

Transferable representations of single-cell transcriptomic data

Ethan Weinberger¹ and Su-In Lee^{1,*}

¹Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle

*corresponding: suinlee@cs.washington.edu

Abstract

Advances in single-cell RNA-seq (scRNA-seq) technologies are enabling the construction of large-scale, human-annotated reference cell atlases, creating unprecedented opportunities to accelerate future research. However, effectively leveraging information from these atlases, such as clustering labels or cell type annotations, remains challenging due to substantial technical noise and sparsity in scRNA-seq measurements. To address this problem, we present HD-AE, a deep autoencoder designed to extract integrated low-dimensional representations of scRNA-seq measurements across datasets from different labs and experimental conditions (<https://github.com/suinleelab/HD-AE>). Unlike previous approaches, HD-AE’s representations successfully transfer to new query datasets without needing to retrain the model. Researchers without substantial computational resources or machine learning expertise can thus leverage the robust representations learned by pretrained HD-AE models to compare embeddings of their own data with previously generated sets of reference embeddings.

Main

New developments in scRNA-seq technologies [1, 2, 3] are dramatically reducing the cost of experiments, facilitating the continual release of new scRNA-seq datasets and enabling the construction of large-scale, annotated reference atlases such as the Human Cell Atlas [4]. Despite this explosion in publicly available data, leveraging knowledge from previously studied scRNA-seq datasets to expedite the analysis of new datasets remains difficult

