

scMerge: Integration of multiple single-cell transcriptomics datasets leveraging stable expression and pseudo-replication

Yingxin Lin¹, Shila Ghazanfar^{1,2†}, Kevin Y. X. Wang^{1†}, Johann A. Gagnon-Bartsch³, Kitty K. Lo¹, Xianbin Su⁵, Ze-Guang Han⁵, John T. Ormerod¹, Terence P. Speed⁴, Pengyi Yang^{1,2*}, Jean Yee Hwa Yang^{1,2*}

¹ School of Mathematics and Statistics, University of Sydney, ² Charles Perkins Centre University of Sydney, ³ Department of Statistics, University of Michigan,

⁴ Bioinformatics Division, Walter and Eliza Hall Institute of Medical Research,

⁵ Key Laboratory of Systems Biomedicine and Collaborative Innovation Center of Systems Biomedicine, Shanghai Center for Systems Biomedicine, Shanghai Jiao Tong University,

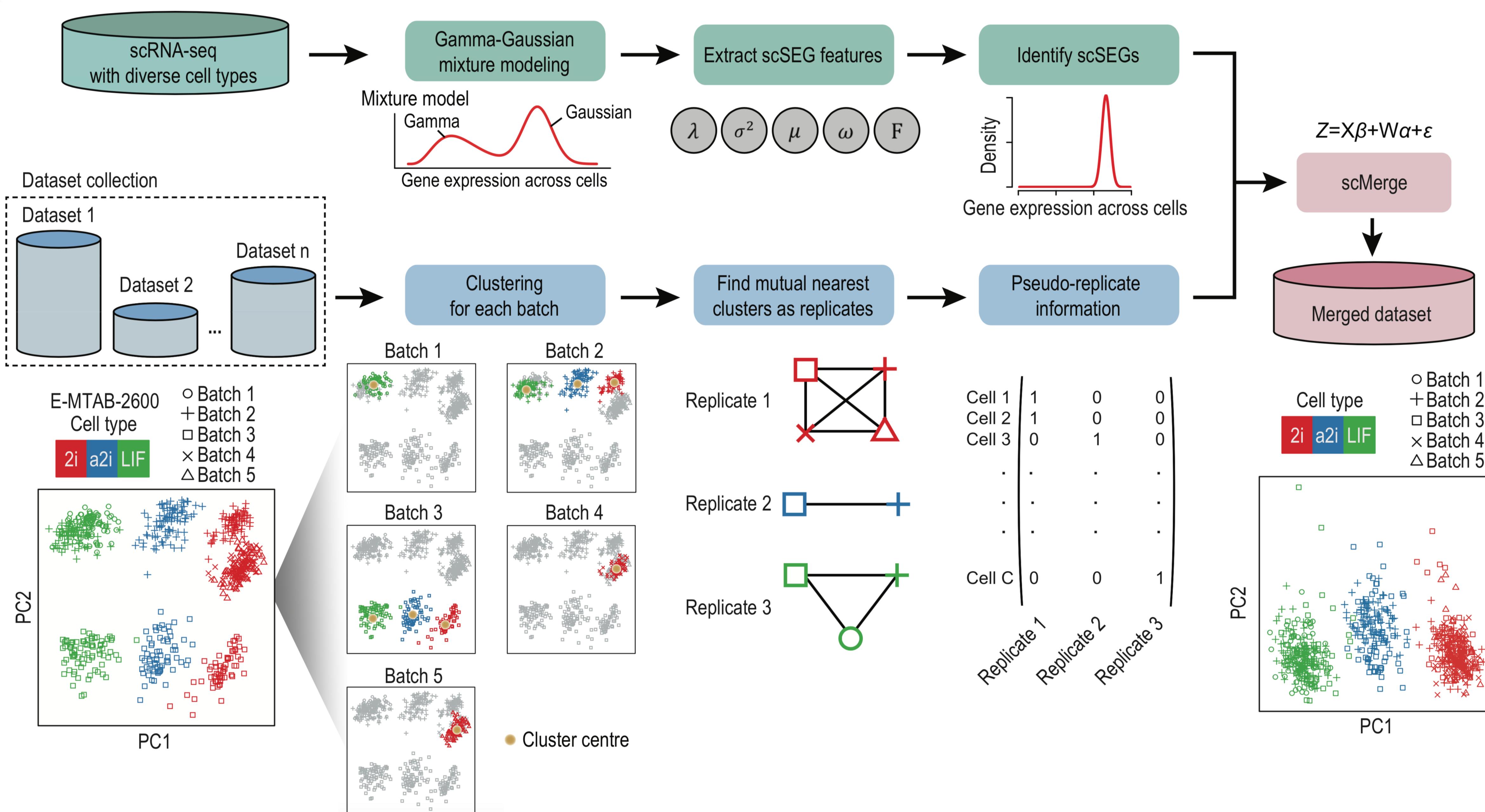


scMerge Framework

Concerted examination of multiple collections of single cell RNA-seq data promises further biological insights that cannot be uncovered with individual datasets.

The scMerge algorithm consists of three key components:

- (i) the identification of **stably expressed genes (scSEGs)** via a Gamma-Gaussian mixture model (Ghazanfar et al. 2016) for use as ‘negative controls’ for estimating unwanted factors;
- (ii) the construction of **pseudo-replicates** to estimate the effects of unwanted factors;
- (iii) the adjustment of the datasets with unwanted variation using a **fastRUVIII** model with the negative control genes and pseudo-replicates, resulting in a single merged dataset.

