

MBC

Medical Bioinformatics Centre

Turku Bioscience

University of Turku and Åbo Akademi University

FI-20520 Turku, Finland

Introduction to integration of scRNA-seq data

The Hitchhiker's Guide to scRNA-seq, iMM

09/07/2024

António Sousa - aggode@utu.fi



**UNIVERSITY
OF TURKU**



Åbo Akademi



InFLAMES
Solution is in Immunity

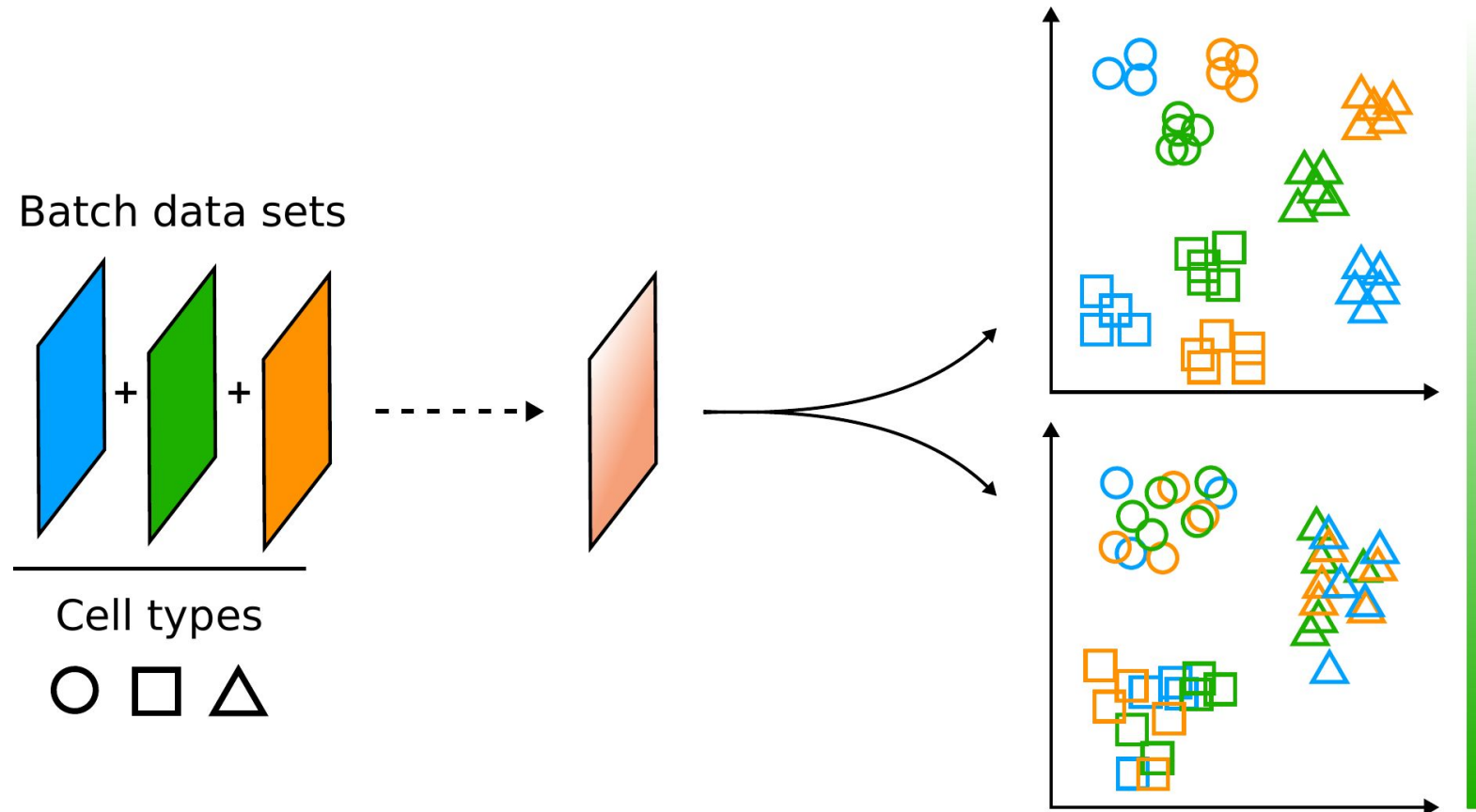


Outline

- Integration: methods & benchmarks
- Strengths & weaknesses of the different methods
- Strategies to assess integration accuracy
- Integration: *de novo* versus reference-based

What is integration?

- Identification of shared cell types/states across heterogeneous scRNA-seq datasets



How do we know if integration is required or not?

- Three main approaches (can be combined):
 - Dimensional reduction techniques
 - Automatic cell annotation
 - Independent sample analysis: clustering → cluster markers → annotation
 - (cluster comparison between samples)

Overview of integration strategies





REVIEW ARTICLE

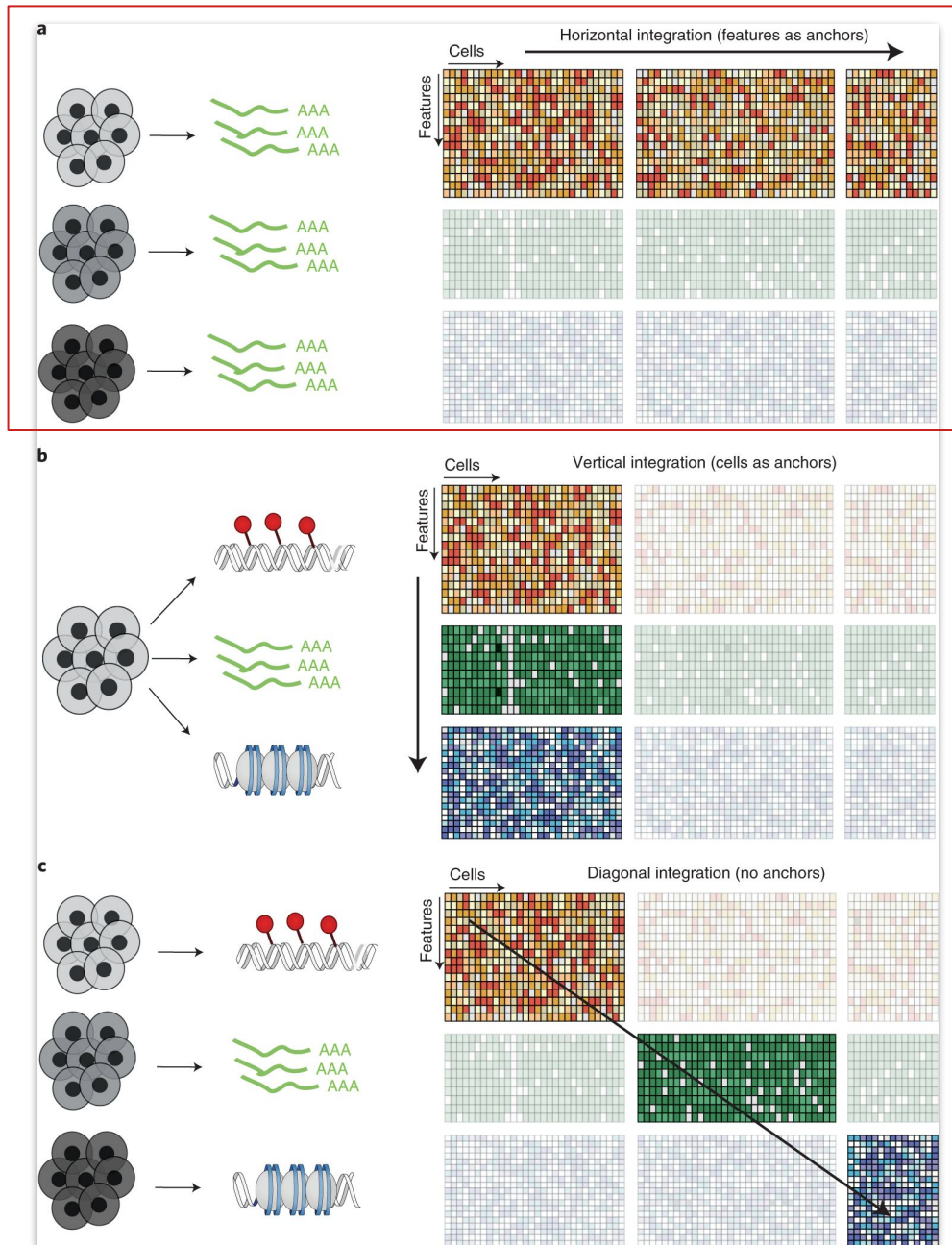
<https://doi.org/10.1038/s41587-021-00895-7>