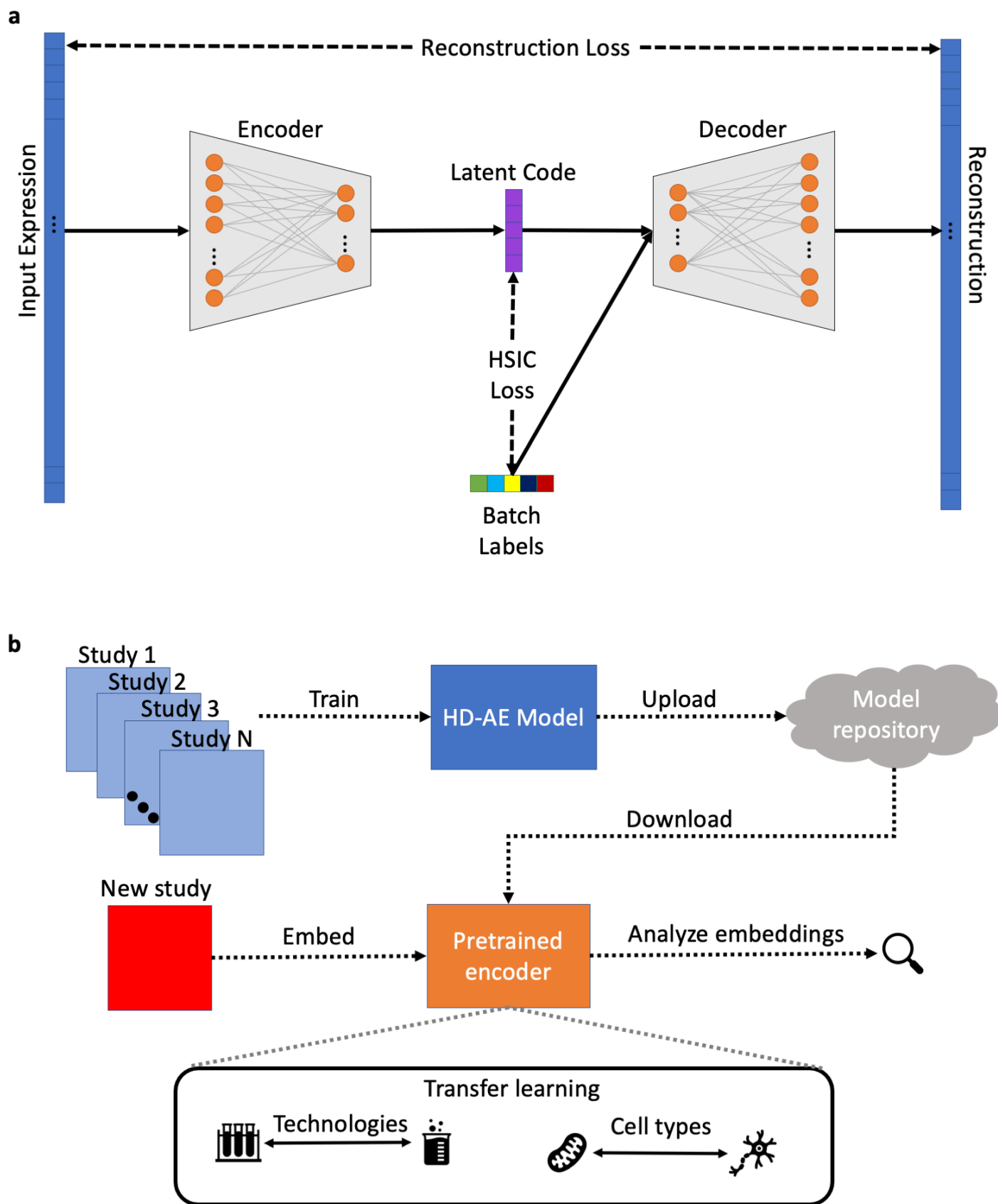


Figure 1: HD-AE models learn unified low-dimensional embeddings of scRNA-seq measurements originating from different experiments. (a) The HD-AE architecture. HD-AE encourages batch effect removal by penalizing the Hilbert-Schmidt Independence Criterion (HSIC) between samples' latent representations and their batch labels (**Methods**). (b) A sample HD-AE transfer learning workflow. Researchers can download pretrained HD-AE models to embed their own data and compare with previously generated sets of reference embeddings.

27 since substantial nuisance factors of variation inherent in scRNA-seq measurements, such
28 as dropout and transcriptional noise, can obscure biological signals of interest [5]. More-
29 over, combining measurements from multiple experiments is complicated by batch effects,
30 i.e., systematic variations between datasets due to differences in experimental conditions or
31 procedures. Batch effects are especially pronounced with scRNA-seq data, since different
32 scRNA-seq protocols have unique sources of bias, sensitivity, and accuracy [6].

33 To address these challenges, several recent works [7, 8, 9, 10, 11, 12, 13] propose data
34 integration methods that produce denoised low-dimensional representations (*embeddings*) of
35 scRNA-seq data. However, these methods are not designed to integrate new query datasets
36 with a previous reference set of embeddings without making users rerun the entire integration
37 pipeline from scratch. This limitation is problematic; raw data from individual scRNA-seq
38 experiments, even those within the same cell atlas, are often stored in different databases and
39 in varying formats, necessitating a time-consuming data collection and preprocessing phase
40 before integration can be performed (**Supplementary Figure 1**). Furthermore, current
41 integration methods scale poorly in terms of computational cost and memory requirements
42 [12] or require specialized hardware (e.g., GPUs), limiting their use to researchers with
43 abundant computational resources.

44 As a response to these limitations, we introduce the Hilbert-Schmidt Deconfounded Au-
45 toencoder (HD-AE, **Figure 1a**), a deep learning approach based on the widely used autoen-
46 coder architecture [14] for learning denoised embeddings of scRNA-seq data. To ensure that
47 samples' embeddings are independent of their batches of origin, we train HD-AE using a loss
48 function that penalizes the Hilbert-Schmidt Independence Criterion (HSIC) [15], a nonpara-
49 metric measure of statistical independence, between samples' embeddings and their batch
50 labels (**Methods**). Removing batch information from the latent space would normally make
51 it difficult for the autoencoder to reconstruct the data faithfully while also preserving true
52 biological structure in the latent space; to mitigate this issue, when training HD-AE, we pass
53 samples' batch labels to the decoder so that batch-specific transformations can be learned
54 to reconstruct the data accurately from the batch-effect-free latent space. *Unlike previously*
55 *proposed scRNA-seq integration methods, pretrained HD-AE models' representations suc-*
56 *cessfully generalize to new batches of data not seen during training, even when those batches*
57 *contain previously unseen cell types.* This lets researchers reuse previously trained HD-AE
58 models off-the-shelf without needing to gather the original data or possess the computational
59 expertise or hardware to train the models themselves from scratch (**Figure 1b**).

60 We first applied HD-AE to construct a reference atlas of pancreas islet cell embeddings
61 using three datasets, each sequenced using a different scRNA-seq protocol. Using UMAP [16]
62 to visualize the raw data, we confirmed that it was clearly separable by batch, even for cells

63 with the same cell type label (**Figure 2a**). After training an HD-AE model and embedding
64 the data into the model’s latent space (**Figure 2b**), we observed that distinctions between
65 batches were removed while cell types remained well-separated, indicating that embedding
66 space variations were due to underlying biological differences rather than technical artifacts.
67 To validate HD-AE’s transfer learning capabilities, we next used our pretrained model to
68 embed a query batch of data collected using the CEL-Seq2 protocol, which was not used
69 to generate any data seen by the model during training. We found that *embeddings of this*
70 *query batch were well-integrated with training batch embeddings* (**Figure 2c**).

