Hanchen Wang^{1,2}, Jure Leskovec^{2,#} and Aviv Regev^{1,#}

¹Research and Early Development, Genentech, ²Department of Computer Science, Stanford University, [#]Equal Correspondence

Although biological studies increasingly rely on embeddings of single cell profiles, the quality of these embeddings can be challenging to assess. Such evaluations are especially important for avoiding misleading biological interpretations, assessing the accuracy of integration methods, and establishing the zero-shot capabilities of foundational models. Here, we posit that current evaluation metrics can be highly misleading. We show this by training a three-layer perceptron, Islander, which outperforms all 11 leading embedding methods on a diverse set of cell atlases, but in fact distorts biological structures, limiting its utility for biological discovery. We then present a metric, scGraph, to flag such distortions. Our work should help learn more robust and reliable cell embeddings.

Embeddings of single cell profiles are now routinely employed as a research tool in biological investigation, to characterize cell types and states, their changes over time, and their distinction between conditions, including diseases, organs, or drug treatments (1, 2). With a dramatic growth in single cell data, including the Human Cell Atlas (3), multiple efforts have focused on universal embeddings for diverse single cell data (4–9). Given their broad utility, it is crucial to scrutinize the embeddings' quality to evaluate the performance of the underlying integration method (10) and zero-shot capabilities of the resulting foundation model (11, 12).

A critical aspect in deriving helpful cell embeddings is the correction of non-biological batch effects that stem from technical variations, such as sample handling and sequencing protocols. These variations can mask biological signals and lead to misleading interpretations. Integration methods thus aim to mitigate batch-specific discrepancies while preserving essential biological information. The effectiveness of these integrated cell embeddings is typically assessed through two lenses: how well the cells from various batches mix together; how closely cells of the same type group together.

Here, we identified an overlooked challenge in the evaluation metrics used to assess embeddings. To demonstrate the limitations of current gold standard metrics for cell profile embeddings (10), we first developed Islander (Fig. 1a), a model that scores best on established metrics, but generates biologically problematic embeddings. Islander is a three-layer perceptron, directly trained on cell type annotations with mixup augmentations (13). We tested Islander across a diverse set of 11 different human tissue cell atlases (brain (14), breast (15), eye (16), fetal gut (17), heart (18), fetal lung (19), pancreas (10), skin (20)), which cover different strengths of batch effects and diverse biological systems, overall comprising more than 3.5 million cells from 10 human organ systems. For each atlas, we trained a Islander model and then compared it with another 12 embedding baselines: three dimension reduction methods (PCA, UMAP, TSNE) (21), eight batch integration methods (Harmony (22), Scanorama (23), BBKNN (24), fastMNN (25), scVI (26), scANVI (27), scGen (28), scPoli (29)), and one foundation model (Geneformer (5)) (Methods). In addition, for each atlas, we compared to the performance of the original authors' integration, if available.

Across all datasets, Islander consistently outperformed all baseline strategies across all 12 metrics (10) (Fig. 1b,c, Extended Data Tables 3-13). This is largely due to the principles underlying the evaluation metrics (10), which focus on assessing the efficiency of cell embeddings in terms of the coherence of cell clustering structures with cell type labels and the blending of batches within clusters. When Islander explicitly aligns these cell embeddings with cell type annotations, it forms well-separated cell 'islands' (Fig. 1d, right), with each island comprising cells annotated as the same type. This

^{© 2024} H. Wang, J. Leskovec & A. Regev.



Figure 1: Drifting Cell Islands highlight limitation of current metrics. a, Islander overview. b, c, Evaluation of cell embeddings. Normalized overall score over 12 metrics (y axis, Methods) for each method (x axis) assessed on using the Fetal Lung Cell Atlas (b) or another 10 cell atlases (c). "Baseline": best baseline results (Methods). d, Fetal Lung Cell Atlas embedding space. Single cell profiles (dot; color-coded by cell type annotation) from the Fetal Lung Cell Atlas embedded by the author's integration method (left, inset zoom in middle) or Islander (right). Annotations: fibroblast cell subsets. e, "Airway fibroblast" cell neighborhood changes across Islander runs. Normalized Euclidean distance (y axis) between the centroids of airway fibroblasts and its five nearest neighbor clusters (x axis) in the 50-dimensional (of author's integration) and 16-dimensional (of Islander's) embedding space. g, Cell islands distort developmental stage structure and cell cell relationships. Cell embeddings as in (d), colored by developmental week (f, color bar) or coarser cell type annotations (g). h, i. scGraph, a metric using learned cell similarity graphs to evaluate cell embeddings. h. Method overview. i. scGraph score (y axis, higher is better) for each method (x axis) assessed using the Fetal Lung Cell Atlas.

alignment significantly boosts the biological variance conservation metrics, leading to top-tier overall performance (Extended Data Tables 3-13).

However, such structure is driven by (and complies well) with the most granular annotation level at the cost of ignoring any higher level relationships and distorting biological structures, potentially obstructing downstream analyses and future discoveries (and would thus not be an advisable for an actual integration method). In particular, when annotated cell subsets follow a continuum, as is the case for fibroblasts, Islander separates its constituent parts (Fig. 1d). In the developing human lung, the original analysis (19) identified multiple sub-types of fibroblasts, each distinguished by different marker genes and spatial locations. While the original embedding preserves a continuum between these fibroblasts (Fig. 1d, left), they are fully separated by Islander (Fig.1d, right). Similarly, the Islander embedding disrupted the developmental continuum, clearly observed in the original study (Fig. 1f, left), but obscured by Islander (Fig. 1f, right).

Moreover, the "cell islands" drifted in different ways across distinct runs, especially for smaller cell subsets. For example, in three separate runs with overall similar scores, the composition of the neighborhood of airway fibroblast cells varied substantially, involving as many as 14 distinct cell types within the five nearest neighbors (Fig. 1e, Extended Data Fig. 1, Extended Data Table 9). Thus, aside from cluster identity, the embedding may be largely arbitrary in all other relationships, and this arbitrariness would carry into downstream analysis or the biologist's interpretation.

Prompted by these limitations of the quality evaluation criteria, we reasoned that focusing solely on the most granular cell relationships in evaluation can pose a substantial limitations, whereas preserving relationships between broader cell types (coarser annotations) is an important additional criterion, and may also be more robust to noise. Indeed, when evaluating the same set of embeddings using broader cell type annotations provided by the authors, Islander now achieved an overall score of 0.523, inferior to PCA (0.557) or the top-performing scVI (26) (0.701) (Extended Data Table 14).

Because hierarchical Cell Ontology annotations are often unavailable (6), we next developed sc-Graph, as a new framework for quality assessment (**Methods**). For each set of cell embeddings, we define an affinity graph to elucidate the similarities between various cell types. scGraph then compares each affinity graph to a consensus graph, derived by aggregating individual graphs from different batches, based on raw reads or PCA loadings. This metrics efficiently highlights the inherent biological structures, emphasizing cell type similarities, while reducing the impact of technical variations across batches. Notably , scGraph does not require that any single batch have cells of all types (Fig. 1h), thus fitting the constraints of real datasets in the domain.

Evaluation by scGraph revealed varied performance across embeddings. Here, Islander had lower scores (Fig. 1i, Extended Data Table 15), while Harmony and scPoli excelled in capturing the complex relationships within functional cellular clusters. Indeed, Islander was the lowest scoring across seven of 11 atlases, underscoring scGraph's ability to detect the "drifting cell islands" artifact. Interestingly, scGraph favored higher-dimensional embeddings, like PCA over a PCA-derived UMAP. Note that scGraph's premise that profiles of functionally-similar cells would be proximate, may not always hold true.

In conclusion, we demonstrated the limitation of current quality metrics by introducing Islander, a three-layer perceptron, an integration approach that outperforms all major methods across diverse cell atlases, but at the cost of "island-like" distortions in the biological structures in cell embedding spaces. To address the inherent limitations of current evaluation metrics, we propose a new approach, scGraph, to helps assess how well the results of an integration approach preserve cell cell relationships at multiple granularities. Our work also highlights the importance of incorporating weaker supervision, as was recently illustrated in approaches where an encoder is regularized with

an additional reconstruction loss (6), or using large-scale unsupervised pre-training (*e.g.*, Universal Cell Embedding (UCE) (7)). These advancements underline the significance of methodological choices in computational biology and offer guiding principles for future research.

Methods

Datasets and pre-processing. Raw sequencing data were downloaded from the respective data providers as of October 1, 2023; details on the datasets and their sources are provided in Extended Data Table 1. The analysis encompasses 11 cell atlases, totaling 3,510,450 cell profiles. A uniform pre-processing protocol was applied across these datasets. Specifically, cell profiles with fewer than 1,000 reads or less than 500 detected genes were filtered out, and genes present in fewer than five cells were also excluded. Normalization was performed using Scanpy (30), scaling each cell's read counts to a total of 10,000 and subsequently applying a log1p transformation.

