

# Adversarial domain translation networks for fast and accurate integration of large-scale atlas-level single-cell datasets

Jia Zhao<sup>1\*</sup>, Gefei Wang<sup>1\*</sup>, Jingsi Ming<sup>2</sup>, Zhixiang Lin<sup>3</sup>, Yang Wang<sup>1</sup>  
The Tabula Microcebus Consortium, Angela Ruohao Wu<sup>4,5†</sup>, Can Yang<sup>1†</sup>

<sup>1</sup>Department of Mathematics, The Hong Kong University of Science and Technology, Hong Kong SAR, China

<sup>2</sup>Academy of Statistics and Interdisciplinary Sciences, KLATASDS-MOE, East China Normal University, Shanghai, China

<sup>3</sup>Department of Statistics, The Chinese University of Hong Kong, Hong Kong SAR, China

<sup>4</sup>Division of Life Science, The Hong Kong University of Science and Technology, Hong Kong SAR, China

<sup>5</sup>Department of Chemical and Biological Engineering, The Hong Kong University of Science and Technology, Hong Kong SAR, China

## Abstract

1 The rapid emergence of large-scale atlas-level single-cell RNA-seq datasets presents remarkable  
2 opportunities for broad and deep biological investigations through integrative analyses. However,  
3 harmonizing such datasets requires integration approaches to be not only computationally  
4 scalable, but also capable of preserving a wide range of fine-grained cell populations. We  
5 created Portal, a unified framework of adversarial domain translation to learn harmonized  
6 representations of datasets. With innovation in model and algorithm designs, Portal achieves  
7 superior performance in preserving biological variation during integration, while achieving  
8 integration of millions of cells in minutes with low memory consumption. We show that Portal  
9 is widely applicable to integrating datasets across samples, platforms and data types (including

---

\*These authors contributed to this work equally.

†Correspondence: [angelawu@ust.hk](mailto:angelawu@ust.hk), [macyang@ust.hk](mailto:macyang@ust.hk).

10 scRNA-seq, snRNA-seq and scATAC-seq). Finally, we demonstrate the power of Portal by  
11 applying it to the integration of cross-species datasets with limited shared information among  
12 them, elucidating biological insights into the similarities and divergences in the spermatogenesis  
13 process among mouse, macaque and human.

## 14 **Introduction**

15 Advances in single-cell sequencing have enabled identification of novel cell types [1, 2], in-  
16 vestigation of gene regulation networks [3, 4], and understanding of cellular differentiation  
17 processes [5, 6]. As single-cell technologies rapidly evolved over recent years, its experimental  
18 throughput substantially increased, allowing researchers to profile increasingly complex and  
19 diverse samples, and accelerating the accumulation of vast numbers of rich datasets over time  
20 [7, 8, 9]. Integrative and comparative analyses of such large-scale datasets originating from  
21 various samples, different platforms and data modalities, as well as across multiple species, offer  
22 unprecedented opportunities to establish a comprehensive picture of diverse cellular behaviors.  
23 Integration is a critical step, to account for heterogeneity of different data sources when taking  
24 advantage of single-cell data from different studies [10]. Thus, integration methods that can  
25 efficiently and accurately harmonize a wide range of data sources are essential for accelerating  
26 life sciences research [11].

27 Although integration methods for single-cell transcriptomics analysis have evolved along  
28 with single-cell sequencing technologies, the rapid accumulation of new and diverse single-cell  
29 datasets has introduced three major challenges to the integration task. First, as the sample size  
30 of each single-cell dataset grows dramatically, numerous extensive datasets with hundreds of  
31 thousands or even millions of cells have been produced [8, 9, 12]. The emergence of large-scale  
32 datasets requires integration methods to be fast, memory-efficient, and scalable to millions  
33 of cells. Second, technology now allows effective, comprehensive characterization of complex  
34 organs, containing rare subpopulations of cells that can now be captured, albeit in small  
35 numbers, thanks to the scale of profiling that is now possible [7, 13]. Investigation into high-  
36 level heterogeneity among cell populations is essential for understanding the mechanism of  
37 complex biological systems. Hence, the ideal integration method needs to carefully preserve fine-  
38 grained cell populations from each atlas-level dataset. Third, the biological origins of datasets  
39 has expanded in diversity, with data now spanning across not only different technological

40 platforms and data types, different individual donors, but even across different species, which  
41 can be especially interesting for evolutionary studies [14, 15, 16]. Integrative analysis of such  
42 diverse datasets would allow researchers to unify resources to address a wider range of biological  
43 questions. Recent single-cell atlasing efforts are a primary example of these challenges – various  
44 human tissue atlases [12, 17], mouse multi-tissue atlases [7, 18], and non-human primate atlases  
45 [19, 20] have been generated, culminating in data from millions of single cells and single  
46 nuclei. Both within and across atlas comparisons are of interest. To perform integrative and  
47 comparative analyses based on such diverse data sources, there is an urgent need for methods  
48 that can flexibly account for heterogeneous dataset-specific effects, while maintaining a high  
49 level of integration accuracy.

50 Many methods have been developed to align single-cell datasets [10], including Harmony  
51 [21], Seurat [22], online iNMF [23], VIPCCA [24], scVI [25], fastMNN [26], Scanorama [27] and  
52 BBKNN [28]. Several of these methods that were designed for large datasets at the time of  
53 publication are now less attractive in terms of scalability in the face of atlas-level dataset sizes.  
54 For instance, a representative category of methods leverages the mutual nearest neighbors  
55 (MNN) to perform data alignment. These MNN-based methods, such as Seurat, fastMNN and  
56 Scanorama, require identification of MNN pairs across datasets, thus the time and memory  
57 costs quickly become unbearably high when the dataset exceeds one million cells. Another  
58 limitation of existing methods is that they are mainly targeted towards integrating datasets of  
59 less complex tissues, utilizing strategies such as MNN, matrix factorization, and soft-clustering  
60 to capture major biological variations. With these strategies, inaccurate mixing of different cell  
61 types can be avoided when clear clustering patterns are present; but when dealing with more  
62 complex tissues, they tend to overcorrect fine-grained cell subpopulations, resulting in the loss  
63 of power in revealing interesting biological variations [29, 30]. Lastly, most existing methods  
64 are designed to correct batch effects caused by technical artifacts. To this end, a number of  
65 methods, like BBKNN and fastMNN, assume that the biological variation is much larger than  
66 the variation of batch effects. This assumption may not be true when applied across data types  
67 and species.

68 To simultaneously address the above three challenges, we created Portal, a machine learning-  
69 based algorithm for aligning atlas-level single-cell datasets with high efficiency, flexibility, and  
70 accuracy. Viewing datasets from different studies as distinct domains with domain-specific

71 effects (including technical variation and other sources of unwanted variation), Portal achieves  
72 extraordinary data alignment performance through a unified framework of domain translation  
73 networks that incorporates an adversarial learning mechanism [31]. To find the correspondence  
74 between two domains, our domain translation network utilizes an encoder to embed cells from  
75 one domain into a latent space where domain-specific effects are removed, and then uses a  
76 generator to map latent codes to another domain. The generator simulates the generation  
77 process of domain-specific effects. In each domain, a discriminator is trained to identify where  
78 poor alignment between the distributions of original cells and transferred cells occurs. The  
79 feedback signal from the discriminator is used to strengthen the domain translation network  
80 for better alignment. The nonlinearity of encoders and generators in the adversarial domain  
81 translation framework enables Portal to account for complex domain-specific effects. In contrast  
82 to existing domain translation methods [32, 33, 34], Portal has the following unique features.  
83 First, Portal has a uniquely designed discriminator which can adaptively distinguish domain-  
84 shared cell types and domain-unique cell types. Therefore, Portal will not force the alignment  
85 of domain-unique cell types, avoiding the risk of overcorrection. Second, without using any  
86 cell type label information, three regularizers of Portal can guide domain translation networks  
87 to find correct correspondence between domains, account for domain-specific effects, and  
88 retain biological variation in the latent space. Third, through a tailored design of lightweight  
89 neural networks and mini-batch optimization accelerated by graphics processing units (GPUs),  
90 Portal can scale up to datasets containing millions of cells in minutes with nearly constant  
91 memory usage. With the above innovations in model and algorithm designs, Portal enables  
92 fast and accurate integration of atlas-level datasets across samples, technological platforms,  
93 data modalities, and species.

94 Through a comprehensive benchmarking study, where integration of heterogeneous collec-  
95 tions of atlas-level single-cell RNA sequencing (scRNA-seq) data are included, Portal shows its  
96 superiority over state-of-the-art alignment algorithms in terms of both computational efficiency  
97 and accuracy. We then show that Portal can accurately align cells from complex tissues profiled  
98 by scRNA-seq and single-nucleus RNA sequencing (snRNA-seq) as well as align scRNA-seq data  
99 and single-cell assay for transposase-accessible chromatin using sequencing (scATAC-seq) data,  
100 even in the presence of highly unbalanced cell type compositions. We also apply Portal to the  
101 integration of cells in differentiation processes, especially the alignment of the gradient of cells

102 in the spermatogenesis process across multiple species (mouse, macaque, and human). Using  
103 these diverse and challenging experiments, we demonstrate Portal’s versatility and power for a  
104 broad range of applications. Comprehensive analyses of real, expert annotated data confirm  
105 that integrated cell embeddings provided by Portal can be reliably used for identification of  
106 rare cell populations via clustering or label transfer, studies of differentiation trajectories,  
107 and transfer learning across data types and across species. Portal is now publicly available  
108 as a Python package (<https://github.com/YangLabHKUST/Portal>), serving as an efficient,  
109 reliable and flexible tool for integrative analyses.