71 To further explore the robustness of our pretrained model, we also embedded a second
72 query batch of data collected using the Smart-seq2 protocol (**Fig 2d**). To simulate a poten-
73 tially more realistic scenario where new batches of data contain cell types not seen by the
74 model during training, we included a cell type (alpha) in this second query batch that was
75 held out during training. For cell types shared between the query and reference batches,
76 we once again found that query batch cell embeddings were well-integrated with reference
77 ones. Moreover, we found that the embeddings of the previously unseen alpha cells formed
78 a distinct cluster well separated from other cell types. This behavior persisted for other
79 choices of held-out cell types (**Supplementary Figure 2**). These results further indicate
80 that *pretrained HD-AE models are able to filter out technical noise between different batches*
81 *of data while preserving meaningful biological variations*.

82 We compared HD-AE’s transfer learning capabilities to those of three previously proposed
83 deep learning methods for producing informative embeddings of scRNA-seq data: SAUCIE
84 [17], scVI [11], and DESC [12]. None of these methods was originally designed for transfer
85 learning. Nevertheless, their shared reliance on autoencoder-based architectures made it
86 straightforward to adapt them for use in the transfer learning setting; we note here, however,
87 that, limitations inherent in scVI forced us to disable its batch effect correction feature to use
88 it for transfer learning (**Supplementary Note 1**). In these experiments, we used the same
89 split between training and query batches as we did with HD-AE. Unsurprisingly, we found
90 that scVI (**Figure 2f**) failed to produce well-integrated embeddings when not explicitly
91 trained to correct for batch effects. Moreover, DESC (**Figure 2e**) and SAUCIE (**Figure**
92 **2g**) did not produce contiguous, well-separated clusters for individual cell types, possibly
93 due to their choices of network architectures and loss functions (**Supplementary Note 2**).
94 Only HD-AE produced clusters that were both well integrated between batches and well
95 separated across cell types.

96 We next assessed HD-AE’s performance when applied to a dataset consisting of nine
97 batches of peripheral blood mononuclear cells (PBMCs) collected using seven different tech-
98 nologies [18]. As we saw with the pancreas islet data, this dataset clearly separated by batch

