# Supplementary Information for

## scMerge leverages factor analysis, stable expression, and pseudoreplication to merge multiple single-cell RNA-seq datasets

**Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang**

**Jean Yee Hwa Yang, Pengyi Yang**
**E-mail: J.Y.H.Y. jean.yang@sydney.edu.au and P.Y. pengyi.yang@sydney.edu.au**

**This PDF file includes:**

      Figs. S1 to S17
      Tables S1 to S2
      Captions for Databases S1 to S2

**Other supplementary materials for this manuscript include the following:**

      Databases S1 to S2
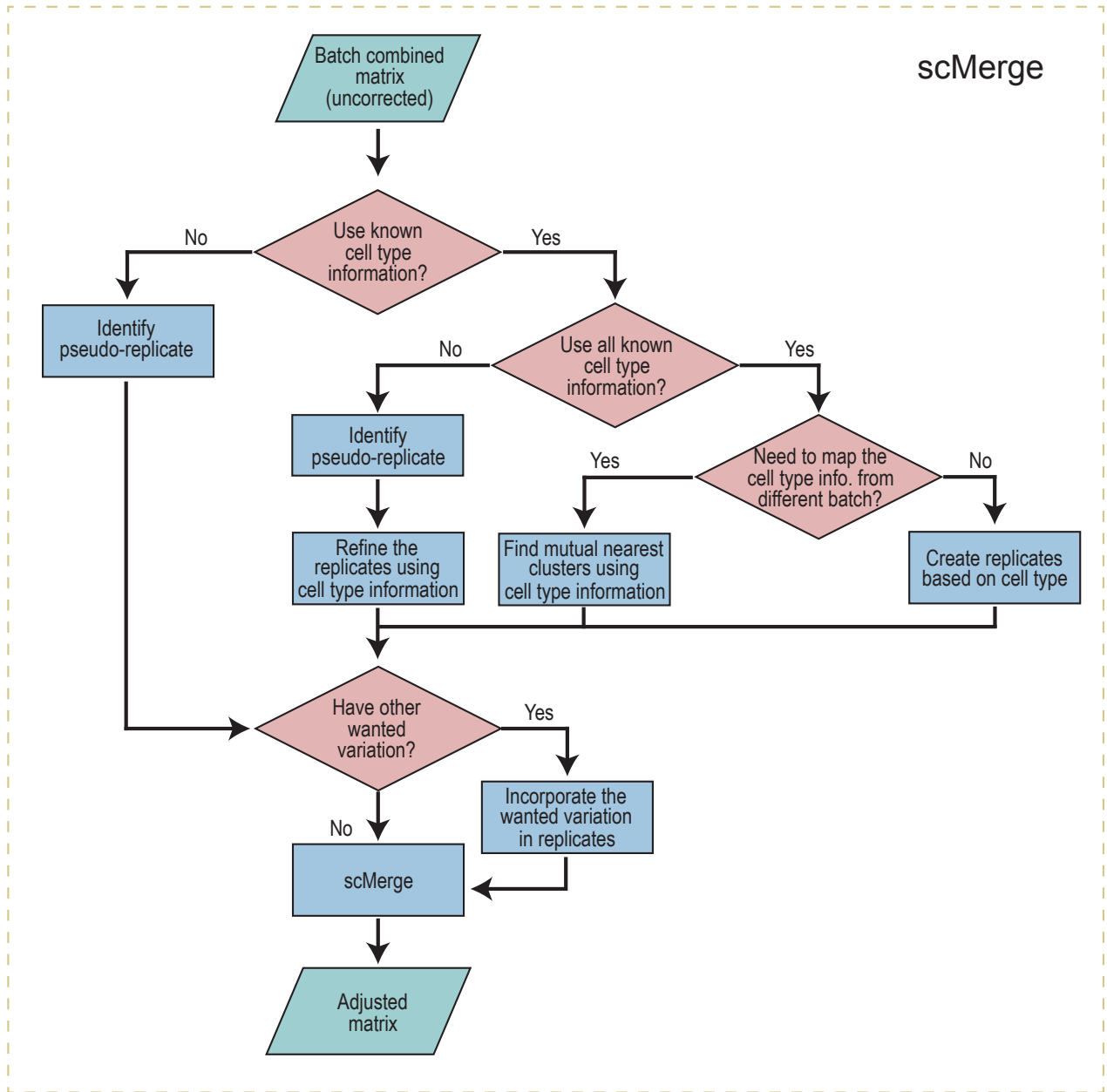
scMerge

Batch combined
matrix
(uncorrected)

Use known
cell type
information?

No

Yes

Identify
pseudo-replicate

Use all known
cell type
information?

No

Yes

Identify
pseudo-replicate

Need to map the
cell type info. from
different batch?

Yes

No

Refine the
replicates using
cell type information

Find mutual nearest
clusters using
cell type information

Create replicates
based on cell type

Have other
wanted
variation?