Baselines. Eleven baseline methods were used for comparison: three dimensionality reduction baselines: PCA, t-SNE (31), and UMAP (32); eight integration methods: Harmony (22), BBKNN (24), Scanorama (23), fastMNN (25), scVI (26), scANVI (27), scGen (28), and scPoli (29); and one pre-trained foundation model: Geneformer (5), for zero-shot embedding extraction. For dimensionality reduction methods, the log1p transformed raw counts from gene-by-cell matrices were provided as input. For each integration method, default hyperparameter settings recommended by the original authors were used. For Geneformer, the largest pre-trained model weights provided by the authors (33) were used. While scANVI, scGen, and scPoli utilize cell type as parts of their computational pipelines, other integration methods do not require such information. The top 1,000 genes were identified as highly variable gene sets.

Assessment metrics. Cell embeddings was assessed using metric as described in previous study (10), and implemented in "scib-metrics" (34). The following evaluation metrics were used (abbreviations noted are used in Extended Data Tables 3-13) . 'I-label' for isolated labels, 'L-NMI' for Leiden nor-malized mutual information, 'L-ARI' for Leiden Averaged Rand Index, 'K-NMI'/'K-ARI' for K-means NMI/ARI, 'S-label/batch' for silhouette label/batch, 'c/i-LISI' for batch-mixing (iLISI) and cell-type separation (cLISI), 'G-Con' for graph connectivity, and 'PCR' for principal component regression. Consistent with previous studies, selection of highly variable genes enhanced the performance of data integration methods.

Islander design. Islander is as a three-layer perceptron with two hidden layers of sizes 128 and 16, respectively, and an output layer matching the total number of cell types as annotated. The first hidden layer incorporates ReLU activation and batch normalization, while cell embeddings are derived from the second hidden layer. The output layer employs a softmax normalization function. In extended experiments, a decoder module was added mirroring the original MLP structure, with output dimensions of 16, 128, and the total number of genes, respectively. Each layer in this extended setup uses ReLU activation and batch normalization, except for the final linear layer.

Training setup. The model was trained in a manner aligned with scvi-tools (35), with mini-batches of 256 randomly sampled cells from all batches, along with their cell type annotations. Islander was trained using cross-entropy loss with mixup (13) augmentations (default setting). The Adam optimizer was used with an initial learning rate of 0.001, over 10 epochs, and a cosine annealing scheduler for learning rate decay. All cells were utilized for training to maximize overfitting.

Neighborhood calculation. Neighborhoods of each cell type were identified by the Euclidean distance between the centroids of cell profiles of each type in the embedding space. To mitigate the effects of batch variation and measurement noise, a trimming strategy was applied. The outermost 20% of data, treated as outliers, were excluded before calculating the centroid coordinates. This ensures a more accurate representation of cell type proximity by focusing on the most representative data points.

scGraph. scGraph quantifies the similarity between two graphs that each represent the closeness between cell types. In these graphs, each entry (x, y) signifies the proximity of cell type x to cell type y. The first graph is derived from the provided embeddings, while the second, serving as a reference, is based on raw counts or PCA loadings from each batch. For the reference, proximity graphs are initially computed from each batch using normalized Euclidean distances between centroids of the cell type profiles. These batch-specific graphs are then amalgamated into a single consensus graph through averaging. The similarity of neighborhoods for each cell type is assessed using Pearson's rank correlation. The final score, reflecting the overall similarity and ranging from -1 to 1 (with higher values indicating greater similarity), is the average across all cell types. The goal is to align the neighborhood graphs from the embeddings with the reference graph derived from the data, indicating that cells with similar profiles are appropriately clustered in the embedding space.

Code availability. The implementation code for the Islander, as well as the tutorial notebooks sufficient to reproduce the results presented in this manuscript, can be accessed via https://github.com/Genentech/Islander. For scIB evaluation pipelines, we use the implementations by Gayso et al from https://github.com/yoseflab/scib-metrics.

Acknowledgment

We thank Romain Lopez, Peng He, Leander Dony, Sara-Jane Dunn, Gocken Eraslan, Adam Gayoso, Graham Heimberg, Kexin Huang, John Marioni, Dana Pe'er, Yusuf Roohani, Yanay Rosen, Andrew Whitehead, and Jiaqi Zhang for invaluable insights, along with other members of the Leskovec and Regev labs and colleagues at the Human Cell Atlas, Chan Zuckerberg Initiative, and Google Deep-Mind for constructive discussions.

References

- [1] Bram Van de Sande, Joon Sang Lee, Euphemia Mutasa-Gottgens, et al. Applications of singlecell rna sequencing in drug discovery and development. *Nature Reviews Drug Discovery*, pages 1–25, 2023.
- [2] Martin Jinye Zhang, Kangcheng Hou, Kushal K Dey, et al. Polygenic enrichment distinguishes disease associations of individual cells in single-cell rna-seq data. *Nature Genetics*, 54(10):1572–1580, 2022.
- [3] Jennifer E Rood, Aidan Maartens, Anna Hupalowska, et al. Impact of the human cell atlas on medicine. *Nature Medicine*, 28(12):2486–2496, 2022.
- [4] Yuhan Hao, Stephanie Hao, et al. Integrated analysis of multimodal single-cell data. Cell, 2021.
- [5] Christina V Theodoris, Ling Xiao, Anant Chopra, et al. Transfer learning enables predictions in network biology. *Nature*, 2023.

- [6] Graham Heimberg, Tony Kuo, et al. Scalable querying of human cell atlases via a foundational model reveals commonalities across fibrosis-associated macrophages. *bioRxiv*, 2023.
- [7] Yanay Rosen, Yusuf Roohani, Ayush Agrawal, et al. Universal cell embeddings: A foundation model for cell biology. *bioRxiv*, 2023.
- [8] Haotian Cui, Chloe Wang, et al. scgpt: Towards building a foundation model for single-cell multi-omics using generative ai. *bioRxiv*, 2023.
- [9] Minsheng Hao et al. Large scale foundation model on single-cell transcriptomics. *bioRxiv*, 2023.
- [10] Malte D Luecken, Maren Büttner, Kridsadakorn Chaichoompu, et al. Benchmarking atlas-level data integration in single-cell genomics. *Nature Methods*, 19(1):41–50, 2022.
- [11] Tianyu Liu, Kexing Li, Yuge Wang, et al. Evaluating the utilities of large language models in single-cell data analysis. *bioRxiv*, 2023.
- [12] Kasia Zofia Kedzierska, Lorin Crawford, Ava Pardis Amini, et al. Assessing the limits of zero-shot foundation models in single-cell biology. *bioRxiv*, 2023.
- [13] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, et al. mixup: Beyond empirical risk minimization. In *ICLR*, 2018.
- [14] Kimberly Siletti, Rebecca Hodge, Alejandro Mossi Albiach, et al. Transcriptomic diversity of cell types across the adult human brain. *Science*, 382(6667):eadd7046, 2023.
- [15] Tapsi Kumar, Kevin Nee, Runmin Wei, et al. A spatially resolved single cell genomic atlas of the adult human breast. *Nature*, 2023.
- [16] Sean K Wang, Surag Nair, Rui Li, et al. Single-cell multiome of the human retina and deep learning nominate causal variants in complex eye diseases. *Cell Genomics*, 2(8), 2022.
- [17] Rasa Elmentaite, Alexander DB Ross, Kenny Roberts, et al. Single-cell sequencing of developing human gut reveals transcriptional links to childhood crohn's disease. *Developmental Cell*, 55(6):771–783, 2020.
- [18] Vincent R Knight-Schrijver, Hongorzul Davaapil, et al. A single-cell comparison of adult and fetal human epicardium defines the age-associated changes in epicardial activity. *Nature Cardiovascular Research*, 1(12):1215–1229, 2022.
- [19] Peng He, Kyungtae Lim, Dawei Sun, et al. A human fetal lung cell atlas uncovers proximal-distal gradients of differentiation and key regulators of epithelial fates. *Cell*, 185(25), 2022.
- [20] Llorenç Solé-Boldo, Günter Raddatz, Sabrina Schütz, et al. Single-cell transcriptomes of the human skin reveal age-related loss of fibroblast priming. *Communications Biology*, 3(1):188, 2020.
- [21] Lukas Heumos, Anna C Schaar, Christopher Lance, et al. Best practices for single-cell analysis across modalities. *Nature Reviews Genetics*, pages 1–23, 2023.
- [22] Ilya Korsunsky, Nghia Millard, Jean Fan, et al. Fast, sensitive and accurate integration of singlecell data with harmony. *Nature Methods*, 16(12):1289–1296, 2019.
- [23] Brian Hie, Bryan Bryson, and Bonnie Berger. Efficient integration of heterogeneous single-cell transcriptomes using scanorama. *Nature Biotechnology*, 37(6):685–691, 2019.
- [24] Krzysztof Polański, Matthew D Young, Zhichao Miao, et al. Bbknn: fast batch alignment of

single cell transcriptomes. Bioinformatics, 36(3):964-965, 2020.

