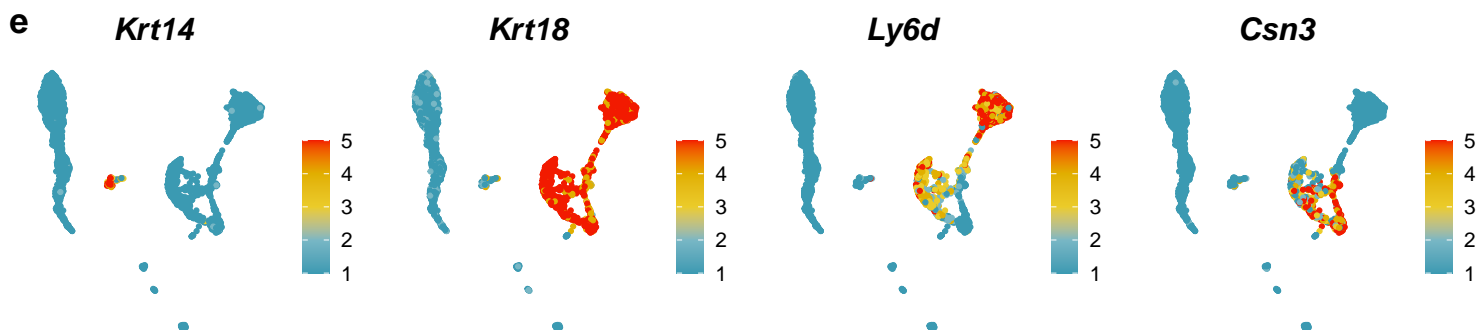
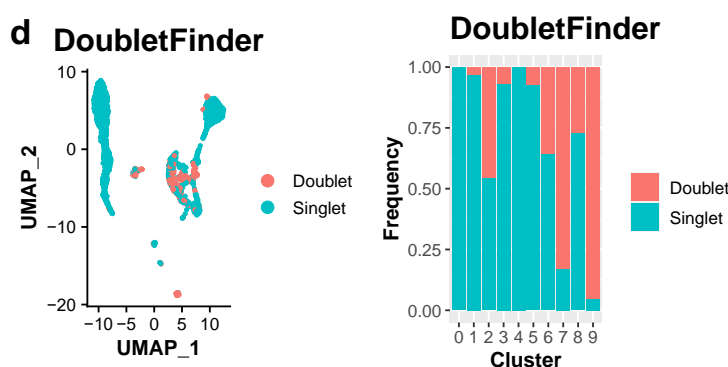
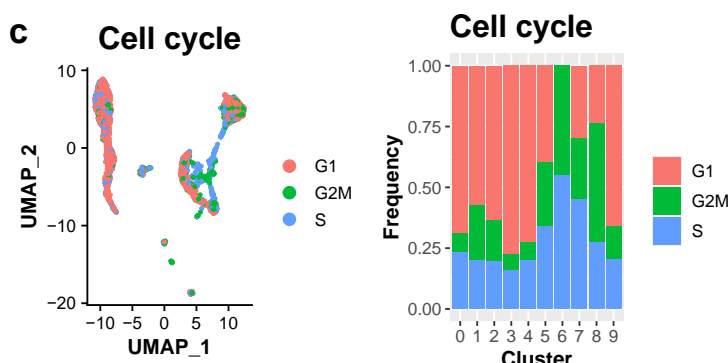
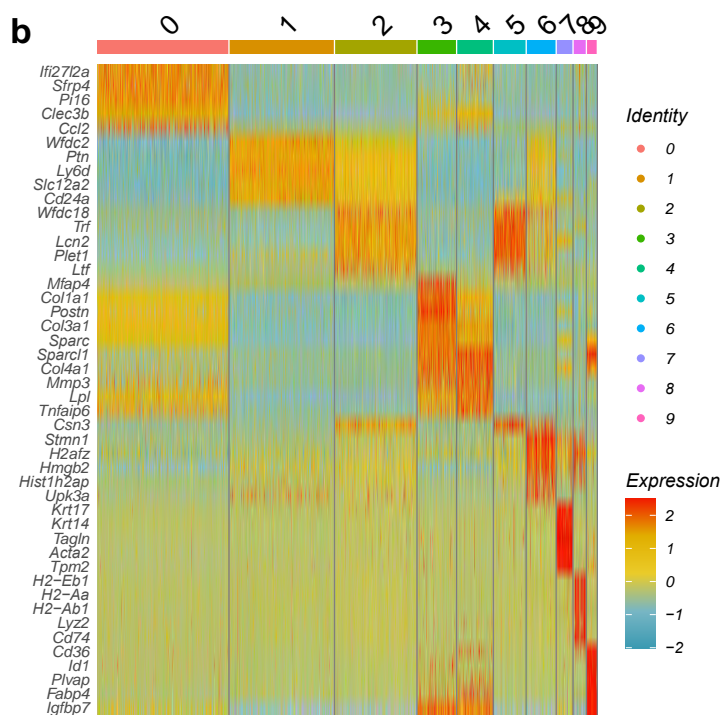
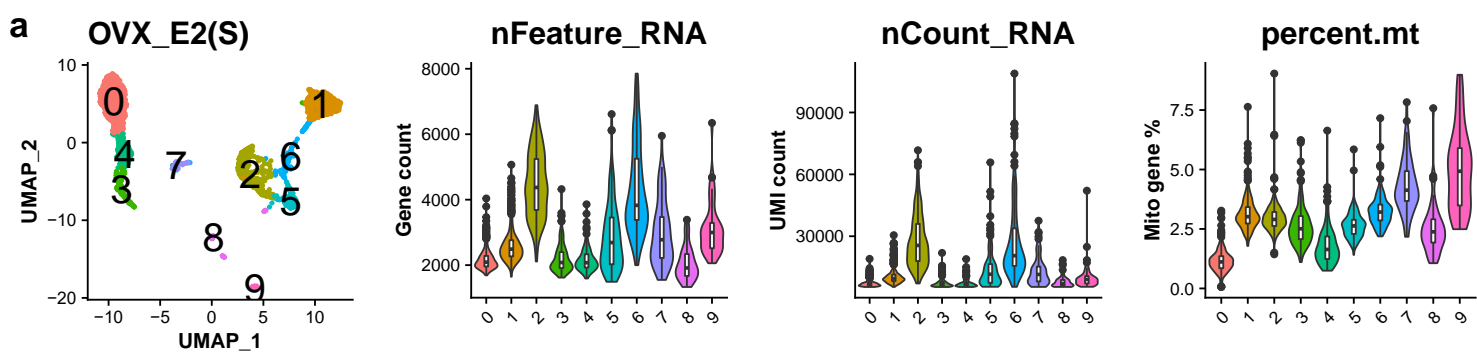


**Supplementary Figure 2 Qualitative characteristics of each dataset and each biological sample. (a-f)** Distribution of the number of detected genes (nFeature\_RNA) and transcripts (nCount\_RNA) and the percentage of the mitochondrial genes (percent.mt) was compared between the five datasets **(a)** and within the Giraddi et al. dataset **(b)**, the Pal et al. dataset **(c)**, the Bach et al. dataset **(d)**, the TabulaMuris dataset **(e)**, and this dataset **(f)**. The n number of each study and sample is summarized in Supplementary Data 2. The box-plot elements were defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. UMI; unique molecular identifier.

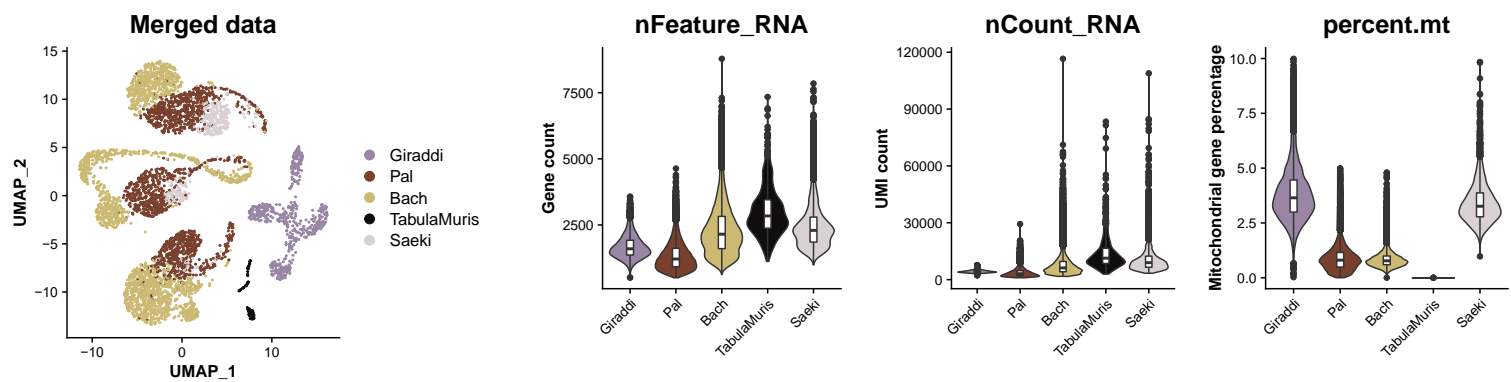
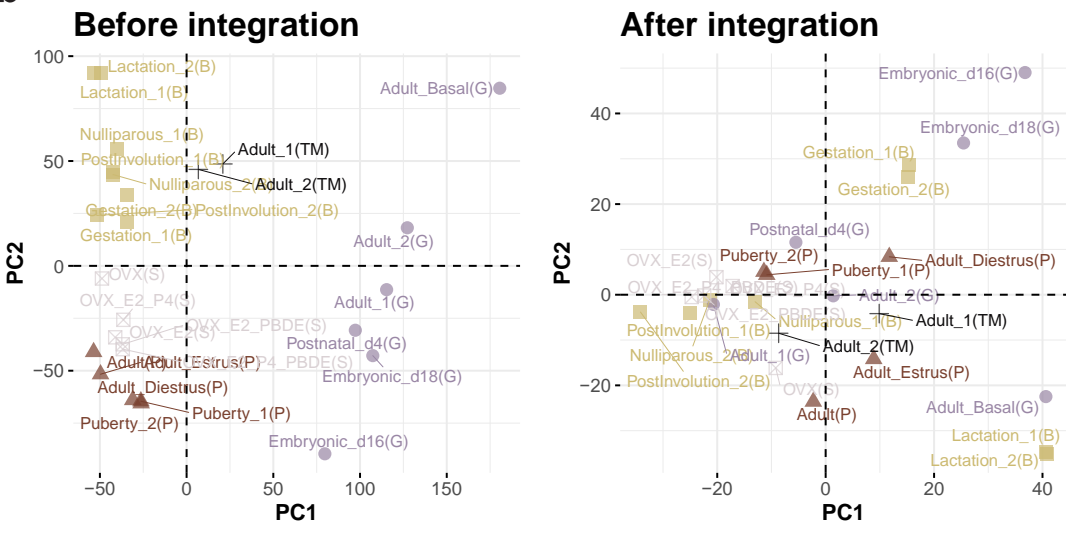
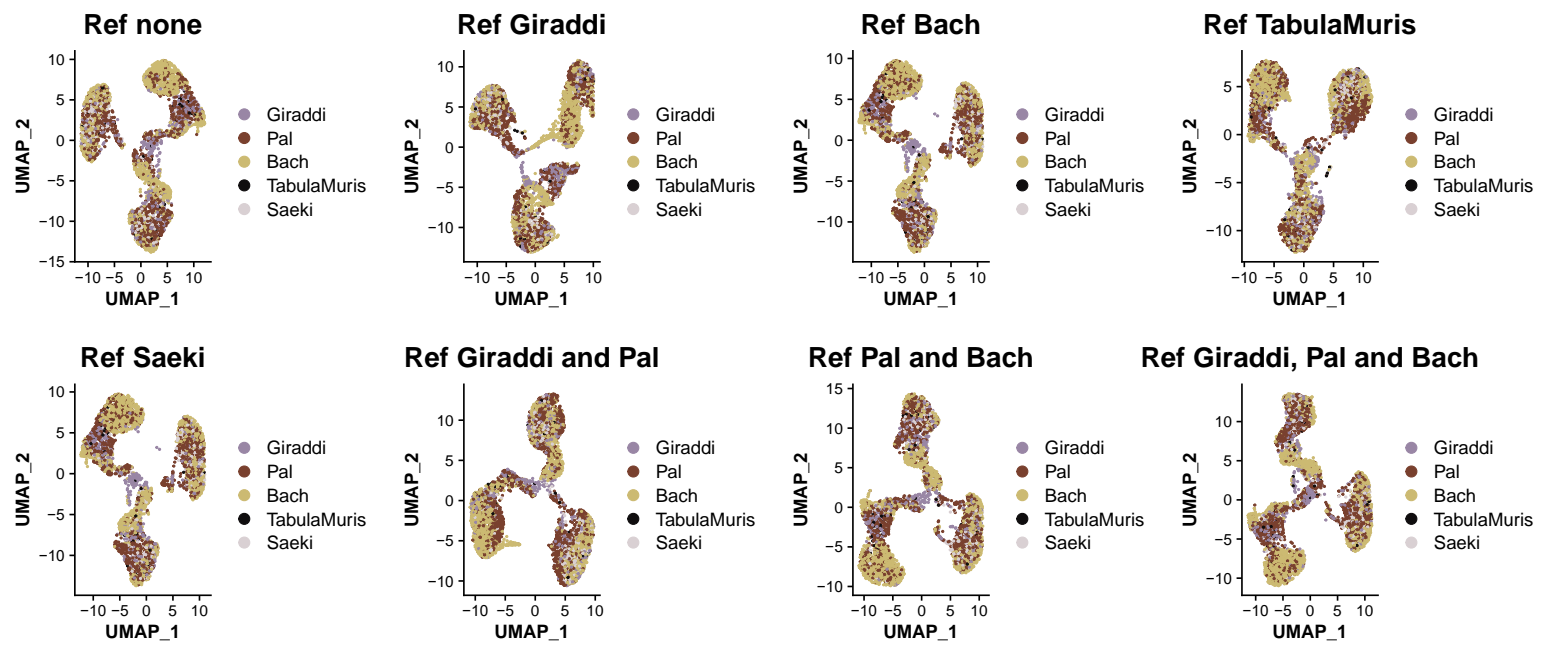


### Supplementary Figure 3 Preprocessing of the scRNAseq data from the adult virgin

**mammary gland in the Pal et al. dataset.** (a) The results of dimension reduction and clustering on a UMAP plot and distribution of nFeature\_RNA, nCount\_RNA, and percent.mt values in each cluster. The n number of each cluster is summarized in Supplementary Data 2. The box-plot elements were defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. (b) The heatmap showing the expression of the top five genes in each cluster. (c) The result of the cell cycle scoring was projected on a UMAP plot and summarized in a stacked bar plot. (d) The results of the *DoubletFinder* analysis was projected on a UMAP plot and summarized in a stacked bar plot. (e) The expression of the marker genes for the mammary epithelium were expressed on UMAP plots (*Krt14* and *Acta2*; Basal cells, *Epcam* and *Krt18*; Luminal cells, *Csn3* and *Elf5*; Luminal alveolar cells, and *Areg* and *Esr1*; Luminal hormone-sensing cells). (f) The results of the correlation analysis of pseudo bulk RNAseq analysis, including a hypothetical doublet cluster between C1 and C2 (Hypo\_C1\_C2). (g) The expression of marker genes for contaminating stromal cells were visualized on UMAP plots (*Coll1a1*; Fibroblasts, *Vim*; Mesenchymal cells, *Cdh5*; Endothelial cells). UMI; unique molecular identifier.



**Supplementary Figure 4 Preprocessing of the scRNAseq data of the mammary gland from the surgically menopausal mouse treated with estrogen for a week in this study.** (a) The results of dimension reduction and clustering on a UMAP plot and distribution of nFeature\_RNA, nCount\_RNA, and percent.mt values in each cluster. The n number of each cluster is summarized in Supplementary Data 2. The box-plot elements were defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. (b) The heatmap showing the expression of the top five genes in each cluster. (c) The results of cell cycle scoring were projected on a UMAP plot and summarized in a stacked bar plot. (d) The results of the *DoubletFinder* analysis was projected on a UMAP plot and summarized in a stacked bar plot. (e) The expression of the marker genes for the mammary epithelium was expressed on UMAP plots (*Krt14*; Basal cells, *Krt18*; Luminal cells, *Csn3*; Luminal alveolar cells, and *Ly6d*; Luminal hormone-sensing cells). (f) The results of the correlation analysis of pseudo bulk RNAseq analysis, including a hypothetical doublet cluster between C1 and C5 (Hypo\_C1\_C5). (g) The expression of marker genes for proliferating cells and contaminating stromal cells were visualized on UMAP plots (*Mki67*; Proliferating cells, *Colla1*; Fibroblasts, *Vim*; Mesenchymal cells, *Ptpnc*; Hematopoietic cells, *Cd52*; Macrophages, *Cdh5*; Endothelial cells). UMI; unique molecular identifier.