Yes

No

Incorporate the
wanted variation
in replicates

scMerge

Adjusted
matrix

**Fig. S1.** Flow chart illustrating the decision-making process associated with scMerge algorithm.

Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang
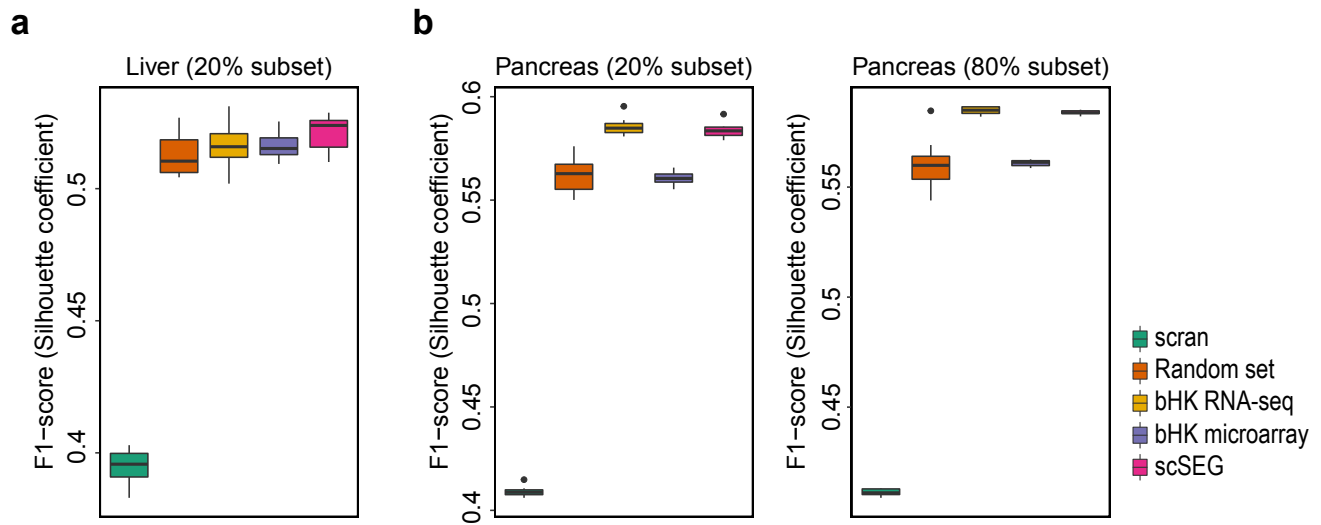
**Fig. S2.** A 1 by 3 panel of boxplots comparing the effect of different types of negative controls for (a) liver datasets and (b) four Pancreas datasets. The y-axis represents the F1score of Silhouette coefficients between cell type mixing and (1 - datasets mixing). Stratified sampling is performed to randomly subset 20% and 80% of cells from the datasets. This procedure is repeated 10 times. The boxplots represent the F1 score results using logcounts, RUVg using random subset of genes, bulk microarray, RNA-Seq data, and scSEG as negative control genes based on subset of the datasets.
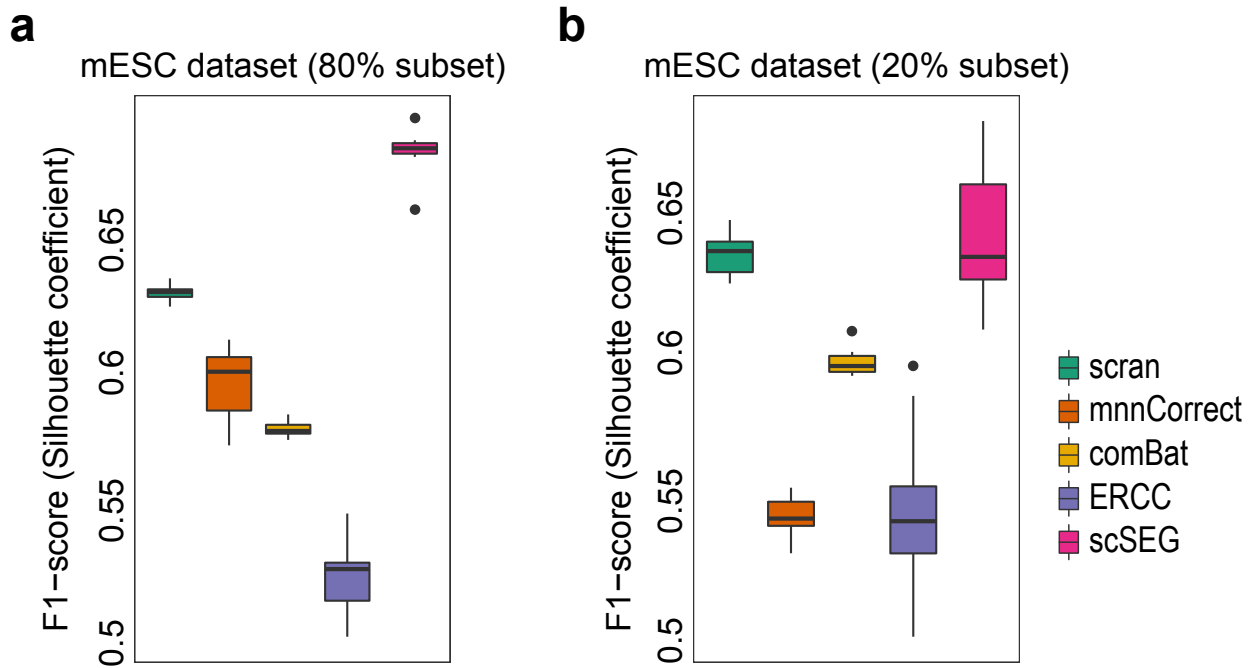
Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**a**