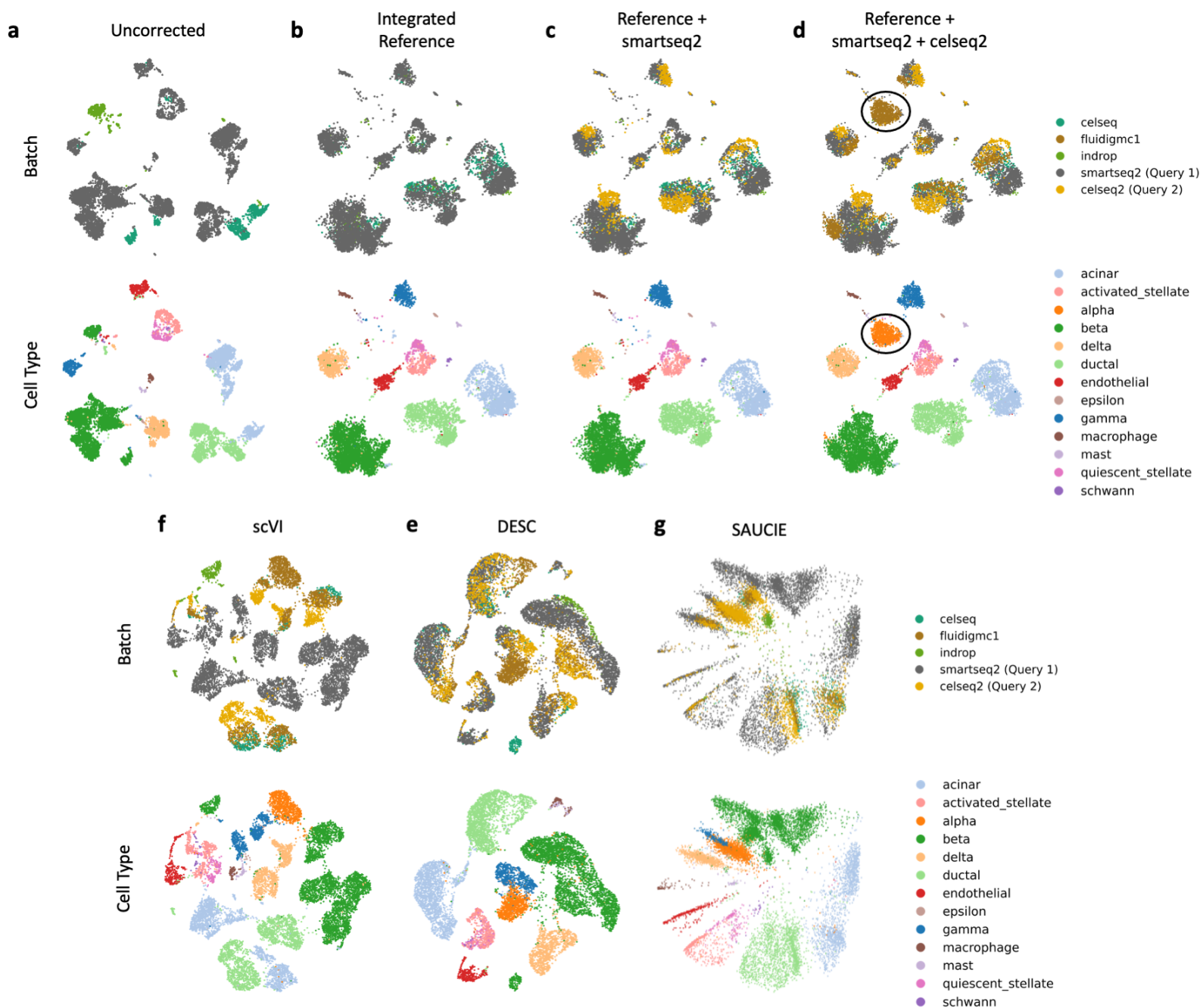


Figure 2: Unlike previously proposed deep learning embedding methods for scRNA-seq data, HD-AE's representations generalize to new datasets at test time. **(a-b)** Three batches of pancreatic islet cells collected using different technologies from different labs before **(a)** and after **(b)** integration with HD-AE. **(c-d)** Querying the reference HD-AE model with two previously unseen batches of data. Black circles indicate cell types held out during training. **(e-g)** Embeddings produced by three previously proposed deep learning methods for scRNA-seq analysis using the same reference and query split as in **(a-d)**.

99 **(Figure 3a)**. In this experiment, HD-AE was trained using seven of the batches; the re-
100 maining two (10x Chromium v3 and Drop-seq) were held out as query batches. Once again,
101 we found that batches were well-mixed and distinctions between cell types were preserved in
102 the HD-AE embedding space **(Figure 3b)**. Moreover, as we found with the pancreas data,
103 our query batch and training batch embeddings were well-integrated **(Figure 3c-d)**.

104 Using this dataset, we also compared the quality of HD-AE embeddings to “full integra-
105 tion” method embeddings (i.e., those from methods designed to integrate all batches of data
106 at once rather than for transfer learning). In particular, we benchmarked HD-AE against
107 seven previously proposed data integration methods: Seurat v3 [8], Conos [10], Harmony [9],
108 scAlign [13], scVI [11], SAUCIE [17], and DESC [12]. Each baseline method was given access
109 to all nine batches of data during training rather than only a set of reference batches. As
110 an additional full integration baseline, we trained an HD-AE model using all batches during
111 training. We report implementation details and hyperparameter choices for all methods in
112 **Supplementary Note 3**. Qualitatively, we found that many baseline methods struggled
113 with this dataset, with only scVI and HD-AE producing well-separated clusters for each cell
114 type **(Supplementary Figure 3)**.

115 We also computed a suite of quantitative metrics to evaluate the quality of each method’s
116 embeddings. To effectively integrate data, a method must balance two potentially opposing
117 goals. First, batches should be well-mixed after integration; that is, the set of k -nearest
118 neighbors around a given cell should be balanced across different batches. To quantify this
119 mixing, we computed the *entropy of batch mixing (EBM)* [7] for each method’s embeddings
120 **(Methods)**. A high EBM can be achieved by randomly mixing the data and disregard-
121 ing biologically meaningful variations. Thus, our evaluation also considered how well local
122 neighborhoods in individual datasets were preserved in the integrated space. To quantify
123 the preservation of this structure, we computed the *k-nearest neighbors purity (kNN purity)*
124 [19] for each method **(Methods)**. A high purity score could trivially be achieved without
125 performing any mixing between batches. Thus, we considered performance on both metrics
126 when evaluating a given method. To compare across metrics, individual metric values were
127 normalized to lie in the range [0, 1].

128 We report our results for individual metrics for a neighborhood size $k = 50$ along with
129 the sums of both metrics to indicate how well each method balances the two **(Figure 3e-g)**.
130 We found that HD-AE outperformed all baseline methods when considering both metrics.
131 This result persisted for varying values of k **(Supplementary Figure 4)**. Remarkably, we
132 found that *HD-AE’s performance in the transfer learning setting was nearly identical to its*
133 *performance when provided all the data during training*.