## 110 Results

### 111 Method Overview: Portal learns a harmonized representation of 112 different datasets with adversarial domain translation.

113 Expression measurements from different datasets fall into different domains due to the existence  
114 of domain-specific effects, including technical variation and other sources of unwanted variation  
115 (Fig. 1a), causing difficulty when performing joint analyses. Without loss of generality, here we  
116 consider two domains,  $\mathcal{X}$  and  $\mathcal{Y}$ . We assume that domain  $\mathcal{X}$  and domain  $\mathcal{Y}$  can be connected  
117 through a low-dimensional shared latent space  $\mathcal{Z}$ , which captures the biological variation and  
118 is not affected by the domain-specific effects. By taking the measurements of cells from  $\mathcal{X}$  and  
119  $\mathcal{Y}$  as inputs, we aim to learn a harmonized representation of cells in latent space  $\mathcal{Z}$  to obtain  
120 data alignment between  $\mathcal{X}$  and  $\mathcal{Y}$ .

121 We achieve the above goal through a unified framework of adversarial domain translation,  
122 namely “Portal”. Domains and the shared latent space are connected by encoders and  
123 generators (Fig. 1b). Encoder  $E_1(\cdot) : \mathcal{X} \rightarrow \mathcal{Z}$  is designed to remove the domain-specific  
124 effects when mapping cells from  $\mathcal{X}$  into  $\mathcal{Z}$ , and generator  $G_1(\cdot) : \mathcal{Z} \rightarrow \mathcal{X}$  is designed to  
125 simulate the domain-specific effects when mapping cells from  $\mathcal{Z}$  into  $\mathcal{X}$ . By symmetry, encoder  
126  $E_2(\cdot) : \mathcal{Y} \rightarrow \mathcal{Z}$  and generator  $G_2(\cdot) : \mathcal{Z} \rightarrow \mathcal{Y}$  are designed with the same role in connecting  
127  $\mathcal{Y}$  and  $\mathcal{Z}$ . To transfer cells between  $\mathcal{Y}$  and  $\mathcal{X}$  through shared latent space  $\mathcal{Z}$  (Fig. 1b),  
128 encoder  $E_2(\cdot)$  and generator  $G_1(\cdot)$  work together to form one domain translation network  
129  $G_1(E_2(\cdot)) : \mathcal{Y} \rightarrow \mathcal{Z} \rightarrow \mathcal{X}$ . Clearly, encoder  $E_1(\cdot)$  and generator  $G_2(\cdot)$  form another domain  
130 translation network  $G_2(E_1(\cdot)) : \mathcal{X} \rightarrow \mathcal{Z} \rightarrow \mathcal{Y}$ . To achieve the mixing of original cells and  
131 transferred cells, discriminators  $D_1(\cdot)$  and  $D_2(\cdot)$  are deployed in domains  $\mathcal{X}$  and  $\mathcal{Y}$  to identify

132 where poor mixing occurs (Fig. 1c). The discriminators’ feedback then guides the domain  
133 translation networks to improve the mixing.

134 However, the well mixing of original cells and transferred cells in each domain does not  
135 imply extraordinary data alignment across domains. First, a domain-unique cell population  
136 should not be mixed with cells from another domain. Second, cell types  $A$  and  $B$  in domain  
137  $\mathcal{X}$  could be incorrectly aligned with cell types  $B$  and  $A$  in domain  $\mathcal{Y}$ , respectively, although  
138 the distributions of original cells and transferred cells are well mixed. To address these issues,  
139 Portal has the following unique features, which distinguishes it from existing adversarial domain  
140 translation frameworks [32, 33]. On one hand, we deploy the tailored design of discriminators  
141  $D_1(\cdot)$  and  $D_2(\cdot)$  such that they can distinguish domain-unique cell types from cell types shared  
142 across different domains. The domain-unique cell types will be treated as outliers and left  
143 in the discriminator’s inactive region (Fig. 1c). In such a way, these cell types will not be  
144 enforced for alignment, avoiding the risk of overcorrection. On the other hand, we design three  
145 regularizers to find correct correspondence across domains and avoid incorrect alignment when  
146 the distributions are well mixed.

147 Specifically, let  $\mathbf{x}$  and  $\mathbf{y}$  be the samples from domains  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. We consider  
148 the following framework of adversarial domain translation,

$$\begin{aligned} \min_{\{E_1, G_1, E_2, G_2\}} \max_{\{D_1, D_2\}} & \mathcal{L}_{\mathcal{X}}(D_1, E_2, G_1) + \mathcal{L}_{\mathcal{Y}}(D_2, E_1, G_2), \\ \text{subject to} & \mathcal{R}_{\text{AE}}(E_1, G_1, E_2, G_2) \leq t_{\text{AE}}, \\ & \mathcal{R}_{\text{LA}}(E_1, G_1, E_2, G_2) \leq t_{\text{LA}}, \\ & \mathcal{R}_{\text{cos}}(E_1, G_1, E_2, G_2) \leq t_{\text{cos}}. \end{aligned} \tag{1}$$

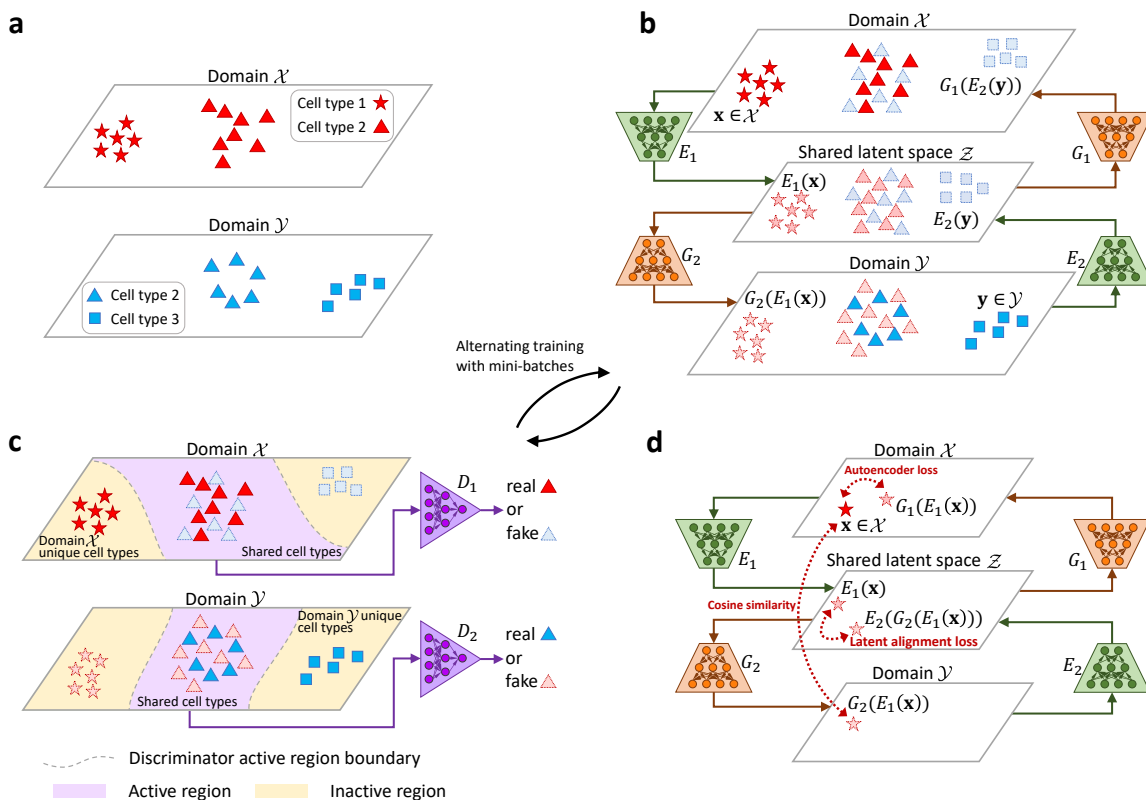
149 In model (1),  $\mathcal{L}_{\mathcal{X}}(D_1, E_2, G_1) := \mathbb{E}[\log D_1(\mathbf{x})] + \mathbb{E}[\log(1 - D_1(G_1(E_2(\mathbf{y}))))]$  and  $\mathcal{L}_{\mathcal{Y}}(D_2, E_1, G_2) :=$   
150  $\mathbb{E}[\log D_2(\mathbf{y})] + \mathbb{E}[\log(1 - D_2(G_2(E_1(\mathbf{x}))))]$  are the objective functions for adversarial learning of  
151 domain translation networks  $G_1(E_2(\cdot))$  and  $G_2(E_1(\cdot))$  in  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. Discriminators  
152  $D_1(\cdot)$  and  $D_2(\cdot)$  are trained to distinguish between “real” cells (i.e. original cells in a domain),  
153 and “fake” cells (i.e. transferred cells generated by domain translation networks) by minimizing  
154  $\mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}}$ , while the domain translation networks are trained against the discriminators by  
155 maximizing  $\mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}}$ . These three regularizers  $\mathcal{R}_{\text{AE}}$ ,  $\mathcal{R}_{\text{LA}}$  and  $\mathcal{R}_{\text{cos}}$  play a critical role in finding  
156 correct correspondence of cells between two domains, accounting for domain-specific effects,  
157 and retaining biological variation in the latent space (Fig. 1d). More specifically, the first  
158 regularizer  $\mathcal{R}_{\text{AE}} := \frac{1}{p} \{\mathbb{E}[\|\mathbf{x} - G_1(E_1(\mathbf{x}))\|_2^2] + \mathbb{E}[\|\mathbf{y} - G_2(E_2(\mathbf{y}))\|_2^2]\}$ , where  $p$  is the dimension-