## mESC dataset (80% subset)

**b**

## mESC dataset (20% subset)

Legend:
- scran
- mnnCorrect
- comBat
- ERCC
- scSEG

**Fig. S3.** A 1 by 2 panel of boxplots showing the effect of using ERCC and scSEG as negative controls with scMerge. Stratified subsampling is performed for the mESC dataset. In each stratified subsampling, we randomly selected 20% (left panel) or 80% (right panel) of the cells from the dataset and perform scMerge with ERCC spike-ins genes and mouse scSEG, and comparing the results with ComBat and mnnCorrect (default settings). This procedure is repeated 10 times.
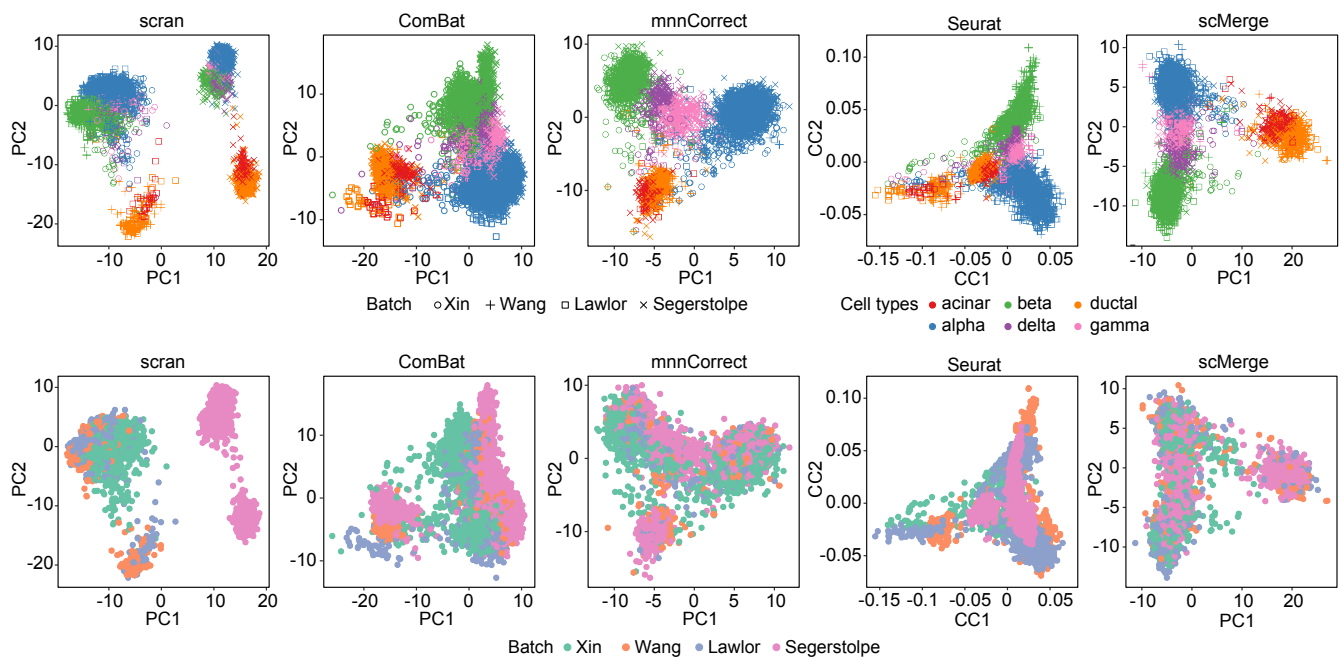
Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Fig. S4.** A 2 by 5 panel of PCA plots of the four Pancreas datasets using the output from scran (logcounts), ComBat, mnnCorrect, Seurat, and scMerge (using scSEG as negative controls). The top row of the panel is color coded by cell types and the second row is color coded by the four different datasets.
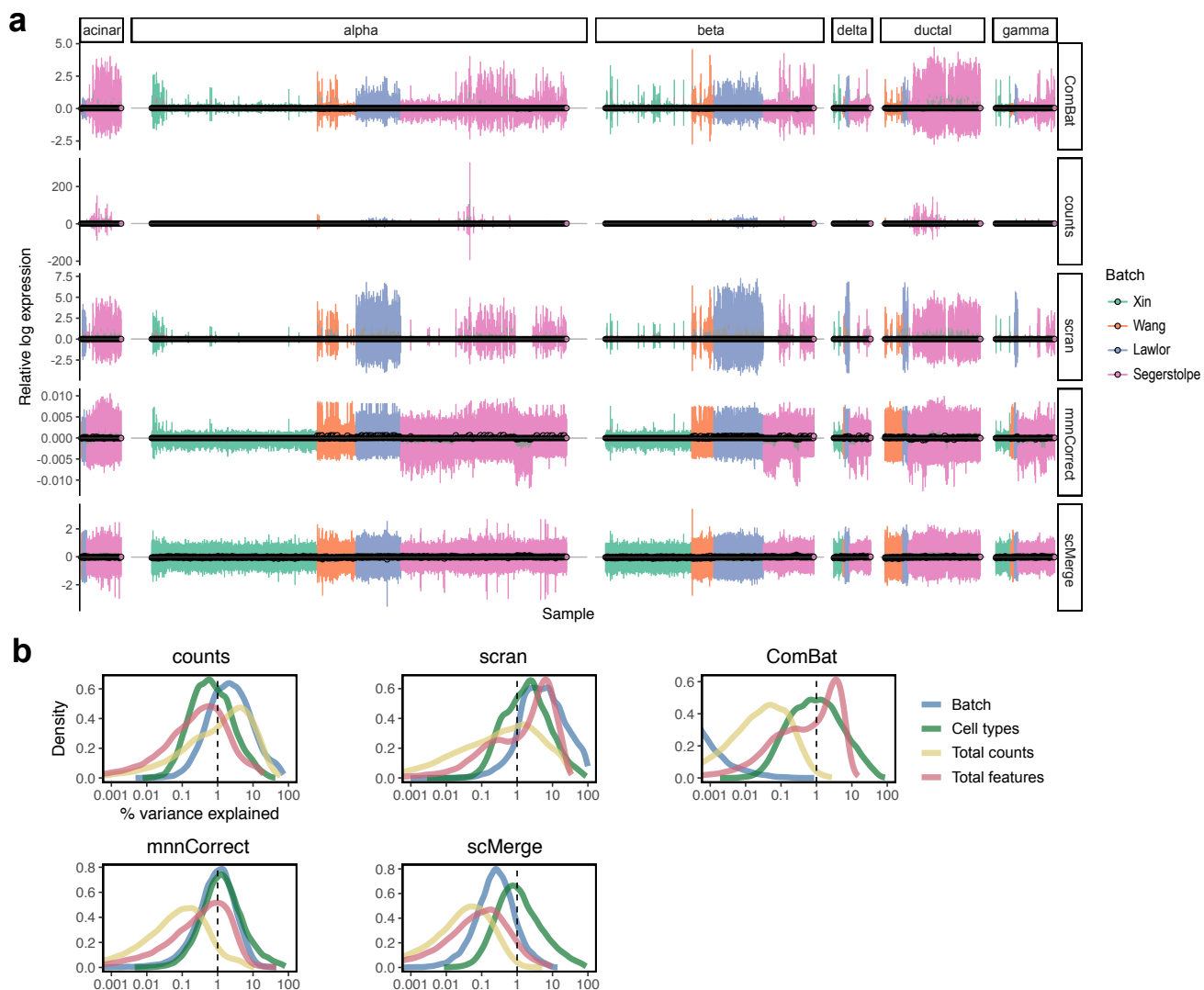
Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Fig. S5.** Diagnostic plots from the Pancreas data collections ("Pancreas4"). A. RLE plots. The boxplot for each cell from the same cell type between different batches of scMerge shares similar inter quantile ranges. B. Percentage of variance explained for each variable. scMerge has cellType explaining the highest percentage of variation for the liver data collection.
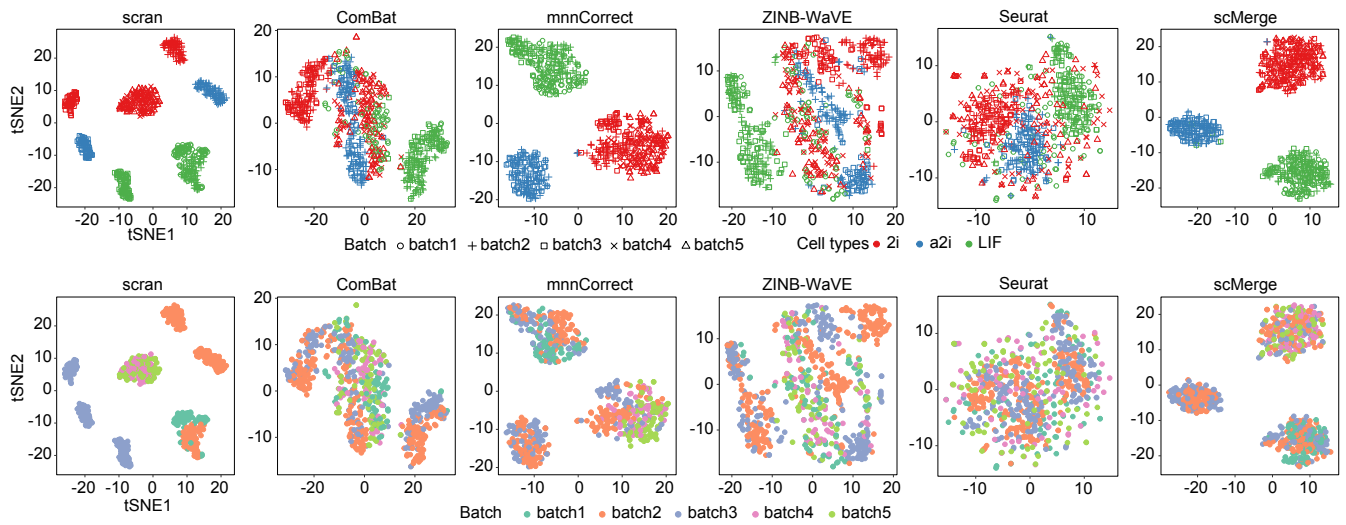
Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Fig. S6.** A 2 by 6 panel of tSNE plots of the mESC datasets using the output from scran (logcounts), ComBat, mnnCorrect, ZINB-WaVE, Seurat, and scMerge (using scSEG as negative controls). The top row of the panel is color coded by three cell types and the second row is color coded by the three batches.
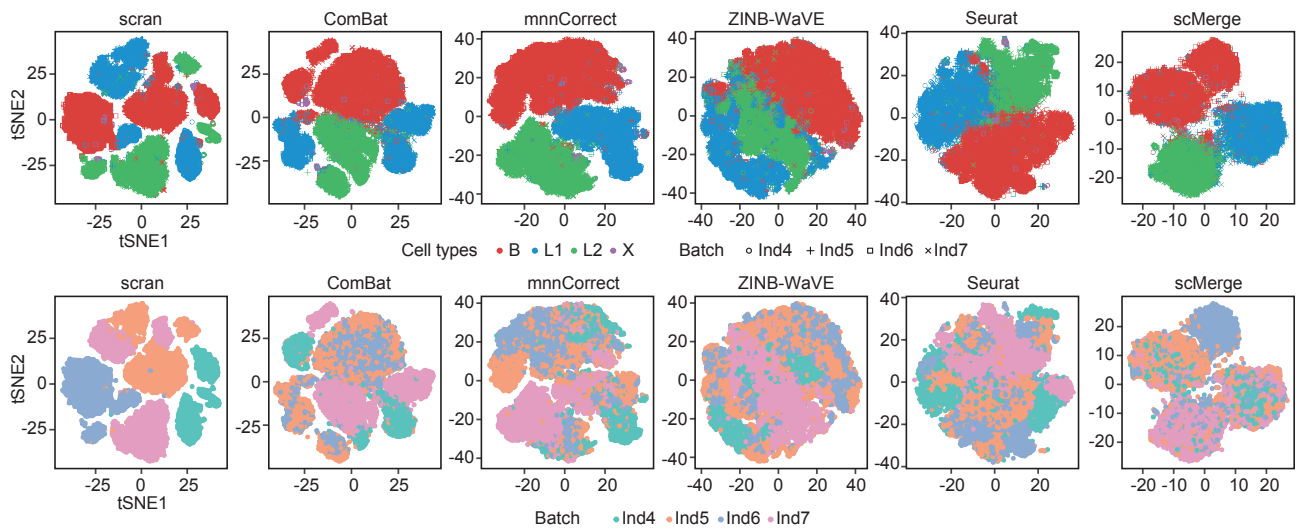
**Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang**