nature
biotechnology

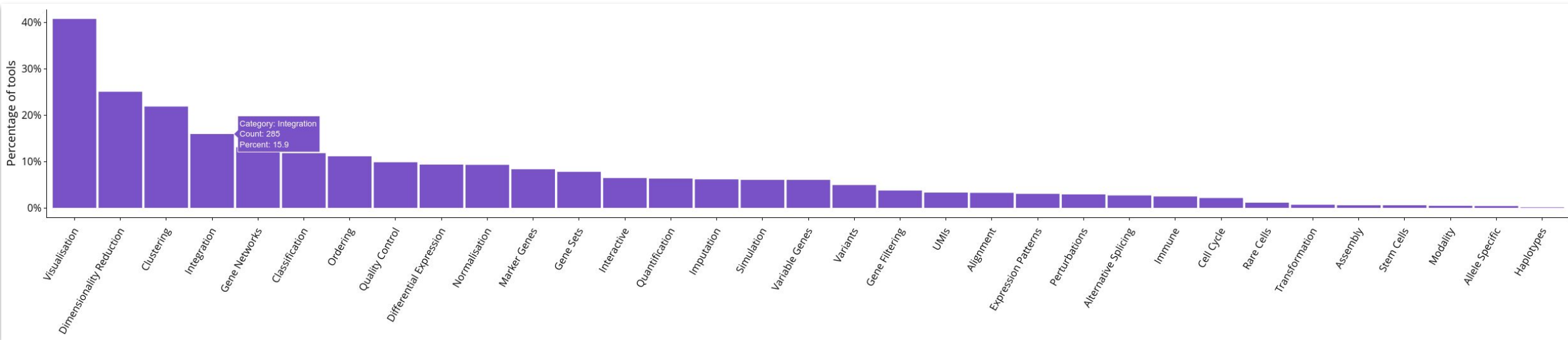
Check for updates

Computational principles and challenges in single-cell data integration

Ricard Argelaguet ^{1,2} ✉, Anna S. E. Cuomo ^{1,3} ✉, Oliver Stegle ^{3,4,5} ✉ and John C. Marioni ^{1,3,6} ✉



Overview of integration methods: www.scrna-tools.org (20/06/24)



285 methods/tools!

Overview of integration methods: www.scrna-tools.org (20/06/24)

Tools

View entries for individual tools

Sorting & Filters

Sort by
Citations

Filter by category
Select multiple categories and click FILTER below

- Cell Cycle
- Classification
- Clustering
- Differential Expression
- Dimensionality Reduction
- Gene Filtering
- Gene Networks
- Gene Sets
- Haplotypes
- Immune
- Imputation
- Integration

FILTER **RESET**

- Seurat CRAN 5.1.0 downloads 65K/month
- Harmony CRAN 1.2.0 downloads 6749/month
- scran in Bioc 8 years rank 87 / 2300
- scvi-tools pypi package 1.1.2 downloads/month 21k
- batchelor in Bioc 5 years rank 99 / 2300
- RaceID CRAN 0.3.3 downloads 447/month
- MIMOSCA
- bseqsc
- MOFA pypi package 1.2 downloads/month 99
- LIGER CRAN 2.0.1 downloads 1092/month
- scanorama pypi package 1.7.4 downloads/month 4k
- scMC
- bbknn pypi package 1.6.0 downloads/month 5k

Seurat

- Shared nearest neighbors in the low dimensional space
- Dimensional reductions: Canonical Correlation Analysis (CCA) or Reciprocal PCA (RPCA)
- Integration or reference-mapping/transfer learning

Comprehensive Integration of Single-Cell Data

Tim Stuart,^{1,4} Andrew Butler,^{1,2,4} Paul Hoffman,¹ Christoph Hafemeister,¹ Ethhymia Papalexi,^{1,2} William M. Mauck III,^{1,2} Yuhao Hao,^{1,2} Marlon Stoeckius,³ Peter Smibert,³ and Rahul Satija^{1,2,5,*}

¹New York Genome Center, New York, NY, USA

²Center for Genomics and Systems Biology, New York University, New York, NY, USA

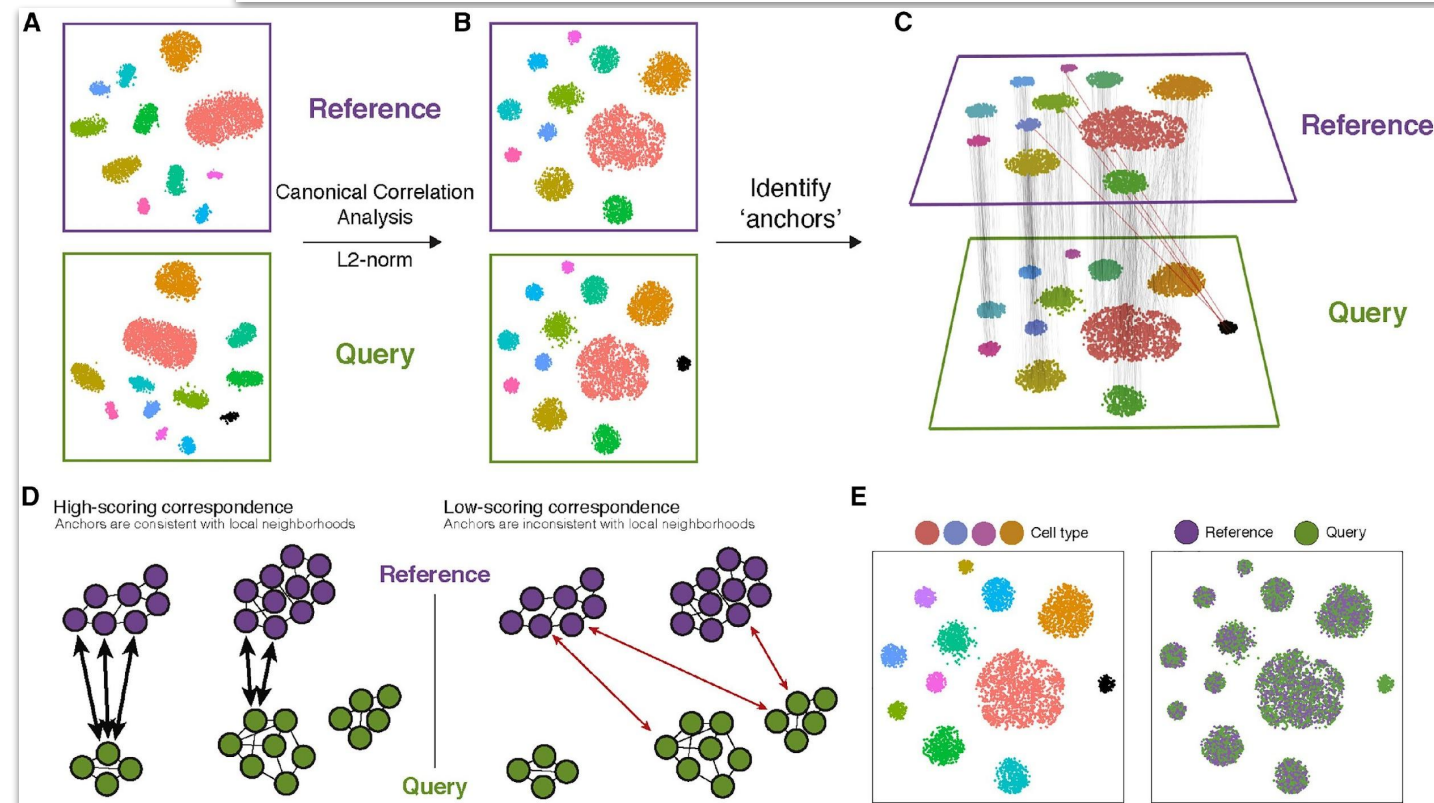
³Technology Innovation Lab, New York Genome Center, New York, NY, USA

⁴These authors contributed equally

⁵Lead Contact

*Correspondence: rsatija@nygenome.org

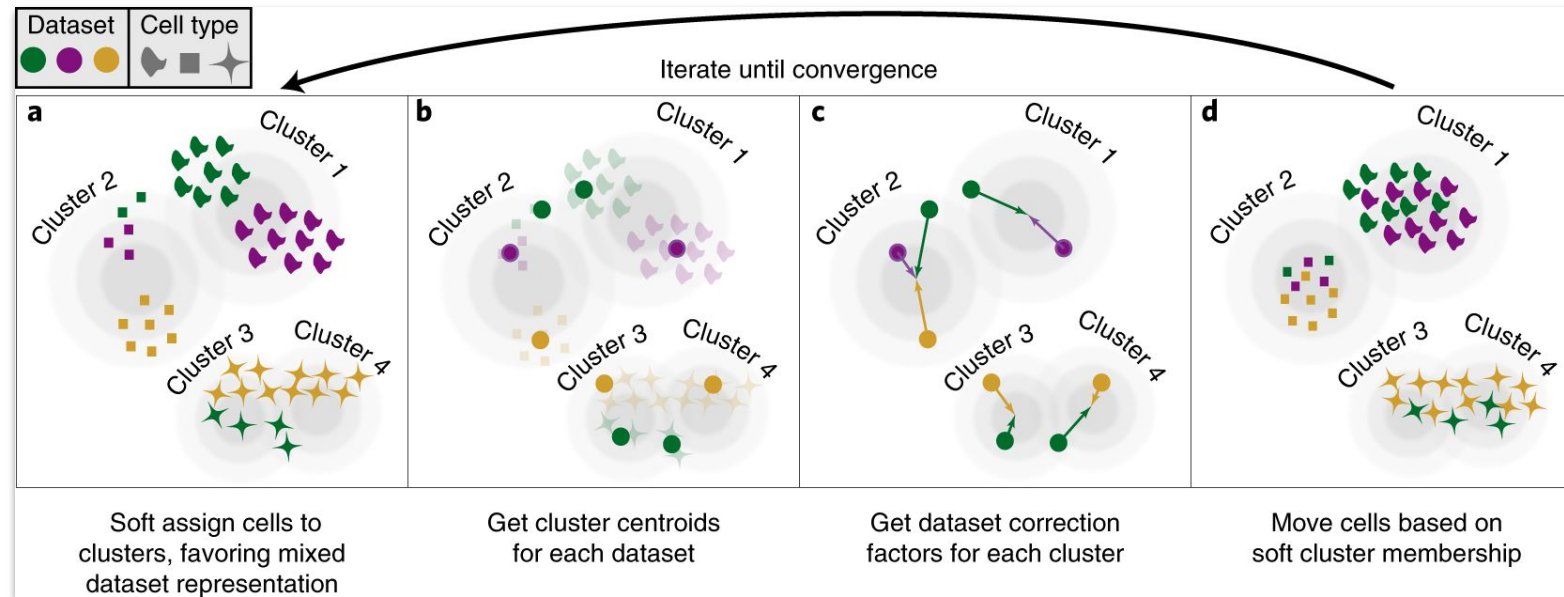
<https://doi.org/10.1016/j.cell.2019.05.031>