- [25] Laleh Haghverdi et al. Batch effects in single-cell rna-sequencing data are corrected by matching mutual nearest neighbors. *Nature Biotechnology*, 36(5):421–427, 2018.
- [26] Romain Lopez, Jeffrey Regier, Michael B Cole, et al. Deep generative modeling for single-cell transcriptomics. *Nature Methods*, 15(12):1053–1058, 2018.
- [27] Chenling Xu, Romain Lopez, et al. Probabilistic harmonization and annotation of single-cell transcriptomics data with deep generative models. *Molecular Systems Biology*, 2021.
- [28] Mohammad Lotfollahi, F Alexander Wolf, and Fabian J Theis. scgen predicts single-cell perturbation responses. *Nature Methods*, 16(8):715–721, 2019.
- [29] Carlo De Donno, Soroor Hediyeh-Zadeh, et al. Population-level integration of single-cell datasets enables multi-scale analysis across samples. *Nature Methods*, 2023.
- [30] F Alexander Wolf, Philipp Angerer, and Fabian J Theis. Scanpy: large-scale single-cell gene expression data analysis. *Genome Biology*, 19:1–5, 2018.
- [31] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [32] Etienne Becht, Leland McInnes, John Healy, et al. Dimensionality reduction for visualizing single-cell data using umap. *Nature Biotechnology*, 37(1):38–44, 2019.
- [33] https://huggingface.co/ctheodoris/Geneformer/tree/main/ geneformer-12L-30M. Accessed: Oct 7, 2023.
- [34] https://scib-metrics.readthedocs.io/en/stable/. Accessed: Oct 1, 2023.
- [35] Adam Gayoso, Romain Lopez, Galen Xing, et al. A python library for probabilistic analysis of single-cell omics data. *Nature Biotechnology*, 40(2):163–166, 2022.
- [36] https://cellxgene.cziscience.com/collections/283d65eb-dd53-496d-adb7-7570c7caa44 Accessed: Oct 1, 2023.
- [37] https://cellxgene.cziscience.com/collections/4195ab4c-20bd-4cd3-8b3d-65601277e73 Accessed: Oct 1, 2023.
- [38] Yapeng Su, Daniel Chen, Christopher Lausted, et al. Multiomic immunophenotyping of covid-19 patients reveals early infection trajectories. *BioRxiv*, 2020.
- [39] https://atlas.fredhutch.org/fredhutch/covid/. Accessed: Oct 1, 2023.
- [40] https://cellxgene.cziscience.com/collections/348da6dc-5bf6-435d-adc5-37747b9ae38 Accessed: Oct 1, 2023.
- [41] https://cellxgene.cziscience.com/collections/17481d16-ee44-49e5-bcf0-28c0780d8c4 Accessed: Oct 1, 2023.
- [42] https://cellxgene.cziscience.com/collections/43b45a20-a969-49ac-a8e8-8c84b211bd0 Accessed: Oct 1, 2023.
- [43] Lisa Sikkema, Daniel C Strobl, Luke Zappia, et al. An integrated cell atlas of the human lung in health and disease. *Nature Medicine*, 2023.
- [44] https://cellxgene.cziscience.com/collections/6f6d381a-7701-4781-935c-db10d30de29 Accessed: Oct 1, 2023.

- [45] https://cellxgene.cziscience.com/collections/2d2e2acd-dade-489f-a2da-6c11aa65402 Accessed: Oct 1, 2023.
- [46] https://doi.org/10.6084/m9.figshare.12420968. Accessed: Oct 1, 2023.
- [47] https://cellxgene.cziscience.com/collections/c353707f-09a4-4f12-92a0-cb741e57e5f Accessed: Oct 1, 2023.

Supplementary

Atlas	# Gene	# Cell	# Class	# Batch	Reference
Brain	59,357	888,263	11	4	Paper (14), Data source (36)
Breast	33,234	703,512	39	126	Paper (15), Data source (37)
COVID	33,537	559,517	31	10	Paper (38), Data source (39)
Eye	36,484	51,645	11	8	Paper (16), Data source (40)
Gut (Fetal)	26,328	62,849	21	9	Paper (17), Data source (41)
Heart	33,234	486,134	27	14	Paper (18), Data source (42)
Lung	28,024	584,444	53	166	Paper (43), Data source (44)
Lung (Fetal, Donor)	26,354	71,752	144	29	Paper (19), Data source (45)
Lung (Fetal, Organoid)	24,653	70,495	28	37	Paper (19), Data source (45)
Pancreas	19,093	16,382	14	9	Paper (10), Data source (46)
Skin	30,933	15,457	13	5	Paper (20), Data source (47)

1. Data availability

Extended Data Table 1: Statistics of cell atlases. "# Class" represents the total number of cell types.

2. Performance

The embedding dimensions of each method is reported in Extended Data Table 2.

Method	PCA	TSNE	UMAP	Harmony	Scanorama	BBKNN	fastMNN
# Dim	50	2	2	50	100	2	50
Method	scVI	scANVI	scGen	scPoli	Geneformer	Islander	
# Dim	30	30	50	10	512	16	

Extended Data Table 2: Embedding dimensions of each method.

The detailed performance of each cell atlas is shown in Extended Data Tables 3-13. The best aggregated scores are highlighted in **bold**, and the term "Author's" denotes the author's integrated embeddings.

ore	Total	0.583	0.614	0.471	0.496	0.531	0.532	0.697	0.693	0.616	0.633	0.603	0.626	0.504	0.592	0.652	0.618	0.709	0.733	0.746	0.555	0.879
egate sc	Bio	0.714	0.752	0.538	0.555	0.574	0.593	0.741	0.741	0.712	0.710	0.599	0.606	0.352	0.688	0.774	0.701	0.787	0.808	0.763	0.624	0.982
Aggr	Batch	0.386	0.407	0.372	0.407	0.466	0.439	0.630	0.620	0.472	0.517	0.610	0.654	0.733	0.447	0.468	0.493	0.590	0.619	0.720	0.451	0.725
	PCR	0.000	0.000	0.000	0.000	0.322	0.000	0.895	0.702	0.318	0.465	0.764	0.848	0.941	0.154	0.000	0.367	0.464	0.589	0.678	0.000	0.879
ion	G-Con	0.966	0.899	0.978	0.988	0.985	0.992	0.957	0.881	0.925	0.728	0.987	0.992	0.400	0.983	0.950	0.985	0.983	0.968	0.988	0.851	0.988
1 correct	KBET	0.079	0.119	0.117	0.155	0.138	0.188	0.240	0.327	0.103	0.293	0.205	0.287	0.889	0.116	0.227	0.127	0.272	0.387	0.687	0.240	0.499
Batch	iLISI	0.005	0.094	0.005	0.121	0.112	0.241	0.142	0.249	0.112	0.174	0.308	0.366	0.470	0.099	0.230	0.123	0.345	0.251	0.398	0.259	0.353
	S-batch	0.882	0.921	0.759	0.773	0.774	0.777	0.914	0.943	0.902	0.926	0.787	0.779	0.963	0.884	0.932	0.865	0.887	0.900	0.850	0.905	0.906
	cLISI	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.975	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
	S-label	0.759	0.769	0.527	0.540	0.691	0.691	0.768	0.770	0.749	0.746	0.754	0.785	0.322	0.638	0.606	0.692	0.801	0.818	0.926	0.545	0.950
tion	K-ARI	0.439	0.442	0.258	0.274	0.302	0.329	0.437	0.443	0.439	0.311	0.349	0.354	0.032	0.326	0.449	0.374	0.597	0.634	0.547	0.352	1.000
conserva	K-NMI	0.745	0.746	0.593	0.618	0.667	0.669	0.749	0.746	0.745	0.674	0.687	0.689	0.152	0.677	0.724	0.709	0.791	0.817	0.814	0.583	0.999
Bio (L-ARI	0.473	0.653	0.046	0.052	0.052	0.054	0.619	0.597	0.489	0.624	0.053	0.055	0.061	0.590	0.960	0.555	0.608	0.670	0.353	0.576	0.995
	I-NMI	0.776	0.841	0.466	0.478	0.479	0.482	0.828	0.830	0.782	0.835	0.480	0.482	0.236	0.814	0.941	0.790	0.827	0.841	0.718	0.768	0.993
	I-label	0.803	0.817	0.873	0.927	0.824	0.925	0.786	0.800	0.781	0.778	0.868	0.880	0.683	0.772	0.738	0.785	0.888	0.878	0.983	0.543	0.940
	HVG		>		>		>		>		>		>	>		<		>	>	>		
	Method	PCA	PCA	TSNE	TSNE	UMAP	UMAP	Harmony	Harmony	Scanorama	Scanorama	BBKNN	BBKNN	fastMNN	scVI	scVI	scANVI	scANVI	scGen	scPoli	Geneformer	Islander

Extended Data Table 3: Benchmarking cell embeddings, on Brain Cell Atlas

Batch Bio Total
1 PCR B
BET G-Con
iLISI KB
S-batch
il cLISI
I S-label
AI K-ARI
RI K-NN
NMI L-A
el L-j
I-lab
HVG I-lab