**Fig. S7.** A 2 by 6 panel of tSNE plots of the Breast Cancer data using the output from scran (logcounts), ComBat, mnnCorrect, ZINB-WaVE, Seurat, and scMerge (using scSEG as negative controls). The top row of the panel is color coded by cell types and the second row is color coded by the four individuals.
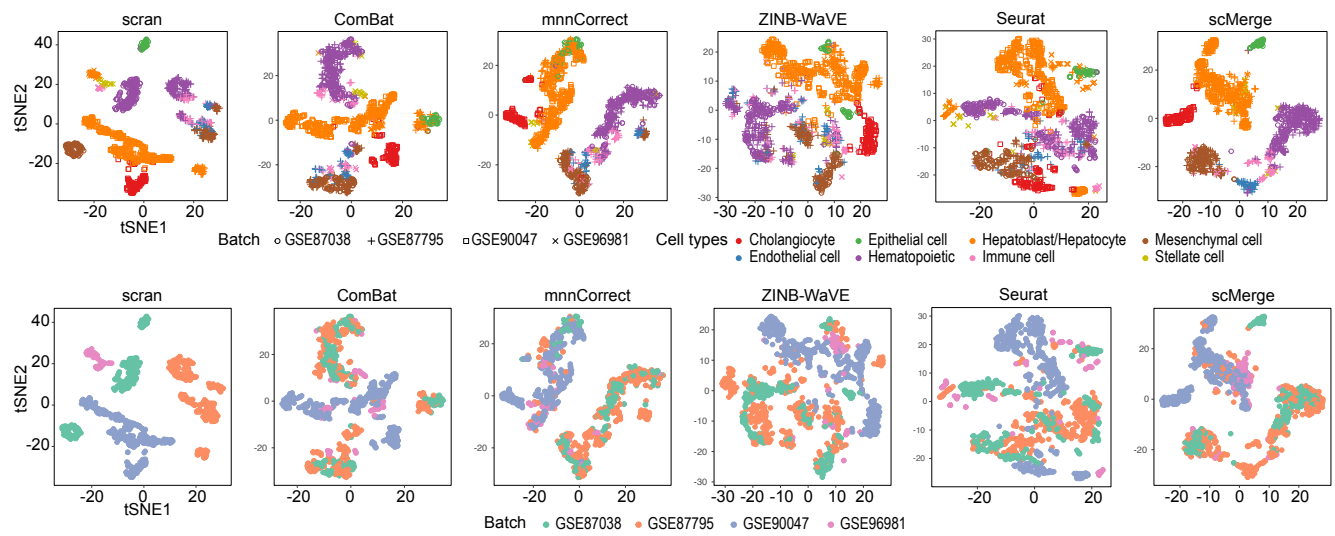
Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Fig. S8.** A 2 by 6 panel of tSNE plots of the Liver data collection using the output from scran (logcounts), ComBat, mnnCorrect, ZINB-WaVE, Seurat, and scMerge (using scSEG as negative controls). The top row of the panel is color coded by cell types and the second row is color coded by the four datasets.

Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T.
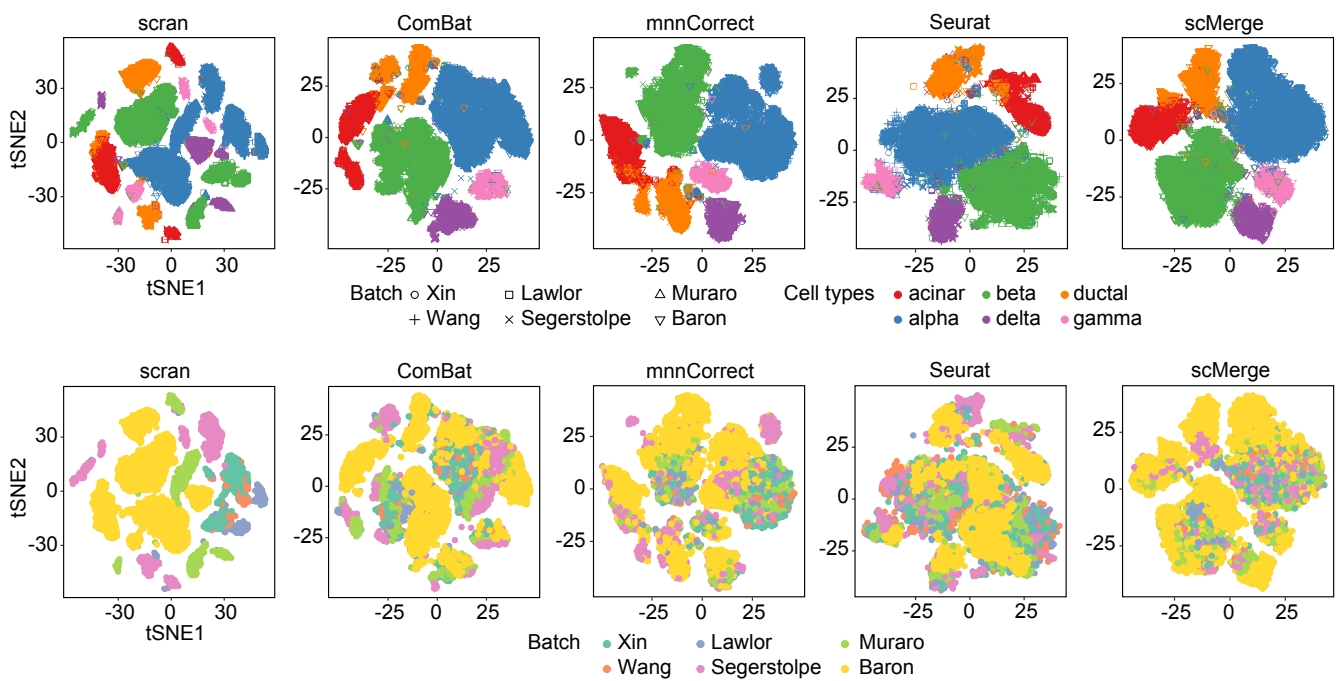Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Fig. S9.** A 2 by 6 panel of tSNE plots of the Olfactory data collection using the output from scran (logcounts), ComBat, mnnCorrect, ZINB-WaVE, Seurat, and scMerge (using scSEG as negative controls). The top row of the panel is color coded by cell types and the second row is color coded by the two datasets.
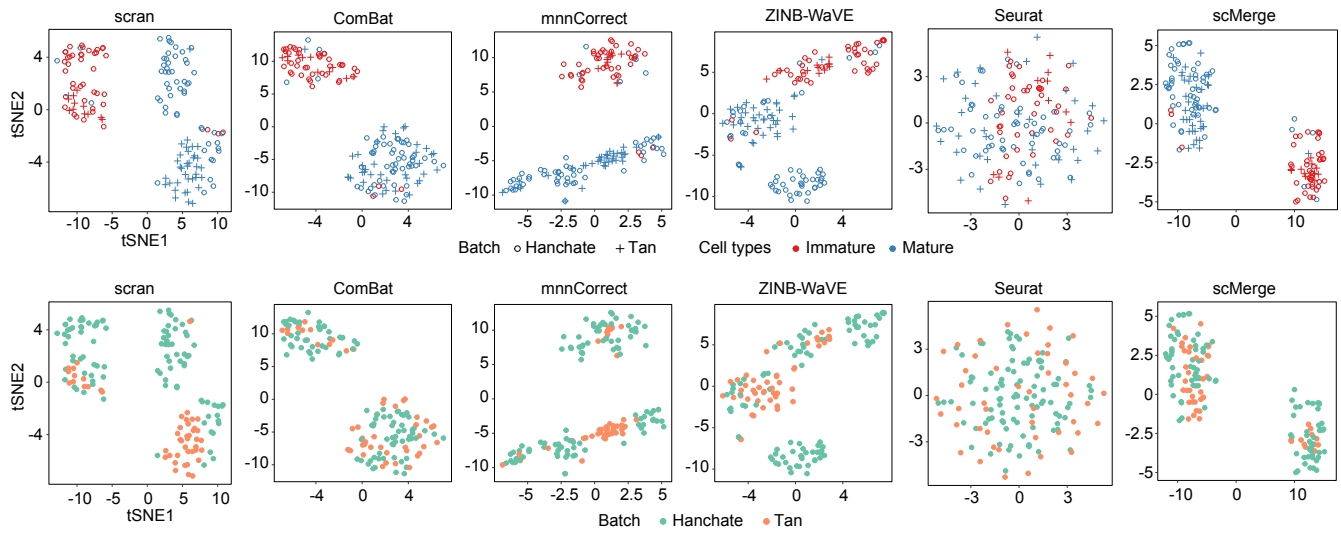
Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Fig. S10.** A 2 by 5 panel of tSNE plots of all six Pancreas related datasets based on the output from scran (logcounts), ComBat, mnnCorrect, Seurat, and scMerge (using scSEG as negative controls). The top row of the panel is color coded by cell types and the second row is color coded by the six different datasets.
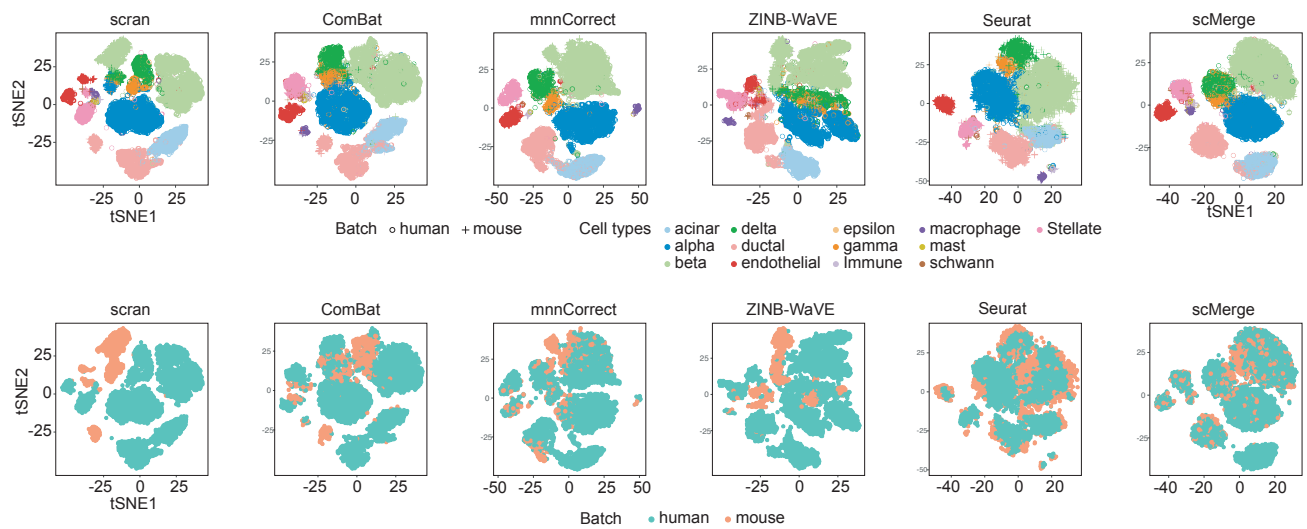
**Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang**