<https://satijalab.org/seurat>

Fast, sensitive and accurate integration of single-cell data with Harmony

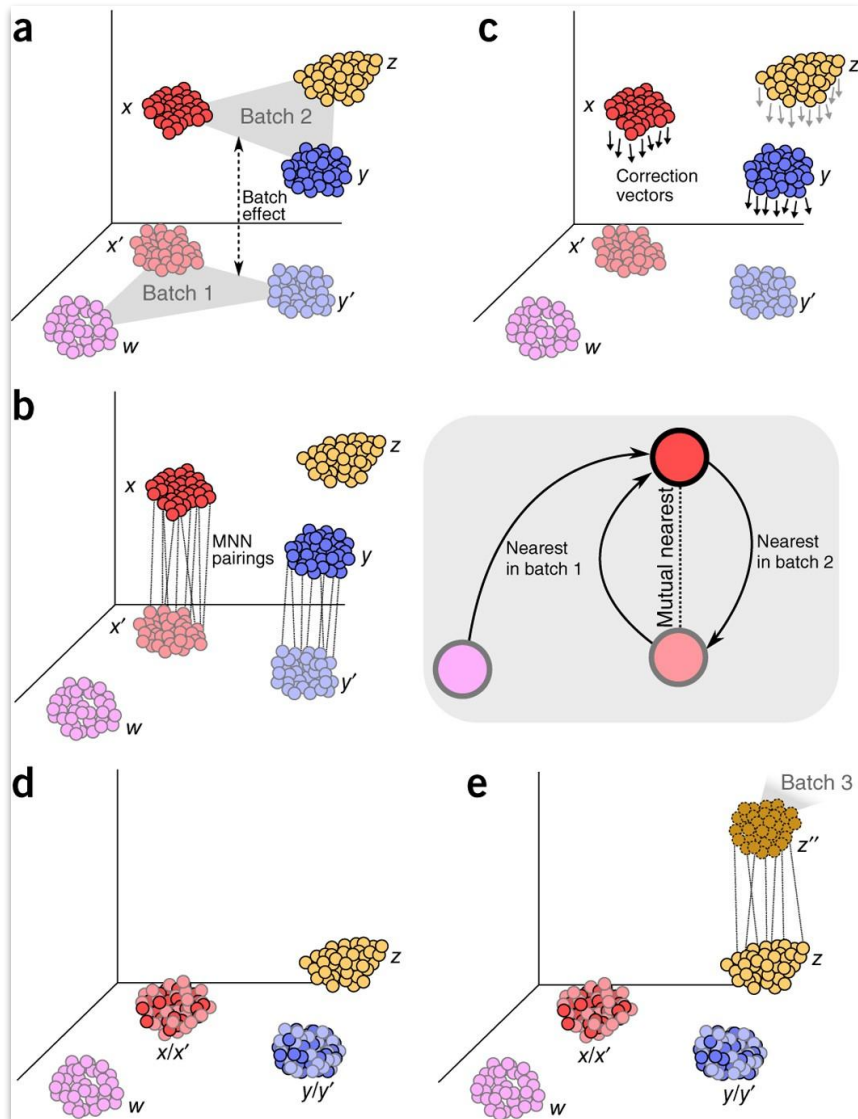
Ilya Korsunsky^{1,2,3,4}, Nghia Millard^{1,2,3,4}, Jean Fan⁵, Kamil Slowikowski^{1,2,3,4},
Fan Zhang^{1,2,3,4}, Kevin Wei², Yuriy Baglaenko^{1,2,3,4}, Michael Brenner², Po-ru Loh^{1,3,4} and
Soumya Raychaudhuri^{1,2,3,4,6*}



- *Approximates dataset-specific cluster centroids to global centroids in the PCA space*

<https://portals.broadinstitute.org/harmony>

fastMNN



Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors

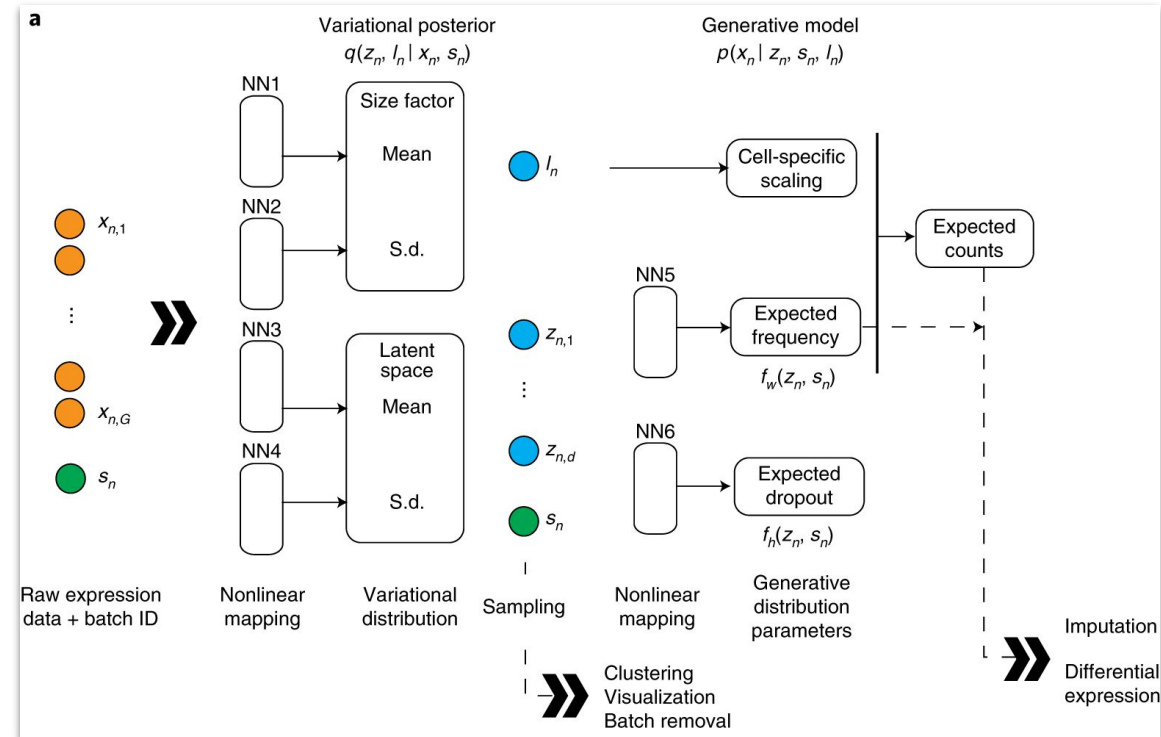
Laleh Haghverdi^{1,2}, Aaron T L Lun³, Michael D Morgan⁴ & John C Marioni^{1,3,4}