Extended Data Table 4: Benchmarking cell embeddings, on Breast Cell Atlas

- 1			-	Bio c	conservat	tion			-	Batch	l correct	ion		Aggi	egate sc	ore
HVG I-label	I-label		L-NMI	L-ARI	K-NMI	K-ARI	S-label	cLISI	S-batch	iLISI	KBET	G-Con	PCR	Batch	Bio	Total
0.49	0.49!	10	0.698	0.669	0.563	0.226	0.520	0.996	0.921	0.122	0.215	0.751	0.000	0.402	0.595	0.518
/ 0.42	0.42		0.732	0.690	0.579	0.223	0.516	0.996	0.932	0.227	0.261	0.630	0.000	0.410	0.595	0.521
0.23	0.23	7	0.522	0.114	0.544	0.187	0.435	0.994	0.724	0.149	0.125	0.608	0.213	0.364	0.433	0.405
< 0.18	0.18	35	0.520	0.112	0.562	0.204	0.442	0.994	0.737	0.292	0.090	0.560	0.152	0.366	0.431	0.405
0.22	0.22	4	0.519	0.125	0.560	0.221	0.487	0.990	0.743	0.353	0.172	0.561	0.088	0.383	0.447	0.421
0.15	0.15	90	0.523	0.123	0.576	0.221	0.498	0.991	0.748	0.407	0.134	0.533	0.000	0.364	0.441	0.410
0.4	0.4	85	0.697	0.674	0.563	0.218	0.522	0.996	0.936	0.319	0.411	0.711	0.609	0.597	0.593	0.595
/ 0.4	0.4	20	0.725	0.668	0.578	0.220	0.517	0.995	0.938	0.320	0.457	0.590	0.514	0.564	0.589	0.579
0.4	0.4	43	0.706	0.669	0.565	0.226	0.517	0.995	0.932	0.270	0.257	0.659	0.345	0.493	0.589	0.550
< 0.4	0.4	-39	0.737	0.690	0.586	0.233	0.515	0.996	0.937	0.285	0.349	0.373	0.330	0.455	0.599	0.542
0.2	0.2	214	0.518	0.119	0.568	0.208	0.497	0.990	0.761	0.414	0.270	0.589	0.252	0.457	0.445	0.450
> 0.0	0.	128	0.533	0.132	0.581	0.216	0.503	0.992	0.758	0.442	0.188	0.546	0.052	0.397	0.441	0.423
< 0.4	0	431	0.732	0.683	0.584	0.238	0.521	0.996	0.929	0.279	0.353	0.584	0.340	0.497	0.598	0.558
0	0	501	0.703	0.657	0.588	0.261	0.519	0.997	0.933	0.228	0.259	0.788	0.404	0.522	0.604	0.571
< 0.	o.	544	0.718	0.660	0.541	0.205	0.502	0.991	0.933	0.307	0.437	0.685	0.466	0.565	0.595	0.583
0	o.	439	0.749	0.719	0.622	0.267	0.544	0.999	0.920	0.256	0.272	0.823	0.307	0.516	0.620	0.578
<. 0.	<u>.</u>	559	0.762	0.739	0.632	0.269	0.554	0.998	0.906	0.304	0.384	0.774	0.278	0.529	0.645	0.598
~ .0	ö	453	0.779	0.768	0.638	0.294	0.553	0.999	0.937	0.327	0.511	0.759	0.292	0.565	0.641	0.611
0	0	231	0.770	0.741	0.666	0.352	0.584	0.998	0.851	0.364	0.494	0.826	0.303	0.568	0.620	0.599
✓ 0.	<u>о</u>	411	0.715	0.567	0.662	0.357	0.572	0.997	0.801	0.418	0.583	0.769	0.242	0.562	0.612	0.592
0	0	.693	0.999	1.000	0.863	0.603	0.713	1.000	0.934	0.420	0.653	0.967	0.693	0.733	0.839	0.797

Extended Data Table 5: Benchmarking cell embeddings, on COVID Cell Atlas

core	Total	0.528	0.563	0.387	0.417	0.436	0.471	0.670	0.698	0.634	0.642	0.506	0.500	0.317	0.656	0.655	0.720	0.667	0.734	0.823	0.743	0.623	0.504	0.899
regate s	Bio	0.633	0.678	0.478	0.494	0.526	0.565	0.697	0.786	0.700	0.733	0.535	0.577	0.266	0.650	0.688	0.759	0.701	0.801	0.924	0.820	0.607	0.541	0.972
Agg	Batch	0.371	0.390	0.250	0.302	0.302	0.331	0.629	0.566	0.535	0.506	0.463	0.386	0.394	0.666	0.607	0.662	0.616	0.634	0.672	0.627	0.647	0.449	0.790
	PCR	0.000	0.000	0.000	0.000	0.000	0.000	0.855	0.626	0.621	0.459	0.571	0.175	0.586	0.861	0.594	0.866	0.725	0.655	0.849	0.570	0.740	0.000	0.849
tion	G-Con	0.751	0.662	0.357	0.538	0.572	0.603	0.802	0.643	0.730	0.504	0.628	0.600	0.050	0.916	0.786	0.918	0.630	0.805	0.930	0.914	0.814	0.789	0.979
1 correct	kBET	060.0	0.207	0.015	0.031	0.031	0.097	0.227	0.273	0.187	0.341	0.079	0.127	0.129	0.252	0.307	0.237	0.382	0.363	0.365	0.525	0.386	0.211	0.707
Batcl	iLISI	0.108	0.175	0.105	0.164	0.121	0.180	0.360	0.389	0.236	0.320	0.275	0.293	0.317	0.403	0.431	0.400	0.459	0.455	0.370	0.360	0.503	0.316	0.492
	S-batch	0.905	0.903	0.770	0.777	0.785	0.773	0.900	0.900	0.903	0.907	0.762	0.734	0.888	0.899	0.916	0.890	0.886	0.893	0.849	0.767	0.791	0.928	0.922
	cLISI	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.844	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
	S-label	0.687	0.705	0.521	0.551	0.594	0.707	0.721	0.718	0.678	0.695	0.706	0.788	0.492	0.565	0.526	0.624	0.650	0.764	0.889	0.915	0.808	0.559	0.920
tion	K-ARI	0.414	0.415	0.177	0.203	0.309	0.389	0.393	0.546	0.551	0.319	0.238	0.378	0.020	0.171	0.237	0.528	0.433	0.547	0.983	0.651	0.400	0.218	0.996
conserva	K-NMI	0.691	0.701	0.481	0.512	0.583	0.625	0.701	0.773	0.754	0.675	0.544	0.643	0.004	0.432	0.619	0.729	0.776	0.777	0.975	0.852	0.710	0.443	0.982
Bio e	L-ARI	0.333	0.550	0.086	0.088	0.087	0.095	0.649	0.967	0.548	0.964	0.094	0.093	0.009	0.968	0.965	0.972	0.591	0.968	0.980	0.632	0.102	0.439	1.000
	IMN-1	0.677	0.742	0.529	0.541	0.538	0.550	0.769	0.855	0.748	0.836	0.545	0.541	0.007	0.850	0.845	0.862	0.775	0.862	0.924	0.819	0.582	0.617	1.000
	I-label	0.629	0.637	0.553	0.563	0.569	0.587	0.644	0.641	0.622	0.639	0.617	0.594	0.483	0.560	0.620	0.597	0.685	0.689	0.716	0.872	0.649	0.513	0.906
	ЫVG		>		>		>		>		>		>	>		>		>	>		>			
	Method	PCA	PCA	TSNE	TSNE	UMAP	UMAP	Harmony	Harmony	Scanorama	Scanorama	BBKNN	BBKNN	fastMNN	scVI	scVI	scANVI	scANVI	scGen	scPoli	scPoli	Author's	Geneformer	Islander