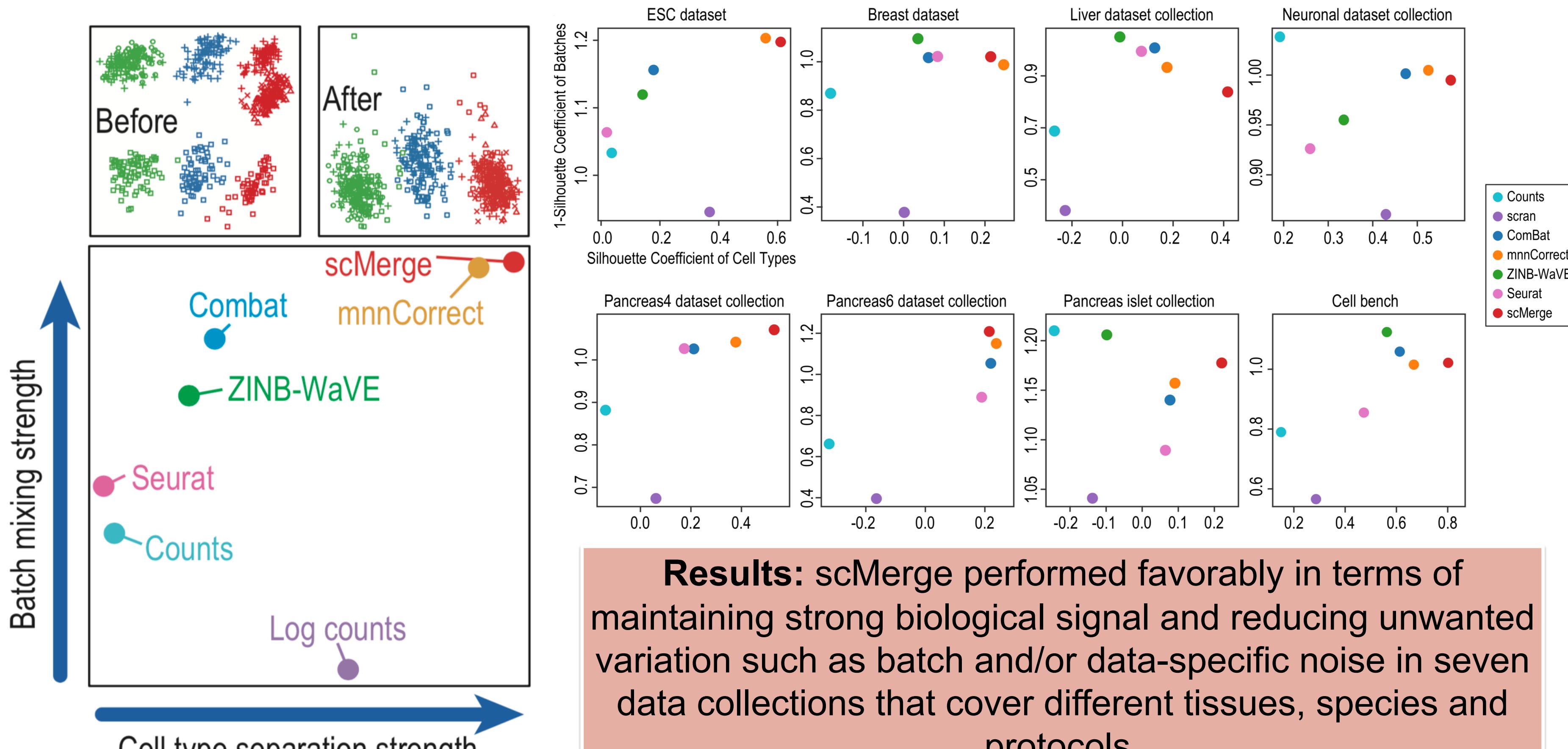


Dataset collections

Name	Accession	Organism	Tissue/organ source	# of cell types	# of cells	Protocol
mESC	E-MTAB-2600	Mouse	mESC	3	704	SMART-seq/C1
Breast	GSE113197	Human	Breast cancer	3	24520	10x Chromium
Liver	GSE87795 GSE90047 GSE87038 GSE96981	Mouse	Liver	8	1249	SMART-seq/C1 SMART-seq2 SMART-seq2 SMART-seq/C1
Neuronal	SRP065920 GSE75413	Mouse	Neuronal	2	145	SMART-seq2 STRT-seq
Pancreas4	GSE81608 GSE83139 GSE86469 E-MTAB-5061	Human	Pancreas	6	4566	SMART-seq/C1 SMART-seq/C1 SMART-seq/C1 SMART-seq2 Cel-seq2 inDrop
Pancreas6	GSE85241 GSE84133	Human	Pancreas Islets	6	1773	
CellBench	NA	Human	Adenocarcinoma cell lines	13	8569	Cel-seq2, Drop-seq, & 10x Chromium
Embryogenesis	GSE45719 GSE57249 E-MTAB-3321 GSE44183 E-MTAB-3929 GSE36552 GEO66507	Mouse Human/mouse Human	Embryogenesis	10+	2144	SMART-seq SMARTer SMART-seq2 Tang et al., 2010 SMART-seq2 Tang et al., 2010 SMARTer

Dataset collections that are used in this study.