- *Correction vectors applied between mutual nearest neighbors pairs*

<https://bioconductor.org/packages/devel/bioc/html/batchelor.html>

Deep generative modeling for single-cell transcriptomics

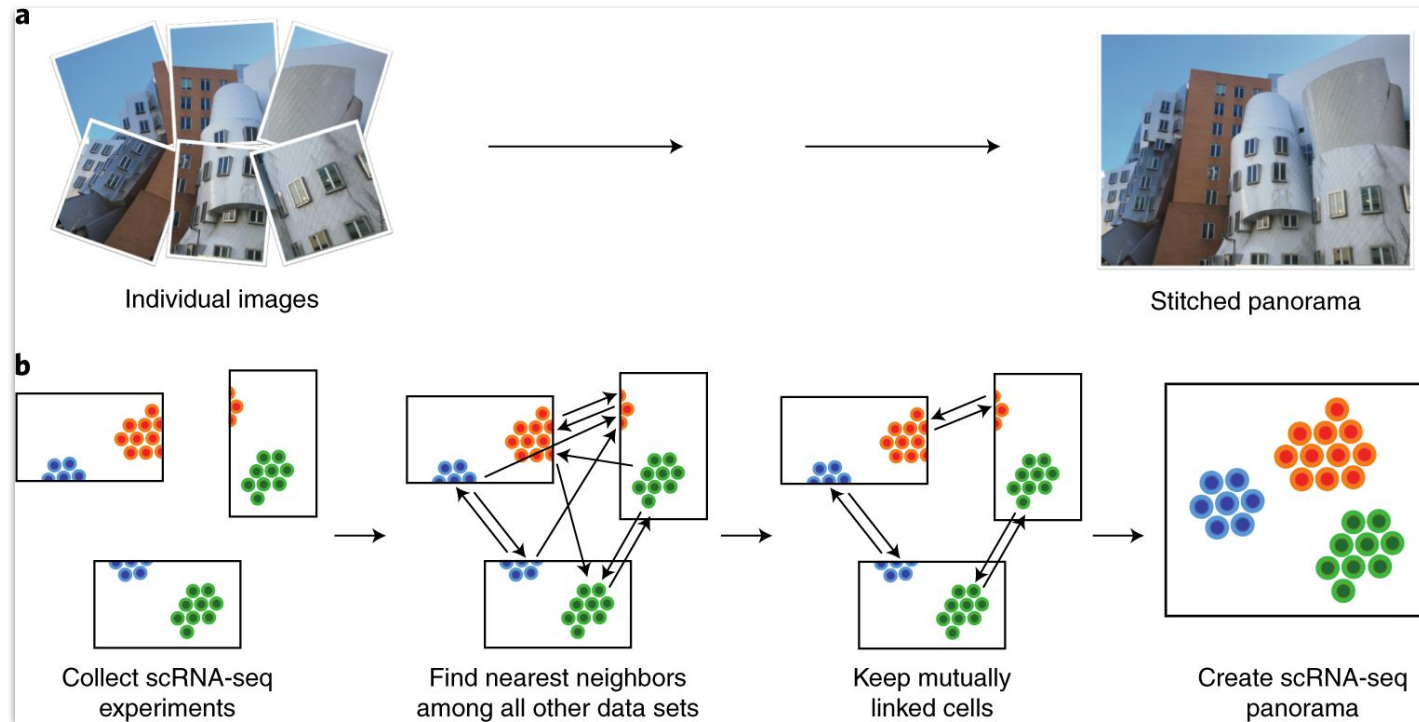
Romain Lopez¹, Jeffrey Regier¹, Michael B. Cole², Michael I. Jordan^{1,3} and Nir Yosef^{1,4,5*}



- “stochastic optimization and deep neural networks to aggregate information across similar cells and genes and to approximate the distributions that underlie observed expression values, while accounting for batch effects”

Efficient integration of heterogeneous single-cell transcriptomes using Scanorama

Brian Hie¹, Bryan Bryson^{2*} and Bonnie Berger^{1,3*}

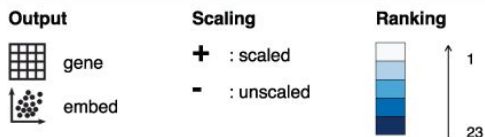


- “Mutually linked cells form matches that can be leveraged to correct for batch effects and merge experiments together, whereby the datasets forming connected components on the basis of these matches become a scRNA-seq ‘panorama’”

<https://github.com/brianhie/scanorama>

Semi-supervised integration

		Output	Features	Scaling	Pancreas	Lung	Immune	T cells
1	ssSTACAS*	embed	HVG	-	2	1	3	1
2	scANVI*	embed	HVG	-	1	3	1	1
3	ssSTACAS*	embed	HVG	+				
4	STACAS	embed	HVG	-				
5	STACAS	embed	HVG	+	1			
6	fastMNN	embed	HVG	-			2	
7	Scanorama	embed	HVG	+			1	
8	Seurat v4 RPCA	embed	HVG	+	3	3		
9	scGen*	gene	HVG	+				
10	Seurat v4 RPCA	embed	HVG	-	2			
11	Harmony	embed	HVG	-				
12	scVI	embed	HVG	-				
13	Scanorama	embed	HVG	-				2
14	scGen*	gene	HVG	-				
15	Harmony	embed	HVG	+				
16	ComBat	gene	HVG	+				
17	fastMNN	embed	HVG	+				
18	Seurat v4 CCA	embed	HVG	-				
19	Scanorama	gene	HVG	-				
20	ComBat	gene	HVG	-				
21	Seurat v4 CCA	embed	HVG	+				
22	Scanorama	gene	HVG	+				
23	Unintegrated	gene	HVG	-				

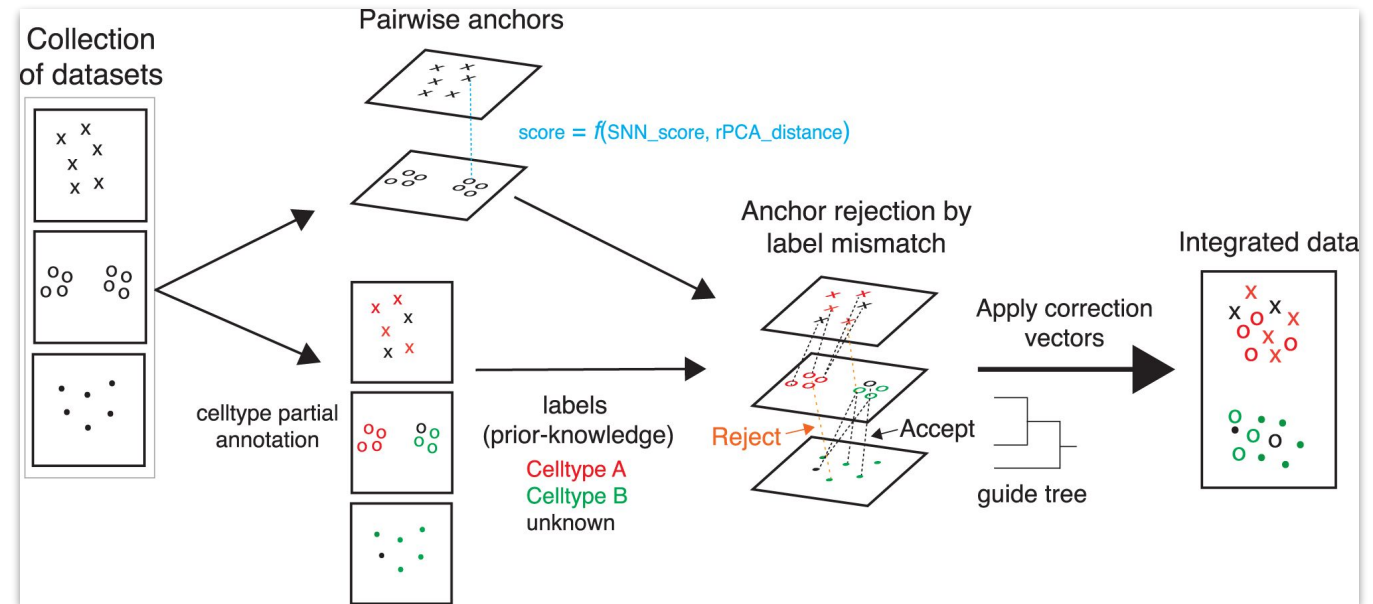


Semi-supervised integration of single-cell transcriptomics data

Received: 25 July 2023

Massimo Andreatta^{1,2,3}, Léonard Héroult^{1,2,3}, Paul Gueguen^{1,2,3}, David Gfeller^{1,2,3}, Ariel J. Berenstein⁴ & Santiago J. Carmona^{1,2,3} ✉

Accepted: 16 January 2024



<https://github.com/carmonalab/STACAS>

Integration scRNA-seq benchmarks

Tran et al. *Genome Biology* (2020) 21:12
<https://doi.org/10.1186/s13059-019-1850-9>

Genome Biology

RESEARCH

Open Access

A benchmark of batch-effect correction methods for single-cell RNA sequencing data

Hoa Thi Nhu Tran[†], Kok Siong Ang[†], Marion Chevrier[†], Xiaomeng Zhang[†], Nicole Yee Shin Lee, Michelle Goh and Jinmiao Chen^{*}