134 To examine how well each method preserved true biological variations, we also quantified

135 how well different cell types clustered after integration. For each method we calculated the
136 adjusted Rand index (ARI) to assess agreement between ground truth cell type annotations
137 and cluster labels assigned by the Leiden community detection algorithm [20] (**Methods**).
138 We found that HD-AE outperformed all baseline methods on this metric (**Figure 3h**). We
139 also plotted the distributions of silhouette scores (**Methods**) for each method (**Figure 3i**).
140 Here, we found that HD-AE, even in the transfer learning setting, was only narrowly bested
141 by scVI in terms of median silhouette score. Moreover, we once again found that HD-AE's
142 performance on these metrics was nearly unchanged between the full integration and transfer
143 learning settings. Taken together, these results further demonstrate *HD-AE's ability to learn*
144 *high-quality transferable representations of scRNA-seq data*.

145 Finally, we used this dataset to explore how the number of training batches affects HD-
146 AE's generalization performance. To do so, we trained HD-AE models with varying numbers
147 of training batches and evaluated the quality of their embeddings of the full dataset as be-
148 fore. Qualitatively, we found that providing more batches during training initially produced
149 more mixing between batches and more compact, well-separated clusters for each cell type
150 (**Supplementary Figure 5**), though this effect appears to have diminishing returns as the
151 number of batches increases. For our quantitative metrics (**Supplementary Figure 6**),
152 we found that HD-AE's silhouette score and ARI performance did not vary considerably for
153 different numbers of training batches; however, we did initially see sharp increases in our
154 combined EBM and kNN purity metric as the number of training batches increased. In par-
155 ticular, HD-AE began to outperform our full integration baseline models for combined EBM
156 and kNN purity when provided with only four of the nine batches during training. These
157 results suggest that *HD-AE achieves most of its generalization potential when it receives only*
158 *a small number of batches during training, potentially enabling the use of pretrained HD-AE*
159 *models even for tissues or organisms that are less well-studied and for which less data is*
160 *publicly available*.

161 This work introduced the HD-AE framework for producing transferable representations of
162 single cell transcriptomic data. In experiments on pancreas islet cell and PBMC scRNA-seq
163 measurements, we found that HD-AE produced well-integrated reference sets of scRNA-seq
164 embeddings and that pretrained HD-AE models successfully generalized to new batches at
165 test time, even when those batches contained previously unseen cell types. This advancement
166 may enable researchers to leverage pretrained deep learning models to obtain embeddings
167 of their own data for use in arbitrary downstream tasks without needing to undertake the
168 burdensome and skill-intensive process of training the models themselves. As part of future
169 work, we envision training HD-AE models on a variety of tissues from various organisms and
170 distributing them for the benefit of the wider scRNA-seq research community.