**Fig. S11.** A 2 by 6 panel of tSNE plots of the Pancreas Islet data collection using the output from scran (logcounts), ComBat, mnnCorrect, ZINB-WaVE, Seurat, and scMerge (using scSEG as negative controls). The top row of the panel color coded by cell types and the second row is color coded by the six datasets.
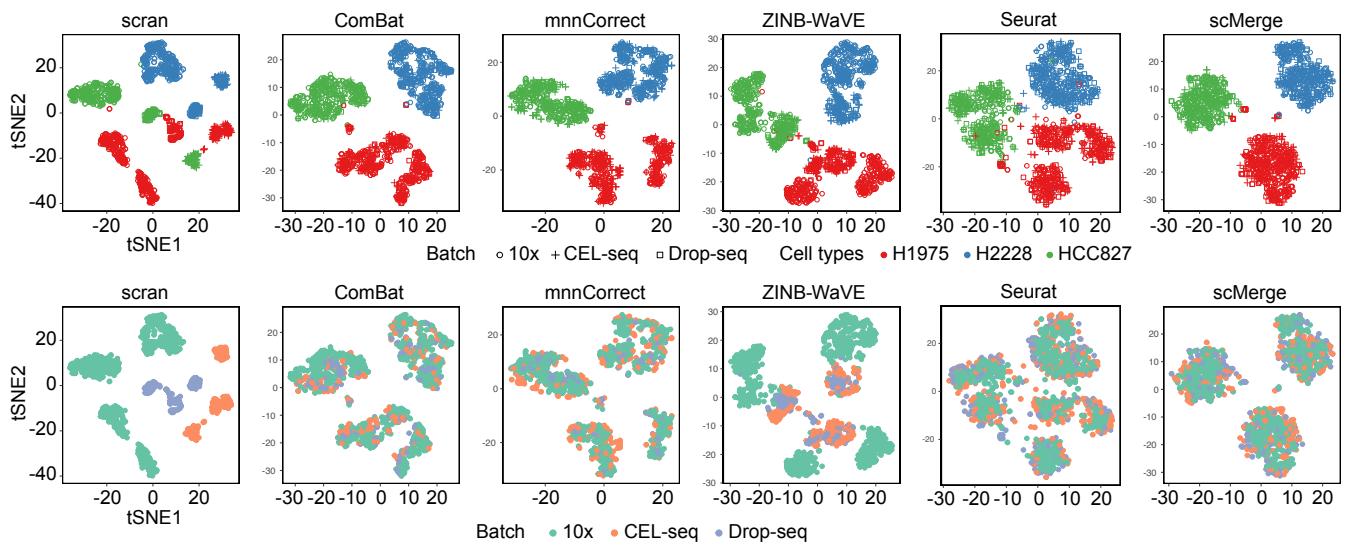
Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Fig. S12.** A 2 by 6 panel of tSNE plots of the CellBench data using the output from scran (logcounts), ComBat, mnnCorrect, ZINB-WaVE, Seurat, and scMerge (using scSEG as negative controls). The top row of the panel is color coded by cell lines and the second row is color coded by the three different platforms.
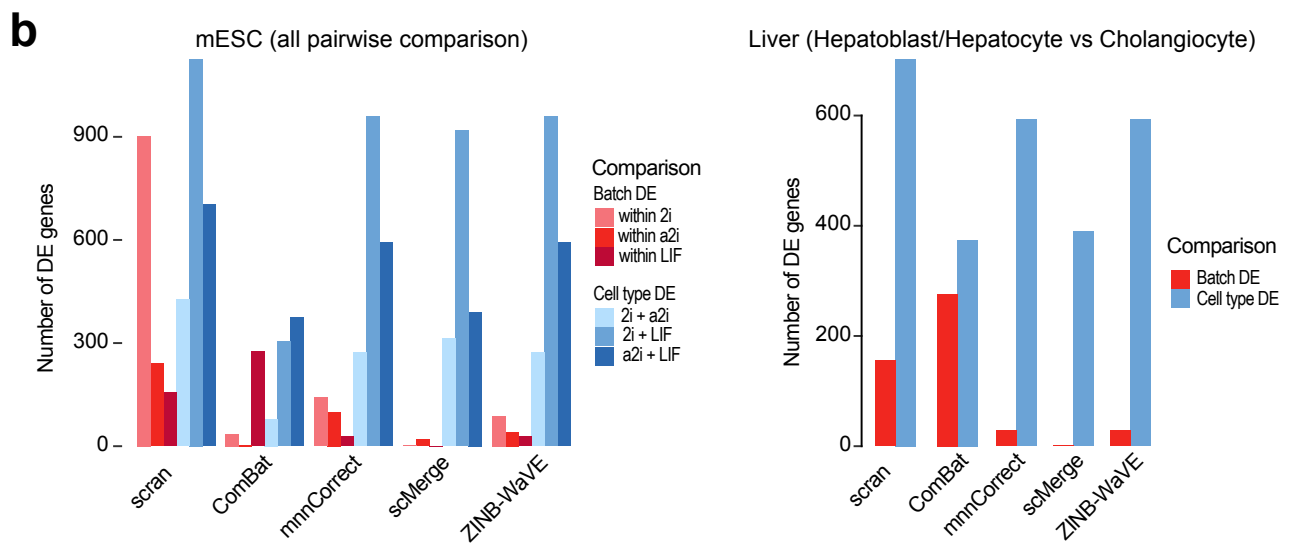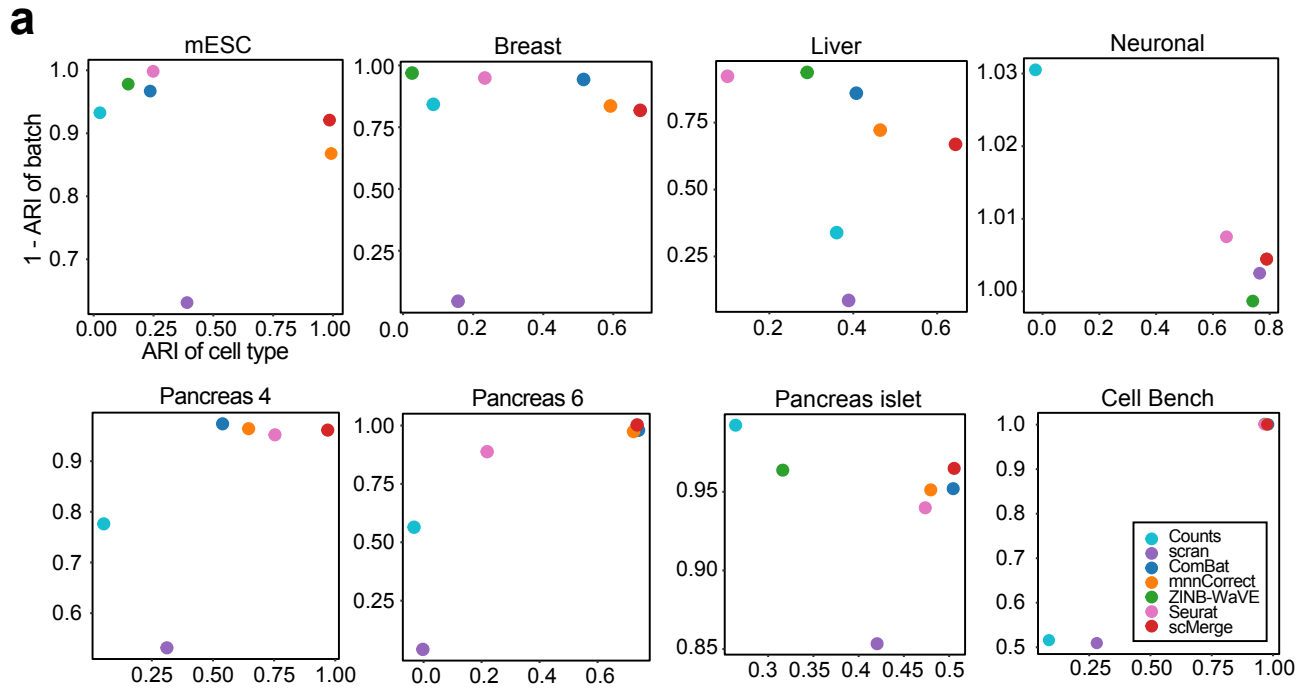
Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Fig. S13.** (a) A 2 by 4 panel of scatter plots of ARI evaluation for no normalization, scran (logcounts), ZINB-WaVE, ComBat, mnnCorrect, Seurat, and scMerge (using scSEG as negative controls) of eight dataset collections (mESC dataset, breast ). x-axes denote the ARI of cell types and y-axes denote the 1 − ARI of batch effects, where desirable outcomes are in the top-right hand corner. (b) A 1 by 2 panel of bar plots of differential expression results of mESC dataset (left panel) and liver dataset collection (right panel) on scran, ComBat, mnnCorrect, scMerge and ZINB-WaVE normalized data. The y-axis indicates the number of DE genes that are selected (adjusted p-value < 0.05 and logFC > 2). Batch DE are selected by performing the DE analysis within the same cell type, while cell type DE are selected by performing the DE analysis between two different cell types.

Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Fig. S14.** A 3 by 4 panel of pseudotime trajectory from Monocle 2 with perturbed data demonstrating the stability of scMerge with output from ComBat, mnnCorrect, ZINB-WaVE and scMerge. The first row of the panel color coded by three datasets; the second row is color coded by the cell types and the third row is color coded by Monocle 2 states.

Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T.
Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