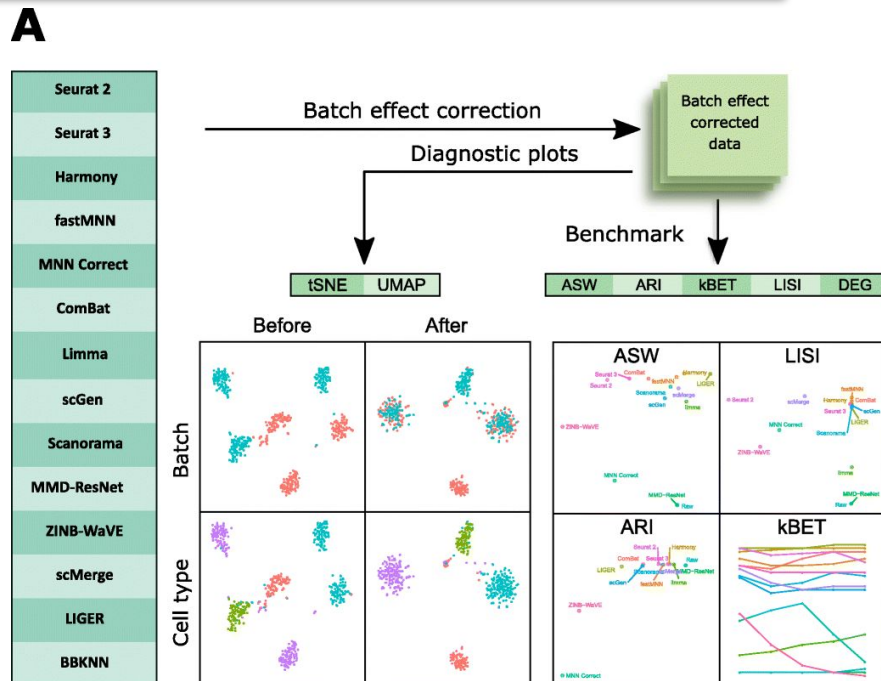
Abstract

Background: Large-scale single-cell transcriptomic datasets generated using different technologies contain batch-specific systematic variations that present a challenge to batch-effect removal and data integration. With continued growth expected in scRNA-seq data, achieving effective batch integration with available computational resources is crucial. Here, we perform an in-depth benchmark study on available batch correction methods to determine the most suitable method for batch-effect removal.

Results: We compare 14 methods in terms of computational runtime, the ability to handle large datasets, and batch-effect correction efficacy while preserving cell type purity. Five scenarios are designed for the study: identical cell types with different technologies, non-identical cell types, multiple batches, big data, and simulated data. Performance is evaluated using four benchmarking metrics including kBET, LISI, ASW, and ARI. We also investigate the use of batch-corrected data to study differential gene expression.

Conclusion: Based on our results, *Harmony*, *LIGER*, and *Seurat 3* are the recommended methods for batch integration. Due to its significantly shorter runtime, *Harmony* is recommended as the first method to try, with the other methods as viable alternatives.

Keywords: Single-cell RNA-seq, Batch correction, Batch effect, Integration, Differential gene expression



B

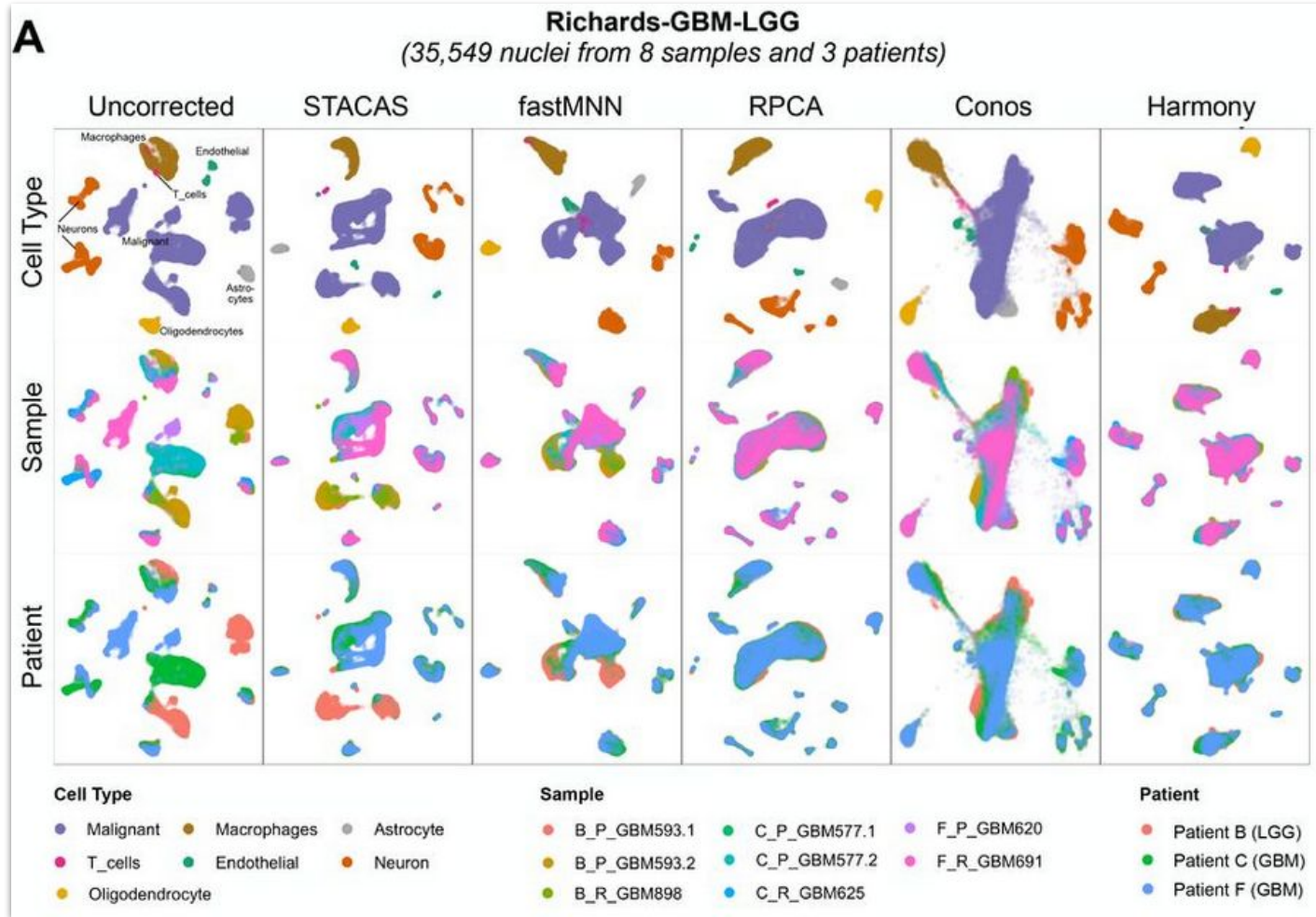
Dataset	Description	Number of batches	Total cell number	Technologies
1	Human Dendritic Cells	2	576	Smart-Seq2
2	Mouse Cell Atlas	2	6,954	Microwell-Seq Smart-Seq2
3	Simulation	Refer to Simulation table		
4	Human Pancreas	5	14,767	inDrop CEL-Seq2 Smart-Seq2 SMARTer SMARTer
5	Human Peripheral Blood Mononuclear Cell	2	15,476	10x 3' 10x 5'
6	Cell line	3	9,530	10x
7	Mouse Retina	2	71,638	Drop-seq
8	Mouse Brain	2	833,206	Drop-seq SPLIT-seq
9	Human Cell Atlas	2	621,466	10x
10	Mouse Haematopoietic Stem and Progenitor Cells	2	4,649	MARS-seq Smart-Seq2

Integration scRNA-seq benchmarks

A comparison of data integration methods for single-cell RNA sequencing of cancer samples

Laura M. Richards^{1,2}, Mazdak Riverin², Suluxan Mohanraj², Shamini Ayyadhury^{2,3}, Danielle C. Croucher^{1,2}, J. Javier Díaz-Mejía², Fiona J. Coutinho⁵, Peter B. Dirks^{4,5,6}, Trevor J. Pugh^{1,2,7,#}

Tumours are routinely profiled with single-cell RNA sequencing (scRNA-seq) to characterize their diverse cellular ecosystems of malignant, immune, and stromal cell types. When combining data from multiple samples or studies, batch-specific technical variation can confound biological signals. However, scRNA-seq batch integration methods are often not designed for, or benchmarked, on datasets containing cancer cells. Here, we compare 5 data integration tools applied to 171,206 cells from 5 tumour scRNA-seq datasets. Based on our results, STACAS and fastMNN are the most suitable methods for integrating tumour datasets, demonstrating robust batch effect correction while preserving relevant biological variability in the malignant compartment. This comparison provides a framework for evaluating how well single-cell integration methods correct for technical variability while preserving biological heterogeneity of malignant and non-malignant cell populations.