Extended Data Table 6: Benchmarking cell embeddings, on Eye Cell Atlas

core	Total	0.484	0.499	0.410	0.457	0.434	0.437	0.612	0.607	0.517	0.557	0.527	0.496	0.432	0.593	0.630	0.612	0.645	0.674	0.602	0.661	0.520	0.475	0.811
regate s	Bio	0.551	0.578	0.515	0.527	0.522	0.532	0.570	0.615	0.570	0.601	0.541	0.563	0.437	0.579	0.625	0.632	0.672	0.734	0.664	0.685	0.581	0.465	0.944
Agg	Batch	0.385	0.381	0.253	0.352	0.301	0.294	0.674	0.595	0.436	0.490	0.506	0.395	0.424	0.613	0.637	0.583	0.605	0.585	0.510	0.625	0.428	0.489	0.612
	PCR	0.000	0.000	0.000	0.459	0.000	0.000	0.746	0.561	0.000	0.333	0.292	0.000	0.343	0.667	0.787	0.526	0.611	0.399	0.002	0.497	0.000	0.362	0.303
tion	G-Con	0.911	0.883	0.463	0.532	0.661	0.683	0.890	0.838	0.897	0.765	0.795	0.763	0.592	0.937	0.923	0.933	0.927	0.959	0.959	0.921	0.920	0.719	0.959
1 correct	kBET	0.111	0.098	0.108	0.069	0.153	0.086	0.551	0.429	0.249	0.269	0.441	0.248	0.187	0.338	0.345	0.339	0.390	0.477	0.489	0.600	0.260	0.349	0.565
Batcl	iLISI	0.009	0.015	0.000	0.003	0.007	0.020	0.282	0.222	0.118	0.157	0.281	0.203	0.107	0.207	0.209	0.209	0.213	0.228	0.231	0.271	0.071	0.133	0.313
	S-batch	0.893	0.908	0.693	0.695	0.686	0.681	0.902	0.925	0.918	0.925	0.720	0.760	0.892	0.915	0.920	0.906	0.885	0.862	0.869	0.835	0.888	0.881	0.922
	cLISI	1.000	1.000	1.000	1.000	0.999	0.999	0.998	0.998	1.000	1.000	1.000	0.999	0.963	0.999	0.998	1.000	1.000	1.000	1.000	1.000	1.000	0.986	1.000
	S-label	0.534	0.539	0.423	0.453	0.394	0.443	0.550	0.535	0.538	0.539	0.504	0.513	0.467	0.531	0.522	0.558	0.590	0.620	0.639	0.596	0.561	0.477	0.833
tion	K-ARI	0.292	0.323	0.226	0.249	0.223	0.238	0.299	0.330	0.322	0.354	0.305	0.283	0.188	0.255	0.352	0.352	0.392	0.424	0.386	0.416	0.358	0.128	0.946
conserva	K-NMI	0.612	0.654	0.576	0.596	0.564	0.570	0.585	0.630	0.619	0.671	0.650	0.648	0.347	0.594	0.668	0.687	0.722	0.753	0.742	0.739	0.638	0.334	0.966
Bio c	L-ARI	0.213	0.290	0.096	0.118	0.163	0.180	0.412	0.509	0.330	0.371	0.255	0.227	0.132	0.422	0.510	0.509	0.573	0.779	0.603	0.526	0.322	0.280	1.000
	L-NMI	0.614	0.644	0.561	0.577	0.584	0.606	0.687	0.718	0.694	0.701	0.662	0.645	0.359	0.713	0.738	0.787	0.791	0.901	0.815	0.759	0.718	0.534	1.000
	I-label	0.593	0.600	0.723	0.699	0.727	0.690	0.460	0.582	0.488	0.574	0.411	0.629	0.605	0.541	0.587	0.528	0.635	0.660	0.464	0.760	0.473	0.518	0.859
	HVG		>		>		>		>		>		>	>		>		>	>		>			
	Method	PCA	PCA	TSNE	TSNE	UMAP	UMAP	Harmony	Harmony	Scanorama	Scanorama	BBKNN	BBKNN	fastMNN	scVI	scVI	scANVI	scANVI	scGen	scPoli	scPoli	Author's	Geneformer	Islander

Extended Data Table 7: Benchmarking cell embeddings, on Fetal Gut Atlas

core	Total	0.478	0.523	0.379	0.404	0.401	0.434	0.638	0.580	0.495	0.555	0.531	0.472	0.281	0.552	0.607	0.576	0.633	0.670	0.623	0.653	0.640	0.477	0.767
regate s	Bio	0.584	0.650	0.493	0.508	0.518	0.517	0.757	0.648	0.598	0.653	0.577	0.543	0.276	0.680	0.649	0.717	0.713	0.748	0.773	0.727	0.771	0.547	0.867
Agg	Batch	0.319	0.332	0.208	0.247	0.224	0.309	0.459	0.478	0.339	0.409	0.463	0.367	0.290	0.360	0.545	0.364	0.514	0.553	0.398	0.541	0.444	0.373	0.616
	PCR	0.000	0.000	0.000	0.078	0.000	0.274	0.000	0.464	0.000	0.364	0.000	0.273	0.416	0.000	0.743	0.000	0.583	0.571	0.000	0.699	0.000	0.094	0.353
tion	G-Con	0.702	0.694	0.396	0.443	0.453	0.494	0.817	0.691	0.682	0.553	0.699	0.556	0.050	0.717	0.743	0.693	0.721	0.913	0.798	0.751	0.825	0.617	0.969
n correct	KBET	0.048	0.077	0.034	0.048	0.042	0.064	0.333	0.189	0.085	0.141	0.464	0.127	0.050	0.121	0.172	0.151	0.203	0.252	0.239	0.317	0.282	0.097	0.529
Batch	iLISI	0.001	0.034	0.001	0.034	0.005	0.089	0.255	0.143	0.041	0.084	0.377	0.184	0.071	0.060	0.129	0.086	0.163	0.190	0.127	0.159	0.239	0.149	0.303
	S-batch	0.841	0.856	0.610	0.631	0.620	0.627	0.891	0.904	0.888	0.903	0.774	0.694	0.863	0.904	0.935	0.892	0.899	0.837	0.824	0.777	0.873	0.906	0.927
	cLISI	0.998	0.996	0.998	0.997	0.994	0.994	0.996	0.995	0.999	0.995	0.994	0.994	0.892	0.997	0.995	0.998	0.998	1.000	0.998	0.997	0.996	0.993	1.000
	S-label	0.532	0.539	0.420	0.443	0.451	0.544	0.647	0.566	0.530	0.546	0.631	0.570	0.465	0.561	0.532	0.589	0.642	0.634	0.693	0.680	0.648	0.476	0.751
tion	K-ARI	0.360	0.422	0.197	0.215	0.244	0.268	0.764	0.412	0.349	0.368	0.286	0.261	0.004	0.385	0.257	0.494	0.459	0.527	0.602	0.436	0.781	0.255	0.658
conserva	K-NMI	0.598	0.616	0.507	0.539	0.544	0.570	0.729	0.617	0.608	0.606	0.648	0.599	0.036	0.620	0.602	0.682	0.674	0.734	0.737	0.704	0.740	0.463	0.885
Bio c	L-ARI	0.417	0.738	0.100	0.110	0.117	0.126	0.782	0.763	0.475	0.761	0.151	0.137	0.011	0.759	0.780	0.785	0.783	0.858	0.781	0.641	0.795	0.547	0.992
	L-NMI	0.611	0.687	0.532	0.542	0.524	0.537	0.797	0.725	0.651	0.697	0.588	0.557	0.066	0.728	0.742	0.754	0.741	0.842	0.754	0.752	0.805	0.594	0.991
	I-label	0.571	0.551	0.700	0.712	0.755	0.581	0.584	0.458	0.578	0.598	0.739	0.680	0.456	0.705	0.638	0.716	0.693	0.644	0.847	0.880	0.630	0.503	0.790
	HVG		>		>		>		>		>		>	>		>		>	>		>			
	Method	PCA	PCA	TSNE	TSNE	UMAP	UMAP	Harmony	Harmony	Scanorama	Scanorama	BBKNN	BBKNN	fastMNN	scVI	scVI	scANVI	scANVI	scGen	scPoli	scPoli	Author's	Geneformer	Islander

Extended Data Table 8: Benchmarking cell embeddings, on Heart Cell Atlas

te score	o Total	64 0.502	96 0.524	05 0.403	25 0.423	03 0.419	31 0.443	15 0.536	37 0.541	76 0.491	36 0.442	67 0.479	67 0.308	61 0.602	61 0.602	16 0.636	16 0.631	02 0.680	56 0.643	55 0.640	93 0.619	50 0.510	84 0.761
Aggrega	atch Bi	407 0.5	416 0.5	249 0.5	271 0.5	293 0.5	311 0.5	566 0.5	546 0.5	364 0.5	301 0.5	345 0.5	370 0.2	514 0.6	513 0.6	515 0.7	503 0.7	499 0.8	475 0.7	468 0.7	509 0.6	451 0.5	577 0.8
	PCR B	0.000 0.	0.000 0.	0.000 0.	0.000 0.	0.000 0.	0.000 0.	0.675 0.	0.625 0.	0.000 0.	0.000 0.	0.000 0.	0.999 0.	0.383 0.	0.375 0.	0.371 0.	0.301 0.	0.183 0.	0.184 0.	0.150 0.	0.354 0.	0.199 0.	0.398 0.
tion	G-Con	0.852	0.867	0.448	0.556	0.562	0.645	0.765	0.673	0.524	0.622	0.728	0.003	0.931	0.931	0.937	0.913	0.894	0.909	0.909	0.950	0.788	0.953
1 correc	kBET	0.317	0.322	0.271	0.275	0.316	0.338	0.561	0.573	0.392	0.402	0.514	0.192	0.415	0.389	0.443	0.508	0.671	0.658	0.658	0.370	0.407	0.724
Batcł	iLISI	0.002	0.003	0.002	0.004	0.007	0.009	0.023	0.024	0.006	0.023	0.027	0.062	0.008	0.008	0.012	0.017	0.029	0.025	0.025	0.006	0.006	0.034
	S-batch	0.867	0.889	0.526	0.519	0.583	0.565	0.805	0.834	0.900	0.459	0.457	0.593	0.835	0.861	0.813	0.778	0.717	0.597	0.597	0.865	0.857	0.774
	cLISI	1.000	1.000	1.000	1.000	0.999	0.999	0.996	0.996	1.000	0.998	0.999	0.877	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.999	1.000
	S-label	0.537	0.545	0.399	0.443	0.400	0.455	0.489	0.499	0.527	0.505	0.554	0.490	0.561	0.550	0.595	0.620	0.673	0.662	0.662	0.553	0.500	0.793
tion	K-ARI	0.313	0.352	0.214	0.225	0.222	0.275	0.199	0.254	0.336	0.272	0.330	0.000	0.339	0.381	0.516	0.487	0.588	0.609	0.633	0.401	0.258	0.659
conserva	K-NMI	0.643	0.690	0.584	0.611	0.603	0.654	0.464	0.509	0.634	0.666	0.702	0.001	0.680	0.679	0.769	0.763	0.842	0.796	0.799	0.715	0.588	0.896
Bio	L-ARI	0.288	0.347	0.146	0.173	0.201	0.225	0.385	0.417	0.271	0.255	0.285	0.000	0.644	0.628	0.742	0.770	0.968	0.760	0.709	0.746	0.294	0.999
	L-NMI	0.653	0.693	0.627	0.648	0.630	0.667	0.609	0.608	0.663	0.667	0.702	0.001	0.799	0.797	0.843	0.845	0.967	0.835	0.830	0.855	0.651	0.999
	I-label	0.516	0.542	0.569	0.575	0.465	0.437	0.465	0.478	0.599	0.388	0.399	0.502	0.602	0.590	0.547	0.530	0.573	0.628	0.655	0.582	0.559	0.842
	HVG		>		>		>		>			>	>		>		>	>		>	>		
	Method	PCA	PCA	TSNE	TSNE	UMAP	UMAP	Harmony	Harmony	Scanorama	BBKNN	BBKNN	fastMNN	scVI	scVI	scANVI	scANVI	scGen	scPoli	scPoli	Author's	Geneformer	Islander