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

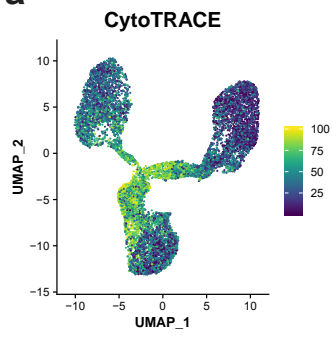
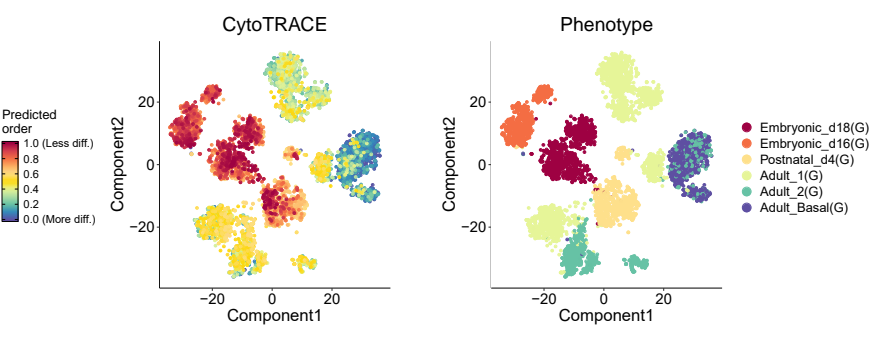
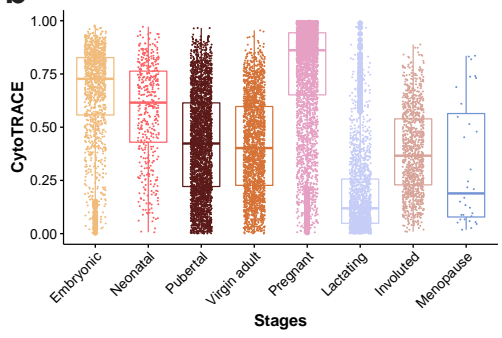
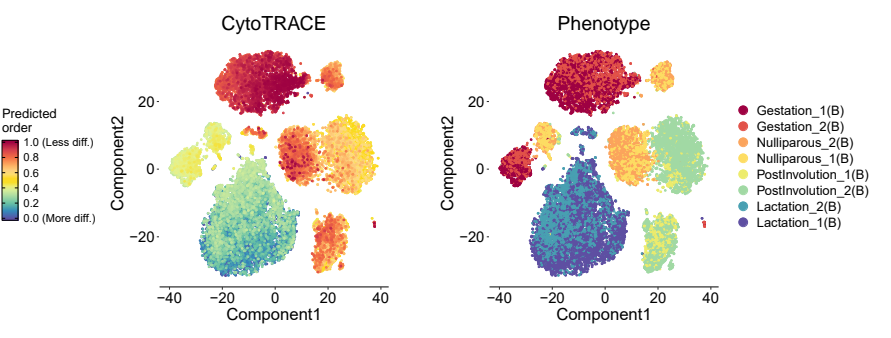
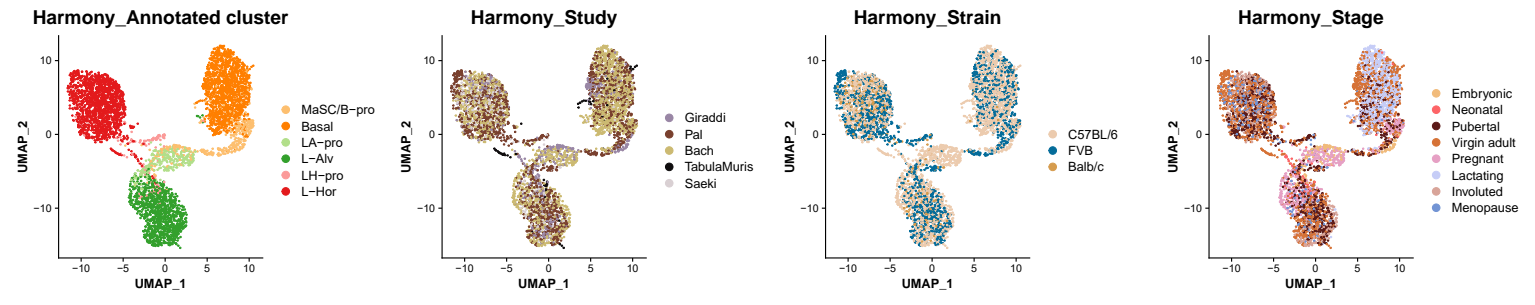
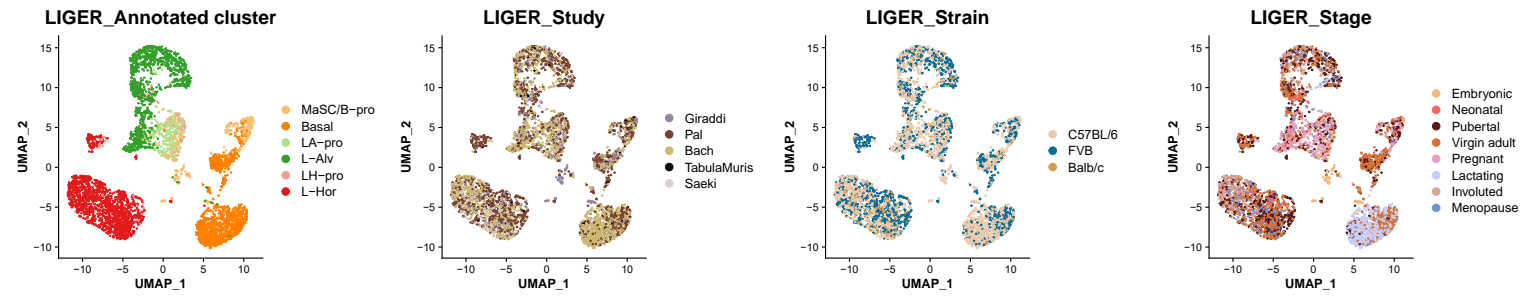
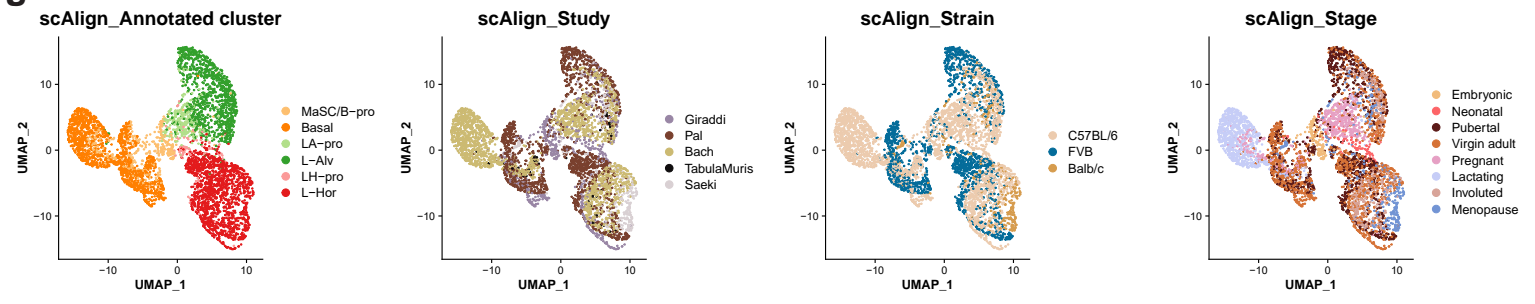
**Supplementary Figure 5 Data integration of the five scRNAseq datasets.** (a) Visualization of the preprocessed and merged datasets on a UMAP plot and distribution of nFeature\_RNA, nCount\_RNA, and percent.mt in each preprocessed dataset. The n number of each study after preprocessing is summarized in Supplementary Data 2. The box-plot elements were defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. Five thousands cells were sampled from the entire data for UMAP visualization. (b) The results of pseudo bulk RNAseq of each sample were color-coded by datasets before and after the anchor-based data integration by *Seurat v3*. (c) Visualization of the data in UMAP dimensionality after the anchor-based data integration using the different dataset(s) as references. Ref none; the datasets were integrated without any reference. Five thousands cells were sampled from the entire data for UMAP visualization. UMI; unique molecular identifier.





**Supplementary Figure 6 Features of the identified clusters in the integrated data. (a)**

Distribution of nFeature\_RNA, nCount\_RNA, and percent.mt values in each cluster. The number of each cluster is as follows: C1 (n=4,217), C2 (n=14,493), C3 (n=3,176), C4 (n=11,856), C5 (n=1,156), and C6 (n=15,509) (biologically independent samples). The box-plot elements were defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. **(b)** Breakdown of each cluster by datasets (Study), mouse strains (Strain), and different developmental stages (Stage). **(c)** The heatmap showing the expression of the top five genes in each cluster. **(d)** The expression of selected genes characteristic to the putative progenitor clusters (C1, C3, and C5; *Birc5*, *Hmgb2*, *Stmn1*), the basal cell lineage (C1 and C2; *Acta2*, *Krt17*, and *Myl9*), the luminal alveolar lineage (C3 and C4; *Lalba*, *Csn2*, and *Spp1*), and the luminal hormone-sensing cells (C5 and C6; *Cited1*, *Ly6d*, and *Prlr*). Five thousands cells were sampled from the entire data for UMAP visualization. UMI; unique molecular identifier.

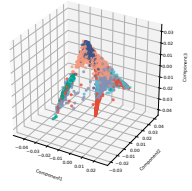
**a****c****b****d****e****f****g**

**Supplementary Figure 7 The results from the CytoTRACE analysis and the different**

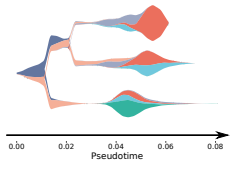
**integration algorithms.** (a) Projection of the results from the Scanorama-based implementation of CytoTRACE on a UMAP plot of the integrated data by *Seurat*. A higher score indicates less differentiated states. (b) The CytoTRACE scores of cells in different stages. The n number of each stage is as follows: Embryonic (n=1,131), Neonatal (n=448), Pubertal (n=3,139), Virgin adult (n=2,139), Pregnant (n=3,375), Lactating (n=1,630), Involved (n=1,068), and Menopause (n=32) (biologically independent samples). The box-plot elements were defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. (c) The results of the CytoTRACE analysis of the Giraddi et al. dataset. The t-Distributed Stochastic Neighbor Embedding (t-SNE) plots of the data were color-coded by the CytoTRACE score and the individual samples. (d) The results of the CytoTRACE analysis of the Bach et al. dataset. The t-SNE plots of the data were color-coded by the CytoTRACE score and the individual samples. (e-g) Data integration by additional algorithms. The UMAP plots after integration were color-coded by their original clustering in *Seurat v3*, datasets (Study), mouse strains (Strain), and different developmental stages (Stage). Five thousands cells were sampled from the entire data for UMAP visualization. (e) The data integrated by *Harmony*. (f) The data integrated by *LIGER*. (g) The data integrated by *scAlign*.

**a**

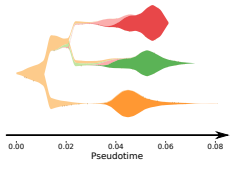
Embryonic\_d18(G) Embryonic\_d18(G) Embryonic\_d18(G) Embryonic\_d18(G)



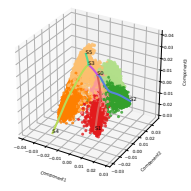
Embryonic\_d16(G) Postnatal\_d4(G) Adult\_2(G)  
Embryonic\_d18(G) Adult\_1(G) Adult\_Basal(G)



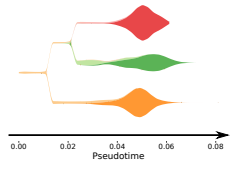
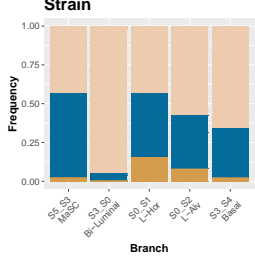
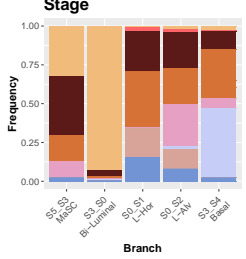
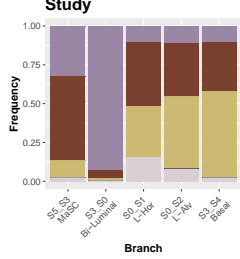
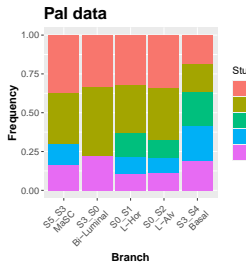
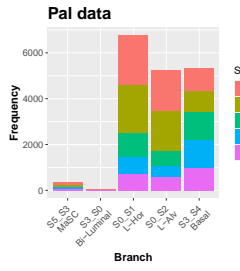
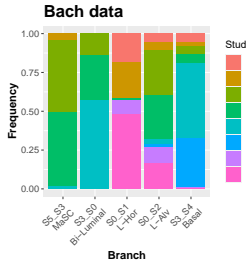
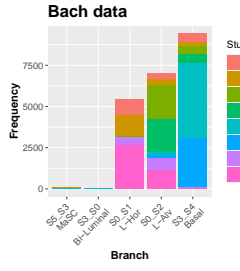
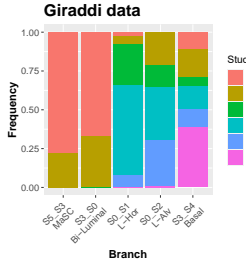
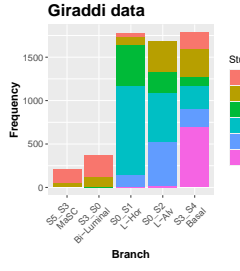
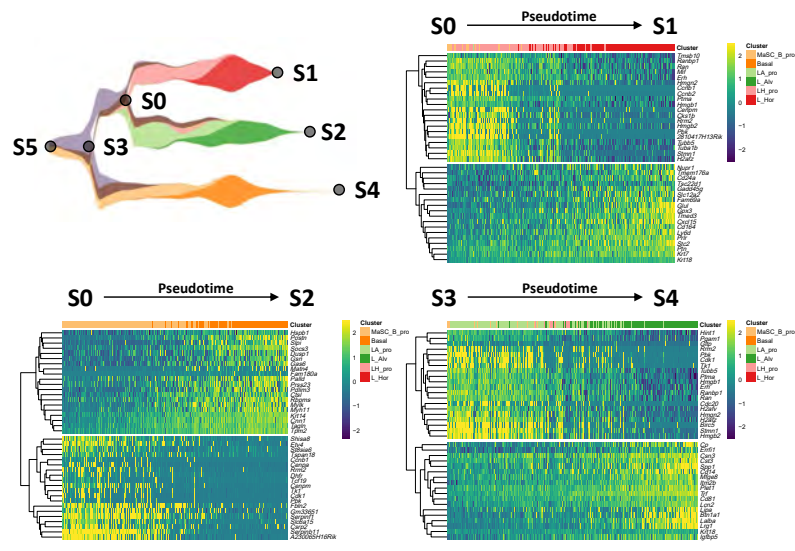
MaSC-B-pro LH-pro LA-iv  
LA-pro LH-pro Basal