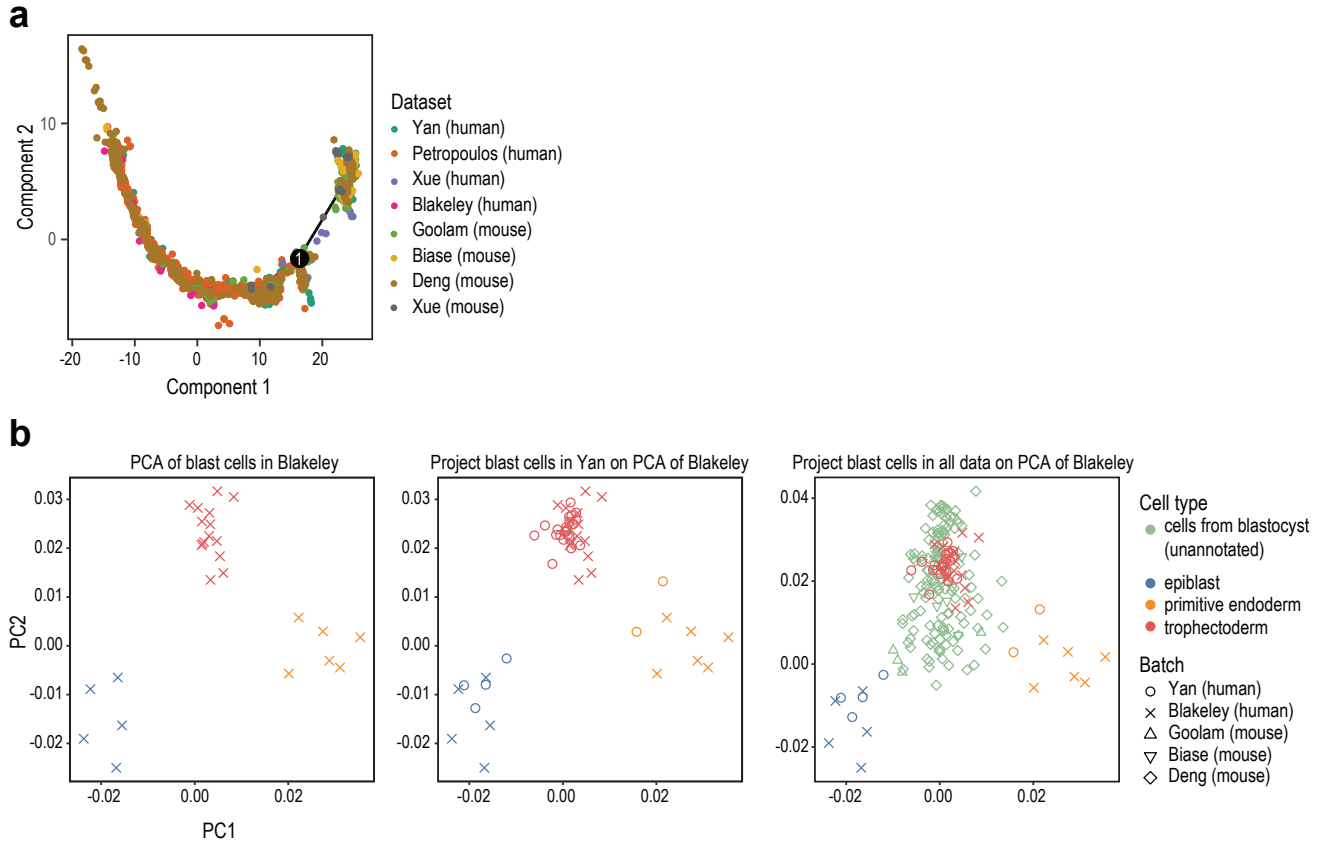**Fig. S15.** (a) A pseudotime trajectory from Monocle 2 with all cells from ESC dataset collections, color coded by dataset. (b) A 1 by 3 panel of PCA plots of blast cells from ESC developmental data collections demonstrating scMerge effectively reproduce their cell type results, colored by the cell types identified from the original paper. The left panel is the PCA plot of blast cells from Blakeley et al. The medium panel is the PCA plot of blast cells from Blakeley et al. where blast cells from Yan et al. are projected. The right panel is the PCA plot of blast cells from Blakeley et al. where blast cells from mouse ESC datasets (Biase et al., Goolam et al. and Deng et al.) are projected.

Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**a** Computation time for different methods
On a scRNA–Seq data with 23699 features and 4566 cells

**b** Gene−wise correlation between RSVD approximation and full SVD results
On a scRNA–Seq data with 23699 features and 4566 cells

**Fig. S16.** (a) Computational time for different rsvd parameters with the Pancreas data collections ("Pancreas4"). This includes 23,699 features and 4566 cells. (b) Gene wise association based on Pearson correlation between rsvd approximation and full SVD calculation.
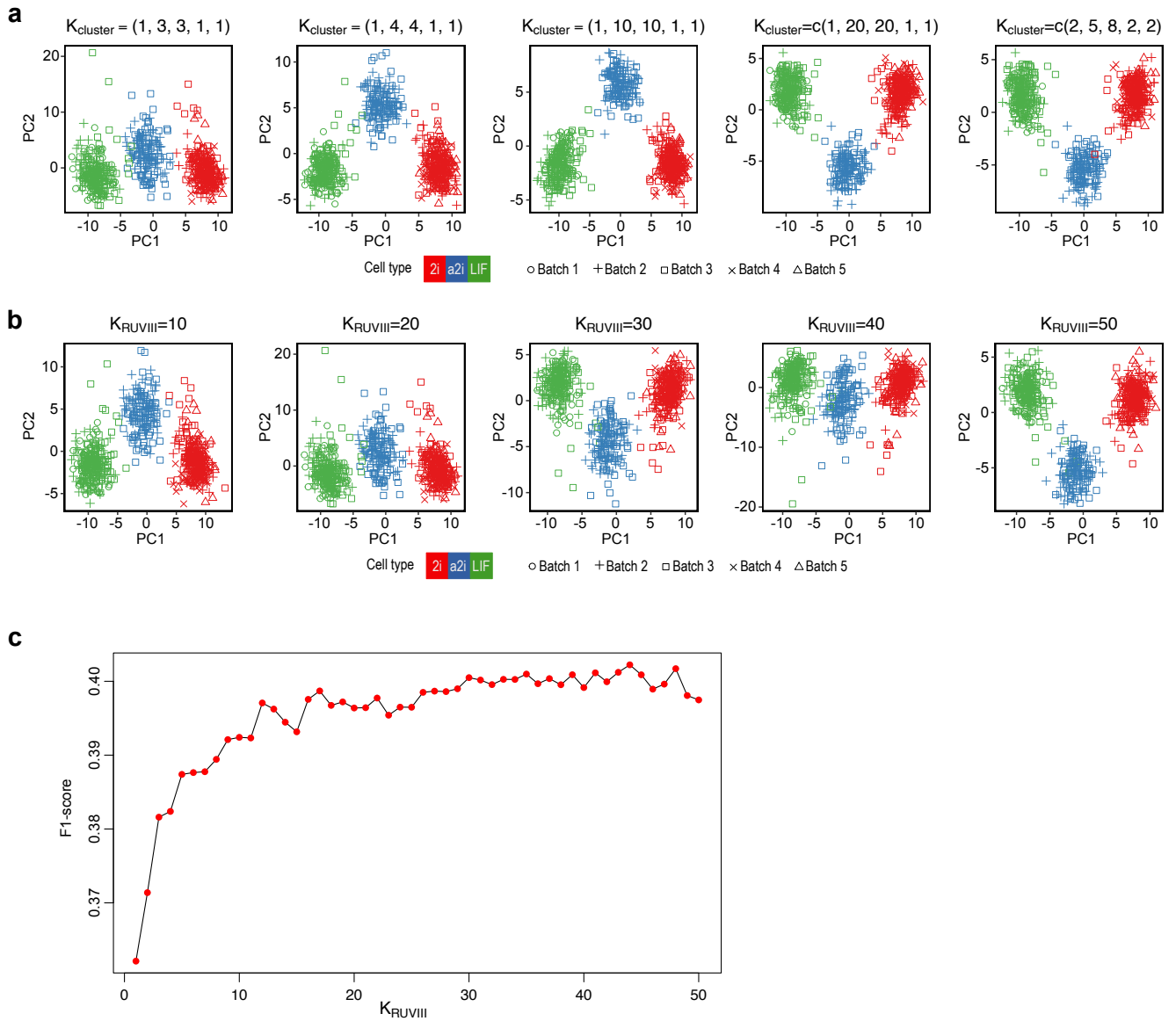
**Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang**

**Fig. S17.** (a) A 1 by 4 panel of PCA plots of mESC datasets using different $k_{Cluster}$ settings. (b) A 1 by 4 panel of PCA plots of mESC datasets using different $k_{RUVIII}$ settings ($k_{RUVIII}$ = 10, 20, 30, 40, 50) (c) A scatter plot showing the performance of scMerge using different $k_{RUVIII}$ from 1 to 50, evaluated by F1 score of Silhouette coefficient. The y-axis represents the F1-score, while the x-axis represents the $k_{RUVIII}$ values.