Evaluation



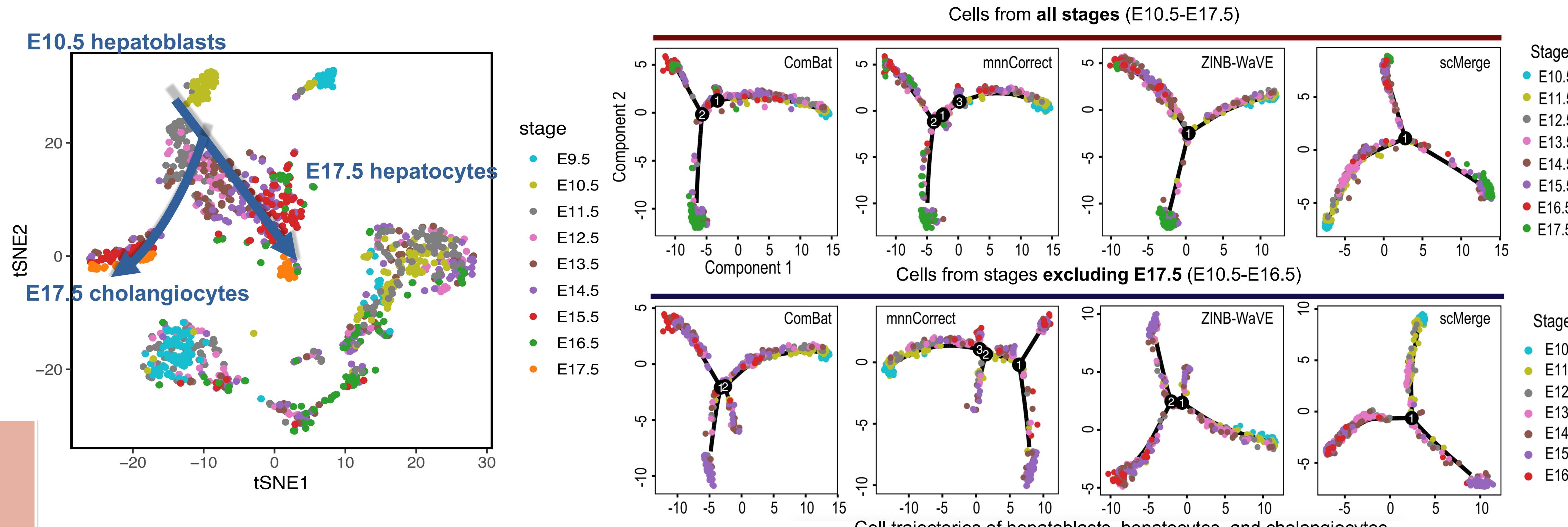
Results: scMerge performed favorably in terms of maintaining strong biological signal and reducing unwanted variation such as batch and/or data-specific noise in seven data collections that cover different tissues, species and protocols.

Case study: Integration of four mouse fetal liver single-cell RNA-seq datasets

Summary: To illustrate the capability of scMerge to enable further downstream analyses, we examined the stability when faced with incomplete data.

Approach: We reconstructed the cell trajectories, using Monocle 2, of hepatoblasts, hepatocytes, and cholangiocytes for both the full Liver data collection and for a subset of the original Liver data collection; with E17.5 time point of GSE90047 removed.

Results: scMerge was most consistent with the full Liver data collection and agrees with current literature.



Acknowledgement

The authors thank all their colleagues, particularly at The University of Sydney, School of Mathematics and Statistics, for their support and intellectual engagement. The following sources of funding for each author, and for the manuscript preparation, are gratefully acknowledged: Australian Research Council Discovery Project grant (DP170100654) to JYH, JT, PY; Discovery Early Career Research Award (DE170100759) to PY; Australia NHMRC Career Development Fellowship (APP1111338) to JYH and the Judith and David Coffey Life Lab at the Charles Perkins Centre, The University of Sydney to SG; Australian Postgraduate Award to KW; Research Training Program Tuition Fee Offset and Stipend Scholarship to YL; NHMRC Program Grant (1054618) to TPS; J.G. was supported by the National Science Foundation under grant no. DMS-1646108. SJTU-USYD Translate Medicine Fund-Systems Biomedicine AF6260003.

Further information

The scMerge R package and more case studies are available at <https://sydneybiox.github.io/scMerge>

Reference:

Ghazanfar et al. BMC Systems Biology 2016, 10(Suppl 5):127 doi: 10.1186/s12918-016-0370-4

Lin et al. BioRxiv 2018, doi: <https://doi.org/10.1101/393280>

Contact: yingxin.lin@sydney.edu.au / jean.yang@sydney.edu.au