**b**

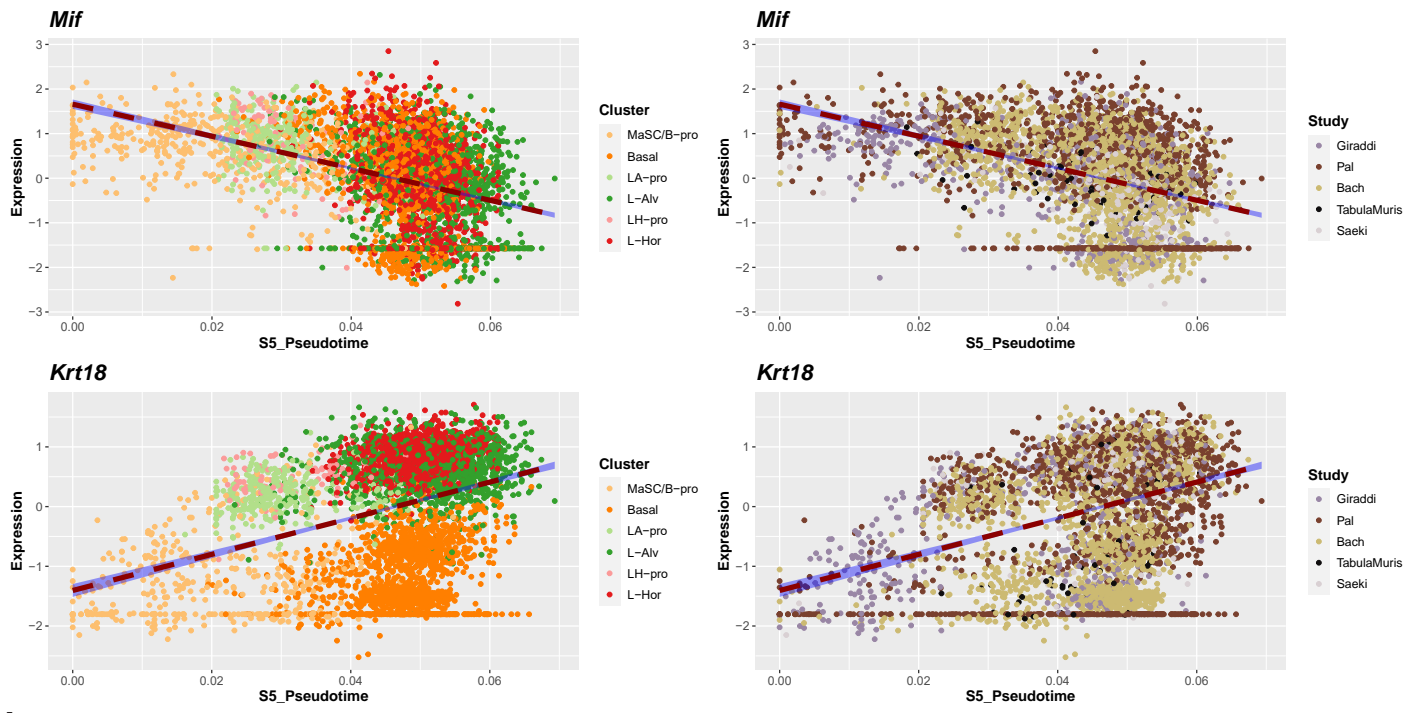
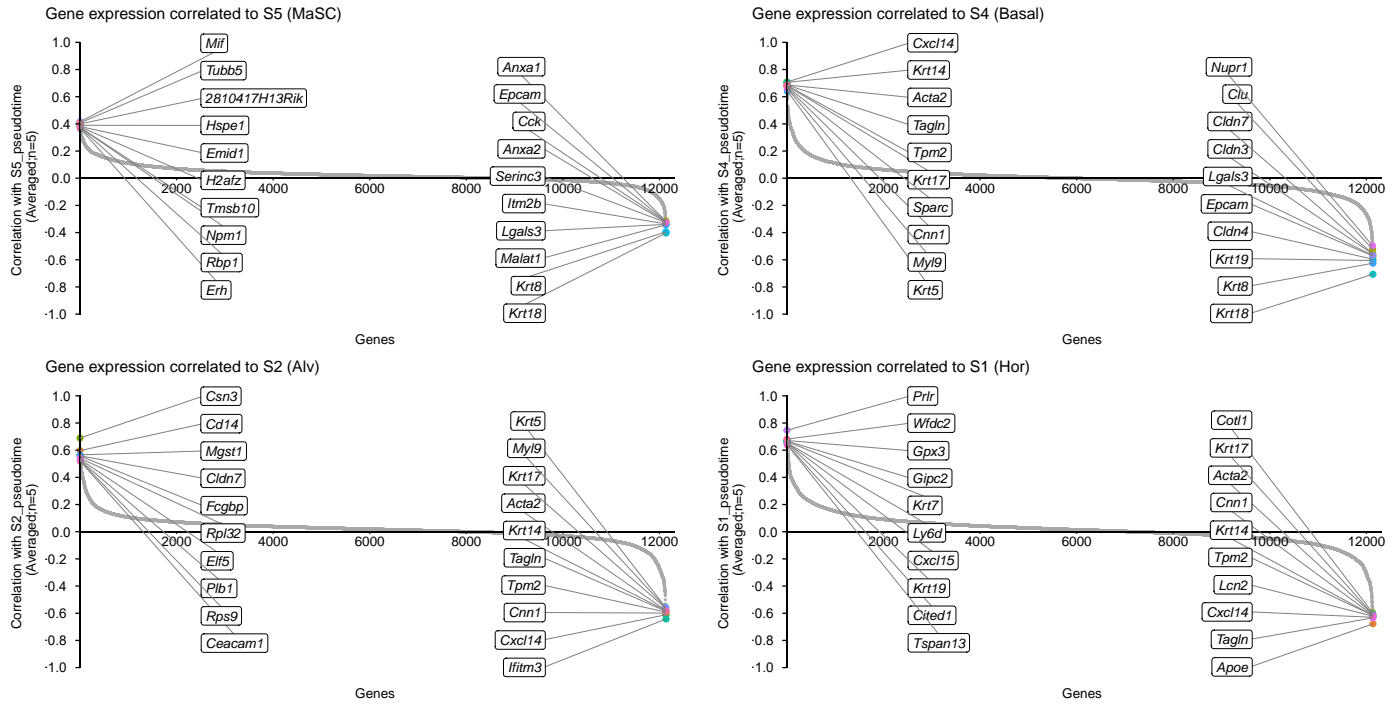
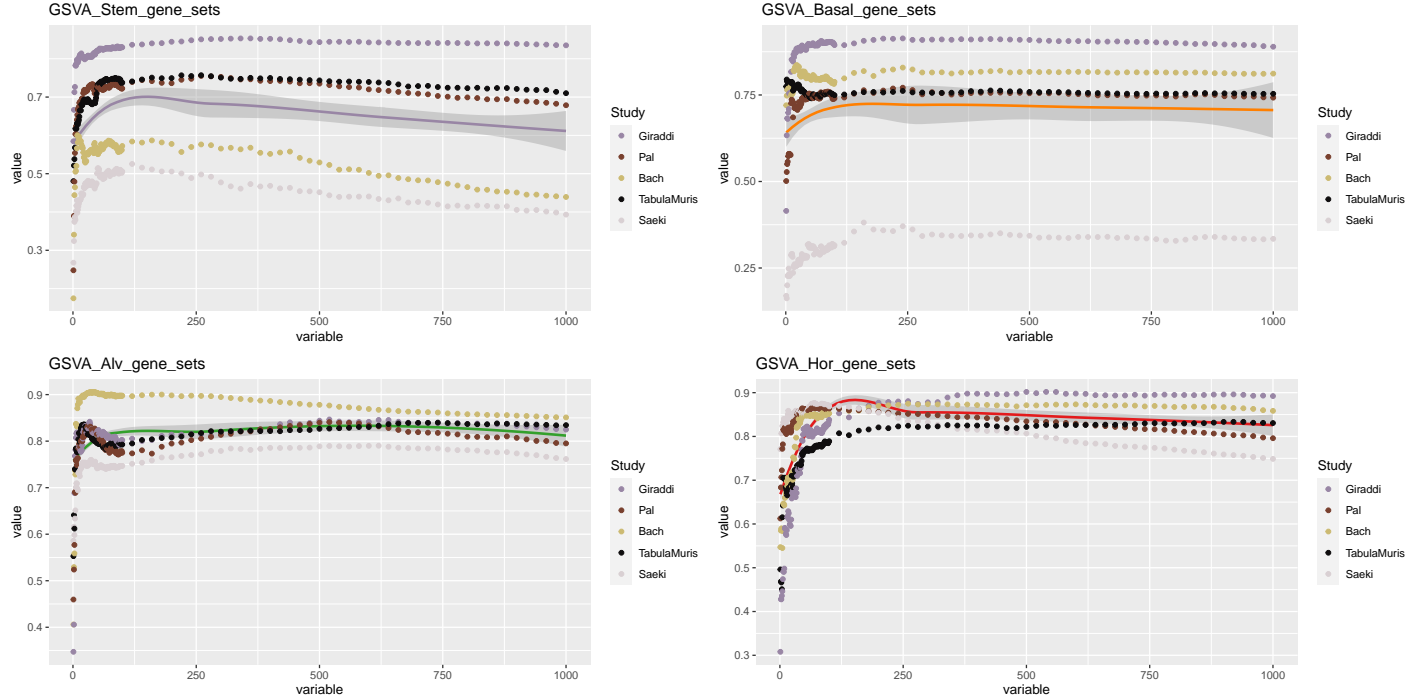
MaSC-B-pro LA-pro LH-pro  
Basal LH-pro LH-pro Basal



MaSC-B-pro LA-pro LH-pro  
LH-pro LA-iv Basal

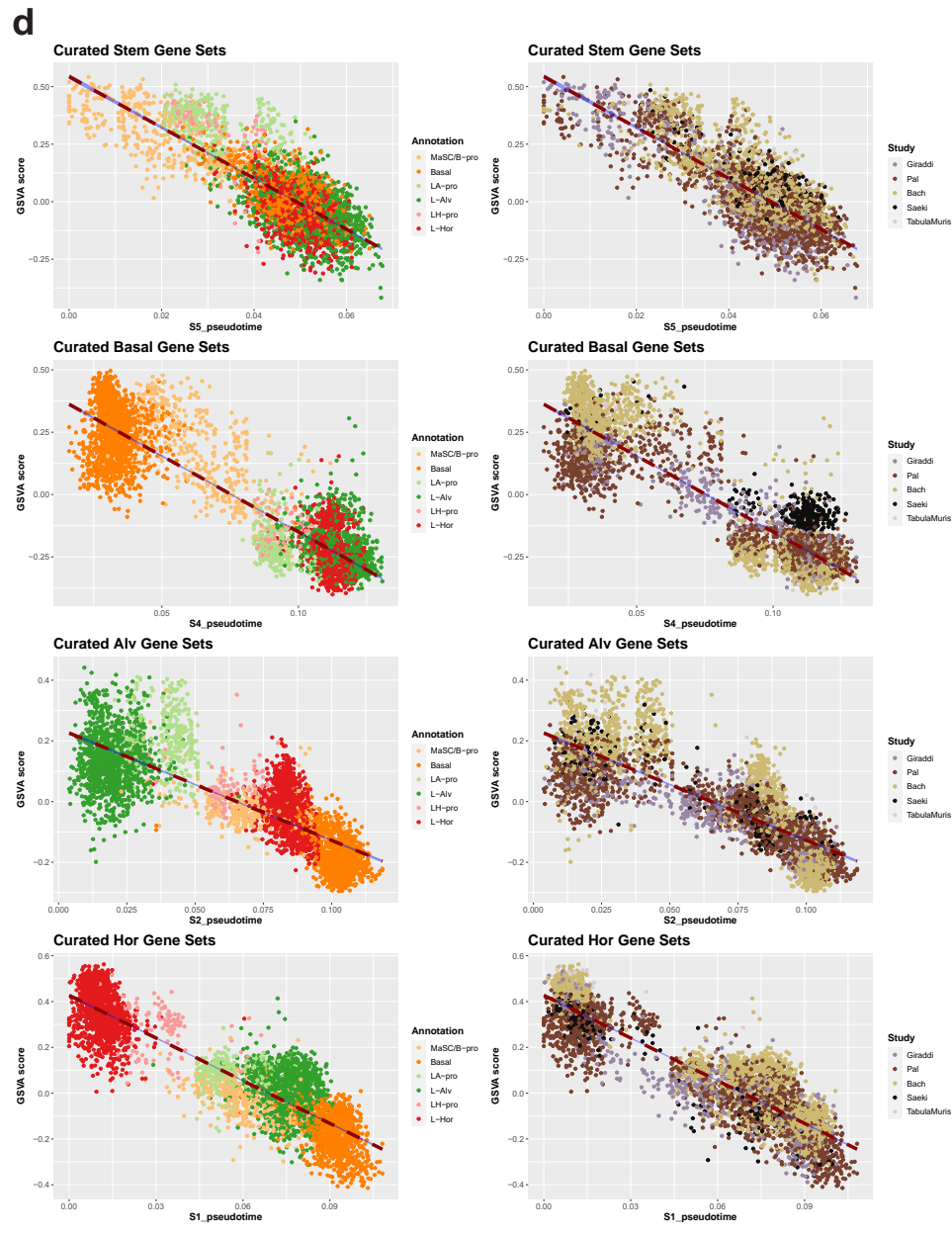
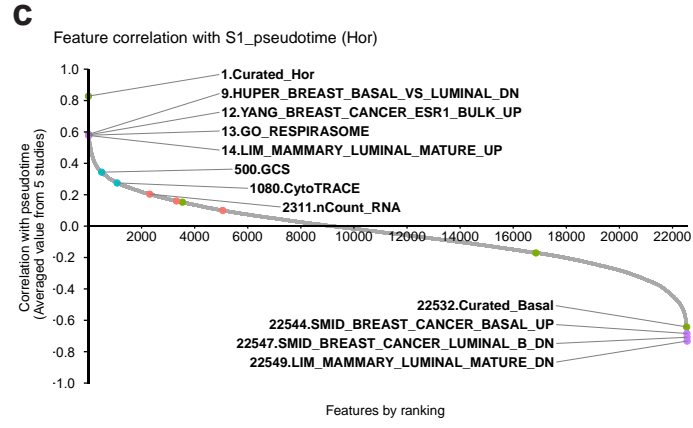
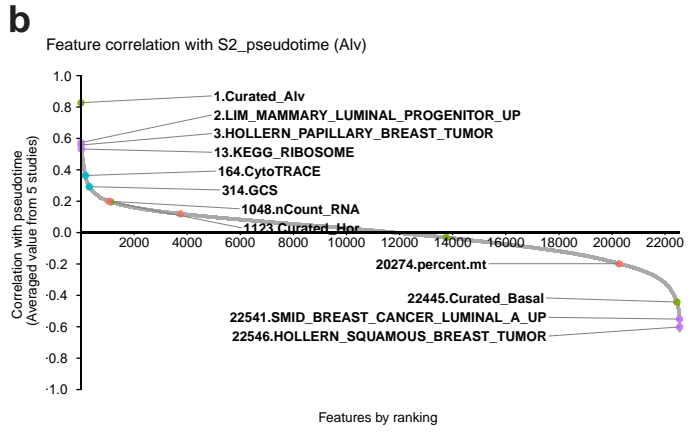
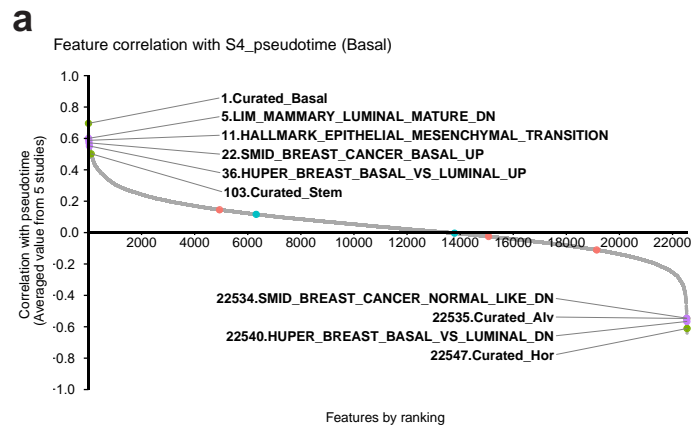
**c****d****e**

**Supplementary Figure 8 STREAM analysis.** (a) The base trajectory was built using the Girardi et al. dataset, projected in Modified Locally Linear Embedding (Mlle) spaces, and plotted on stream plots, color-coded by individual samples and clusters. (b) The 50K mouse mammary epithelial cells were mapped onto the base trajectory, projected in the Mlle space, and plotted on a stream plot. (c) Breakdown of each branch by datasets (Study), mouse strains (Strain), and different developmental stages (Stage). (d) Breakdowns of Girardi et al., Bach et al., and Pal et al. datasets shown in stacked bar plots. The absolute number and proportion of cells in each branch were presented and color-coded by samples. (e) Expression of the top 20 transitionally expressed genes in S0\_S1 (L-Hor), S0\_S2 (L-Alv), and S3\_S4 (Basal) branches were shown in heatmaps. Cells were sorted in the order of the relevant pseudotime.