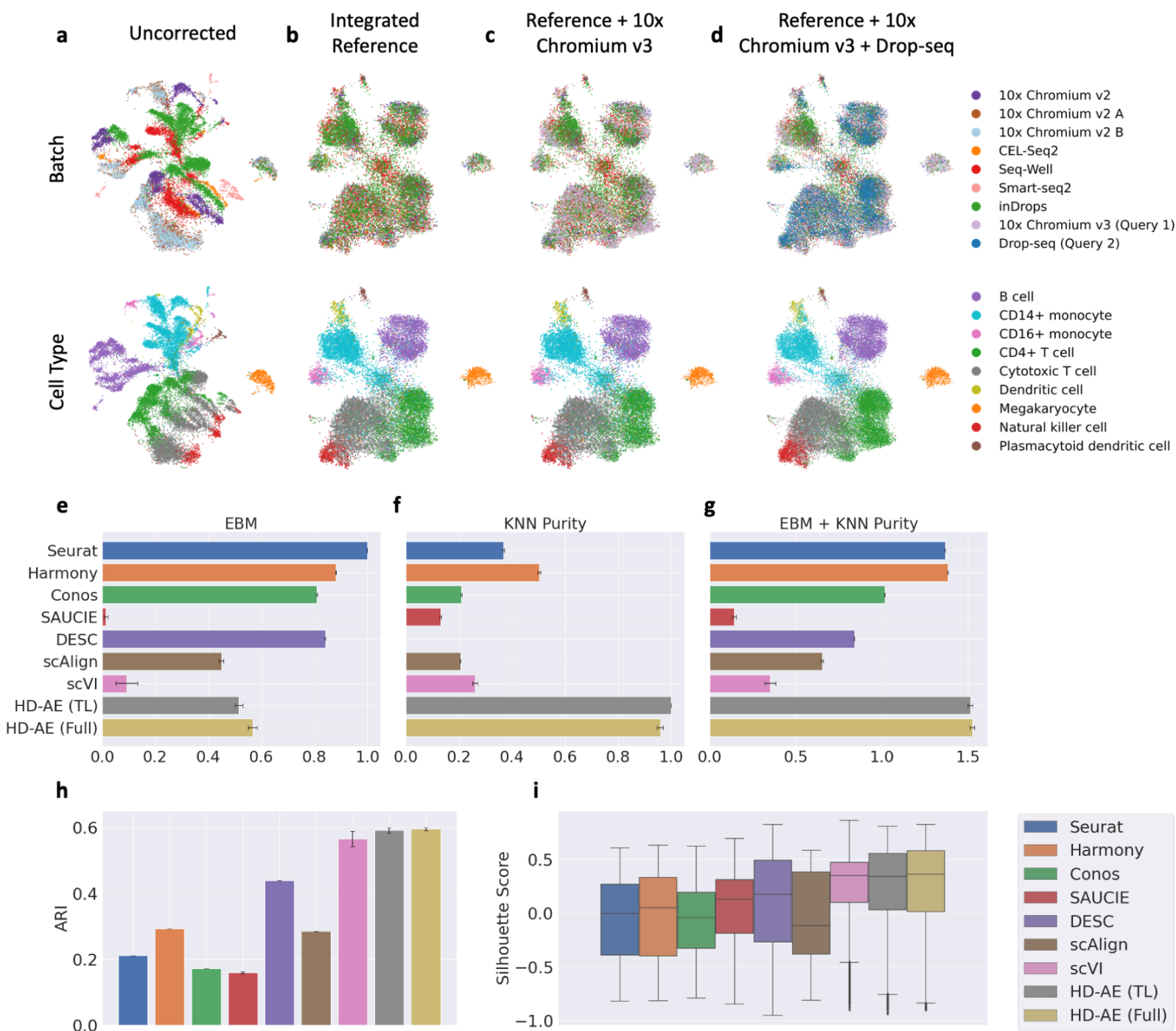


Figure 3: Transfer learning with HD-AE produces higher quality embeddings than previously proposed full integration workflow embeddings. **(a-b)** UMAP plots of seven batches of PBMC data before **(a)** and after **(b)** integration using HD-AE. **(c-d)** Querying the reference HD-AE model with two previously unseen batches of data. **(e-g)** Metrics for evaluating integration performance. Entropy of batch mixing (EBM, **(e)**) quantifies cell mixing across batches, while k -nearest neighbors purity (kNN purity, **(f)**) quantifies the preservation of within-batch local structure. We also report the sum of these metrics (right) to evaluate how well each method balances the two properties. HD-AE (TL) refers to HD-AE trained with seven of the nine batches and embedding the remaining batches via transfer learning, while HD-AE (Full) refers to HD-AE trained with all nine batches. kNN purity and EBM were normalized to lie in the range $[0, 1]$ to enable comparison across the metrics. For each method we report the mean and standard error across five random subsamples of the data. **(h-i)** Separation between cell types in the embedding space as quantified by the adjusted Rand index **(h)** and silhouette scores **(i)**.

171 Methods

172 Autoencoder Model

173 HD-AE extends the standard autoencoder architecture. An autoencoder consists of two
174 networks: (1) an encoder network $f_\phi: X \rightarrow Z$ parameterized by ϕ , which maps from an
175 input space $X \in \mathbb{R}^M$ to a latent space $Z \in \mathbb{R}^D$, and (2) a decoder network $g_\psi: Z \rightarrow X$
176 parameterized by ψ , which maps the latent space representation of a sample back to the
177 original input space. The goal of the encoder network is to learn to map a given sample to the
178 latent space Z so that the decoder network can simultaneously learn to faithfully reconstruct
179 the original sample from its latent space representation. Moreover, we assume that $M \gg D$;
180 therefore, the latent space Z acts like an information bottleneck, capturing the strongest
181 sources of variation in the original data in order to perform accurate reconstructions. In
182 our implementations, both subnetworks consist of fully connected layers with rectified linear
183 unit (ReLU) activations between them. We measure reconstruction loss using mean-squared
184 error, so we train our two networks to solve the optimization problem

$$185 \min_{\phi, \psi} \mathbb{E} \|x_i - g_\psi(f_\phi(x_i))\|_2^2,$$

186 where the expectation is taken over our training data.

187 The Hilbert Schmidt Independence Criterion (HSIC)

188 For random variables X and Y with probability distribution p_{XY} , the HSIC measures the
189 statistical dependence between the two. In particular, an HSIC of zero between X and Y is
190 zero if and only if X and Y are independent, while a higher HSIC corresponds to a stronger
191 level of dependence.

192 If $\{(x_i, y_i)\}_{i=1}^n$ are independently and identically distributed samples drawn from p_{XY} ,
193 the HSIC can be empirically estimated via

$$194 \widehat{\text{HSIC}}(\{(x_i, y_i)\}_{i=1}^n) = \frac{1}{(n-1)^2} \text{Tr}(KHLH).$$

195 Here, $K_{ij} = k(x_i, x_j)$ and $L_{ij} = l(y_i, y_j)$ are Gram matrices for kernel functions k and l ,
196 respectively, where k and l must be universal kernels, a class that includes the widely used
197 Gaussian and Laplacian kernels [15]. Moreover, H is a centering matrix, and Tr denotes the
198 trace operator. See **Supplementary Note 4** for further details on the HSIC.