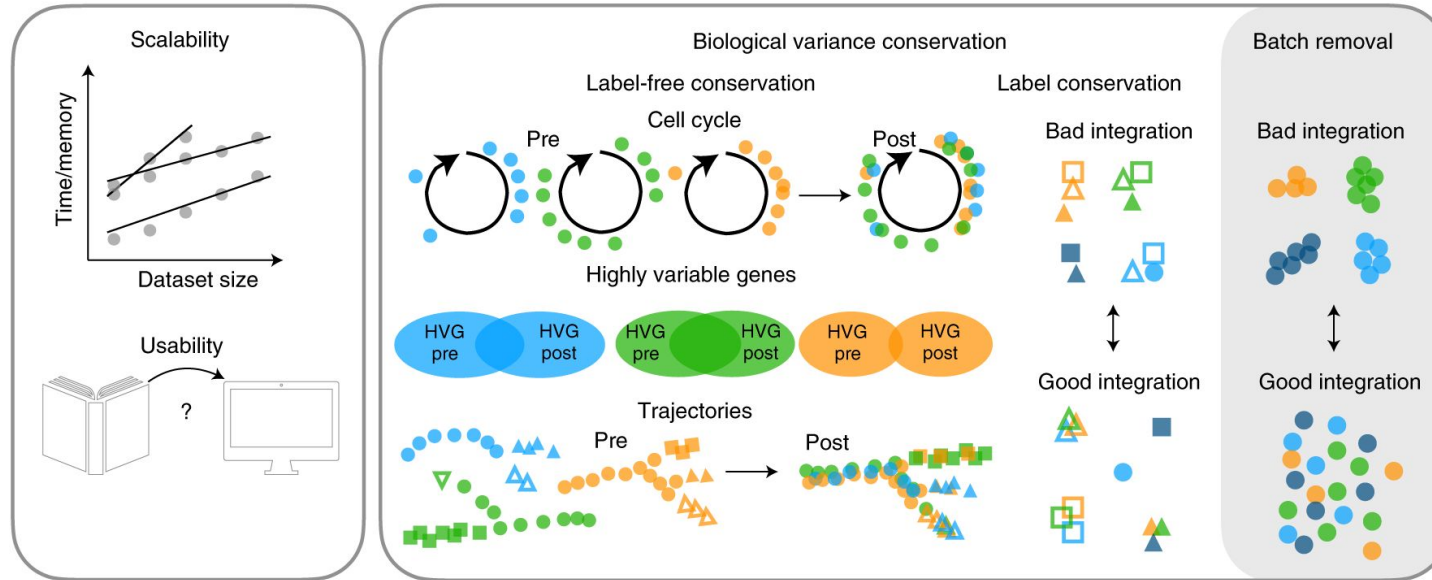
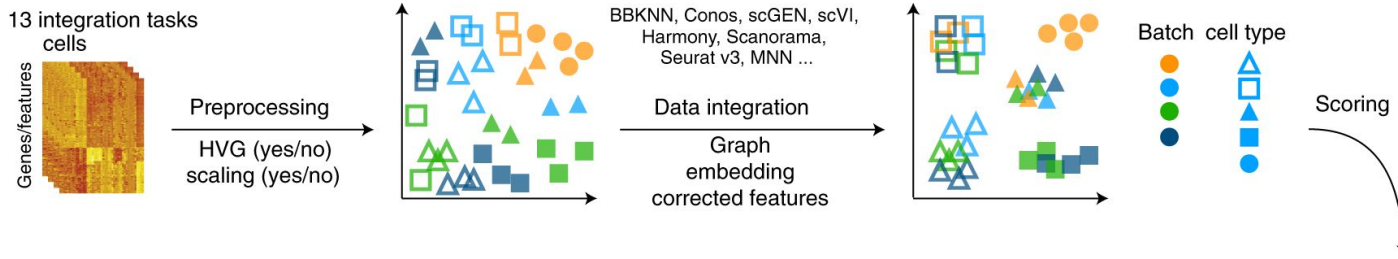


Integration scRNA-seq benchmarks

OPEN

Benchmarking atlas-level data integration in single-cell genomics

Malte D. Luecken¹, M. Büttner¹, K. Chaichoompu¹, A. Danese¹, M. Interlandi², M. F. Mueller¹, D. C. Strobl¹, L. Zappia^{1,3}, M. Dugas⁴, M. Colomé-Tatché^{1,5,6} and Fabian J. Theis^{1,3,5}



Method		RNA	Simulations	Usability	Scalability
1	scANVI*	HVG -	1 2	1	3
2	Scanorama	HVG +	1 2	1	1
3	scVI	HVG -	3 3	1	1
4	fastMNN	HVG -	2	1	1
5	scGen*	HVG -	3 1	1	1
6	Harmony	HVG -	1	1	1
7	fastMNN	HVG -	1	1	1
8	Seurat v3 RPCA	HVG +	2	1	1
9	BBKNN	HVG -	2	3	2
10	Scanorama	HVG +	1	1	1
11	ComBat	HVG -	3	3	1
12	MNN	HVG +	1	1	1
13	Seurat v3 CCA	HVG -	1	1	1
14	trVAE	HVG -	1	1	1
15	Conos	HVG -	1	1	1
16	DESC	FULL -	3	1	1
17	LIGER	HVG -	1	1	1
18	SAUCIE	HVG +	1	3	1
19	Unintegrated	FULL -	1	1	1
20	SAUCIE	HVG +	1	3	1

Rank Name Output Features Scaling Pancreas Lung Immune (human) Immune (human/mouse) Mouse brain Sim 1 Sim 2 Package Paper Time Memory

Output: Genes (grid), Embedding (up arrow), Graph (network icon).
 Scaling: + Scaled, - Unscaled.
 Ranking: 1 (white), 20 (black).

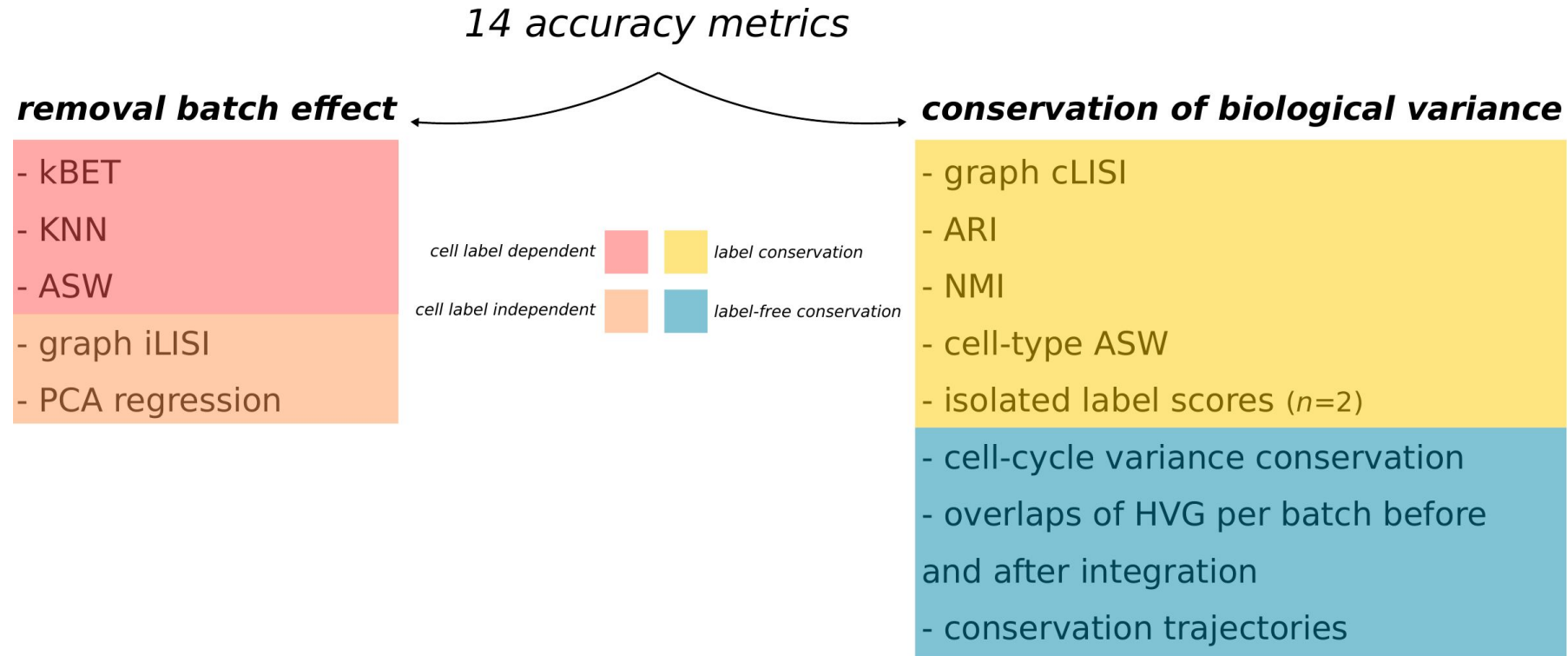
Conclusions: strengths & weaknesses of the different methods

- Integration: trade-off between biological signal preservation and batch correction
- Inconsistent results across benchmarks: different methods, integration tasks, performance metrics
- Benchmarks: showed that performance of integration is inherently dependent on the complexity of the integration task, highlighting that there is no single method fitting all the tasks
- Methods more suitable to correct strong batch effects (prioritize batch correction): Seurat CCA, Harmony and fastMNN
- Methods more suitable to correct mild to moderate batch effects (prioritize biological conservation): Seurat RPCA, Scanorama and scVI

Strategies to assess integration accuracy

- Two types of assessment:
 - qualitative (dimensional reduction visualizations)
 - quantitative (objective metrics)

Measuring integration performance (quantitative): scib package



Input: *integrated joint embedding or feature corr. matrix*

Paper: <https://www.nature.com/articles/s41592-021-01336-8>

scIB github repo: <https://github.com/theislab/scib>

scib-pipeline github repo: <https://github.com/theislab/scib-pipeline>

Measuring integration performance

Published Online: 23 March, 2021 | Supp Info: <http://doi.org/10.26508/lsa.202001004>
 Downloaded from life-science-alliance.org on 20 June, 2024