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

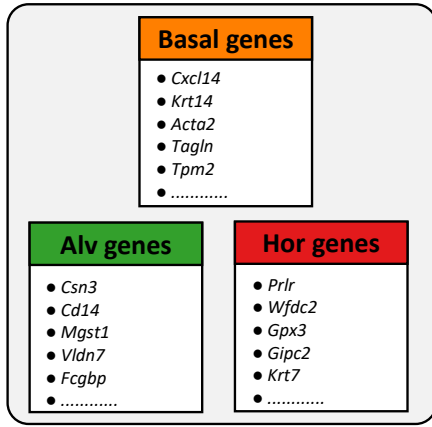
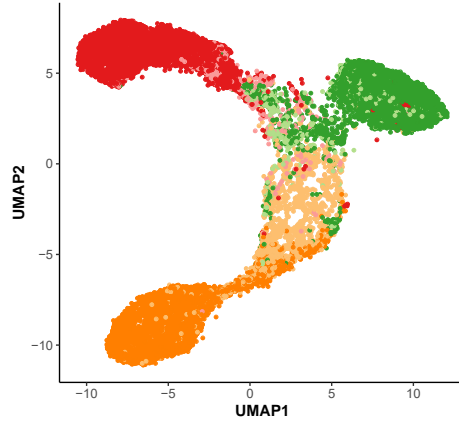
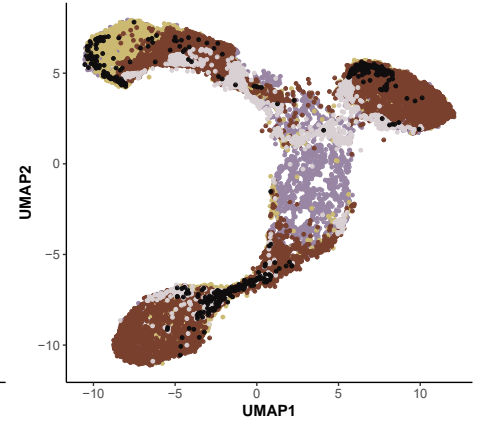
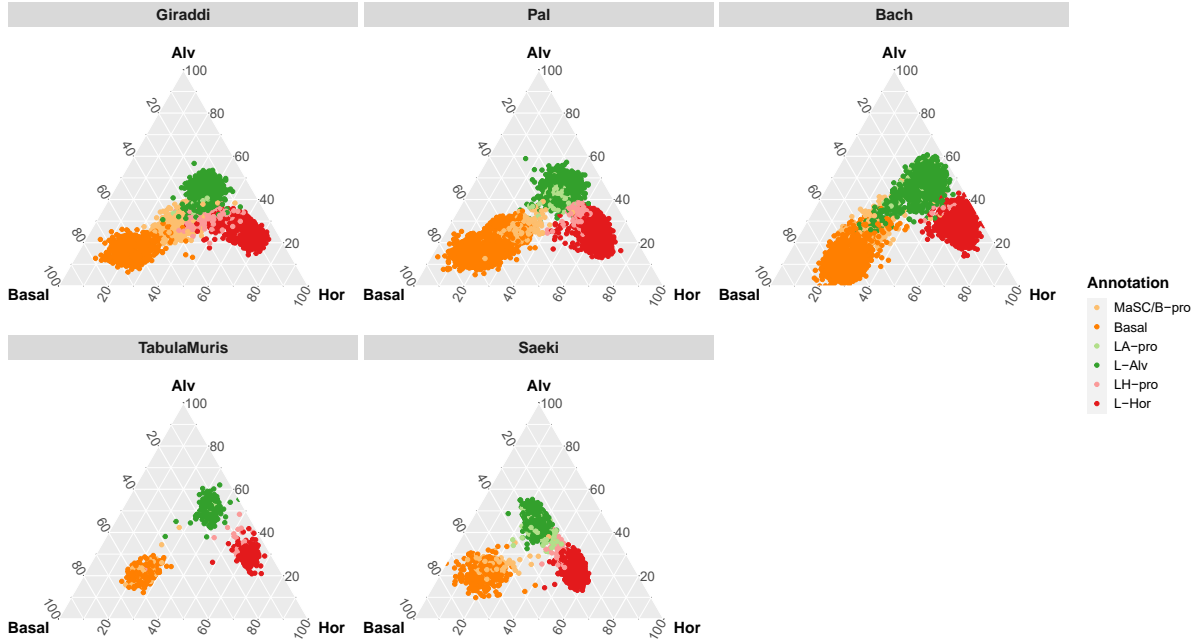
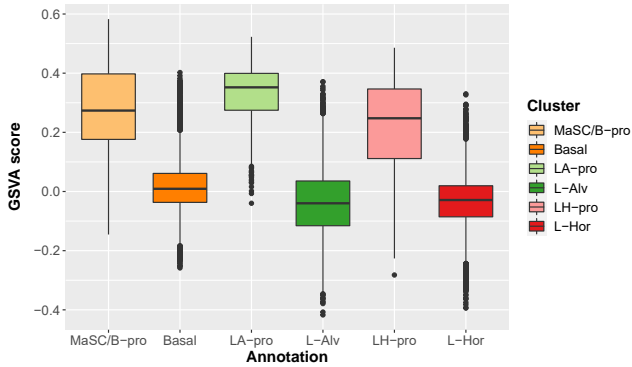
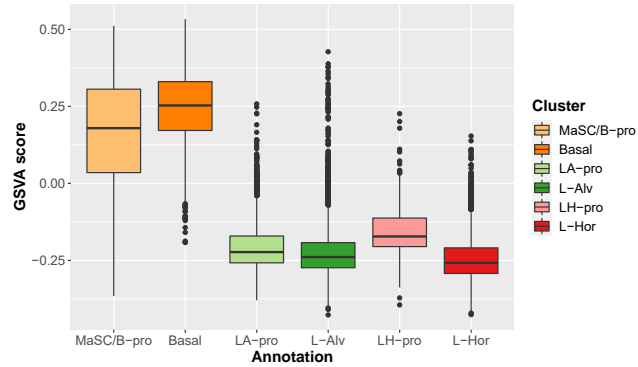
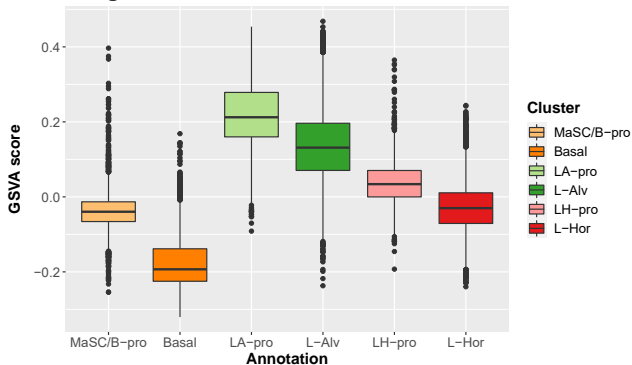
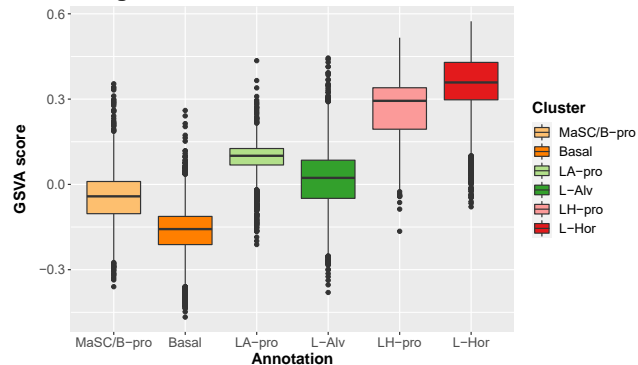
**Supplementary Figure 9 Identification of the lineage-specific genes and gene sets.** (a) The expression of the top positively (*Mif*) and negatively (*Krt18*) correlated genes with S5 pseudotime (“Stem” state). The X-axis represents the passage of pseudotime from the S5 node and the Y-axis represents gene expression that was individually scaled in each dataset. Each point represents a single cell. The points were color-coded by clusters and datasets. The dashed lines indicate regression lines. Five thousands cells were sampled from the entire data for UMAP visualization. (b) The results of the comprehensive correlation analysis between the gene expression and the pseudotimes. The combinations of the 12,319 genes and the four different pseudotimes [S5; Stem, S4; Basal, S2; Alv, and S1; Hor] were considered. The genes were ordered according to the correlation coefficients and the top 10 positively and negatively correlated genes were labeled. The X-axis represents the rank order and the Y-axis represents the averaged correlation coefficients from the five datasets. (c) The change in the performance of the gene sets according to the number of the top correlated genes in a gene set. The X-axis represents the number of gene(s) in a gene set and the Y-axis represents the correlation between the scGSVA score and the pseudotime [S5; Stem, S4; Basal, S2; Alv, and S1; Hor]. The correlation was calculated on a dataset basis. The solid lines represent regression curves in a LOESS model and the surrounding grey areas represent the confidence intervals.



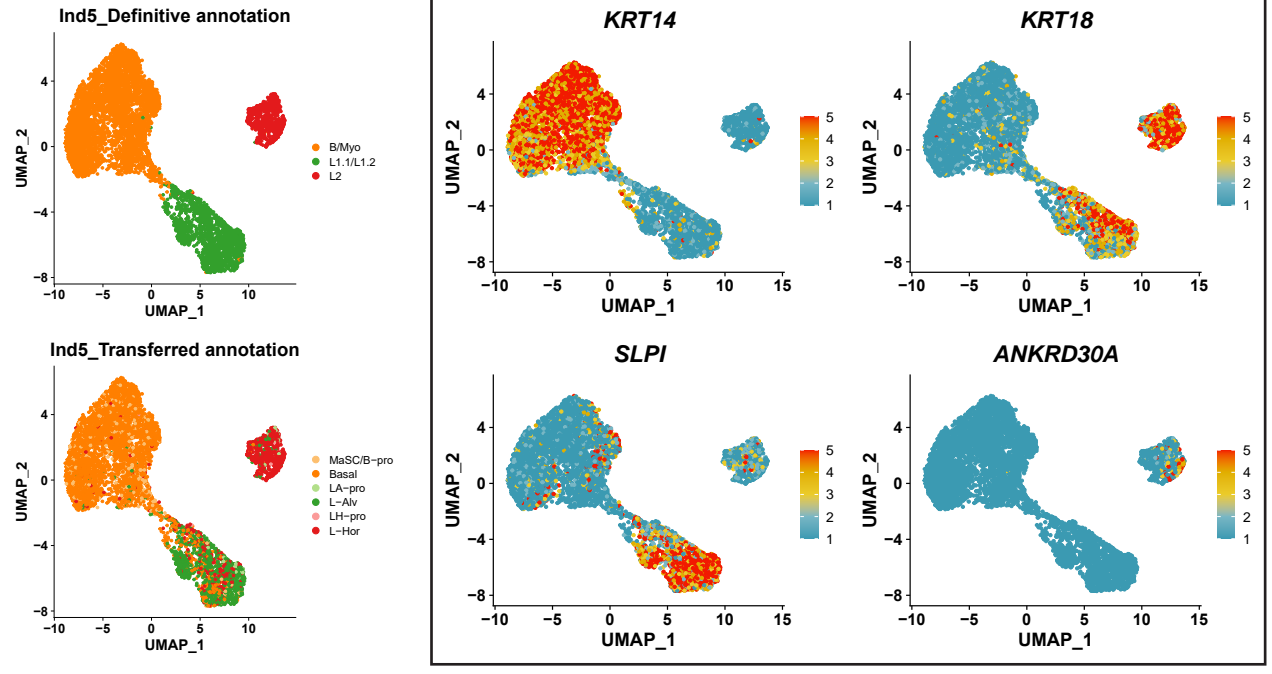
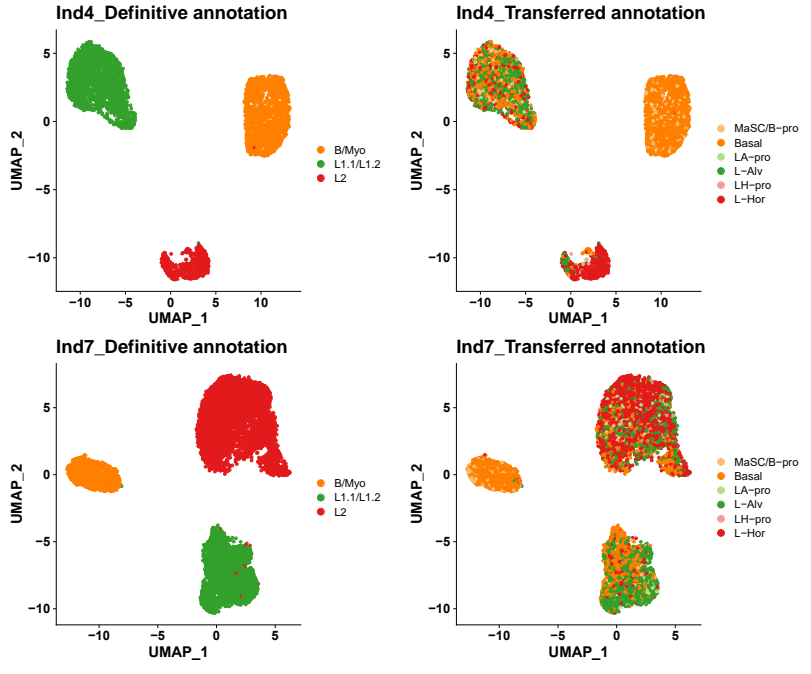
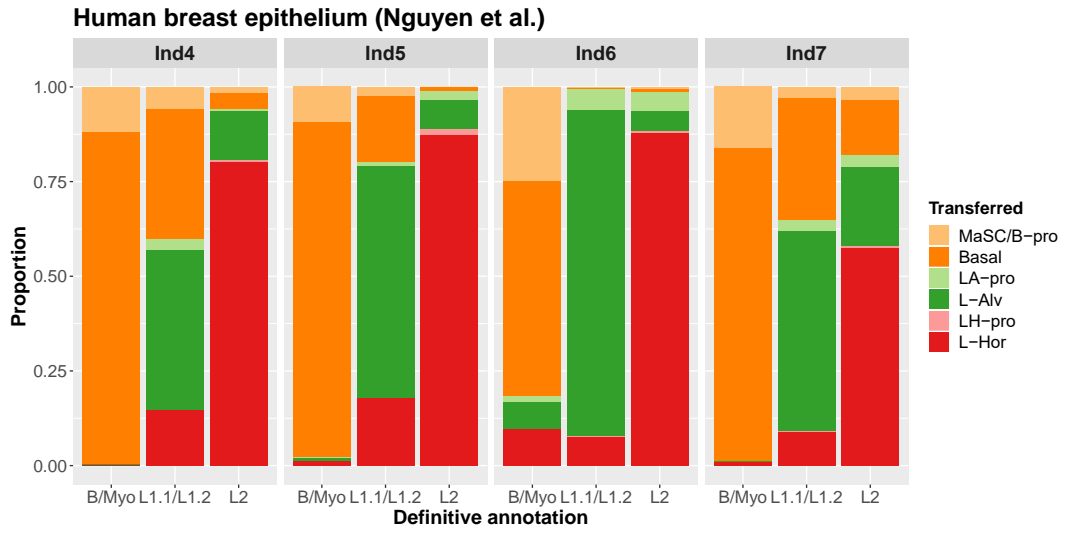