Extended Data Table 9: Benchmarking cell embeddings, on Lung Cell Atlas

				Bio	conserva	tion				Batc	n correct	ion		Aggr	egate sc	ore
Method	HVG	I-label	IMN-T	L-ARI	K-NMI	K-ARI	S-label	cLISI	S-batch	iLISI	KBET	G-Con	PCR	Batch	Bio	Total
PCA		0.571	0.809	0.581	0.723	0.237	0.532	1.000	0.834	0.043	0.620	0.821	0.000	0.464	0.636	0.567
PCA	>	0.581	0.799	0.619	0.716	0.191	0.535	0.999	0.868	0.059	0.640	0.797	0.000	0.473	0.634	0.570
TSNE		0.583	0.762	0.318	0.720	0.160	0.499	1.000	0.542	0.042	0.480	0.663	0.000	0.345	0.577	0.484
TSNE	>	0.585	0.767	0.350	0.716	0.157	0.504	1.000	0.568	0.059	0.509	0.693	0.000	0.366	0.582	0.496
UMAP		0.580	0.765	0.387	0.713	0.162	0.489	0.999	0.576	0.068	0.572	0.708	0.000	0.385	0.585	0.505
UMAP	>	0.548	0.771	0.404	0.713	0.162	0.524	0.999	0.586	0.085	0.560	0.707	0.000	0.388	0.589	0.508
Harmony		0.542	0.706	0.324	0.640	0.188	0.482	0.996	0.823	0.142	0.808	0.764	0.281	0.564	0.554	0.558
Harmony	>	0.566	0.652	0.300	0.577	0.098	0.483	0.995	0.859	0.131	0.672	0.613	0.329	0.521	0.524	0.523
Scanorama		0.491	0.814	0.675	0.708	0.221	0.529	1.000	0.852	0.083	0.728	0.805	0.000	0.494	0.634	0.578
Scanorama	>	0.565	0.778	0.657	0.719	0.195	0.528	0.999	0.896	0.070	0.697	0.646	0.000	0.462	0.635	0.565
BBKNN		0.391	0.740	0.340	0.692	0.148	0.527	0.998	0.578	0.155	0.771	0.715	0.000	0.444	0.548	0.506
BBKNN	>	0.541	0.746	0.382	0.689	0.147	0.509	0.998	0.598	0.124	0.625	0.651	0.000	0.400	0.573	0.504
fastMNN	>	0.397	0.181	0.036	0.180	0.015	0.375	0.968	0.781	0.157	0.095	0.048	0.557	0.328	0.307	0.316
scVI		0.549	0.706	0.351	0.629	0.125	0.513	0.996	0.826	0.133	0.835	0.843	0.399	0.607	0.553	0.574
scVI	>	0.579	0.709	0.335	0.640	0.124	0.505	0.996	0.855	0.129	0.774	0.849	0.646	0.650	0.555	0.593
scANVI		0.525	0.773	0.544	0.671	0.169	0.526	0.998	0.813	0.123	0.839	0.846	0.143	0.553	0.601	0.582
scANVI	>	0.574	0.826	0.706	0.732	0.215	0.552	1.000	0.820	0.112	0.799	0.854	0.300	0.577	0.658	0.625
scGen	>	0.464	0.851	0.699	0.726	0.437	0.553	1.000	0.736	0.134	0.694	0.838	0.000	0.480	0.676	0.598
scPoli		0.453	0.858	0.647	0.787	0.285	0.618	1.000	0.735	0.137	0.861	0.900	0.000	0.526	0.664	0.609
scPoli	>	0.616	0.791	0.530	0.742	0.311	0.545	0.998	0.681	0.141	0.789	0.773	0.274	0.532	0.648	0.601
Author's		0.575	0.844	0.561	0.774	0.347	0.567	1.000	0.834	0.070	0.780	0.897	0.000	0.516	0.667	0.607
Geneformer		0.492	0.640	0.304	0.520	0.107	0.475	0.996	0.829	0.114	0.672	0.624	0.410	0.530	0.505	0.515
Islander (Run1)		0.818	0.999	1.000	0.901	0.449	0.793	1.000	0.854	0.124	0.889	0.972	0.240	0.616	0.851	0.757
Islander (Run2)		0.824	0.999	1.000	0.891	0.406	0.793	1.000	0.853	0.123	0.883	0.970	0.217	0.609	0.845	0.751
Islander (Run3)		0.817	0.999	1.000	0.894	0.440	0.794	1.000	0.854	0.123	0.888	0.970	0.249	0.617	0.849	0.756
		Extend	led Data	Table 1	0: Bencl	ımarkin	g cell en	nbeddin	gs, on (L	ung, Fet	cal, Don	or) Cell	Atlas			

17

				Bio	conserva	tion				Batch	i correct	ion		Aggi	regate sc	ore
Method	HVG	I-label	IMN-1	L-ARI	K-NMI	K-ARI	S-label	cLISI	S-batch	iLISI	kBET	G-Con	PCR	Batch	Bio	Total
PCA		0.466	0.441	0.264	0.399	0.126	0.512	0.992	0.791	0.000	0.486	0.728	0.000	0.401	0.457	0.435
PCA	>	0.439	0.451	0.235	0.402	0.125	0.506	0.992	0.821	0.000	0.521	0.738	0.000	0.416	0.450	0.436
TSNE		0.281	0.378	0.061	0.381	0.093	0.430	0.991	0.522	0.000	0.341	0.513	0.000	0.275	0.374	0.334
TSNE	>	0.313	0.378	0.068	0.391	0.095	0.420	0.990	0.564	0.000	0.332	0.558	0.000	0.291	0.379	0.344
UMAP		0.294	0.396	0.105	0.396	0.102	0.429	0.985	0.546	0.000	0.402	0.587	0.000	0.307	0.387	0.355
UMAP	>	0.382	0.393	0.097	0.398	0.111	0.468	0.986	0.578	0.000	0.491	0.647	0.000	0.343	0.405	0.380
Harmony		0.546	0.324	0.177	0.267	0.084	0.465	0.982	0.793	0.034	0.607	0.466	0.276	0.435	0.406	0.418
Harmony	>	0.484	0.336	0.148	0.215	0.045	0.428	0.982	0.870	0.017	0.621	0.608	0.712	0.565	0.377	0.452
Scanorama		0.532	0.441	0.298	0.359	0.105	0.484	0.991	0.837	0.000	0.477	0.698	0.000	0.403	0.459	0.436
Scanorama	>	0.442	0.430	0.249	0.386	0.119	0.488	0.991	0.824	0.000	0.548	0.595	0.081	0.410	0.444	0.430
BBKNN		0.434	0.374	0.094	0.363	0.091	0.438	0.981	0.593	0.036	0.462	0.632	0.000	0.345	0.396	0.376
BBKNN	>	0.396	0.346	0.079	0.357	0.090	0.367	0.980	0.626	0.017	0.323	0.336	0.000	0.260	0.374	0.328
fastMNN	>	0.440	0.254	0.111	0.161	0.035	0.444	0.955	0.846	0.006	0.106	0.131	0.463	0.310	0.343	0.330
scVI		0.485	0.344	0.251	0.164	0.032	0.457	0.977	0.870	0.024	0.581	0.706	0.821	0.600	0.387	0.473
scVI	>	0.517	0.283	0.136	0.183	0.040	0.452	0.970	0.880	0.031	0.557	0.722	0.934	0.625	0.369	0.471
scANVI		0.473	0.471	0.360	0.320	0.095	0.458	0.992	0.854	0.018	0.576	0.708	0.622	0.555	0.453	0.494
scANVI	>	0.508	0.455	0.347	0.366	0.111	0.453	0.990	0.832	0.013	0.588	0.751	0.737	0.584	0.462	0.511
scGen	>	0.572	0.908	0.890	0.649	0.243	0.642	1.000	0.714	0.019	0.663	0.895	0.364	0.531	0.701	0.633
scPoli		0.631	0.477	0.301	0.415	0.115	0.485	0.993	0.734	0.010	0.581	0.845	0.000	0.434	0.488	0.466
scPoli	>	0.479	0.484	0.298	0.427	0.149	0.511	0.989	0.666	0.010	0.610	0.759	0.113	0.432	0.477	0.459
Author's		0.647	0.480	0.272	0.442	0.192	0.555	0.995	0.796	0.002	0.527	0.769	0.000	0.419	0.512	0.475
Geneformer		0.436	0.337	0.208	0.242	0.074	0.477	0.980	0.828	0.012	0.466	0.457	0.598	0.472	0.394	0.425
Islander		0.785	0.955	0.890	0.736	0.313	0.752	1.000	0.844	0.027	0.687	0.978	0.000	0.507	0.776	0.668
		Exten	ded Dat	a Table [11: Benc	hmarkin	g cell em	ibedding	gs, on (Lı	ıng, Feta	ıl, Orgaı	oid) Ce	ll Atlas			