199 HD-AE

200 HD-AE, an extension of the standard autoencoder model, was specifically designed to learn
201 batch-effect-free latent representations. To do so, we added a regularization term to the
202 autoencoder objective to minimize the empirical HSIC between latent representations and
203 batch labels. Removing batch information from the latent space would usually complicate
204 an autoencoder’s efforts to reconstruct the data faithfully while preserving true biological
205 structure in the latent space; to mitigate this issue, when training HD-AE, we passed batch
206 labels to the decoder so that batch-specific transformations could be learned to reconstruct
207 the data accurately from the batch-effect-free latent space.

208 Suppose we have n total gene expression samples. Let b_i denote the batch label for the
209 i th sample x_i , let B denote the vector of batch labels for all samples, and let (by an abuse
210 of notation) $f_\phi(X) \in \mathbb{R}^{n \times d}$ denote the matrix of latent representations of the dataset X . We
211 then have the HD-AE objective function

$$212 \min_{\phi, \psi} \underbrace{\mathbb{E} \|x_i - g_\psi(f_\phi(x_i), b_i)\|_2^2}_{\text{reconstruction error}} + \lambda \cdot \underbrace{\widehat{\text{HSIC}}(f_\phi(X), B)}_{\text{batch effect penalty}},$$

213 which we can optimize via stochastic gradient descent. In all our experiments Gaussian
214 kernel functions with $\sigma^2 = 1$ were used to compute $\widehat{\text{HSIC}}$.

215 Datasets and Preprocessing

216 Pancreas Data

217 Our pancreas data came from the `panc8` dataset [21] provided in the `SeuratData` R pack-
218 age available at <https://github.com/satijalab/seurat-data>. We used data from the
219 `celseq`, `fluidigm1`, and `indrop` batches, as indicated by the `tech` field in the R object for
220 training, and we used cells from the `smartseq2` and `celseq2` batches for testing generaliza-
221 tion performance. We preprocessed the data by first filtering down the collection of datasets
222 to the top 2000 highly variable genes as determined by the `Seurat` R package. For `scVI`, we
223 used this filtered data directly; for other models, we normalized the data using the `Seurat`
224 normalization workflow. For `DESC`, we also scaled the data after normalization using the
225 `scale_bygroup` function from the `DESC` Python package.

226 PBMC Data

227 Our PBMC data came from the `pbmcsca` dataset [18] available in the `SeuratData` R package.
228 For our transfer learning HD-AE model, we used data from the `CEL-Seq2`, `10x Chromium`

229 (v2), 10x Chromium (v2) A, 10x Chromium (v2) B, Drop-seq, Seq-Well, and inDrops
230 batches, as indicated by the `Method` field in the R object for training. During preprocessing,
231 we removed any cells with a cell type label of `Unassigned`; otherwise, our preprocessing
232 workflow was the same as for the pancreas data.

233 Evaluation Metrics

234 Entropy of Batch Mixing

235 Letting c denote the number of batches, the entropy of batch mixing (EBM) is defined as

$$236 \quad EBM = \sum_{i=1}^c x_i \log(x_i),$$

237 where c is the number of batches, x_i denotes the proportion of cells originating from a batch
238 i in a given region, and $\sum_{i=1}^c x_i = 1$. To assess the EBM for a given method, we followed
239 a standard [11] estimation procedure: we randomly chose 100 cells, calculated “regional”
240 EBM values for each cell using the batch proportions from the cell’s 50 nearest neighbors in
241 the integrated space, and then averaged over the 100 regional EBMs.

242 kNN Purity

243 For a given batch, two similarity matrices were constructed. The first was computed using
244 that batch’s cells’ gene expression values pre-integration; the second was computed using
245 that batch’s cells’ representations in the integrated space. We then computed the ratio
246 of the intersection of these matrices’ corresponding k nearest neighbors graphs over their
247 union. We repeated this procedure for each batch in a dataset and reported the average of
248 this statistic.

249 Silhouette Score

250 For a given cell i , the silhouette score $s(i)$ is defined as follows. Let $a(i)$ be the average
251 distance between i and the other cells in i ’s cluster, and let $b(i)$ be the smallest average
252 distance between i and all other cells in a different cluster. The silhouette score $s(i)$ is then

$$253 \quad s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}.$$

254 A silhouette score close to one indicates that i is tightly clustered with cells with the
255 same ground truth label. A score close to -1 indicates that a cell has been grouped with cells
256 with a different label.