**Supplementary Figure 10 Performance evaluation of the curated lineage-specific gene sets.**

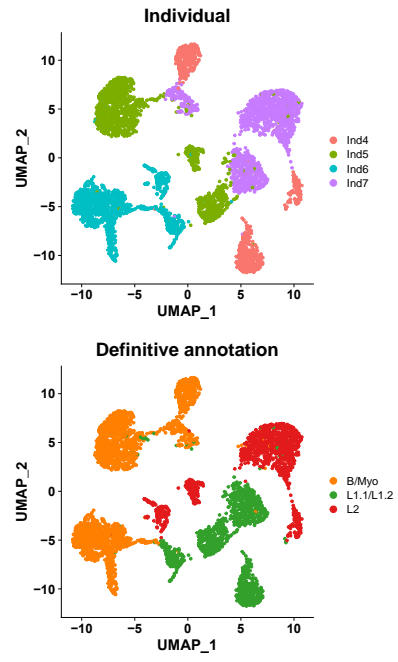
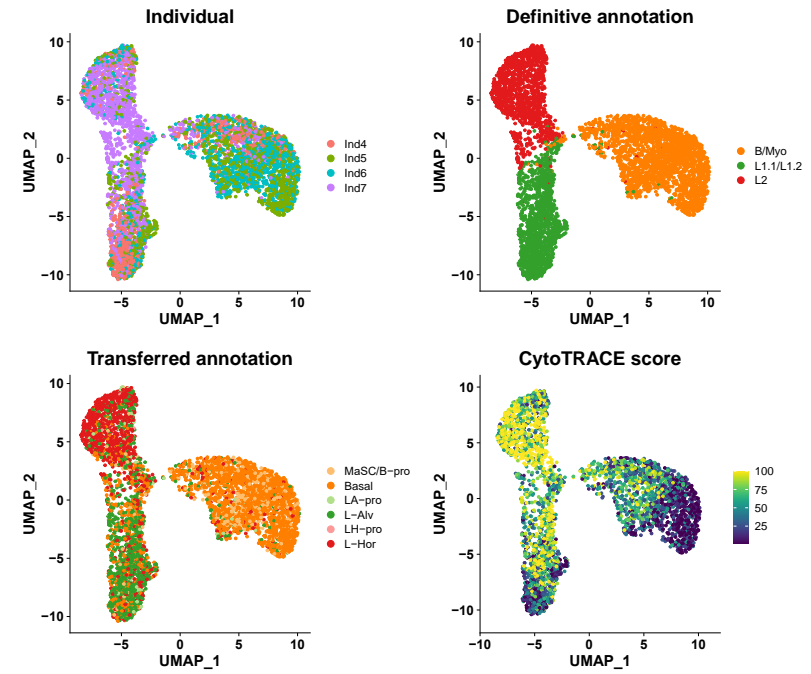
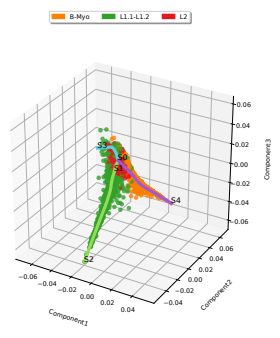
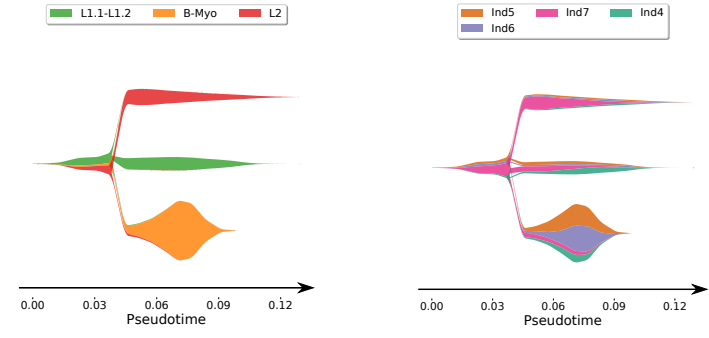
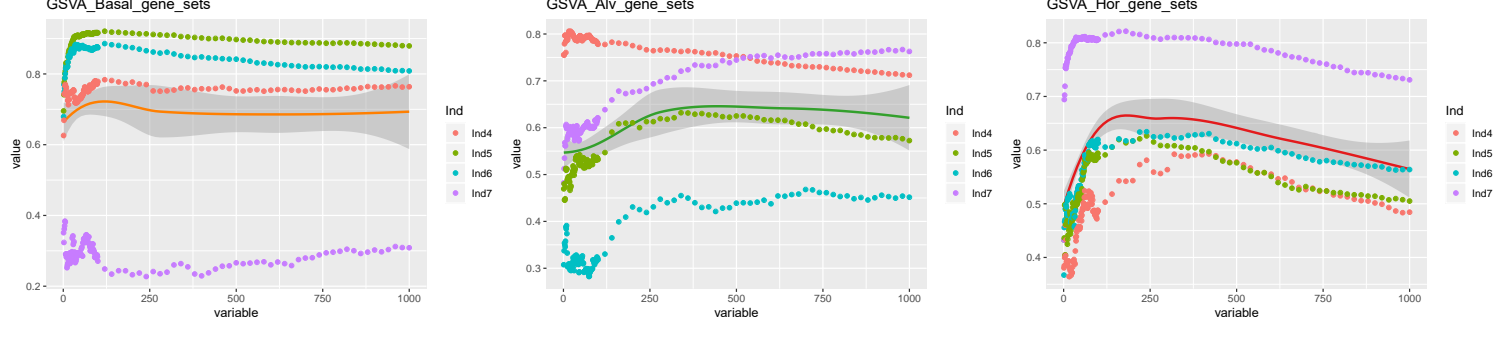
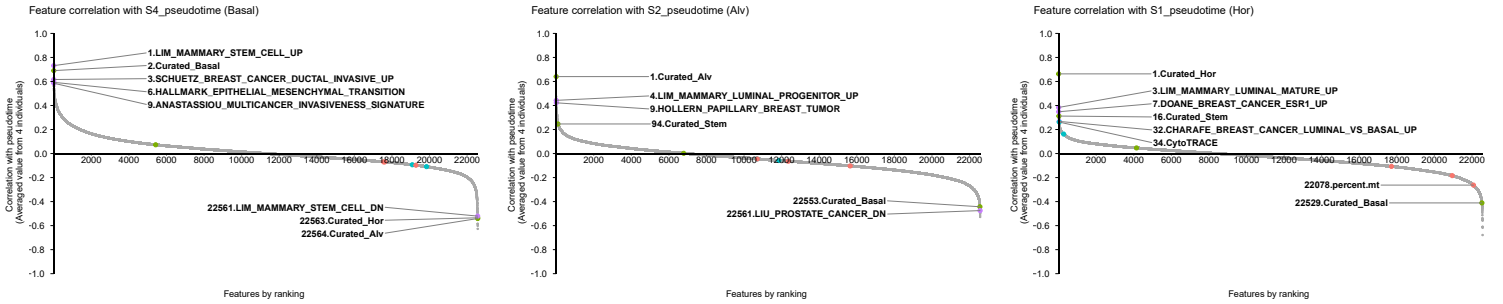
**(a-c)** Correlation of the scGSVA scores for the examined features with S4 pseudotime (Basal, **a**), S2 pseudotime (Alv, **b**), and S1 pseudotime (Hor, **c**). The X-axis represents the rank of the features in terms of the correlation coefficients and the Y-axis represents the correlation coefficients of the scGSVA scores with the pseudotime of interest. The selected features are highlighted and labeled on the plot. The results for the curated gene sets, algorithms, selected RNA-based features, and the cellular characteristics are colored in green, blue, purple, and red, respectively. The numbers on the text labels indicate their rankings. **(d)** The performance of the curated gene sets. The X-axis represents the passage of pseudotime from the node of interest and the Y-axis represents scGSVA scores for the curated gene sets (Stem; 160 genes, Basal; 240 genes, Alv; 500 genes, and Hor; 200 genes) that were individually calculated in each dataset. Each point represents a single cell. The points were color-coded by clusters and datasets. The dashed lines indicate regression lines. Five thousands cells were sampled from the entire data for UMAP visualization.

**a****GSVA analysis****All data (by cluster)****All data (by study)****b****Visualization by ternary plot****c****Stem gene set****Basal gene set****Alv gene set****Hor gene set**

**Supplementary Figure 11 Summarization and lineage inference of the scRNAseq data with the curated gene sets.** (a) The data were summarized with three differentiation-associated gene sets (Basal, Alv, and Hor), and projected on the UMAP plot, color-coded by clusters and datasets. (b) Visualization of single cells on ternary plots using the scGSVA scores for the three differentiation-associated gene sets (Basal, Alv, and Hor). (c) Distribution of scGSVA scores for the four curated lineage-specific gene sets in each cluster. The n number of each cluster is as follows: MaSC/B-pro (n=4,217), Basal (n=14,493), LA-pro (n=3,176), L-Alv (n=11,856), LH-pro (n=1,156), and L-Hor (n=15,509) (biologically independent samples). The box-plot elements were defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers.

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

**Supplementary Figure 12 scRNAseq data of the human breast epithelium.** (a) The definitive and transferred annotation of the cells from individual #five (Ind5). According to the original publication, the expression of marker genes for the definitive annotation was visualized on UMAP plots (*KRT14*; Basal and Myoepithelial cells, *KRT18*; Luminal cells, *SLPI*; L1.1 and L1.2 cells, and *ANKRD30A*; L2). (b) The definitive and transferred annotation of the cells from individual #four (Ind4) and individual #seven (Ind7) were projected on UMAP dimensionality. (c) The stacked bar plots represent the agreement regarding the lineage inference between the definitive and the transferred annotation, broken down by individuals.

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

**Supplementary Figure 13 Data integration and identification of the lineage-specific gene**

**sets in the human breast epithelium scRNAseq data. (a)** The merged data, before integration,

visualized on UMAP plots and color-coded by individuals and definitive annotation. Five

thousands cells were sampled from the entire data for UMAP visualization. **(b)** The integrated

data visualized on UMAP plots, colored by individuals, definitive annotation, transferred

annotation, and CytoTRACE scores. Five thousands cells were sampled from the entire data for

UMAP visualization. **(c)** The trajectory learned in the Mlle space using 1,000 human breast

epithelial cells from each individual. **(d)** The stream plots of the 24K human breast epithelial

cells, color-coded by definitive clusters and individuals. **(e)** The change in the performance of the

gene sets according to the number of the top correlated genes in a gene set. The X-axis represents

the number of gene(s) in a gene set and the Y-axis represents the correlation between the

scGSVA score and the pseudotime [S4; Basal, S2; Alv, and S1; Hor]. The correlation was

calculated on an individual basis. The solid lines represent regression curves in a LOESS model

and the surrounding grey areas represent the confidence intervals. **(f)** Correlation of the scGSVA

scores for the examined features with S4 pseudotime (Basal), S2 pseudotime (Alv), and S1

pseudotime (Hor). The X-axis represents the rank of the features in terms of the correlation

coefficients and the Y-axis represents the correlation coefficients of the scGSVA scores with the

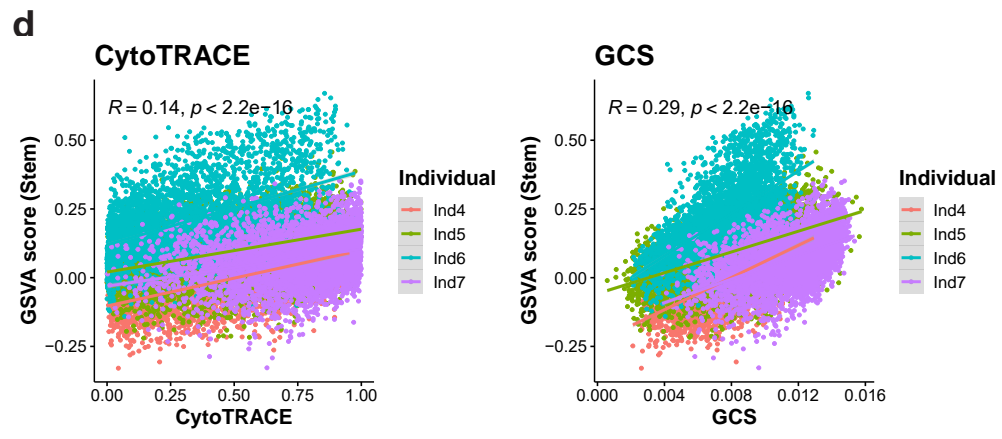
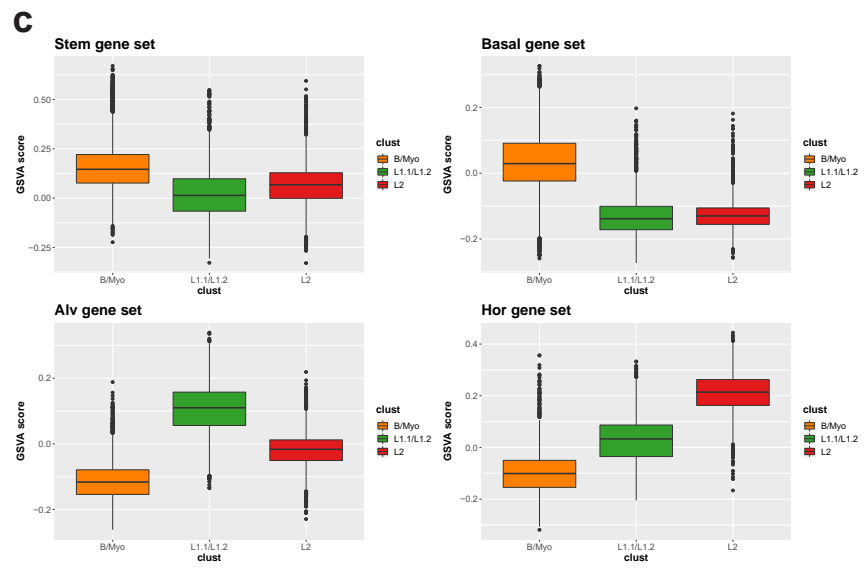
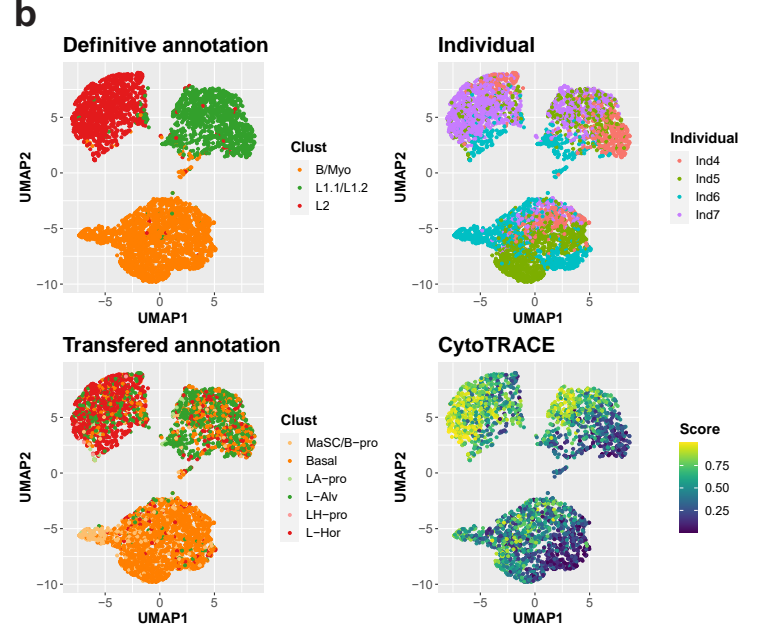
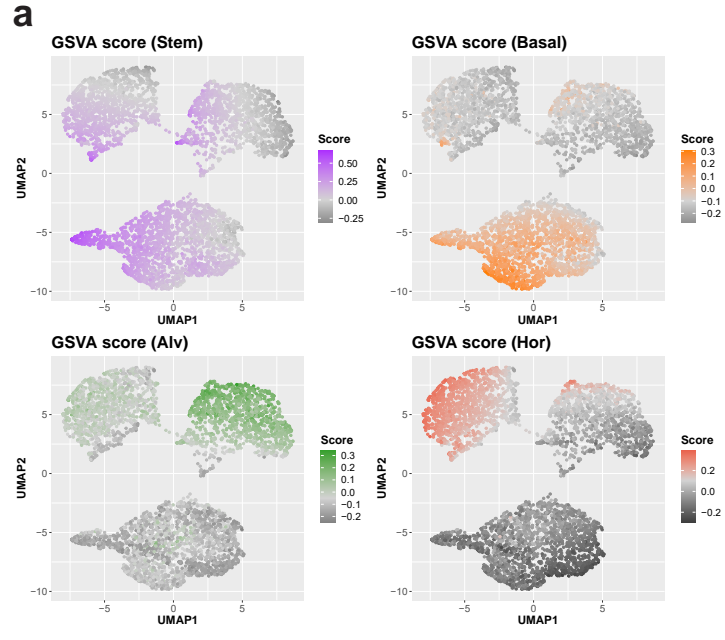
pseudotime of interest. The selected features are highlighted and labeled on the plot. The results

for the curated gene sets, algorithms, selected RNA-based features, and cellular characteristics

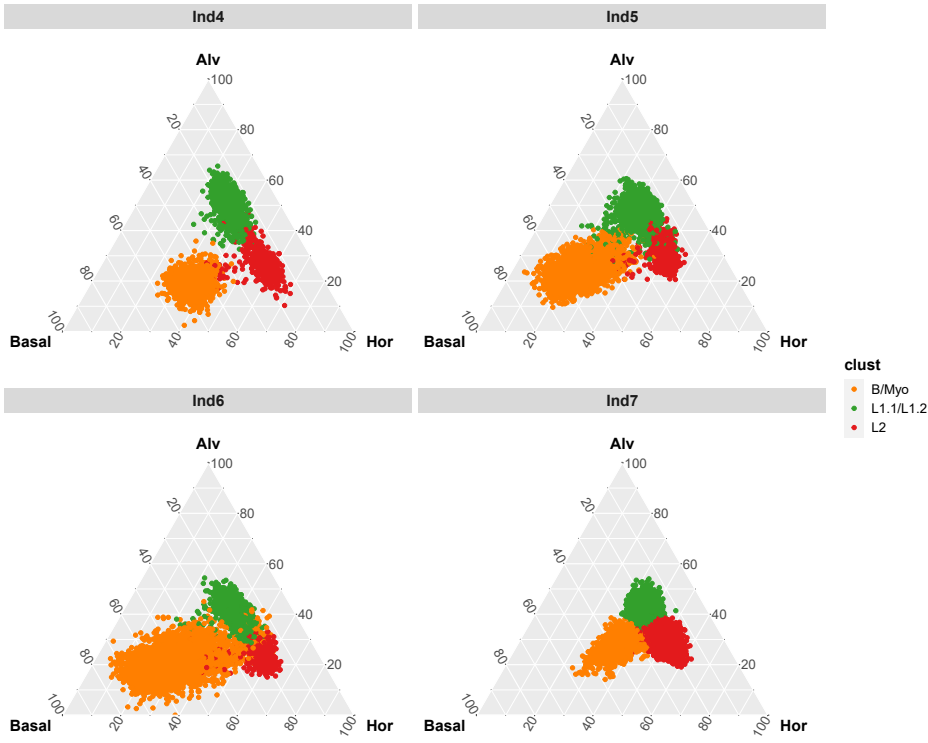
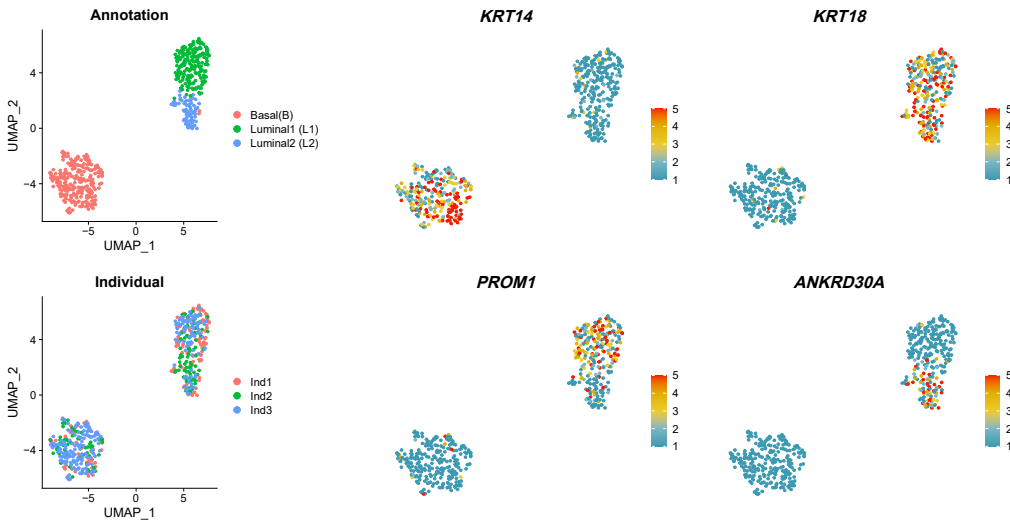
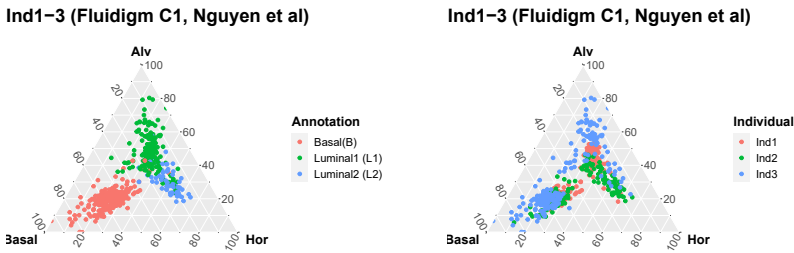
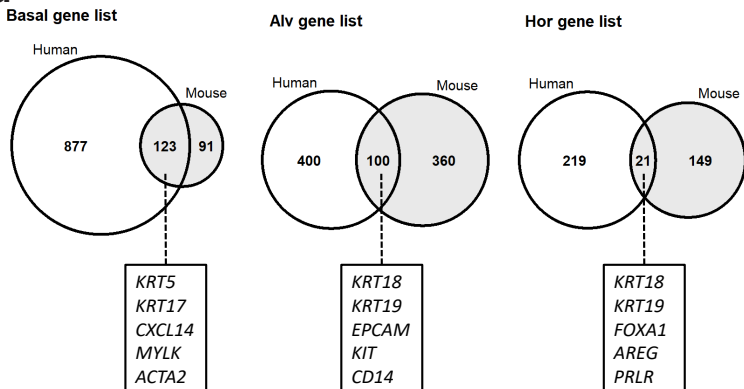
are colored in green, blue, purple, and red, respectively. The numbers on the text labels indicate