159 ality of domains  $\mathcal{X}$  and  $\mathcal{Y}$ , requires the autoencoder consistency in domains  $\mathcal{X}$  and  $\mathcal{Y}$ ; the  
160 second regularizer  $\mathcal{R}_{\text{LA}} := \frac{1}{q} \{ \mathbb{E} [\|E_1(\mathbf{x}) - E_2(G_2(E_1(\mathbf{x})))\|_2^2] + \mathbb{E} [\|E_2(\mathbf{y}) - E_1(G_1(E_2(\mathbf{y})))\|_2^2] \}$ ,  
161 where  $q$  is the dimensionality of  $\mathcal{Z}$ , imposes the consistency constraint in the latent space;  
162 and the third regularizer  $\mathcal{R}_{\text{cos}} := \mathbb{E} \left[ 1 - \frac{\langle \mathbf{x}, G_2(E_1(\mathbf{x})) \rangle}{\|\mathbf{x}\|_2 \|G_2(E_1(\mathbf{x}))\|_2} \right] + \mathbb{E} \left[ 1 - \frac{\langle \mathbf{y}, G_1(E_2(\mathbf{y})) \rangle}{\|\mathbf{y}\|_2 \|G_1(E_2(\mathbf{y}))\|_2} \right]$  introduces  
163 the cross-domain correspondence by preserving the cosine similarity between a sample and  
164 its transferred version;  $t_{\text{AE}}$ ,  $t_{\text{LA}}$  and  $t_{\text{cos}}$  are their corresponding constraint parameters. More  
165 detailed explanation can be found in the Methods section.

166 We solve the above optimization problem via alternating updates by stochastic gradient  
167 descent. The algorithm is extremely computationally efficient with the support of stochastic  
168 optimization accelerated by GPUs. After the training process, Portal learns a harmonized  
169 representation of different domains in shared latent space  $\mathcal{Z}$ . Samples from  $\mathcal{X}$  and  $\mathcal{Y}$  can  
170 be transferred into latent space  $\mathcal{Z}$  to form an integrated dataset  $\{E_1(\mathbf{x})\}_{\mathbf{x} \in \mathcal{X}} \cup \{E_2(\mathbf{y})\}_{\mathbf{y} \in \mathcal{Y}}$   
171 using encoders  $E_1(\cdot)$  and  $E_2(\cdot)$ , facilitating the downstream integrative analysis of cross-domain  
172 single-cell datasets.

## 173 **Accurate integration of atlas-level datasets within minutes and re-** 174 **quiring lower memory consumption compared to other methods.**

175 The rapid accumulation of large-scale single-cell datasets requires integration algorithms  
176 to efficiently handle datasets containing millions of cells without loss of accuracy. For a  
177 comprehensive comparison, we first benchmarked Portal and existing representative methods,  
178 including Harmony [21], Seurat v3 [22], online iNMF [23], VIPCCA [24], scVI [25], fastMNN  
179 [26], Scanorama [27] and BBKNN [28], in terms of integration performance following a recent  
180 benchmarking study [30]. Using a number of scRNA-seq datasets from diverse tissue types  
181 with curated cell cluster annotations, including mouse spleen, marrow, and bladder [7], we  
182 quantitatively evaluated the integration performance of each method. We first evaluated  
183 alignment performance, which can sometimes be interpreted as batch correction performance, of  
184 all compared methods. The score for batch correction was computed by leveraging a collection  
185 of batch correction metrics designed in existing studies, including k-nearest neighbor batch-effect  
186 test (kBET) [35], principal component regression of the batch covariate (PCR batch) [35],  
187 average silhouette width across batches (batch ASW) [35], graph integration local inverse  
188 Simpson’s Index (graph iLISI) [30, 21] and graph connectivity [30]. The higher the batch



**Figure 1: Overview of Portal.** **a.** Portal regards different single-cell datasets as different domains. Joint analyses of these datasets are confounded by domain-specific effects, representing the unwanted technical variation. **b.** Portal employs encoders  $E_1(\cdot), E_2(\cdot)$  to embed the biological variation of domains  $\mathcal{X}$  and  $\mathcal{Y}$  into a shared latent space  $\mathcal{Z}$ , where domain-specific effects are removed. The generating process of domain-specific effects are captured by two generators  $G_1(\cdot)$  and  $G_2(\cdot)$ . Encoder  $E_1(\cdot)$  and generator  $G_2(\cdot)$  form a domain translation network  $G_2(E_1(\cdot))$  mapping from  $\mathcal{X}$  to  $\mathcal{Y}$ ; Encoder  $E_2(\cdot)$  and generator  $G_1(\cdot)$  form another domain translation network mapping from  $\mathcal{Y}$  to  $\mathcal{X}$ . **c.** Encoders and generators are trained by competing against specially designed discriminators  $D_1(\cdot)$  and  $D_2(\cdot)$ . In each domain, a discriminator is trained to distinguish between original cells in this domain and cells transferred from another domain, providing feedback signals to assist alignment. To prevent overcorrection of domain-unique cell types, the discriminators in Portal with the tailored design are also able to distinguish between domain-unique cell types and domain-shared cell types. With this design, Portal can focus only on merging cells of high probability to be of domain-shared cell types, while it remains inactive on cells of domain-unique cell types. **d.** Portal leverages three regularizers to help it find correct and consistent correspondence across domains, including the autoencoder regularizer, the latent alignment regularizer and the cosine similarity regularizer.

189 correction score, the higher the degree of mixing across datasets. We also assessed the score  
 190 for conservation of biological variation using different metrics, including adjusted rand index  
 191 (ARI) [36], normalized mutual information (NMI) [37], cell type ASW, graph cell type local  
 192 inverse Simpson's Index (graph cLISI) [30, 21], isolated label F1 [30], isolated label silhouette



193 [30] and cell cycle conservation [30]. By jointly accounting for these metrics, the score can be  
194 used to evaluate different methods' ability to preserve information such as cell type identities.  
195 Inappropriate merging of cell types during integration will result in a low score of biological  
196 variation conservation. Finally, we computed the overall score as a 40:60 weighted average of  
197 the batch correction score and the conservation of biological variation score to indicate the  
198 overall integration performance. Based on our benchmarking results, we found that in general,  
199 BBKNN, Scanorama, fastMNN, scVI and VIPCCA had less satisfactory overall integration  
200 performance compared to the other four methods (Fig. 2a, the first three columns; Figs. S5, S7,  
201 S9 and S11). As indicated by the relatively low batch correction scores of BBKNN, Scanorama,  
202 fastMNN and scVI, we found that observable batch effects still exist in the integration results  
203 that they produced (Figs. S6, S8 and S10). Although VIPCCA showed reasonable performance  
204 in terms of removing batch effects, incorrect mixing of distinct cell types was often observed in  
205 VIPCCA's integration results (Fig. S6). Therefore, its overall scores are relatively low due to  
206 the loss of biological variation (Figs. S5, S7).

207 Among those methods with high user popularity, Harmony, Seurat, and online iNMF also  
208 showed the best overall integration performance results (Fig. 2a, the first three columns; Figs.  
209 S6, S8 and S10). To offer precise and robust integration performance, Seurat [22] utilizes  
210 the detection of mutual nearest neighbors (MNN) to build correspondence between datasets  
211 in the shared embedding space obtained by applying canonical correlation analysis (CCA).  
212 Harmony [21] learns a simple linear correction for dataset-specific effects by running an iterative  
213 soft clustering algorithm, enabling fast computation on large datasets. Online iNMF [23] is a  
214 recently developed approach based on widely used integration method LIGER [38]. It extends  
215 LIGER's non-negative matrix factorization to an iterative and incremental version to improve  
216 its scalability, while it has nearly the same performance as LIGER. For the remainder of this  
217 study, we focus our discussion on comparisons between Portal and these three high-performing  
218 and popular methods in the main text. The comparisons with other methods are provided in  
219 Fig. 2a (the last three columns), and Supplementary Information (Figs. S12 - S18).

220 Next, we evaluated the speed, memory usage, alignment quality, and integration accuracy  
221 using a more challenging integration task. We used two mouse brain atlases [8, 9] as bench-  
222 marking datasets for a more in-depth comparison of Portal and three other methods. One atlas  
223 contains Drop-seq data of 939,489 cells, and another one contains 10X Genomics (10X) data