Research Article



Check for updates

CellMixS: quantifying and visualizing batch effects in single-cell RNA-seq data

Almut Lütge^{1,2}, Joanna Zyprych-Walczak³, Urszula Brykczynska Kunzmann⁴, Helena L Crowell^{1,2}, Daniela Calini⁵, Dheeraj Malhotra⁵, Charlotte Soneson^{2,4}, Mark D Robinson^{1,2}

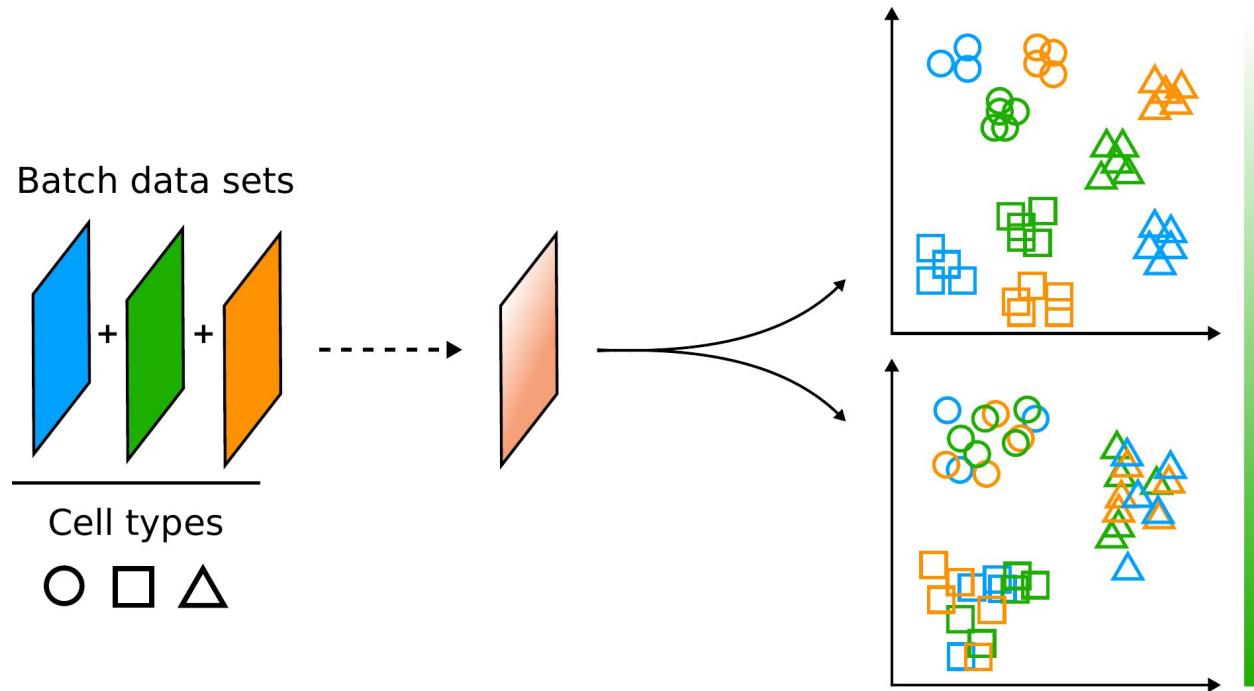
Table 1. Batch mixing metrics: short summaries of metrics included in the benchmark.

Metric	Level	Basis	Short description	Interpretation
Cell-specific Mixing Score (cms)	Cell	knn, pca	Test for whether distance distributions from a neighbourhood are batch specific	P-value: Probability to observe as large differences in distance distributions assuming the same underlying distribution
Local Inverse Simpson Index (lisi)	Cell	knn	Inverse of the sum of batch probabilities within weighted knn	Effective number of batches in neighbourhood
Entropy	Cell	knn	Sum of the products of the batch probabilities and their log within each cell's knn	Randomness in the data according to the batch variable
Mixing metric (mm)	Cell	knn	Median position of the fifth cell from each batch within its knn	Number of cells within knn until each batch is represented by five cells
Graph connectivity (graph)	Cell type	knn-graph	Fraction of directly connected cells within cell type graphs	Proportion of non-distorted cell type relationships
k-nearest neighbour Batch effect test (kBet)	Cell type	knn	Test for equal batch proportions within a random cell's knn	P-value: Probability to observe as large differences in batch proportions assuming the same underlying proportions
Average silhouette width (asw)	Cell type	pca	Average relationship of within and between batch-cluster distances for each cell type	Indication of how well clusters are separated
Principal component regression (pcr)	Global	pca	Correlation of the batch variable with principal components weighted by their variance attributes	Proportion of variance attributed to batch

<https://github.com/almutlue/CellMixS>

Integration (*de novo*) versus reference-mapping

- Integration (*de novo*): integration → clustering → cluster markers → annotation

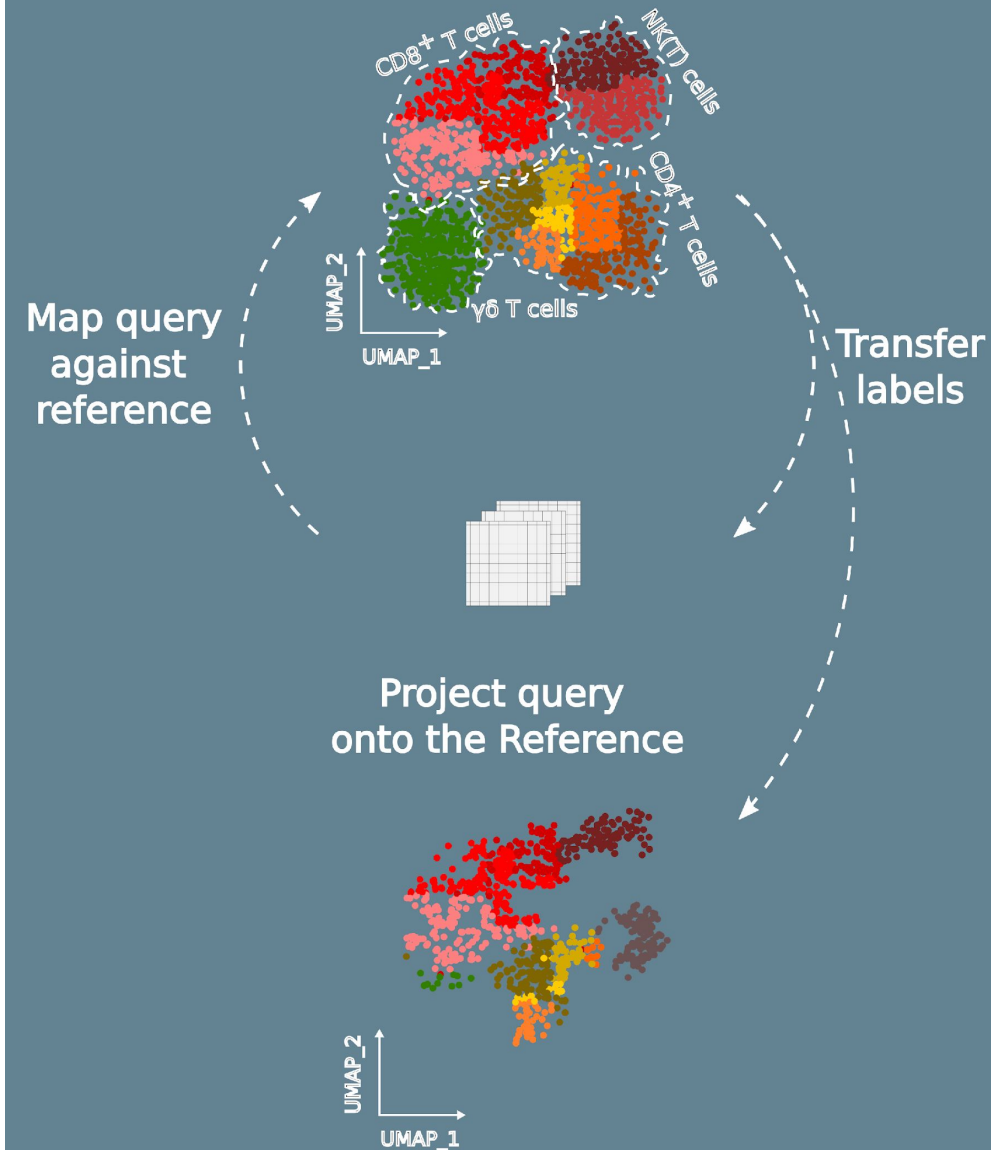


- Application: newly poorly described single-cell transcriptomic data

Integration (*de novo*) versus reference-mapping

- **Reference-mapping:** projecting & transferring labels from a previously annotated reference to a new query data set
- **Application:** experiments generating single-cell transcriptomic data with well-known identities for which are references available

Reference Mapping



Reference-mapping

Cell

Resource

Comprehensive Integration of Single-Cell Data

Tim Stuart,^{1,4} Andrew Butler,^{1,2,4} Paul Hoffman,¹ Christoph Hafemeister,¹ Efthymia Papalexi,^{1,2} William M. Mauck III,^{1,2} Yuhao Hao,^{1,2} Marlon Stoeckius,³ Peter Smibert,³ and Rahul Satija^{1,2,5,*}