257 Adjusted Rand Index

258 The adjusted Rand index (ARI) measures agreement between reference clustering labels and
259 labels assigned by a clustering algorithm. Given a set of n cells and two sets of clustering
260 labels describing those cells, the overlap between clustering labels can be described using a
261 contingency table, where each entry indicates the number of cells in common between the
262 two sets of labels. Mathematically, the ARI is calculated as

$$263 \text{ARI} = \frac{\sum_{ij} \binom{n_{ij}}{2} - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}},$$

264 where n_{ij} is the number of cells assigned to cluster i based on the reference labels and
265 cluster j based on a clustering algorithm, a_i is the number of cells assigned to cluster i in the
266 reference set, and b_j is the number of cells assigned to cluster j by the clustering algorithm.

267 In our experiments, we assigned cells to clusters using the Leiden community detection
268 algorithm. Because the results of this algorithm depend heavily on its resolution hyperpa-
269 rameter, for each method we tried a number of resolution values in the range $[0.5, 1.0]$ and
270 reported the best resulting ARI for each method.

271 References

- 272 [1] Gierahn, T. M. *et al.* Seq-well: portable, low-cost RNA sequencing of single cells at
273 high throughput. *Nat. Methods* **14**, 395–398 (2017).
- 274 [2] Hashimshony, T. *et al.* Cel-seq2: sensitive highly-multiplexed single-cell RNA-seq.
275 *Genome Biology* **17**, 77 (2016).
- 276 [3] Xin, Y. *et al.* Use of the Fluidigm C1 platform for RNA sequencing of single mouse
277 pancreatic islet cells. *Proceedings of the National Academy of Sciences* **113**, 3293–3298
278 (2016).
- 279 [4] Regev, A. *et al.* Science forum: the human cell atlas. *Elife* **6**, e27041 (2017).
- 280 [5] Kharchenko, P. V., Silberstein, L. & Scadden, D. T. Bayesian approach to single-cell
281 differential expression analysis. *Nat. Methods* **11**, 740–742 (2014).
- 282 [6] Svensson, V. *et al.* Power analysis of single-cell RNA-sequencing experiments. *Nat.*
283 *Methods* **14**, 381–387 (2017).

- 284 [7] Haghverdi, L., Lun, A. T., Morgan, M. D. & Marioni, J. C. Batch effects in single-
285 cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat.*
286 *Biotechnol.* **36**, 421–427 (2018).
- 287 [8] Stuart, T. *et al.* Comprehensive integration of single-cell data. *Cell* **177**, 1888–1902
288 (2019).
- 289 [9] Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with
290 harmony. *Nat. Methods* 1–8 (2019).
- 291 [10] Barkas, N. *et al.* Joint analysis of heterogeneous single-cell RNA-seq dataset collections.
292 *Nat. Methods* **16**, 695–698 (2019).
- 293 [11] Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling
294 for single-cell transcriptomics. *Nat. Methods* **15**, 1053–1058 (2018).
- 295 [12] Li, X. *et al.* Deep learning enables accurate clustering with batch effect removal in
296 single-cell RNA-seq analysis. *Nat. Commun.* **11**, 1–14 (2020).
- 297 [13] Johansen, N. & Quon, G. scAlign: a tool for alignment, integration, and rare cell
298 identification from scRNA-seq data. *Genome Biology* **20**, 1–21 (2019).
- 299 [14] Hinton, G. E. & Salakhutdinov, R. R. Reducing the dimensionality of data with neural
300 networks. *Science* **313**, 504–507 (2006).
- 301 [15] Gretton, A., Bousquet, O., Smola, A. & Schölkopf, B. Measuring statistical dependence
302 with hilbert-schmidt norms. In *International Conference on Algorithmic Learning The-*
303 *ory*, 63–77 (Springer, 2005).
- 304 [16] McInnes, L., Healy, J. & Melville, J. UMAP: Uniform manifold approximation and
305 projection for dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).
- 306 [17] Amodio, M. *et al.* Exploring single-cell data with deep multitasking neural networks.
307 *Nat. Methods* **16**, 1139–1145 (2019).
- 308 [18] Ding, J. *et al.* Systematic comparison of single-cell and single-nucleus RNA-sequencing
309 methods. *Nat. Biotechnol.* **38**, 737–746 (2020).
- 310 [19] Xu, C. *et al.* Probabilistic harmonization and annotation of single-cell transcriptomics
311 data with deep generative models. *Molecular Systems Biology* **17**, e9620 (2021).

- 312 [20] Traag, V. A., Waltman, L. & Van Eck, N. J. From Louvain to Leiden: guaranteeing
313 well-connected communities. *Sci. Rep.* **9**, 1–12 (2019).
- 314 [21] Lab, S. *panc8.SeuratData: Eight Pancreas Datasets Across Five Technologies* (2019).
315 R package version 3.0.2.