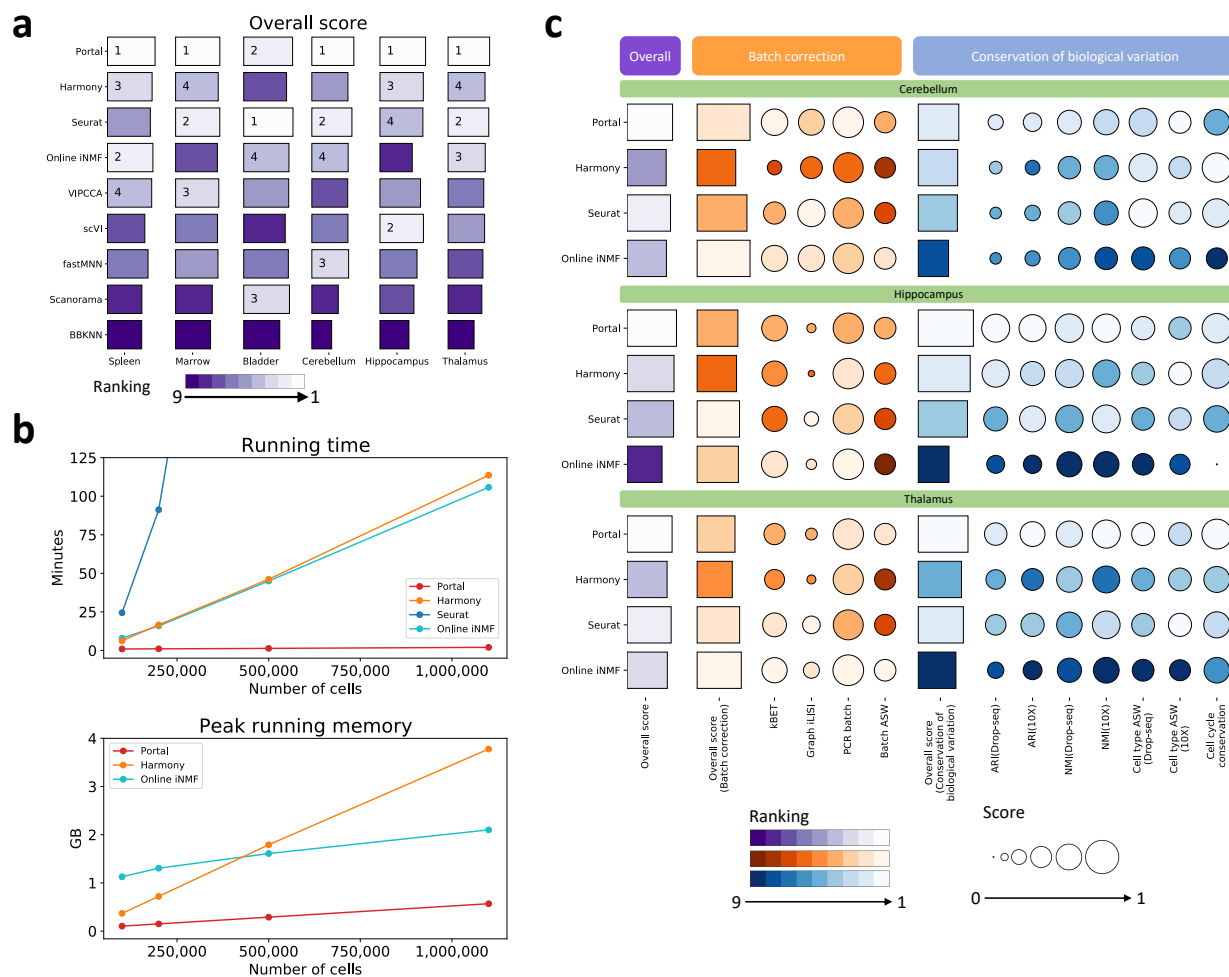


Figure 2: **Benchmarking of Portal and other state-of-the-art integration methods.**

**a.** Overall scores of the compared methods evaluated on mouse spleen, marrow, bladder, cerebellum, hippocampus and thalamus datasets. The ranking was visualized by color gradient, where lighter color indicates better performance. **b.** The running time and the peak running memory required by the benchmarked methods. The datasets were sampled from two mouse brain atlas datasets ( $n = 100,000, 250,000, 500,000,$  and  $1,100,167$ ). Seurat required 24.52 GB on the dataset with 100,000 cells, which was not comparable to the other three benchmarked methods in terms of the peak running memory usage. **c.** Batch correction and biological variation conservation evaluated using three shared tissues from two mouse brain atlases (profiled by Drop-seq and 10X), including cerebellum, hippocampus, and thalamus. Biological variation conservation performance was assessed based on fine-grained annotations provided by the original publications [8, 9].

224 of 160,678 cells. These two mouse brain atlases have data from three shared brain regions:  
 225 cerebellum, hippocampus, and thalamus. There are many small clusters of neuron subtypes  
 226 in these datasets, where gene expressions between subclusters could have a relatively small  
 227 difference. Thus, these datasets are more challenging to integrate compared to data with clear  
 228 clustering patterns.

229 First, Portal has superior integration accuracy even when handling datasets which contain  
230 many subclusters with small difference. The score of biological variation conservation shows  
231 that Portal outperforms other state-of-the-art methods in cluster identity preservation, as the  
232 scores were assessed based on fine-grained cell type and subtype annotations. In particular, for  
233 all three brain regions tested, Portal has the highest ARI and NMI scores among the compared  
234 methods (Fig. 2c).

235 Second, Portal also outperforms the other three methods on scalability, in terms of time  
236 and memory consumption. For this benchmarking test, we obtained datasets from the original  
237 full-sized datasets by combining the two atlases and subsampling proportionally from each  
238 atlas, with each dataset having increasing sample size ranging from 100,000 to 1,100,167 (full  
239 dataset). The running time and the peak running memory of all methods were recorded using  
240 these datasets on the same GPU server. The results show that Portal's running time and peak  
241 running memory remained almost constant even when the sample size increased dramatically  
242 (Fig. 2b). Compared to the other three methods, the running time required by Portal was also  
243 substantially less. On the dataset containing 500,000 cells, Portal's running time was 80 seconds;  
244 when number of cells grew to 1,100,167, Portal's running time only increased to 120 seconds. In  
245 comparison, Harmony and online iNMF both needed more than 40 minutes to integrate 500,000  
246 cells and more than 100 minutes to complete the integration of the full dataset. The running  
247 time of Seurat increased most rapidly among the compared methods. It took as much as 511  
248 minutes (over 8.5 hours) to integrate the 500,000-cell dataset. The computational efficiency  
249 of Portal is owing to two important factors in its design: 1) its algorithm takes advantage of  
250 GPU-accelerated stochastic optimization, such that Portal reads data in mini-batches from the  
251 disk rather than having to load the entire dataset at once, which enables fast integration of  
252 large single-cell datasets using small amounts of memory; and 2) lightweight neural networks  
253 are adopted in Portal to further improve computational efficiency. As such, Portal is also the  
254 most memory-efficient approach among the benchmarked methods (Fig. 2b). Peak running  
255 memory required by Portal ranged from 0.29 GB on 500,000-cell dataset to 0.57 GB on the  
256 full million-cell dataset. Notably, Portal's lightweight networks and mini-batch stochastic  
257 optimization algorithm enable us to control the GPU peak running memory usage at a constant  
258 level of 0.06 GB. Among compared methods, online iNMF used less memory than Harmony  
259 and Seurat when the sample size became larger than 500,000, because it is also trained in

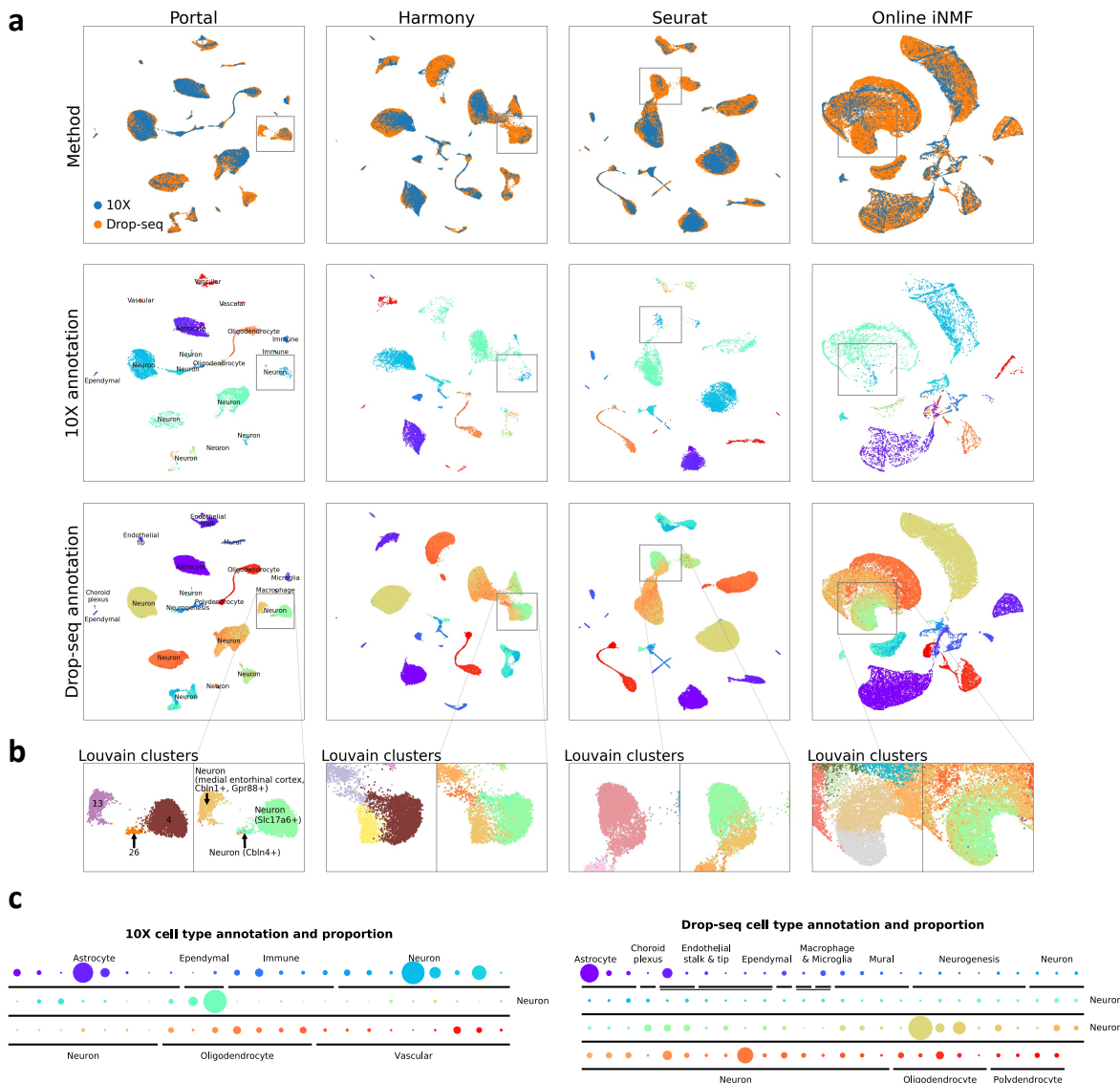
260 mini-batches. However, its peak running memory was 2.10 GB on the million-cell dataset,  
261 which is 2.7 times more than Portal's. Seurat required remarkably more memory usage than  
262 the other three methods. For clarity of visualization, we did not display the peak running  
263 memory required by Seurat as it ranged from 24.52 GB on the 100,000-cell dataset to 276.41  
264 GB on the 500,000-cell dataset.