¹New York Genome Center, New York, NY, USA

²Center for Genomics and Systems Biology, New York University, New York, NY, USA

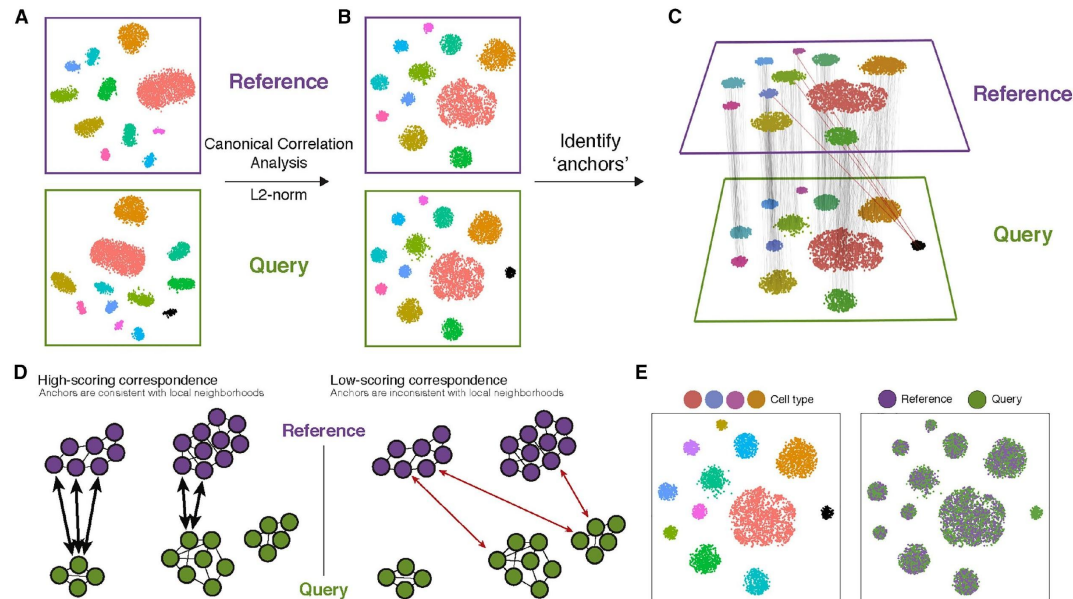
³Technology Innovation Lab, New York Genome Center, New York, NY, USA

⁴These authors contributed equally

⁵Lead Contact

*Correspondence: rsatija@nygenome.org

<https://doi.org/10.1016/j.cell.2019.05.031>



nature
biotechnology

ANALYSIS

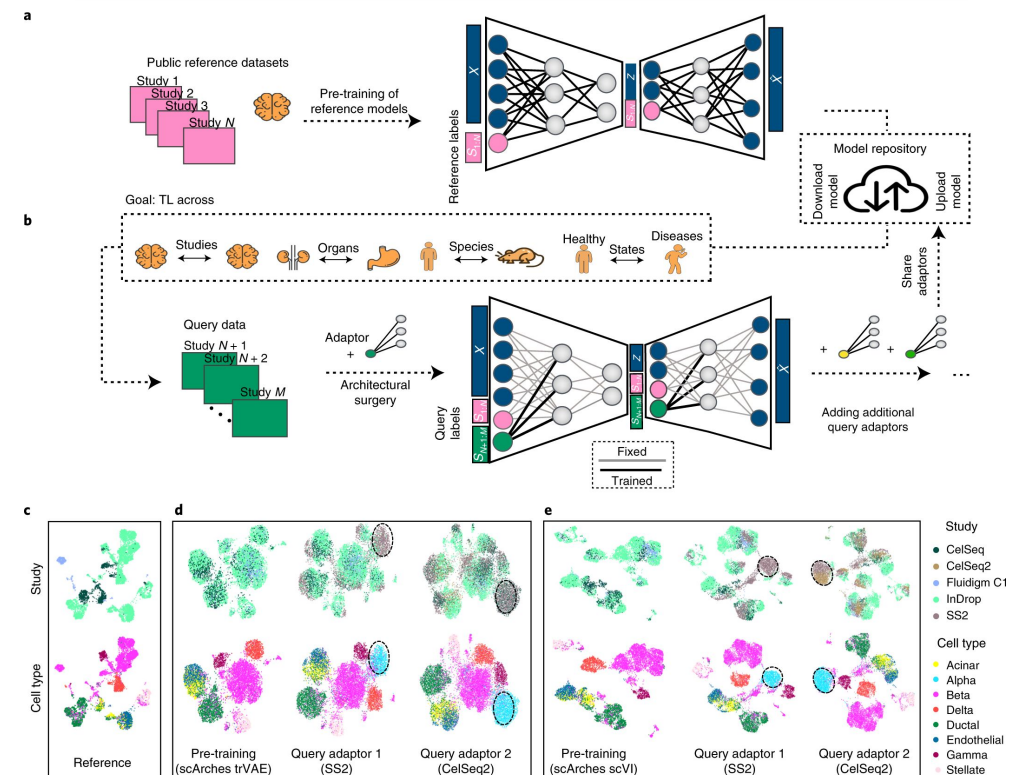
<https://doi.org/10.1038/s41587-021-01001-7>

Check for updates

OPEN

Mapping single-cell data to reference atlases by transfer learning

Mohammad Lotfollahi^{1,2}, Mohsen Naghipourfar¹, Malte D. Luecken¹, Matin Khajavi¹, Maren Büttner¹, Marco Wagenstetter¹, Žiga Avsec³, Adam Gayoso⁴, Nir Yosef^{4,5,6,7}, Marta Interlandi⁸, Sergei Rybakov^{1,9}, Alexander V. Misharin¹⁰ and Fabian J. Theis^{1,2,9} ✉



Additional resources about integration

- ***Analysis of single cell RNA-seq data*** Sanger course:

<https://www.singlecellcourse.org/scrna-seq-dataset-integration.html>

- ***Single-cell best practices*** website:

https://www.sc-best-practices.org/cellular_structure/integration.html



This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No.: 955321



scIB: kBET, kNN & ASW (dependent on cell label identity)

kBET (*k*-nearest-neighbor batch effect test):

- compares the batch label distribution of a local neighborhood of a (*k*NN) graph for a subset of a given cell label against the global neighborhood – Pearson’s Chi-squared test. Null hypothesis is that the distributions are similar, *i.e.*, batches are well mixed. K-BET result is the average rejection rate. Lower value means no/weak batch effect.

kNN (*k*-nearest-neighbor graph connectivity):

- measures if all the cells from a given cell identity label are all connected (1).

ASW (*average silhouette width*):

- measures the within- and between-cluster distances (to the closest cluster). The average of all widths corresponds to the ASW. The ASW ranges between -1 - 0 - 1 (misclassification - overlapping clusters - well-separated dense clusters). Computed on the integrated embedding (or PCA). Used for assessing batch effects as well as cell label conservation (some differences).

$$GC = \frac{1}{|C|} \sum_{c \in C} \frac{|LCC(G(N_c; E_c))|}{|N_c|}.$$

Paper: <https://www.nature.com/articles/s41592-021-01336-8>

scIB github repo: <https://github.com/theislab/scib>

scib-pipeline github repo: <https://github.com/theislab/scib-pipeline>

scIB: iLISI & PCA (*independent on cell label identity*)

Graph iLISI (*graph integration local inverse Simpson's Index*):

- “*the inverse Simpson's index is used to determine the number of cells that can be drawn from a neighbor list before one batch is observed twice*”. It is also used for cell-identity, known as cLISI. A value of 0 corresponds to low batch integration or cell-type separation.

PCA regression:

- the R -squared of the linear regression between each principal component and the batch variable (covariate) is calculated. The variance contribution of the batch effect by each PC is calculated as the product of the variance and the R -squared. The sum of the individual variances gives the total variance explained by the batch variable.

$$\text{Var}(C|B) = \sum_{i=1}^G \text{Var}(C|PC_i) \times R^2(PC_i|B),$$

Paper: <https://www.nature.com/articles/s41592-021-01336-8>

scIB github repo: <https://github.com/theislab/scib>

scib-pipeline github repo: <https://github.com/theislab/scib-pipeline>

scIB: ARI & NMI (*label conservation*)

ARI (*Adjusted Rand Index*):

- it compares two clustering results, i.e., cell-labels vs louvain clustering results on the integrated data set, accounting for both, matches and mismatches. Values of 0 or 1 correspond to a random classification or to a total match.

NMI (*Normalized Mutual Information*):

- it compares the overlap between cell-type labels with the louvain clustering results on the integrated data set. The overlap is scaled (0-1) “using the mean of the entropy terms for cell-type and cluster labels”. A value of 0 means uncorrelated clustering and a value of 1 perfect match.

Paper: <https://www.nature.com/articles/s41592-021-01336-8>

scIB github repo: <https://github.com/theislab/scib>

scib-pipeline github repo: <https://github.com/theislab/scib-pipeline>