the ranking of features.

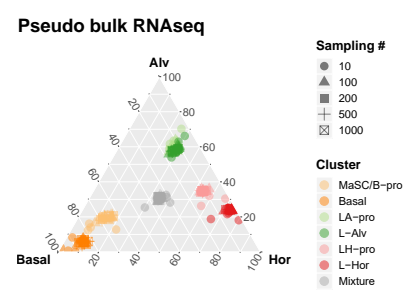
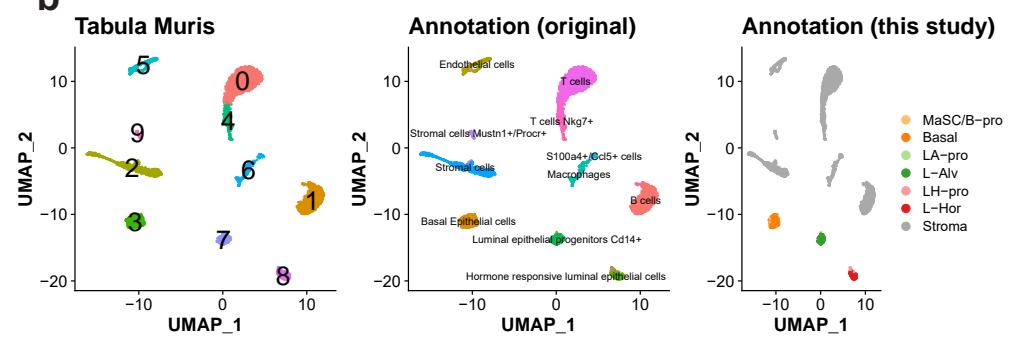
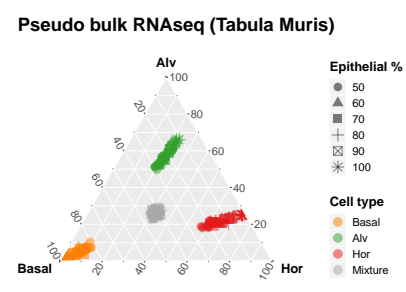
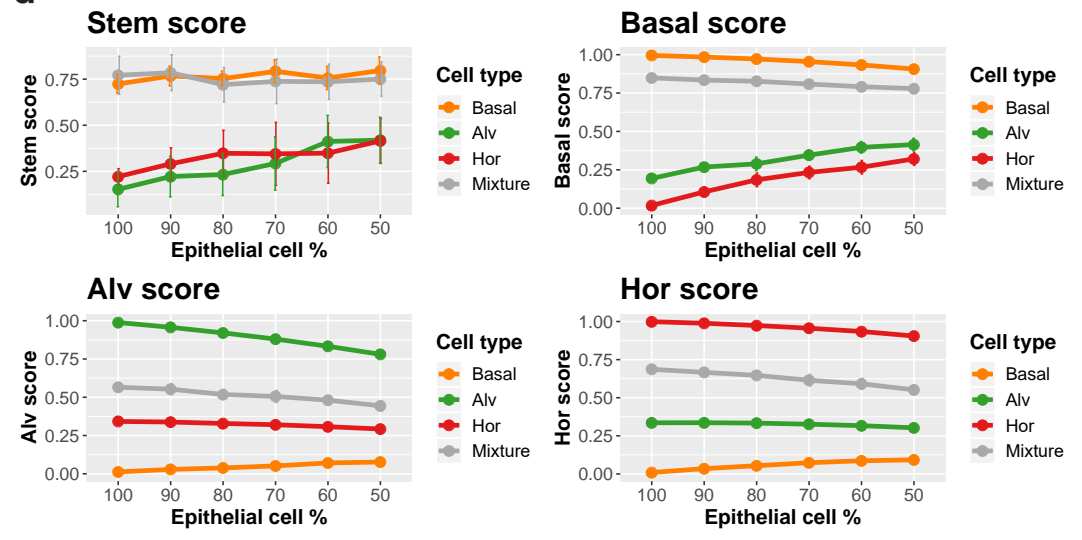
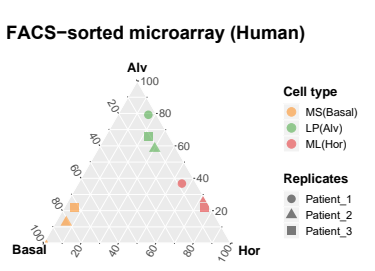
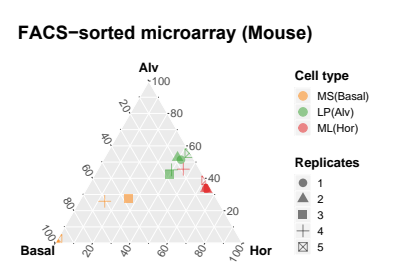




**Supplementary Figure 14 Lineage-specific gene sets-based summarization of the human breast epithelium scRNAseq data.** (a) The scGSVA scores for the four curated gene sets used for summarization were projected on UMAP plots. Five thousands cells were sampled from the entire data for UMAP visualization. (b) UMAP plots colored by definitive annotation, individuals, transferred annotation, and CytoTRACE scores. Five thousands cells were sampled from the entire data for UMAP visualization. (c) Distribution of the scGSVA scores for the four curated lineage-specific gene sets in each cluster. The n number of each cluster is as follows: B/Myo (n=11,205), L1.1/L1.2 (n=6,687), and L2 (n=6,388) (biologically independent samples). The box-plot elements were defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. (d) The correlations of the scGSVA scores for the curated Stem gene set with the CytoTRACE score and Gene Count Signature (GCS). n = 14,967 cells, Spearman's rank correlation.

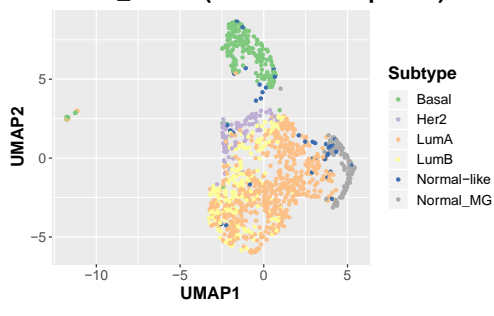
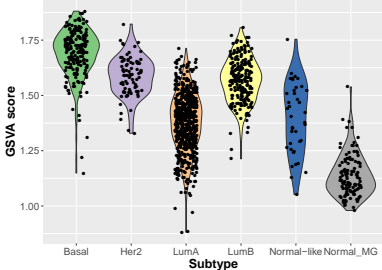
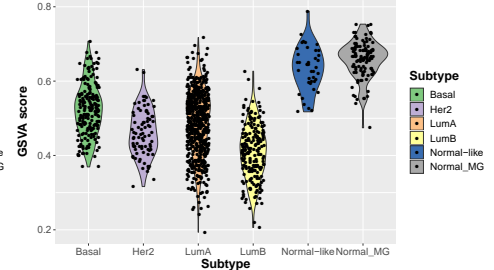
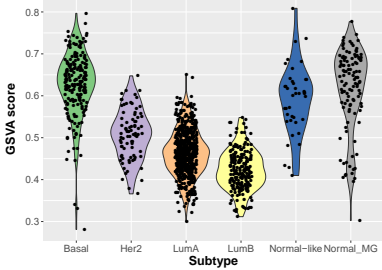
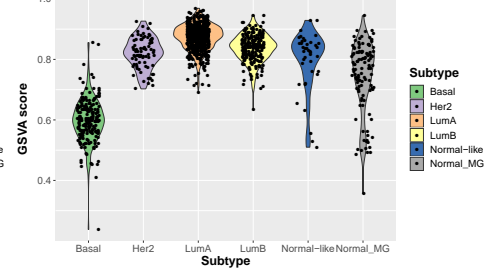
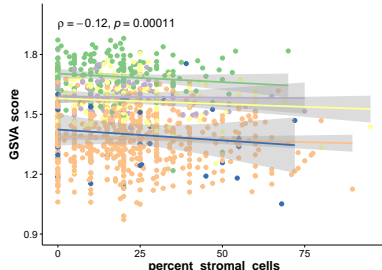
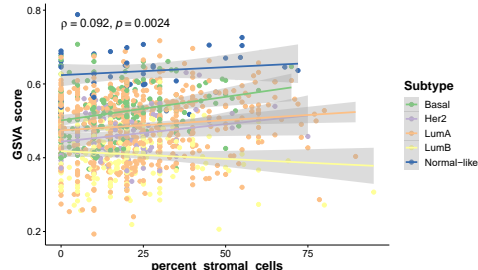
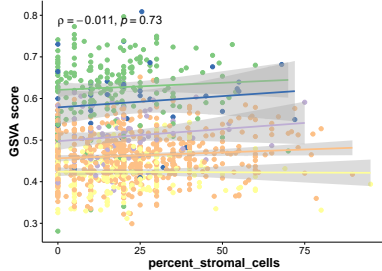
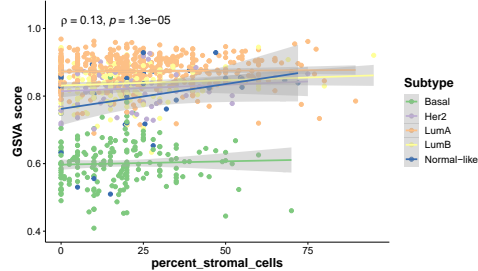
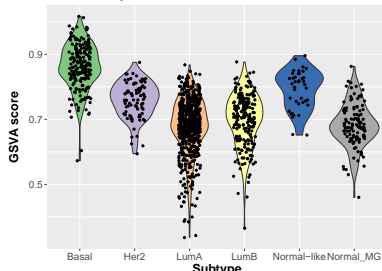
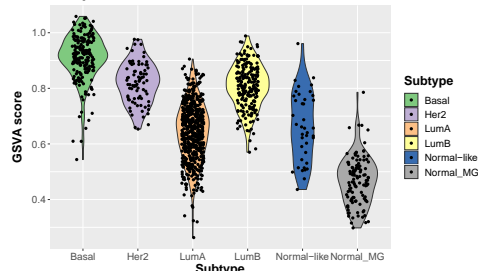
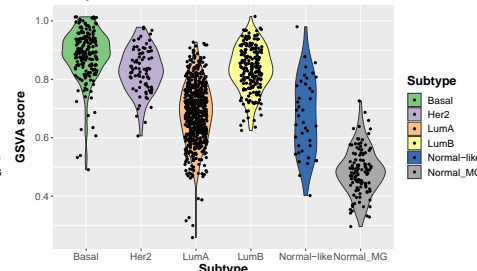
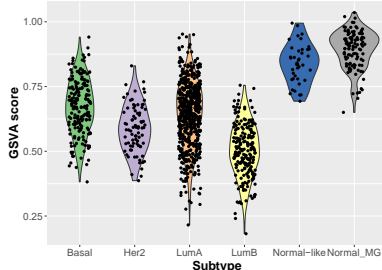
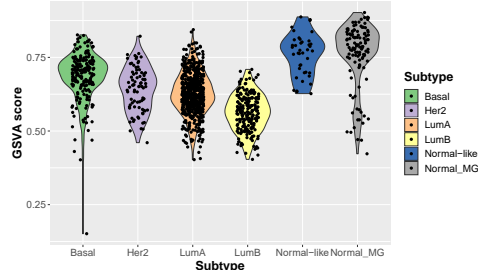
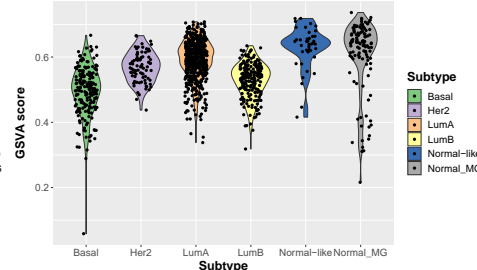
**a****Human breast epithelium****b****c****d**

**Supplementary Figure 15 Application of the human lineage gene sets to the additional datasets and the comparison between the species.** (a) Visualization of the single cells on ternary plots using the scGSVA scores for the three differentiation-associated gene sets (Basal, Alv, and Hor). The four individuals from the 10X data were shown. (b) The additional three individuals from the Fluidigm C1 data visualized on UMAP plots and color-coded by definitive annotation and individuals. According to the original publication, the expression of marker genes for the definitive annotation was visualized on UMAP plots (*KRT14*; Basal and Myoepithelial cells, *KRT18*; Luminal cells, *PROM1*; L1 cells, and *ANKRD30A*; L2 cells). (c) Visualization of additional data on ternary plots based on the scGSVA scores for the three differentiation-associated gene sets (Basal, Alv, and Hor). The plots were color-coded by definitive annotation and individuals. (d) Comparisons of the lineage gene sets between mice and humans. The representative common genes were shown.