265 Finally, and importantly, Portal's high performance in speed and memory consumption  
266 does not compromise its ability to align cell type clusters. The batch correction score shows  
267 that Portal's alignment ability is comparable to, if not better than, the other state-of-the-art  
268 methods, indicating that Portal is capable to effectively remove domain-specific effects.

### 269 **Portal preserves subcluster and small cluster identities in complex** 270 **tissues thereby facilitating identification of rare subpopulations.**

271 When integrating complex tissues, one problem that can arise is the inadvertent loss of small  
272 cell populations and subpopulations. Due to more nuanced differences between clusters, or due  
273 to the imbalance in cell numbers for very small cell populations, these "fine-grained" groups of  
274 cells may become inappropriately combined with other groups after integration. In the brain,  
275 for example, there are many subpopulations of neurons which are distinguished from each other  
276 using a few key gene markers while still all bearing the neuron signature; furthermore, some of  
277 these neuronal subtypes could be rare compared to other subtypes. To demonstrate that Portal  
278 can preserve the nuanced information of such small cell populations and subpopulations, we  
279 performed further analysis on the mouse hippocampus tissue integration results. Both mouse  
280 brain atlas datasets contain extensive data for this brain region (Fig. 3), and both studies  
281 identified a wide range of transcriptionally distinct cell subpopulations, including a variety  
282 of neuron subtypes, whose nuanced transcriptional differences should ideally be preserved by  
283 integration methods.

284 After applying Portal and the other three benchmarked methods to integrate the data, we  
285 used the integrated cell representations to perform clustering. Using the Louvain method [40]  
286 with default resolution, we obtained 29 (Portal), 29 (Harmony), 25 (Seurat) and 30 (online  
287 iNMF) clusters, respectively (Fig. S19). Particularly, we focused on one region where the  
288 cell proportions between two datasets were highly unbalanced, as marked in Fig. 3a. Only a  
289 few of cells in this region are from the 10X dataset, making it challenging to build alignment



**Figure 3: Preservation of fine-grained neuron subpopulations in the integration of hippocampus datasets.** **a.** We visualized integration results from Portal, Harmony, Seurat and online iNMF of hippocampus datasets profiled by Drop-seq and 10X with UMAP [39]. Top panels are UMAP plots colored by profiling methods. Middle and bottom panels are UMAP plots of cells from the 10X dataset, the Drop-seq dataset after integration respectively, colored by fine-grained annotations (c). **b.** We marked a region containing three distinct neuron subpopulations. Results from Louvain clustering algorithm were presented for a comparison of cluster identity preservation performance. **c.** Cell type annotations and proportions of the two datasets from their original publications [8, 9]. The comparison among proportions of subpopulations was visualized by the sizes of corresponding dots.

290 between datasets while preserving subpopulations from the Drop-seq dataset. In the original  
 291 publication [8], cells from the Drop-seq dataset within the marked region were all annotated  
 292 as neurons but further classified into three transcriptionally distinct subpopulations, namely:

293 *Cbln1*+/*Grp88*+ medial entorhinal cortex neurons; *Slc17a6*+ neurons; and *Cbln4*+ neurons.  
294 Among the benchmarked methods, Portal was the only method that clearly clustered these cells  
295 into three coherent groups in the integrated embedding space. Specifically, clusters 4, 13, and  
296 26 identified by the Louvain method recovered the *Slc17a6*+ neuron; *Cbln1*+/*Grp88*+ medial  
297 entorhinal cortex neuron; and the *Cbln4*+ neuron subpopulations, respectively (Fig. 3b). Each  
298 cluster was confirmed by the high expression level of the annotated marker genes (Fig. S20a).  
299 Notably, these three groups only accounted for 4.79%, 1.76% and 0.32% of the total sample  
300 size, respectively, demonstrating Portal’s ability to preserve identities of rare subpopulations.  
301 However, the differences among these three subpopulations were not well preserved by the other  
302 three methods, making it difficult to detect them each distinctly using the Louvain clustering  
303 method (Fig. 3a, b). As shown in Fig. S20c, we also identified eight protein coding genes  
304 that were the most significantly differentially expressed among clusters, indicating the different  
305 functions of each of the three neuron subtypes. Cluster 4 showed high expression levels of  
306 *Camk2n1*, *Map1b*, *Nrgn*, *Syt1*, and no detectable expression of *Camk2d*, *Igfbp5*, *Nr4a2* and  
307 *Ntng1*. A different pattern was observed in cluster 13: high expression of *Camk2d*, *Camk2n1*,  
308 *Map1b* and *Syt1*, and no detectable expression of the other four genes. Cluster 26, meanwhile,  
309 showed moderate levels of expression of all eight genes. In the marked region, cells from the 10X  
310 dataset were mostly concentrated in clusters 4 and 13. The alignment by Portal was confirmed  
311 by the consistent gene expression levels in clusters 4 and 13 between the two datasets (Fig.  
312 S20b). Besides the eight differentially expressed genes, we also examined a larger set of genes,  
313 and computed the cross correlation of these genes pairwise between cells from all three groups.  
314 This analysis showed that cells within each cluster had higher similarity in gene expression  
315 than cells from other clusters, further showing the biological difference between these three  
316 clusters that should not be mixed after integration (Fig. S20d). The above results highlight  
317 Portal’s power to preserve rare cell types.

318 The integrative analysis on the hippocampus tissue demonstrates Portal’s ability to maintain  
319 nuanced transcriptional differences for small subpopulations. This means that Portal can also  
320 be used to “call out” rare subpopulations in one dataset based on integration with another  
321 dataset via label transfer. To illustrate this feature, we take 10X and SMART-seq2 (SS2) data  
322 generated for a mouse lung scRNA-seq atlas [7] as an example: the typically larger sample size  
323 of the 10X dataset facilitates powerful clustering analyses for identification of cell types; while

324 the greater sequencing depth and sensitivity of SS2 enables deeper investigation into cell biology  
325 [41]. To leverage the different strengths of the two technologies, we used Portal to perform  
326 integrated analysis on 1,676 SS2 cells and 5,404 10X cells (Fig. S21a). Specifically, we defined  
327 the 10X dataset annotations from the original publication [7] as reference labels (Fig. S21b),  
328 then made use of the Portal's integration results to identify cell types for the SS2 dataset based  
329 on these reference labels. After integration, for each SS2 cell, label transfer was performed  
330 by detecting its nearest neighbors among 10X cells. From this analysis, we identified four  
331 subpopulations of myeloid cells for the SS2 dataset, namely alveolar macrophages, dendritic  
332 cell and interstitial macrophages, classical monocytes, and non-classical monocytes (Fig. S21d).  
333 Transferred labels of these four subpopulations were validated by known marker gene expression  
334 levels [42]. For example, compared to classical monocytes, non-classical monocytes showed  
335 lower expression of *Ccr2* and higher expressions of *Trem14* (Fig. S22). Consistent with the gene  
336 expression pattern of alveolar macrophages in the 10X dataset, alveolar macrophages annotated  
337 by Portal in the SS2 dataset had high expression levels of marker genes *Mrc1* and *Siglec5*.  
338 Notably, in the SS2 dataset, the alveolar macrophage subpopulation only accounted for 0.78%  
339 of total sample size, and could not be distinguished from the other SS2-profiled macrophages in  
340 the original publication [7]. Based on the original labels, alveolar macrophages were unidentified  
341 as they were labeled in a more general group named “dendritic cell, alveolar macrophage,  
342 and interstitial macrophage” (Fig. S21c). Making good use of the larger 10X dataset, Portal  
343 successfully identified extremely rare subpopulations within the SS2 dataset. We then used  
344 the mouse lemur bladder scRNA-seq datasets from the Tabula Microcebus Consortium [43] as  
345 another example to demonstrate Portal's ability for discovering rare subpopulations via label  
346 transfer. In this example, mouse lemur bladder tissue was also profiled by both SS2 and 10X.  
347 When we integrated these datasets and transferred labels from the 10X dataset to the SS2  
348 dataset using Portal, we were able to distinguish a very small myofibroblast subpopulation of  
349 just 11 cells in the SS2 dataset from the rest of the fibroblasts (Fig. S23a). We verified their  
350 myofibroblast identity based on their high expressions of known marker genes *ACTA2*, *MYH11*,  
351 *TAGLN* [44] (Fig. S23b).