Bio conservation IVG I-label L-NMI K-ARI K-ARI	Bio conservation I-label L-NMI L-ARI K-NMI K-ARI	Bio conservation L-NMI L-ARI K-NMI K-ARI	Bio conservation	conservation K-NMI K-ARI	tion K-ARI		S-label	cLISI	S-batch	Batch	correct kBET	ion G-Con	PCR	Agg1 Batch	regate sc Bio	tore Total
					TTATAT_NT_NT		ט-זמטכו		ט-טמוכוו	IULI	NULLI		ז כוו	חמורוו		TOLAT
0.496 0.594 0.214 0	0.496 0.594 0.214 0	0.594 0.214 0	0.214 0	-	0.459	0.313	0.499	1.000	0.770	0.000	0.057	0.553	0.000	0.276	0.511	0.417
\[\lambda 0.634 0.685 0.421 0 \]	0.634 0.685 0.421 0	0.685 0.421 0	0.421 0	0	.620	0.408	0.572	1.000	0.849	0.002	0.206	0.656	0.000	0.343	0.620	0.509
0.363 0.573 0.161 0.	0.363 0.573 0.161 0.	0.573 0.161 0.	0.161 0.	0.	438	0.203	0.417	1.000	0.482	0.000	0.030	0.424	0.000	0.187	0.451	0.345
\[\lambda 0.954 0.634 0.236 0.1 \]	0.954 0.634 0.236 0.1	0.634 0.236 0.1	0.236 0.1	0	589	0.340	0.451	1.000	0.546	0.000	0.060	0.550	0.000	0.231	0.601	0.453
0.372 0.580 0.211 0.5	0.372 0.580 0.211 0.5	0.580 0.211 0.3	0.211 0.3	0	387	0.167	0.386	1.000	0.511	0.000	0.036	0.555	0.000	0.220	0.443	0.354
V 0.932 0.630 0.240 0.5	0.932 0.630 0.240 0.5	0.630 0.240 0.5	0.240 0.5	0.5	89	0.343	0.496	1.000	0.570	0.000	0.114	0.587	0.000	0.254	0.604	0.464
0.553 0.869 0.908 0.6	0.553 0.869 0.908 0.6	0.869 0.908 0.6	0.908 0.6	0.6	23	0.426	0.570	1.000	0.885	0.198	0.421	0.853	0.918	0.655	0.707	0.686
V 0.586 0.904 0.945 0.79	0.586 0.904 0.945 0.79	0.904 0.945 0.79	0.945 0.79	0.7	92	0.647	0.619	1.000	0.868	0.220	0.538	0.765	0.790	0.636	0.785	0.725
0.555 0.754 0.564 0.61	0.555 0.754 0.564 0.61	0.754 0.564 0.61	0.564 0.61	0.61	4	0.392	0.556	1.000	0.905	0.037	0.278	0.863	0.674	0.551	0.633	0.601
\[\lambda 0.609 0.878 0.916 0.76 \]	0.609 0.878 0.916 0.76	0.878 0.916 0.76	0.916 0.76	0.76	2	0.598	0.597	1.000	0.934	0.144	0.384	0.694	0.673	0.566	0.766	0.686
0.473 0.623 0.271 0.62	0.473 0.623 0.271 0.62	0.623 0.271 0.62	0.271 0.62	0.62	9	0.367	0.547	1.000	0.660	0.023	0.190	0.676	0.671	0.444	0.558	0.512
V 0.923 0.681 0.319 0.723	0.923 0.681 0.319 0.723	0.681 0.319 0.723	0.319 0.723	0.723	~	0.461	0.695	1.000	0.669	0.062	0.217	0.850	0.628	0.485	0.686	0.606
✓ 0.429 0.099 0.023 0.03!	0.429 0.099 0.023 0.03	0.099 0.023 0.03	0.023 0.035	0.035	10	0.015	0.475	0.841	0.792	0.190	0.004	0.150	0.853	0.398	0.274	0.323
0.700 0.880 0.916 0.63	0.700 0.880 0.916 0.63	0.880 0.916 0.63	0.916 0.63	0.63	ы	0.442	0.558	0.999	0.851	0.209	0.424	0.913	0.964	0.672	0.733	0.709
V 0.664 0.920 0.953 0.74:	0.664 0.920 0.953 0.74	0.920 0.953 0.74	0.953 0.74:	0.74;	~	0.451	0.573	1.000	0.881	0.277	0.527	0.910	0.856	0.690	0.757	0.731
0.720 0.917 0.949 0.73	0.720 0.917 0.949 0.73	0.917 0.949 0.73	0.949 0.73	0.73	2	0.712	0.572	1.000	0.844	0.207	0.411	0.923	0.955	0.668	0.800	0.747
 0.664 0.920 0.953 0.74 	0.664 0.920 0.953 0.74	0.920 0.953 0.74	0.953 0.74	0.7	43	0.451	0.573	1.000	0.881	0.277	0.527	0.910	0.856	0.690	0.757	0.731
✓ 0.813 0.960 0.981 0.83	0.813 0.960 0.981 0.83	0.960 0.981 0.83	0.981 0.83	0.83	33	0.638	0.700	1.000	0.777	0.286	0.637	0.938	0.754	0.678	0.847	0.779
0.786 0.919 0.949 0.78	0.786 0.919 0.949 0.78	0.919 0.949 0.78	0.949 0.78	0.78	33	0.541	0.701	1.000	0.822	0.192	0.485	0.931	0.871	0.660	0.811	0.751
V 0.682 0.912 0.945 0.8	0.682 0.912 0.945 0.8	0.912 0.945 0.8	0.945 0.8	0.8	29	0.733	0.760	1.000	0.758	0.344	0.677	0.939	0.757	0.695	0.837	0.780
0.893 0.999 0.999 0.9	0.893 0.999 0.999 0.9	0.999 0.999 0.9	0.999 0.9	0.9	80	0.984	0.914	1.000	0.873	0.254	0.558	0.977	0.897	0.712	0.967	0.865

Extended Data Table 12: Benchmarking cell embeddings, on Pancreas Cell Atlas

core	Total	0.542	0.551	0.494	0.522	0.480	0.523	0.708	0.718	0.597	0.631	0.682	0.615	0.657	0.668	0.733	0.693	0.776	0.738	0.715	0.754	0.700	0.516	0.812	
regate s	Bio	0.665	0.683	0.620	0.650	0.623	0.669	0.784	0.754	0.696	0.712	0.731	0.726	0.727	0.687	0.768	0.747	0.839	0.773	0.826	0.812	0.794	0.585	0.950	
Agg	Batch	0.359	0.352	0.304	0.330	0.265	0.305	0.595	0.664	0.447	0.510	0.607	0.448	0.552	0.639	0.681	0.612	0.682	0.685	0.548	0.665	0.561	0.412	0.605	
	PCR	0.000	0.000	0.313	0.364	0.000	0.000	0.000	0.753	0.000	0.562	0.000	0.244	0.682	0.683	0.833	0.537	0.784	0.598	0.000	0.491	0.000	0.055	0.000	
tion	G-Con	0.874	0.814	0.521	0.592	0.684	0.782	0.891	0.766	0.829	0.516	0.954	0.895	0.784	0.901	0.926	0.910	0.940	0.938	0.963	0.941	0.904	0.835	0.966	
1 correct	kBET	0.029	0.021	0.075	0.053	0.004	0.018	0.721	0.545	0.246	0.287	0.686	0.081	0.187	0.345	0.335	0.355	0.375	0.556	0.483	0.589	0.563	0.141	0.666	
Batcl	iLISI	0.022	0.052	0.005	0.018	0.016	0.070	0.437	0.334	0.246	0.264	0.548	0.269	0.210	0.386	0.402	0.385	0.404	0.423	0.394	0.455	0.424	0.128	0.472	
	S-batch	0.868	0.875	0.606	0.624	0.620	0.655	0.928	0.923	0.916	0.919	0.848	0.752	0.897	0.881	0.910	0.871	0.910	0.910	0.900	0.851	0.913	0.901	0.923	
	cLISI	0.999	0.999	1.000	1.000	0.998	0.998	1.000	0.999	0.999	0.999	1.000	0.999	0.999	0.995	0.998	0.999	1.000	1.000	1.000	1.000	1.000	0.976	1.000	
	S-label	0.581	0.583	0.566	0.612	0.610	0.626	0.646	0.588	0.576	0.575	0.718	0.700	0.587	0.544	0.575	0.559	0.621	0.650	0.750	0.739	0.647	0.521	0.885	
tion	K-ARI	0.548	0.591	0.488	0.520	0.458	0.566	0.708	0.726	0.534	0.612	0.683	0.666	0.672	0.541	0.770	0.590	0.873	0.644	0.791	0.758	0.721	0.340	0.919	
conserva	K-NMI	0.699	0.703	0.670	0.688	0.649	0.701	0.809	0.784	0.697	0.746	0.807	0.771	0.766	0.633	0.793	0.698	0.886	0.796	0.874	0.840	0.807	0.525	0.961	
Bio (L-ARI	0.555	0.585	0.344	0.412	0.387	0.471	0.873	0.776	0.741	0.684	0.582	0.540	0.686	0.750	0.807	0.932	0.936	0.779	0.892	0.738	0.907	0.485	1.000	
	IMNI-T	0.693	0.744	0.659	0.678	0.654	0.693	0.895	0.825	0.804	0.795	0.797	0.737	0.795	0.820	0.839	0.922	0.926	0.886	0.900	0.849	0.918	0.737	1.000	
	I-label	0.580	0.578	0.612	0.641	0.604	0.625	0.554	0.579	0.522	0.571	0.533	0.668	0.583	0.525	0.592	0.528	0.631	0.659	0.573	0.763	0.557	0.513	0.886	
	HVG		>		>		\mathbf{i}		>		>		>	>		<		>	>		>				
	Method	PCA	PCA	TSNE	TSNE	UMAP	UMAP	Harmony	Harmony	Scanorama	Scanorama	BBKNN	BBKNN	fastMNN	scVI	scVI	scANVI	scANVI	scGen	scPoli	scPoli	Author's	Geneformer	Islander	