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

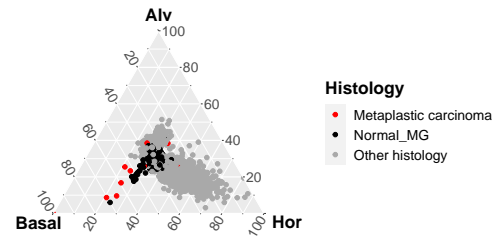
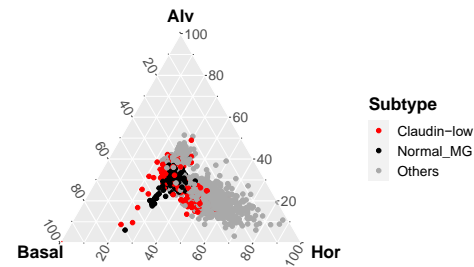
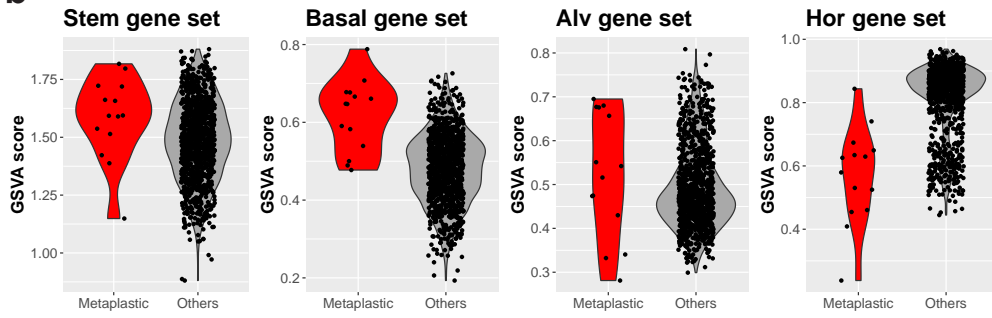
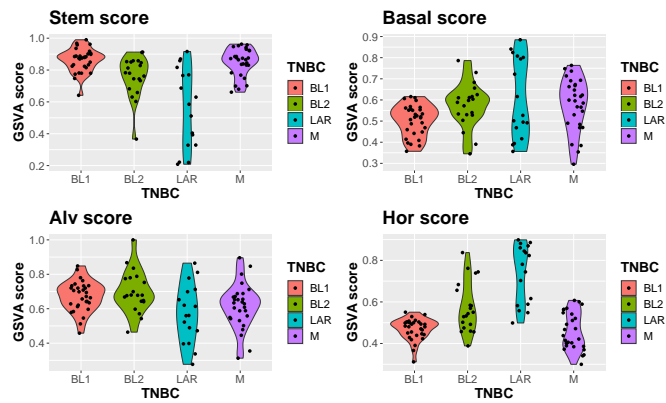
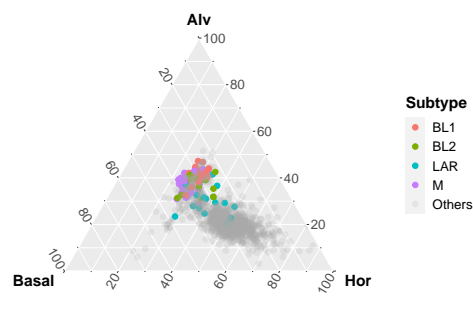
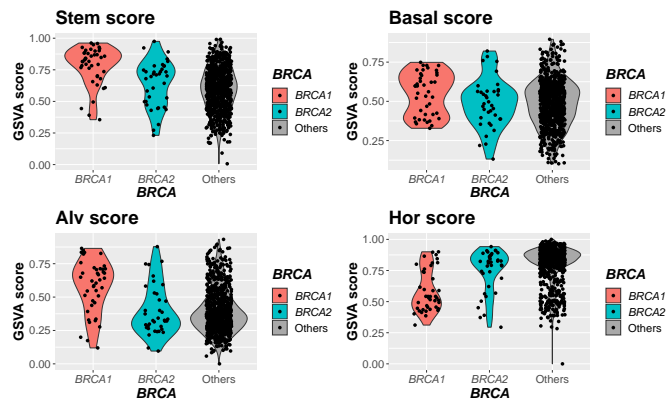
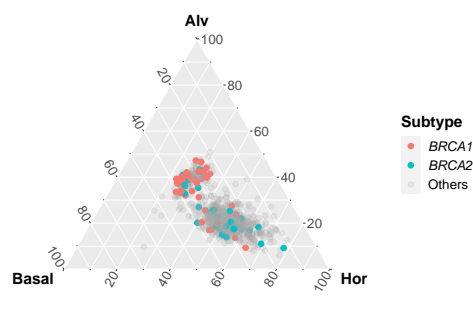
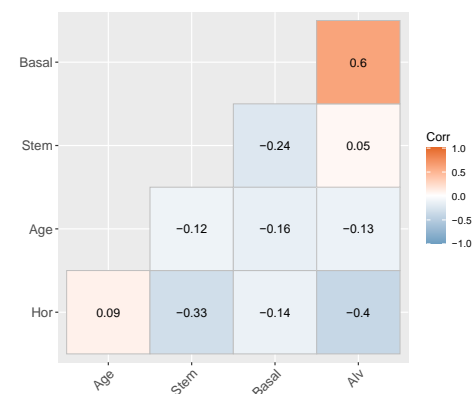
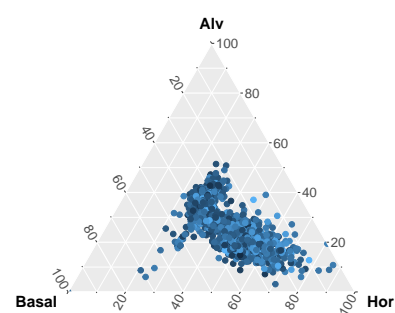
**Supplementary Figure 16 Extention of the lineage inference to the bulk RNAseq data. (a)**

The pseudo bulk RNA seq of the different cell types in the mouse mammary epithelium with different sampling numbers from the scRNAseq data. Mixture indicates that cells were evenly sampled from the six different cell types. **(b)** The Tabula Muris data used for the pseudo tissue bulk RNAseq. The Louvain clustering, original annotations in the project, and clustering of the epithelial cells obtained from the integrative analysis were visualized on UMAP plots. **(c)** The pseudo bulk RNA seq of the three different cell lineages in the mouse mammary epithelium mixed with the varying number of stromal cells. Mixture indicates that cells were evenly sampled from the three different lineages. **(d)** The changes in the scGSVA scores for the curated lineage-specific gene sets, with the increasing proportion of the stromal cells in the pseudo bulk RNAseq. Each data was generated by averaging the results from 10 independent simulations. Error bars represent standard deviations. **(e)** The FACS-sorted transcriptome analysis of the three major cell types in the mouse and human mammary epithelium. The results were scored using the scGSVA analysis with the lineage-specific gene sets

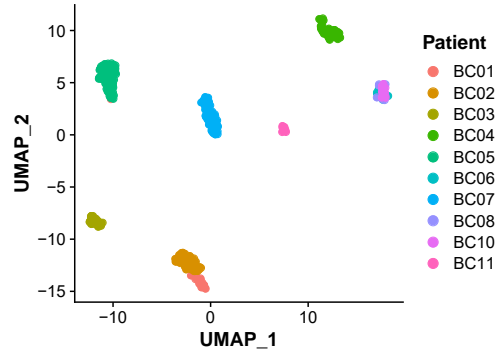
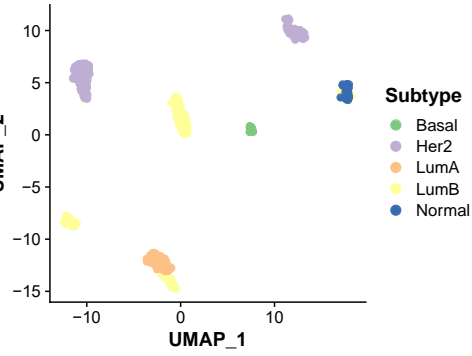
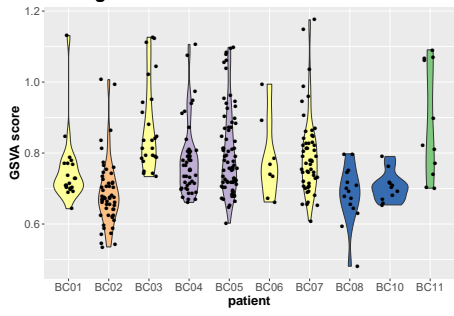
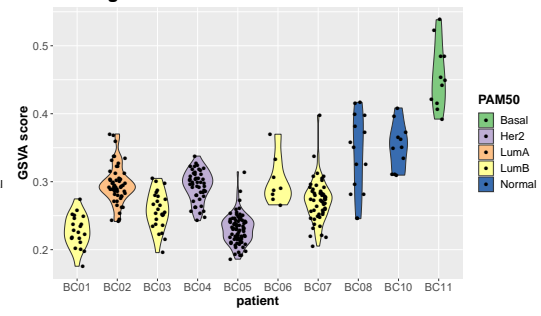
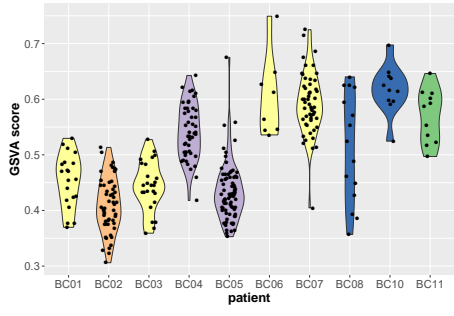
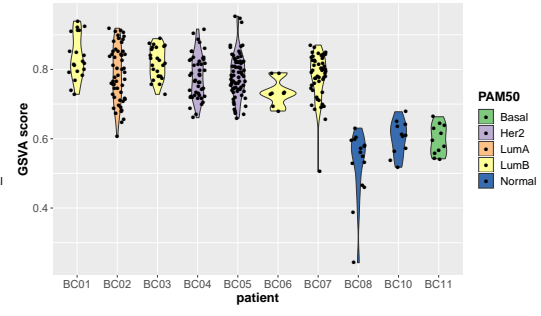
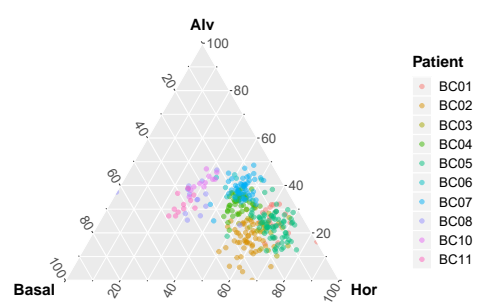
**a****TCGA\_BRCA (whole transcriptome)****b****Stem gene set****Basal gene set****Alv gene set****Hor gene set****c****Stem gene set****Basal gene set****Alv gene set****Hor gene set****d****MaSC/B-pro score****LA-pro score****LH-pro score****Basal score****L-Alv score****L-Hor score**

**Supplementary Figure 17 scGSVA analysis of the TCGA BRCA data.** (a) The summarization of the data in UMAP dimensionality, considering the expression of all detected transcripts. (b) Distribution of the scGSVA scores for the four curated lineage-specific gene sets in each subtype. (c) The correlation of the scGSVA scores for the lineage-specific gene sets with the stromal proportion in the tumor tissue (n=1,083). Spearman's rank correlation. (d) Distribution of the scGSVA scores for the DEGs of the six clusters identified in the mouse integrated data in each subtype. The n number of each subtype is as follows: Basal (n=194), Her2 (n=82), LumA (n=567), LumB (n=207), Normal-like (n=40), and Normal\_MG (n=113) (biologically independent samples).

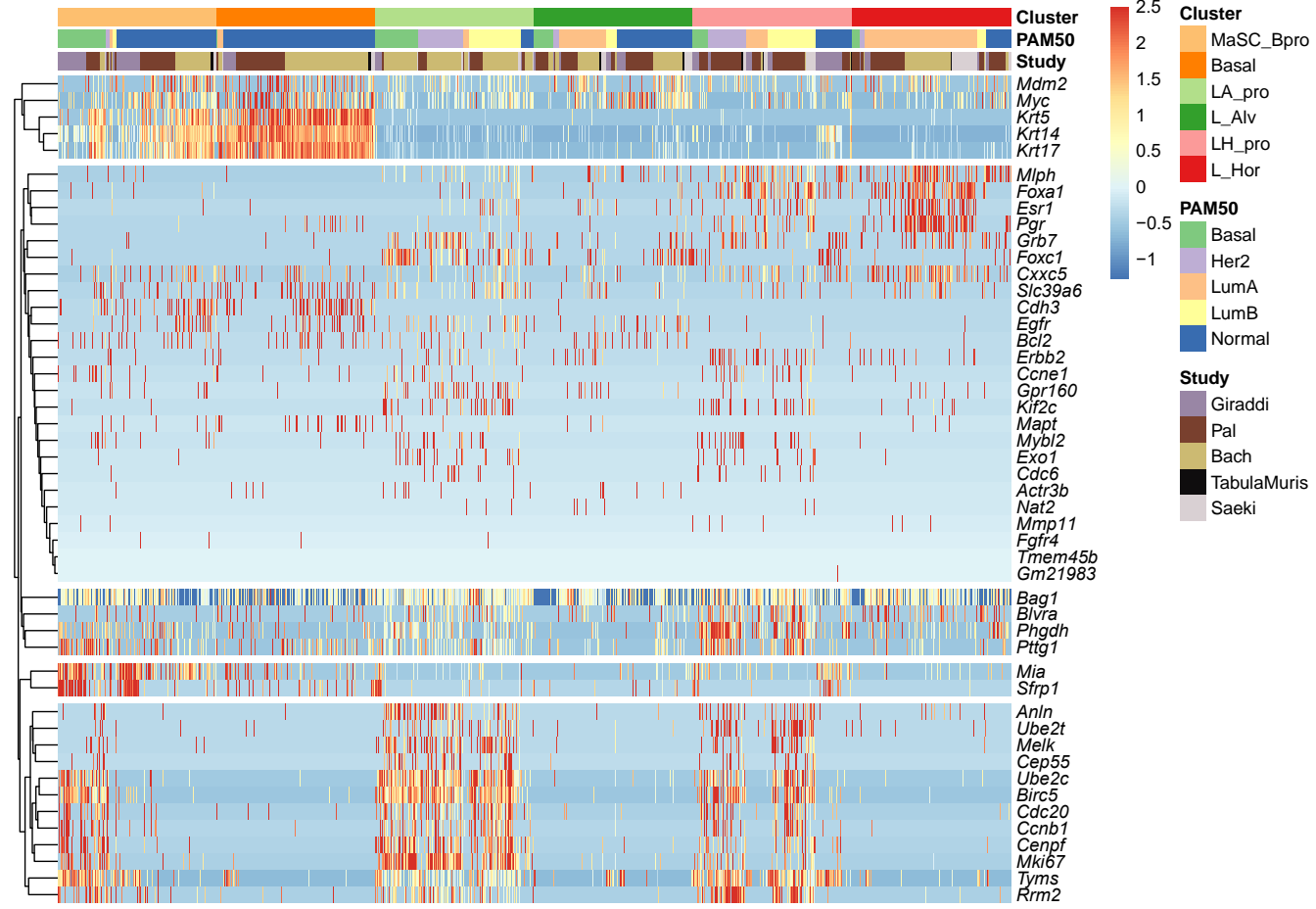
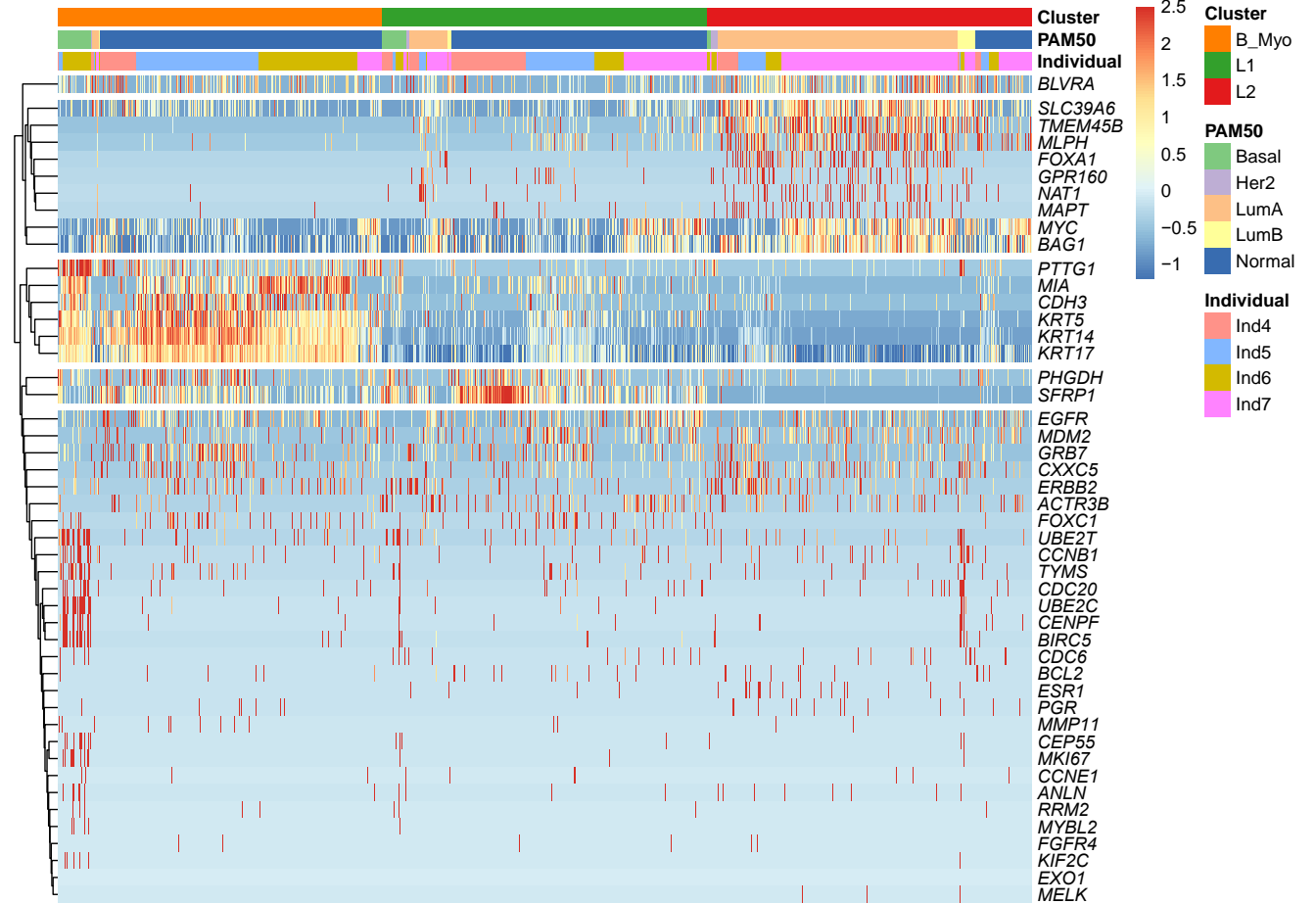