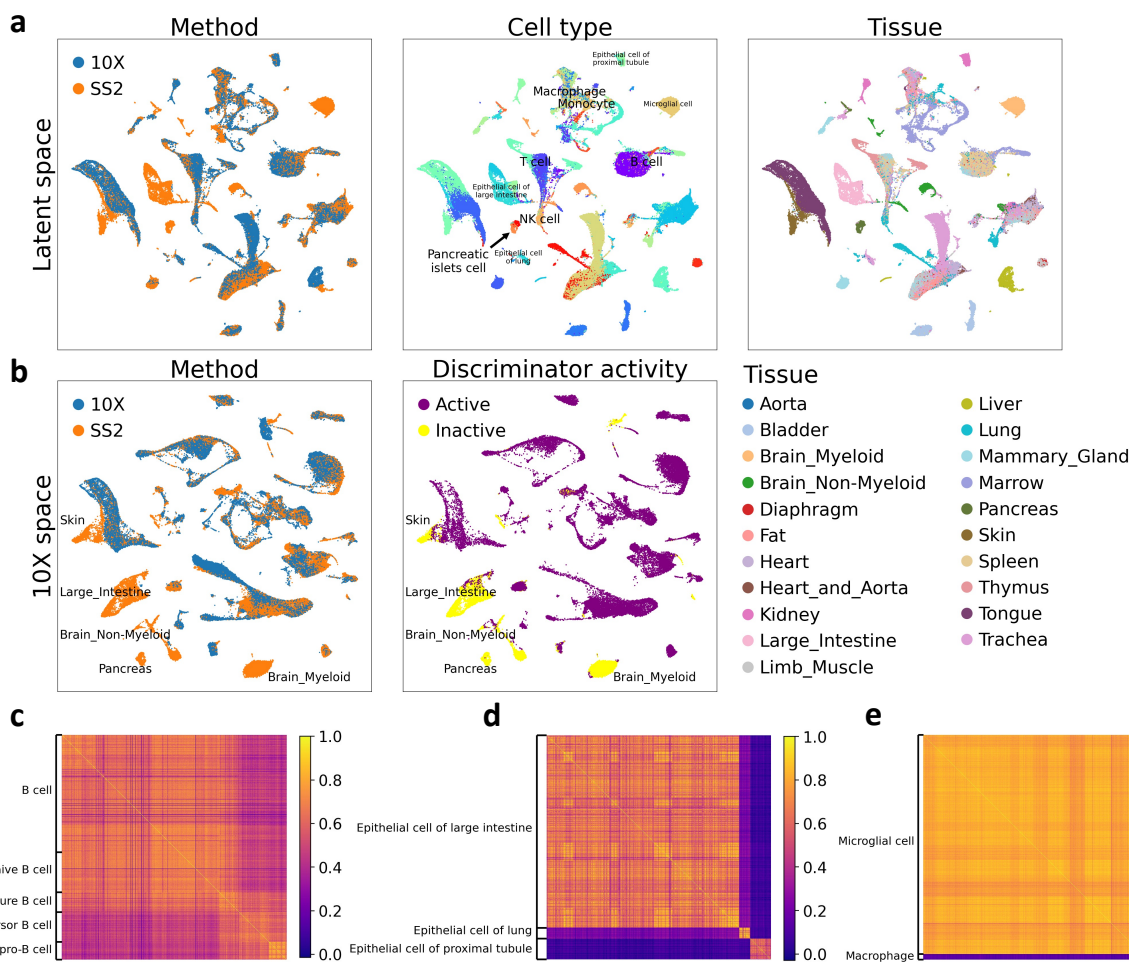


Figure 4: **Construction of mouse cell atlas across the entire organism by integrating atlas datasets from the Tabula Muris project.** We applied Portal to integrate the datasets obtained by 10X and SS2. There were cells from unique tissues presented in the SS2 dataset. **a.** UMAP plots of Portal’s integration results in the shared latent space, colored by profiling methods, cell types and tissues. **b.** Portal also transferred cells from the space of SS2 dataset to the space of the 10X dataset (10X space). In the 10X space, 10X cells were fixed as reference. Portal only aligned SS2 cells of shared cell types between datasets to 10X cells, while maintaining the identities of SS2 cells belonging to tissue-unique cell types. This was achieved by the special design of discriminator activity in Portal. **c, d.** Correlations among cells from subpopulations of B cells (**c**) and epithelial cells (**d**). **e.** Transcriptional distinction between macrophage and microglial cells.

### 352 Integration of comprehensive whole-organism cell atlases.

353 So far, Portal has shown impressive performance in aligning tissue-level atlases where nuanced  
 354 transcriptional differences among subpopulations can be maintained after integration. We  
 355 next assess Portal’s capabilities under another challenging scenario: integrating two atlases  
 356 across an entire organism, where one of the atlases includes many more organs and tissue



357 types than the other. This can be very challenging for some integration algorithms due to  
358 having “missing cell types” in one of the datasets [10]. In contrast to these approaches, Portal  
359 uses discriminators with tailored design in the adversarial domain translation framework to  
360 distinguish domain-specific cell types from cell types shared across domains automatically, and  
361 is thus robust to non-overlapped tissue samples.

362 To build a foundation for extensive study of cell populations across the whole organism, the  
363 Tabula Muris Consortium [7] profiled cells from 20 tissues using a combination of SS2 (44,779  
364 cells) and 10X (54,865 cells) (Fig. 4). Notably, seven of these 20 tissues were only profiled by  
365 SS2 but not 10X: brain (myeloid and non-myeloid), diaphragm, fat, large intestine, pancreas  
366 and skin. We used Portal to build a comprehensive integrated mouse atlas that merges all  
367 the cells, and we found Portal to show extraordinary accuracy in aligning cells of the same  
368 cell type from the two datasets profiled by different platforms, not only in the shared latent  
369 space but also in both domains (Figs. 4a, b and S24). After Portal integration, tissue-specific  
370 cell types of SS2-only tissues, such as microglial cells in brain (myeloid), cell types in large  
371 intestine, and pancreatic islets cells, were all successfully and correctly remained separated  
372 from other cell types. The other three benchmarked methods, however, failed to retain many  
373 tissue-specific cell types unmixed with other cell types. For instance, they mixed microglial  
374 cells together with other macrophage cells, even though the data from these two cell types were  
375 clearly transcriptionally different (Figs. 4e and S24).

376 Using this construction of a mouse cell atlas across organs, we also confirmed that the  
377 designed boundaries for discriminator active region in Portal (Fig. 1c) indeed helped to  
378 maintain the biological variation. By looking into the domain of 10X data (10X space), the  
379 discriminator in the 10X domain was found inactive for tissue-specific cell types that were only  
380 in the SS2 dataset (Fig. 4b). For these cells, Portal ensured that their identities were preserved  
381 by making the adversarial learning objective inactive on them automatically. Portal’s ability  
382 to conserve information of cell populations indicates its reliability for integrating atlas-level  
383 single-cell datasets across entire organisms.

384 Besides the alignment between datasets, Portal’s integration result could characterize the  
385 similarities and differences among cell types. For example, immune cells such as B cells, T cells,  
386 natural killer cells (NK cells), monocytes and macrophages were profiled by both platforms  
387 and contained in multiple tissues including brain (myeloid), diaphragm, fat, kidney, limb

388 muscle, liver, lung, mammary gland, marrow, spleen, and thymus. Portal correctly kept the  
389 subpopulations belonging to the same type of immune cells close to each other, revealing the  
390 resemblance of immune cells across different tissues. For instance, the transcriptional correlation  
391 of all types of B cells, containing B cells, naive B cells, immature B cells, precursor B cells, and  
392 late pro-B cells confirmed such similarity (Fig. 4c). In addition, the epithelial cells of different  
393 tissues were identified by Portal as disjoint clusters, which was consistent with the biological  
394 distinction among these cell types (Fig. 4d).

## 395 **Portal successfully and efficiently aligns datasets across different data** 396 **types.**

397 As most of existing methods were developed only for integrating scRNA-seq datasets, aligning  
398 datasets with different data types could be problematic for these approaches. Here we illustrate  
399 that Portal can flexibly account for the distinction between different data types and yield  
400 accurate integration results.

401 We first examined integration of scRNA-seq data and snRNA-seq data. For frozen samples  
402 such as biobanked tissues, and for tissue types that have unique morphology or phenotypes, such  
403 as brain, fat, or bone, it can be challenging or sometimes even impossible to extract intact cells  
404 for scRNA-seq profiling [45, 46]. To bypass this issue, snRNA-seq has been developed. Although  
405 nuclear transcriptomes are shown to be representative of the whole cell [47], distinctions between  
406 the whole cell and nucleus in terms of the transcript type and composition make scRNA-seq  
407 data and snRNA-seq data intrinsically different [45]. Aligning these two types of data is  
408 desirable, as the combined dataset enables joint analysis that can take advantages of both  
409 techniques, and help to improve statistical power for the analysis. Especially for comparing  
410 multiple complex tissues, with some cell types being shared and others being non-overlapping,  
411 researchers could benefit from such integrated joint analysis – one example being the integration  
412 of brain snRNA-seq data with scRNA-seq data of blood to examine similarities and differences  
413 between immune cells in each tissue milieu. However, due to the inherent difference in these two  
414 data types, aligning scRNA-seq and snRNA-seq data is not the same as batch effects correction.  
415 Compared to batch effects among scRNA-seq datasets, technical noise and unwanted variation  
416 arising from different data types are often more complex and have higher strength [45, 48].  
417 Thus, using standard batch effects correction to integrate across data types may result in loss

418 of alignment accuracy or important biological signals.