Extended Data Table 13: Benchmarking cell embeddings, on Skin Cell Atlas

			Bio c	conservat	tion				Batch	correct	ion		Aggı	regate sc	ore
I-label L-NN	L-NN	Į	L-ARI	K-NMI	K-ARI	S-label	cLISI	S-batch	iLISI	KBET	G-Con	PCR	Batch	Bio	Total
0.575 0.77	0.7	02	0.480	0.743	0.453	0.594	1.000	0.849	0.043	0.258	0.872	0.000	0.404	0.659	0.557
0.291 0.5	0	594	0.099	0.636	0.270	0.510	1.000	0.625	0.042	0.205	0.633	0.000	0.301	0.486	0.412
0.399 0.0	ö.	550	0.167	0.684	0.327	0.531	1.000	0.627	0.068	0.305	0.798	0.000	0.360	0.537	0.466
0.588 0.	o.	783	0.556	0.778	0.695	0.613	1.000	0.748	0.142	0.615	0.794	0.603	0.581	0.716	0.662
0.569 0.	<u>.</u>	790	0.477	0.732	0.464	0.589	1.000	0.869	0.083	0.396	0.942	0.132	0.484	0.660	0.590
0.901 0.	0.	682	0.186	0.756	0.436	0.569	1.000	0.570	0.155	0.513	0.895	0.079	0.442	0.647	0.565
0.620 0.8	0.	886	0.896	0.780	0.567	0.571	1.000	0.851	0.133	0.478	0.934	0.668	0.613	0.760	0.701
0.639 0.	o.	848	0.618	0.778	0.581	0.606	1.000	0.830	0.123	0.470	0.909	0.527	0.572	0.724	0.663
0.691 0.	o.	751	0.390	0.882	0.841	0.733	1.000	0.674	0.137	0.495	0.775	0.331	0.482	0.755	0.646
0.713 0.	o.	663	0.174	0.673	0.358	0.553	1.000	0.590	0.107	0.392	0.829	0.000	0.384	0.591	0.508
0.506 0.	o.	759	0.550	0.547	0.319	0.527	1.000	0.848	0.113	0.403	0.834	0.405	0.521	0.601	0.569
0.650 0.	o.	686	0.321	0.586	0.329	0.582	1.000	0.748	0.124	0.474	0.365	0.240	0.390	0.593	0.512
0.678 0.	<u>.</u>	686	0.321	0.616	0.369	0.589	1.000	0.746	0.123	0.469	0.365	0.217	0.384	0.608	0.519
0.695 0.	0	687	0.321	0.616	0.365	0.585	1.000	0.747	0.123	0.485	0.364	0.249	0.394	0.610	0.523
0.328 0.	0.	637	0.151	0.445	0.271	0.442	1.000	0.485	0.149	0.462	0.385	0.360	0.368	0.468	0.428
0.700 0.	<u>.</u>	635	0.147	0.430	0.294	0.515	1.000	0.475	0.149	0.461	0.375	0.538	0.400	0.532	0.479
0.354 0	0	.634	0.147	0.405	0.259	0.495	1.000	0.494	0.152	0.473	0.375	0.312	0.361	0.471	0.427

ypes.
cell t
road
ng bi
s, usi
Atla
) Cell
Donor
Fetal,
(Lung,
uo ,
embeddings
cell
narking
3ench
[4: I
Table 1
Data
Extended]

Method	HVG	Brain	Breast	COVID	Eye	Gut (F)	Heart	Lung (F,D)	Lung (F,O)	Lung	Pancreas	Skin
PCA		0.932	0.885	0.808	0.941	0.968	0.746	0.769	0.707	0.827	0.909	0.914
PCA	>	0.949	0.938	0.962	0.979	0.979	0.713	0.769	0.682	0.910	0.952	0.953
TSNE		0.620	0.504	0.599	0.745	0.724	0.634	0.448	0.373	0.481	0.730	0.708
TSNE	>	0.633	0.632	0.616	0.694	0.699	0.657	0.448	0.426	0.655	0.608	0.718
UMAP		0.580	0.592	0.649	0.695	0.652	0.532	0.563	0.462	0.635	0.621	0.637
UMAP	>	0.684	0.617	0.726	0.762	0.634	0.618	0.563	0.354	0.652	0.624	0.703
Harmony		0.924	0.908	0.837	0.946	0.636	0.835	0.421	0.491	0.746	0.928	0.721
Harmony	>	0.941	0.957	0.971	0.981	0.978	0.766	0.731	0.703	0.890	0.954	0.973
Scanorama		0.929	0.872	0.813	0.949	0.934	0.735	0.669	0.689	0.738	0.894	0.922
Scanorama	>	0.915	0.916	0.930	0.924	0.946	0.727	0.669	0.718	0.774	0.913	0.936
BBKNN		0.544	0.600	0.772	0.674	0.645	0.677	0.462	0.453	0.654	0.667	0.733
BBKNN	>	0.562	0.601	0.765	0.618	0.739	0.608	0.613	0.588	0.669	0.583	0.768
scVI		0.820	0.731	0.622	0.782	0.726	0.718	0.440	0.678	0.632	0.633	0.677
scVI	>	0.863	0.732	0.704	0.795	0.710	0.757	0.438	0.668	0.683	0.643	0.696
scANVI		0.857	0.774	0.773	0.818	0.798	0.818	0.518	0.708	0.712	0.634	0.651
scANVI	>	0.927	0.848	0.848	0.820	0.820	0.888	0.508	0.695	0.829	0.696	0.761
scGen	>	0.949	0.880	0.868	0.954	0.946	0.688	0.694	0.607	0.872	0.901	0.899
scPoli		nan	nan	0.629	0.824	0.857	0.829	0.732	0.645	0.826	0.836	0.690
scPoli	>	0.869	0.818	0.741	0.933	0.849	0.858	0.729	0.647	0.826	0.733	0.887
Author's		0.644	0.884	nan	0.593	0.676	0.829	0.379	0.337	0.634	nan	0.706
Geneformer		0.799	0.631	nan	0.693	0.821	0.693	0.597	0.625	0.668	nan	0.800
Islander		0.549	0.464	0.480	0.753	0.633	0.577	0.256	0.697	0.513	0.676	0.654

Extended Data Table 15: Benchmarking cell embeddings, using scGraph. "F", "D" and "O" represents fetal, donor and organoid, respectively. "nan" means the embeddings are not available, due to memory limitations (>500G in RAM) or unavailability of raw counts (Geneformer). We underline the worst methods among integration methods, Geneformer and Islander, for each atlas.

3. Visualization



Extended Data Figure 1: Drifting Cell Islands, different runs of Islander on fetal lung atlas.



Breast Cell Atlas



- CD4-positive, alpha-beta T cell
- basal cell
- pericyte
- myeloid cell
- neutrophil

- non-classical monocyte
- alternatively activated macrophage
- activated CD4-positive, alpha-beta T cell
- effector memory CD4-positive, alpha-beta T cell
- activated CD8-positive, alpha-beta T cell

- effector memory CD8-positive, alpha-beta T cell
- unswitched memory B cell
- class switched memory B cell
- 🗕 lgG plasma cell
- IgA plasma cell
- conventional dendritic cell
- endothelial cell of lymphatic vessel
- capillary endothelial cell
- luminal epithelial cell of mammary gland
- mammary gland epithelial cell
- vein endothelial cell
- endothelial cell of artery

Extended Data Figure 2: Drifting Cell Islands, Part I

COVID Cell Atlas



Extended Data Figure 3: Drifting Cell Islands, Part II





Extended Data Figure 5: Drifting Cell Islands, Part IV