**a****TCGA\_BRCA (Histology)****TCGA\_BRCA (Claudin-low subtype)****b****c****TCGA\_TNBC\_subtyping****d****TCGA\_BRCA\_mutation****e****Relationship\_origin\_and\_age**

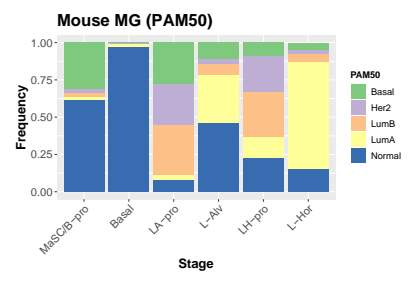
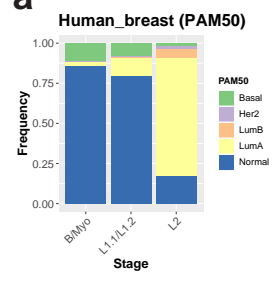
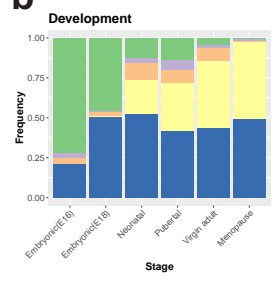
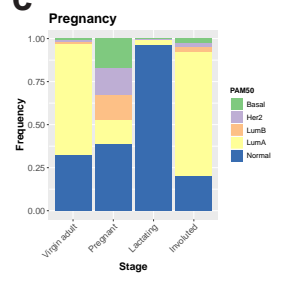
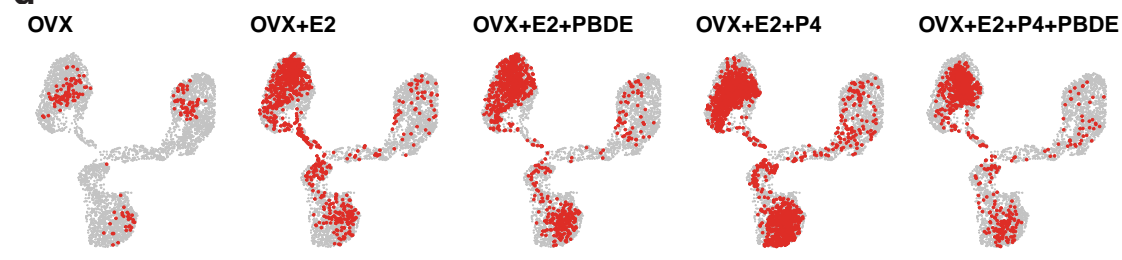
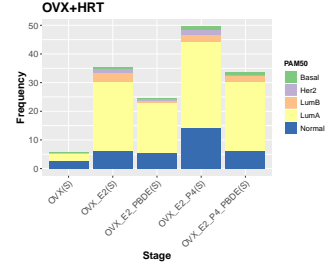
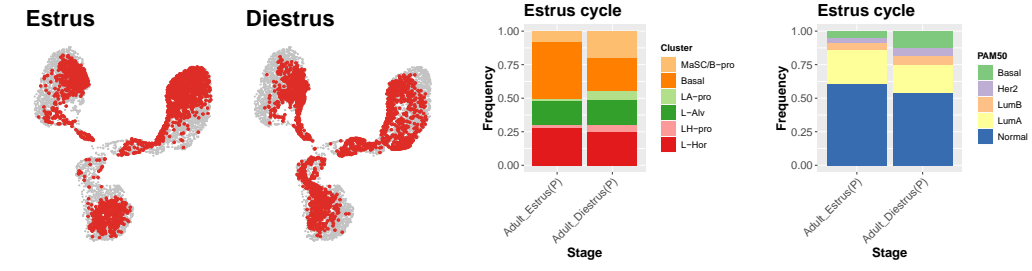
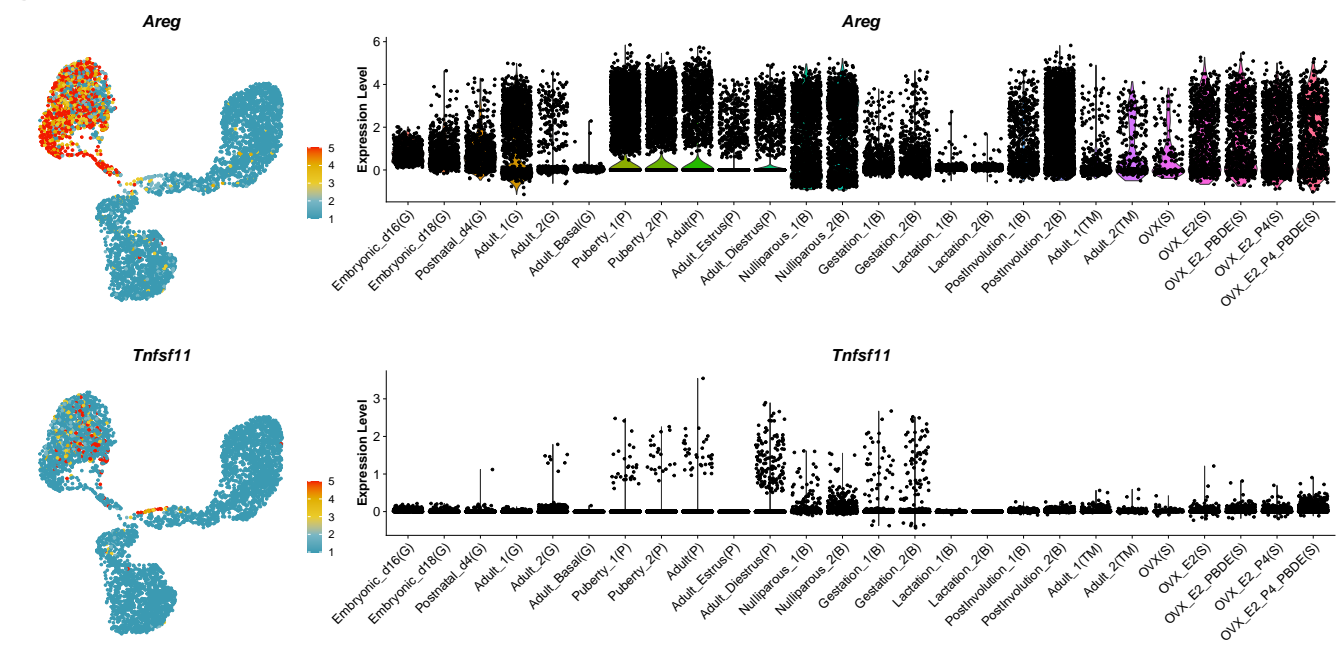
**Supplementary Figure 18 Breakdown of the TCGA data by additional features.** (a) The ternary plots of the TCGA data for the three differentiation gene sets. Metaplastic carcinomas and Claudin-low tumors were highlighted in red. (b) The comparisons of the scores for the four lineage-specific gene sets between metaplastic carcinomas (n=14) and the other histologies (n=1,075) (biologically independent samples). (c) The TNBCtype-4 subtypes. Distribution of the subtypes were overlaid on the ternary plot and the scGSVA scores were broken down by subtypes. The n number of each subtype is as follows: BL1 (n=30), BL2 (n=20), LAR (n=17), and M (n=28) (biologically independent samples). (d) The *BRCA* status. The *BRCA* status were overlaid on the ternary plot and the scGSVA scores were broken down by status. The n number of each group is as follows: BRCA1 (n=42), BRCA2 (n=37), and Others (n=767) (biologically independent samples). (e) The relationship of scGSVA scores with age at diagnosis. Age at diagnosis was overlaid on the ternary plot and correlation analysis was performed with each scGSVA score.

**a****BRCA (scRNAseq)****BRCA (scRNAseq)****b****Stem gene set****Basal gene set****Alv gene set****Hor gene set****c****BRCA(scRNAseq)**

**Supplementary Figure 19 scRNAseq data of human breast cancer.** (a) UMAP plots of the human breast cancer scRNAseq data colored by patients and subtypes. (b) Distribution of the scGSVA scores for the four curated lineage-specific gene sets in each patient. The cell number of each patient is as follows: BC01 (n=20), BC02 (n=53), BC03 (n=25), BC04 (n=47), BC05 (n=75), BC06 (n=8), BC07 (n=52), BC08 (n=15), BC10 (n=11), and BC11 (n=11) (biologically independent samples). (c) The ternary plot of the human breast cancer scRNAseq data for the three differentiation gene sets, color-coded by patients.

**a****b**

**Supplementary Figure 20 Reflection of the PAM50 molecular subtyping on the scRNAseq data of the mammary epithelium.** The heatmaps visualized the expression of the genes used for the PAM50 classification grouped by clusters, subtypes, and datasets. **(a)** Mouse normal mammary epithelium. **(b)** Human normal breast epithelium.

**a****b****c****d****e****f****g**

**Supplementary Figure 21 Reorganization of the gland during development, pregnancy, menopausal HRT with exposure to PBDEs, and estrus cycle and its implication for the risk of specific breast cancer subtypes.** (a) The results of the PAM50 subtyping of each cluster in the normal human and mouse mammary epithelium. (b) Distribution of PAM50 classification at different life stages. (c) Distribution of PAM50 classification during pregnancy. (d) The reorganization of the mammary gland by the HRTs and the PBDEs exposure in the surgically menopausal mice. The cells from each sample were highlighted in red on the UMAP plots. Five thousands cells were sampled from the entire data for background UMAP visualization. (e) The stacked bar plot represents the proportion of the epithelial cells in the entire gland, color-coded by the PAM50 classification. (f) The mammary gland reorganization during the estrus cycle. The cells from each sample were highlighted on the UMAP dimensionality. The proportion of each cluster and PAM50 subtype in the mammary epithelium is summarized in the stacked bar plots. Five thousands cells were sampled from the entire data for background UMAP visualization. (g) Expression of the soluble RNA messengers in the mammary gland and their fluctuations. The expression of amphiregulin (*Areg*) and RANKL (*Tnfsf11*) were visualized on UMAP plots and violin plots, grouped by samples. The n number of each sample is summarized in Supplementary Data 2. HRT; hormone replacement therapy. HRT: hormone replacement therapy.



## Supplementary References

1. Anaconda Software Distribution. Anaconda. <https://anaconda.com> (2020).
2. Huber, W. *et al.* Orchestrating high-throughput genomic analysis with Bioconductor. *Nat. Methods* 12, 115–121 (2015).
3. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping Identifiers for the Integration of Genomic Datasets with the R/Bioconductor package biomaRt. *Nat. Protoc.* 4, 1184–1191 (2009).
4. Yu, G., Wang, L.-G., Han, Y. & He, Q.-Y. clusterProfiler: an R Package for Comparing Biological Themes Among Gene Clusters. *OMICS* 16, 284–287 (2012).
5. Wilke, C. cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'. R package version 1.0.0. <https://CRAN.R-project.org/package=cowplot> (2019).
6. Gulati, G. S. *et al.* Single-cell transcriptional diversity is a hallmark of developmental potential. *Science* 367, 405–411 (2020).
7. McGinnis, C. S., Murrow, L. M. & Gartner, Z. J. DoubletFinder: Doublet Detection in Single-Cell RNA Sequencing Data Using Artificial Nearest Neighbors. *Cell Syst.* 8, 329–337.e4; 10.1016/j.cels.2019.03.003 (2019).
8. Wickham, H., François, R., Henry, L. & Müller, K. dplyr: A Grammar of Data Manipulation. R package version 0.8.3. <https://CRAN.R-project.org/package=dplyr> (2019).
9. Lê, S., Josse, J. & Husson, F. FactoMineR: An R Package for Multivariate Analysis. *J. Stat. Softw.* 25, 1–18 (2008).
10. Kassambara, A. & Mundt, F. factoextra: Extract and Visualize the Results of Multivariate Data Analyses. R package version 1.0.5. <https://CRAN.R-project.org/package=factoextra> (2017).

11. Gendoo, D. M. A. *et al.* Genefu: an R/Bioconductor package for computation of gene expression-based signatures in breast cancer. *Bioinformatics* 32, 1097–1099 (2016).
12. Kassambara, A. ggcorrplot: Visualization of a Correlation Matrix using 'ggplot2'. R package version 0.1.3. <https://CRAN.R-project.org/package=ggcorrplot> (2019).
13. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag, 2009).
14. Slowikowski, K. ggrepel: Automatically Position Non-Overlapping Text Labels with 'ggplot2'. R package version 0.8.1. <https://CRAN.R-project.org/package=ggrepel> (2019)..
15. Kassambara, A. ggpubr: 'ggplot2' Based Publication Ready Plots. R package version 0.2.3. <https://CRAN.R-project.org/package=ggpubr> (2019).
16. Hamilton, N. E. & Ferry, M. ggtern: Ternary Diagrams Using ggplot2. *J. Stat. Softw.* 87, 1–17 (2018).
17. Hänzelmann, S., Castelo, R. & Guinney, J. GSVA: gene set variation analysis for microarray and RNA-Seq data. *BMC Bioinformatics* 14, 7; 10.1186/1471-2105-14-7 (2013).
18. Thomas, K. *et al.* Jupyter Notebooks - a publishing format for reproducible computational workflows in *Positioning and Power in Academic Publishing: Players, Agents and Agendas* (eds. Loizides, F. and Schmidt, B.) 87–90 (IOS Press, 2016).
19. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat. Methods* 16, 1289–1296 (2019).
20. Liu, J. *et al.* Jointly defining cell types from multiple single-cell datasets using LIGER. *Nat. Protoc.* 15, 3632–3662 (2020).
21. Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* 9, 90–95 (2007).
22. Dolgalev, I. msigdb: MSigDB Gene Sets for Multiple Organisms in a Tidy Data Format. R package version 7.0.1. <https://CRAN.R-project.org/package=msigdb> (2019).

23. Oliphant, T. E. *Guide to NumPy: 2nd Edition*. (Continuum Press, 2015).
24. McKinney, W. Data Structures for Statistical Computing in Python. In: Proceedings of the 9th Python in Science Conference, Volume 445, 2010.
25. Kolde, R. pheatmap: Pretty Heatmaps. R package version 1.0.12. <https://CRAN.R-project.org/package=pheatmap> (2019).
26. Sievert, C. *Interactive Web-Based Data Visualization with R, plotly, and shiny*. (CRC Press, 2020).
27. Rossum, G. V. & Drake, F. L. *Python 3 Reference Manual*. (CreateSpace Independent Publishing Platform, 2009).
28. R Core Team. R: A language and environment for statistical computing. <http://www.R-project.org/> (2013).
29. Wickham, H. Reshaping Data with the reshape Package. *J. Stat. Softw.* 21, 1–20 (2007).
30. RStudio Team. RStudio: Integrated Development for R. <http://www.rstudio.com/> (2020).
31. Johansen, N. & Quon, G. scAlign: a tool for alignment, integration, and rare cell identification from scRNA-seq data. *Genome Biol.* 20, 166; 10.1186/s13059-019-1766-4 (2019).
32. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* 177, 1888-1902.e21 (2019).
33. Chen, H. *et al.* Single-cell trajectories reconstruction, exploration and mapping of omics data with STREAM. *Nat. Commun.* 10, 1903; 10.1038/s41467-019-09670-4 (2019).
34. Wickham, H. stringr: Simple, Consistent Wrappers for Common String Operations. R package version 1.4.0. <https://CRAN.R-project.org/package=stringr> (2019).

35. Morgan M., Obenchain, V., Hester, J. & Pagès, H. SummarizedExperiment: SummarizedExperiment container. R package version 1.14.1.  
<https://bioconductor.org/packages/SummarizedExperiment> (2019).
36. Colaprico, A. *et al.* TCGAAbiolinks: an R/Bioconductor package for integrative analysis of TCGA data. *Nucleic Acids Res.* **44**, e71; 10.1093/nar/gkv1507 (2016).
37. Santamaría, L. P. & del Peso, L. TFEA. ChIP: Analyze Transcription Factor Enrichment. R package version 1.4.2. <https://bioconductor.org/packages/release/bioc/html/TFEA.ChIP.html> (2019).
38. Konopka, T. umap: Uniform Manifold Approximation and Projection. R package version 0.2.5.0. <https://CRAN.R-project.org/package=umap> (2020).
39. Garnier, S. viridis: Default Color Maps from 'matplotlib'. R package version 0.5.1.  
<https://CRAN.R-project.org/package=viridis> (2018).
40. Ram, K. & Wickham, H. wesanderson: A Wes Anderson Palette Generator. R package version 0.3.6. <https://CRAN.R-project.org/package=wesanderson> (2018).