18 of 21 Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Table S1. More detailed summary of datasets and data collections used in this study.**

**Supplementary Table 1**

| Type of merge | | | Name | ID | Author | DOI or URL | Protocol | Organism | Tissue | # of cell types | # of cells | # of batches |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Within experiment | | | mESC | E-MTAB-2600 | Kolodziejczyk | 10.1016/j.stem.2015.09.011 | SMARTer/C1 | Mouse | Mouse ESC | 3 | 704 | 5 |
| | | | Breast | GSE113197 | Nguyen | 10.1038/s41467-018-04334-1 | 10x Chromium | Human | Breast cancer | 3 | 24520 | 4 |
| Across platforms with significant depth difference | Across data experiments | | Liver | GSE87795 | Su | 10.1186/s12864-017-4342-x | SMARTer/C1 | Mouse | Liver | 8 | 1236 | NA |
| | | | | GSE90047 | Yang | 10.1002/hep.29353 | Smart-Seq2 | | | | | |
| | | | | GSE87038 | Dong | 10.1186/s13059-018-1416-2 | STRT-seq | | | | | |
| | | | | GSE96981 | Camp | 10.1038/nature22796 | SMARTer/C1 | | | | | |
| | | | Neuronal | SRP065920 | Tan | 10.15252/msb.20156639 | Smart-Seq2 | Mouse | Neuronal | 2 | 145 | NA |
| | | | | GSE75413 | Hanchate | 10.1126/science.aad2456 | STRT-seq | | | | | |
| | | Pancreas6 | Pancreas4 | GSE81608 | Xin | 10.1016/j.cmet.2016.08.018 | SMARTer/C1 | Human | Pancreas | 6 | 4566 | NA |
| | | | | GSE83139 | Wang | 10.2337/db16-0405 | Smart-Seq | | | | | |
| | | | | GSE86469 | Lawlor | 10.1101/gr.212720.116 | SMARTer/C1 | | | | | |
| | | | | E-MTAB-5061 | Segerstolpe | 10.1016/j.cmet.2016.08.020 | Smart-Seq2 | | | | | |
| | | | | GSE85241 | Muraro | 10.1016/j.cels.2016.09.002 | Cel-seq2 | Human | Pancreas | 6 | 1773 | NA |
| | | | | GSE84133 | Baron | 10.1016/j.cels.2016.08.011 | inDrop | Human+ mouse | Pancreas Islets | 13 | 8569 | 2 (human & mouse) |
| | | | CellBench | cellBench | | https://github.com/LuyiTian/CellBench_data | Cel-seq2, Drop-seq, 10x Chromium | Human | Adenocarcinoma cell lines | 3 | 1401 | 3 per cell types |
| | Across organisms | | ESC | GSE45719 | Deng | 10.1126/science.1245316 | Smart-Seq | Mouse | ESC | 10+ | 2144 | NA |
| | | | | GSE57249 | Biase | 10.1101/gr.177725.114 | SMARTer | | | | | |
| | | | | E-MTAB-3321 | Goolam | 10.1016/j.cell.2016.01.047 | Smart-Seq2 | | | | | |
| | | | | GSE44183 | Xue | 10.1038/nature12364 | Tang et al., 2010* | Human + mouse | | | | |
| | | | | E-MTAB-3929 | Petropoulos | 10.1016/j.cell.2016.03.023 | Smart-Seq2 | | | | | |
| | | | | GEO66507 | Blakeley | 10.1242/dev.123547 | SMARTer | human | | | | |
| | | | | GSE36552 | Yan | 10.1038/nsmb.2660 | Than et al., 2010* | | | | | |

Within experiment
Across data experiments
Across platforms with significant depth difference
Across organisms
* Tang et al, Nature Protocol, 2010 (DOI: 10.1038/nprot.2009.236)

**Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang**

**Table S2.** $F1_{sil}$ and $F1_{ARI}$ of all methods (counts, scran (logcounts), ComBat, mnnCorrect, ZINB-WaVE, Seurat, and scMerge) across all datasets.

## Supplementary Table 2

| Silhouette Coefficient | counts | scran | ComBat | mnnCorrect | ZINB-WaVE | Seurat | scMerge |
|---|---|---|---|---|---|---|---|
| mESC | 0.52 | 0.56 | 0.58 | 0.68 | 0.57 | 0.52 | 0.69 |
| Breast | 0.42 | 0.27 | 0.52 | 0.55 | 0.53 | 0.53 | 0.59 |
| Liver | 0.35 | 0.26 | 0.53 | 0.52 | 0.51 | 0.52 | 0.53 |
| Neuronal | 0.56 | 0.54 | 0.60 | 0.61 | 0.56 | 0.53 | 0.61 |
| Pancreas 4 | 0.44 | 0.41 | 0.56 | 0.59 | NA | 0.55 | 0.63 |
| Pacnreas 6 | 0.33 | 0.27 | 0.57 | 0.60 | NA | 0.51 | 0.61 |
| Pancreas Islet | 0.47 | 0.47 | 0.55 | 0.56 | 0.52 | 0.54 | 0.60 |
| cellBench | 0.47 | 0.39 | 0.64 | 0.63 | 0.65 | 0.54 | 0.65 |

| ARI | counts | scran | ComBat | mnnCorrect | ZINB-WaVE | Seurat | scMerge |
|---|---|---|---|---|---|---|---|
| mESC | 0.49 | 0.43 | 0.54 | 0.60 | 0.53 | 0.55 | 0.63 |
| Breast | 0.47 | 0.04 | 0.58 | 0.55 | 0.50 | 0.54 | 0.55 |
| Liver | 0.27 | 0.08 | 0.53 | 0.48 | 0.54 | 0.50 | 0.48 |
| Neuronal | 0.50 | 0.64 | 0.64 | 0.64 | 0.63 | 0.63 | 0.64 |
| Pancreas 4 | 0.45 | 0.38 | 0.60 | 0.61 | NA | 0.61 | 0.65 |
| Pacnreas 6 | 0.36 | 0.04 | 0.63 | 0.62 | NA | 0.51 | 0.64 |
| Pancreas Islet | 0.56 | 0.53 | 0.58 | 0.58 | 0.56 | 0.57 | 0.59 |
| cellBench | 0.35 | 0.36 | 0.66 | 0.66 | 0.66 | 0.66 | 0.66 |

Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang

**Additional data table S1 (SupplementaryFile1.xlsx)**

Excel spreadsheet with six sheets listing scSEG and housekeeping genes derived from bulk microarray and bulk RNA-Seq, for human and mouse each. Mouse bulk microarray and bulk RNA-Seq derived housekeeping gene lists are homologues of the human gene lists given.

**Additional data table S2 (SupplementaryFile2.xlsx)**

Excel spreadsheet with pseudotime estimation for all cells from ESC data collections.

**Yingxin Lin, Shila Ghazanfar, Kevin Y. X. Wang, Johann A. Gagnon-Bartsch, Kitty K. Lo, Xianbin Su, Ze-Guang Han, John T. Ormerod, Terence P. Speed, Pengyi Yang and Jean Yee Hwa Yang**