419 We evaluated Portal's ability to integrate snRNA-seq data and scRNA-seq data using three  
420 mouse brain atlas datasets, including one snRNA-seq dataset profiled by SPLiT-seq [49], and  
421 two scRNA-seq datasets profiled by Drop-seq and 10X [8, 9]. In this task, we applied integration  
422 methods to harmonize these three atlases across all brain regions. To test the accuracy of  
423 integration results, we only used cells that had annotations provided by the authors in each  
424 atlas project. After selecting cells with cell type annotations, 319,359 cells in the Drop-seq  
425 dataset, 160,678 cells in the 10X dataset, and 74,159 nuclei in the SPLiT-seq remained for  
426 integration.

427 Prior to any integration, the raw datasets were clustered by the experimental methods  
428 rather than the cell types (Fig. S25a), and shared cell types between the three datasets did not  
429 align well, indicating the initial discrepancy between the three large datasets. After integration,  
430 UMAP visualizations showed that the different alignment methods gave varying results. Portal  
431 (Fig. S25b) and Seurat (Fig. S25d) achieved the best alignment of data among different  
432 methods, showing good mixing of cells annotated with the same cell type label, while also  
433 preserving subcluster data structure in the integrated results. In particular, the alignment of  
434 scRNA-seq (10X, Drop-seq) and snRNA-seq (SPLiT-seq) datasets was comparably good as that  
435 of the two scRNA-seq datasets, indicating successful alignment between the two data types  
436 without loss of biologically important variations between clusters. Online iNMF (Fig. S25e),  
437 although it successfully clustered and aligned the same cell types together, within each cluster  
438 the streaky pattern suggested potential numerical artefacts in the integrated data. Furthermore,  
439 online iNMF alignment resulted in loss of biological variation, which was most easily observable  
440 in the coalescence of the previously distinct neuron subpopulations (Fig. S25a) into one large  
441 amorphous cluster (Fig. S25e). Harmony, however, showed poor mixing of the snRNA-seq  
442 data in some of the cell types, such as the astrocytes, where the scRNA-seq datasets were  
443 well-mixed after alignment, but the snRNA-seq data were not mixed well with the rest (Fig.  
444 S25c). Similar to online iNMF, some of the neurons' subcluster structure appeared to be lost  
445 after the integration by Harmony. Overall, Portal and Seurat presented the best scRNA-seq  
446 and snRNA-seq data alignment performance; however, not including data preprocessing time,  
447 Seurat took over 17 hours to complete the task, while Portal only took 87 seconds (details of  
448 the procedure are included in the Methods section).

449 We further assessed Portal’s ability to integrate across data types when no cell type, or  
450 very few cell types are shared. In this scenario, we applied Portal to integrate one human  
451 PBMC scRNA-seq dataset [50] with two human brain snRNA-seq datasets [51, 52], respectively,  
452 as two examples. In the first example (Fig. S26), where no cell type was shared between  
453 datasets, Portal did not mixed any two populations of cells together, showing its robustness.  
454 More importantly, it embedded monocytes and dendritic cells from the PBMC dataset close to  
455 macrophages from the brain dataset, indicating the similarities among these immune cells across  
456 tissue types. In comparison, overcorrection was observed in results from other state-of-the-art  
457 methods. For example, Seurat mixed T cells from blood with excitatory neurons and inhibitory  
458 neurons from brain inappropriately (Fig. S26). The reliability of Portal was also demonstrated  
459 in the second example (Fig. S27). It correctly aligned cells of the only shared cell type (T cell)  
460 between datasets, while it did not mix other distinct cell types (Fig. S27).

461 Besides integration of scRNA-seq data and snRNA-seq data, we then applied Portal to  
462 align scRNA-seq data and scATAC-seq data. As an epigenomic profiling method, scATAC-seq  
463 measures chromatin accessibility, providing a complementary view to scRNA-seq. Integrative  
464 analyses of scRNA-seq and scATAC-seq data are very helpful to leverage and unify information  
465 from the both aspects [53, 22]. For this task, we used one scRNA-seq PBMC dataset profiled  
466 by CITE-seq and one scATAC-seq PBMC dataset profiled by ASAP-seq [54]. For a better  
467 evaluation, we compared Portal with Seurat, online iNMF and VIPCCA, which had shown  
468 their ability of cross-omics integration in the original publications. A recent state-of-the-art  
469 method, scJoint [55], was also included in the comparison, as it was designed specifically for  
470 scRNA-seq and scATAC-seq data alignment.

471 As shown in the UMAP visualizations (Fig. 5), Portal, scJoint and VIPCCA were able  
472 to align the two datasets correctly, while online iNMF and Seurat did not align some cell  
473 clusters: for example, monocytes in online iNMF’s integration, and a cluster of mixed cell types  
474 from ASAP-seq in Seurat’s result. Among the benchmarked methods, Portal showed superior  
475 performance on the preservation of biological signals. After Portal’s integration, B cells and T  
476 cells were kept as disjoint clusters, while subpopulations of T cells were remained to be close  
477 to each other. In comparison, the coalescence of the previously distinct cell type clusters in  
478 VIPCCA’s result indicates the loss of information. Unlike other methods, scJoint requires cell  
479 type label information of scRNA-seq datasets as its input. It utilizes the cell type annotations

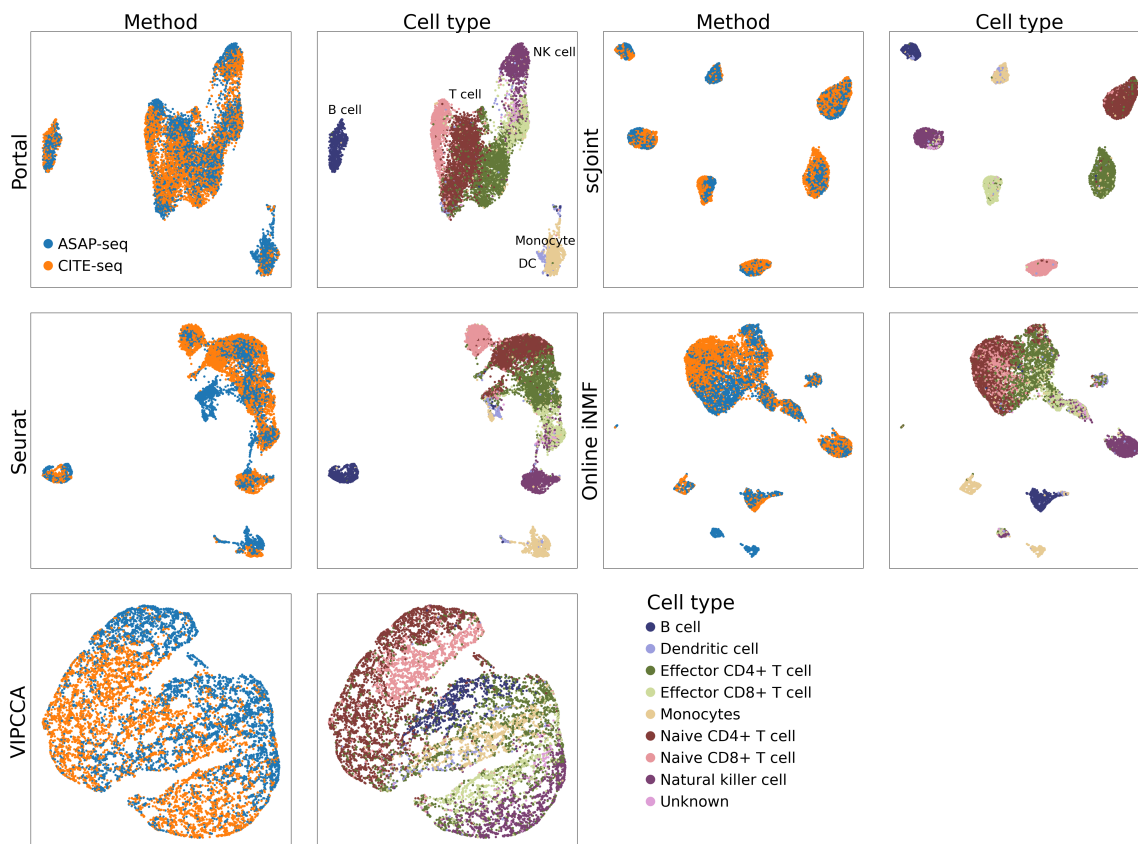


Figure 5: **Comparison of Portal and other cross-omics integration methods on the alignment of scRNA-seq and scATAC-seq data.** We applied Portal, scJoint, Seurat, online iNMF, and VIPCCA to align the scRNA-seq dataset (profiled by CITE-seq) and the scATACseq dataset (profiled by ASAP-seq) of peripheral blood mononuclear cells (PBMCs) [54]. UMAP plots were colored by profiling methods and cell types, respectively.

480 to construct embedding of cells. As a result, cells from the scRNA-seq dataset with different  
 481 cell type labels are forced to form disjoint clusters. Biological information was largely lost in  
 482 scJoint’s integration of PBMC data: the subpopulations of T cells (naive CD4+ T cells, naive  
 483 CD8+ T cells, effector CD4+ T cells, effector CD8+ T cells) lost their similarity and became  
 484 far apart from each other (Fig. 5). Portal and scJoint were also benchmarked with a more  
 485 challenging task: we manually removed B cells from the CITE-seq dataset such that B cells  
 486 became a dataset-specific population. The results further demonstrated Portal’s robustness  
 487 to unbalanced cell type compositions even in cross-omics integration, while scJoint showed  
 488 comparatively inferior performance (Fig. S28).

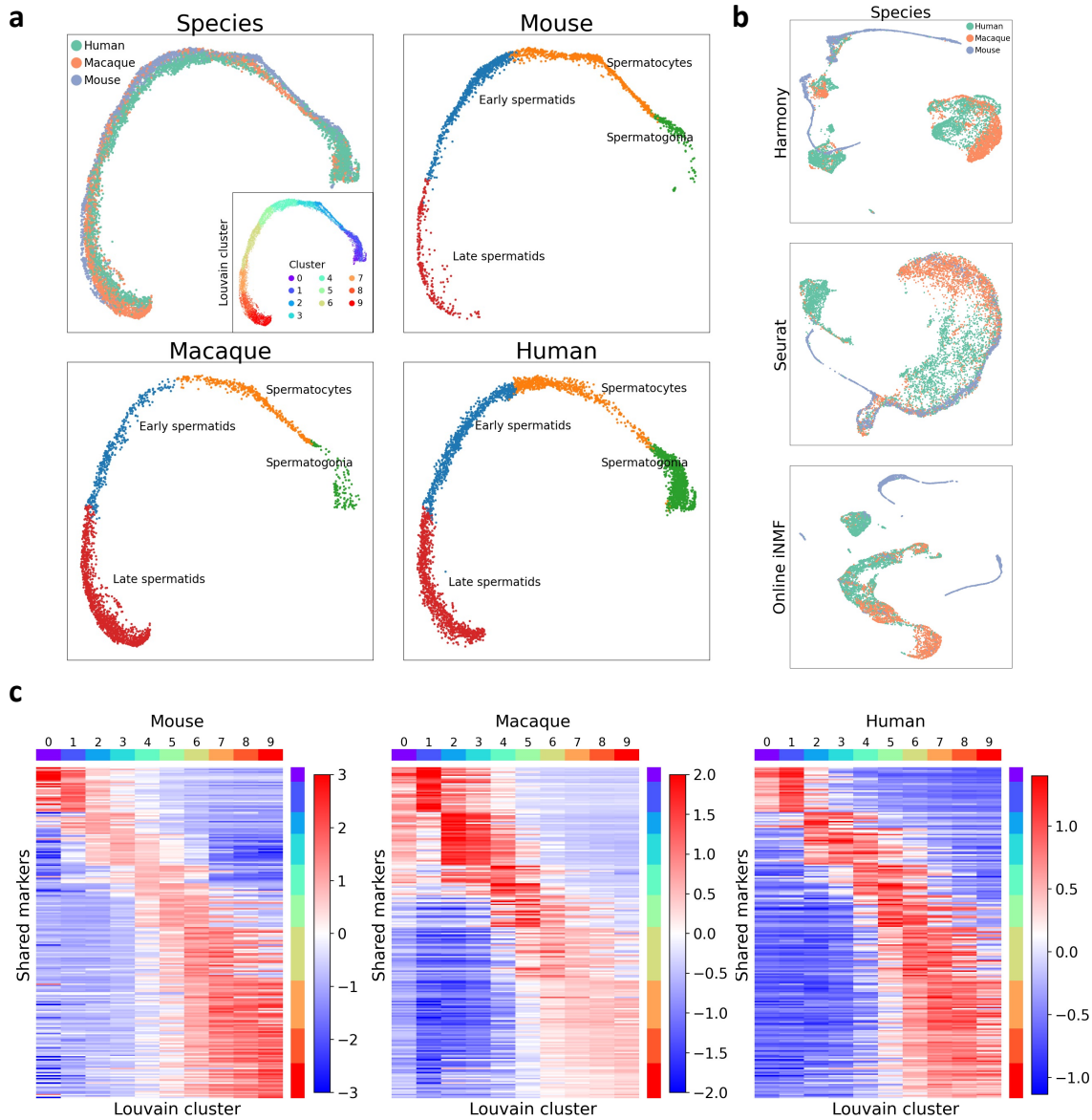


Figure 6: **Integration of spermatogenesis datasets across different species, including mouse, macaque and human.** **a.** The UMAP plot of Portal's result colored by species, as well as UMAP plots of integrated mouse, macaque, human datasets visualized separately. Ten clusters were obtained by applying the Louvain clustering algorithm, facilitating detailed comparative analysis across species. **b.** Integration results of Harmony, Seurat and online iNMF. **c.** Portal identified 228 highly variable genes that are shared in the spermatogenesis process across all three mammalian species.

489 **Portal aligns spermatogenesis differentiation process across multiple**  
490 **species.**

491 Portal does not need to specify the structure and the strength of unwanted variation when  
492 integrating datasets. Instead, it can flexibly account for general difference between datasets,

493 including batch effects, technical noises, and other sources of unwanted variation, by nonlinear  
494 encoders and generators in the adversarial domain translation framework. Therefore, Portal is  
495 also applicable for merging datasets with intrinsic biological divergence, revealing biologically  
496 meaningful connections among these datasets. In this section, we demonstrate that Portal  
497 can successfully align scRNA-seq datasets of the testes from different species including mouse,  
498 macaque and human (Fig. 6).

499 Compared to merging datasets from the same species, cross-species integration poses  
500 additional unique challenges. Although the transcriptomes of different species may share  
501 expression of homologous or orthologous genes, the number of shared genes varies between  
502 different species and is limited. Furthermore, two species may have genes with very similar  
503 sequence and be annotated in the transcriptome by the same name, but have altered function,  
504 which means that expression of the same gene in different species can denote different cell  
505 function [56]. In other words, the amount of information one can utilize for integration becomes  
506 limited and fuzzier while the variation across datasets becomes far larger, with limited number  
507 of shared genes and even fewer shared highly variable genes across different species. Nonetheless,  
508 cross-species integration can be very meaningful despite its challenges, as it can generate quick  
509 draft annotations of new or less-studied species' atlases and cell types via label transfer from  
510 well-studied species. This saves time in the manual annotation process of single-cell tissue atlas  
511 generation for new species. Such integration can also enable detailed comparisons between  
512 species, such as comparisons of cell type composition, discovery of cell types unique to a  
513 particular species, or cross-species comparisons of the same cell types.

514 Mammalian spermatogenesis is a continuous and irreversible differentiation process from  
515 spermatogonial stem cells (SSCs) to sperm cells [57, 58, 59, 60, 16]. Due to the unique  
516 degenerate nature of the Y chromosome (Y-chr), Y-chr gene expression is intricately and  
517 tightly regulated in the spermatogenesis process through meiotic sex chromosome inactivation  
518 (MSCI) [61, 62, 63, 64, 65]. Interestingly, Y-linked genes are highly divergent between different  
519 species, including between closely related primates such as the chimpanzee, macaque, and  
520 human [61, 66, 67]; yet MSCI as a process is conserved across many species and is required for  
521 male fertility [64, 68]. This evidence suggests that while the evolution of genes on the Y-chr  
522 generated diverse species-specific genetic combinations, the tight control of gene expression  
523 through MSCI is required to ensure genetic stability [61]. Recently, cross-species comparisons of

524 “escape genes” that are able to maintain or re-activate their expression despite MSCI repression  
525 during spermatogenesis have generated fascinating insights on evolutionary biology, and on  
526 sex chromosome evolution [63, 65, 69, 16]. In this biological context, integrating datasets  
527 with continuous and gradient developmental trajectories, such as for spermatogenesis data,  
528 requires integration methods to preserve the continuous structure of each dataset, while still  
529 providing high accuracy of cell type alignment between datasets. This is more difficult when,  
530 like in spermatogenesis data, there are no distinct clusters, making integration of such data a  
531 particularly difficult task. After confirming Portal’s capability of preserving the gradual change  
532 of cells based on two examples (Figs. S29 and S30), we perform cross-species integration of  
533 testes datasets from three species, including one mouse [59], one macaque and one human  
534 [16], aligning the different stages of spermatogenesis across species thereby highlighting unique  
535 features of each. The successful integration of these spermatogenesis trajectories serves as a  
536 demonstration of the power of Portal in complex and low-information data alignment, and how  
537 it can facilitate the annotation and discovery process for new single-cell tissue atlases.

538 We first annotated the mouse sample according to the pattern of marker genes (Sper-  
539 matogonia: *Sycp1*, *Uchl1*, *Crabp1*, *Stra8*; Spermatocytes: *Piwil1*, *Pttg1*, *Insl6*, *Spag6*; Early  
540 spermatids: *Tssk1*, *Acrv1*, *Spaca1*, *Tsga8*; Late spermatids: *Prm1*, *Prm2*, *Tnp1*, *Tnp2*) [57, 58].  
541 Then we used Portal to harmonize the three samples, where the integration was accomplished  
542 in the mouse sample domain: The cells from the mouse sample were used as reference, and  
543 cells from the other two species were mapped to the mouse sample domain by Portal. Based on  
544 our annotation of the mouse sample, we transferred the broad cell type labels to cells from the  
545 macaque and human samples according to the nearest neighbors, using the alignment given by  
546 Portal (Fig. 6a). To check whether the alignments were correct for broad cell type identities, we  
547 visualized the UMAPs for cells from each species labeled by their original published annotations  
548 [16], and we confirmed concordant cell type integration across species (Fig. S31). Then, we used  
549 Louvain clustering algorithm to cluster the cells from all three species based on integrated cell  
550 representations. Ten clusters were found, and the cluster names were relabeled by their order  
551 of progression from the spermatogonia along the developmental trajectory (Fig. 6a). We then  
552 visualized the expression of known spermatogenesis markers [57, 58, 16] in each Louvain cluster  
553 and found that the Louvain clusters generated by Portal’s alignment clearly captured the key  
554 transcriptomic features for each stage of spermatogenesis, and correctly identified cells from



555 each stage for all three species (Fig. S32, S33). Furthermore, each Louvain cluster represented  
556 a more fine-grained classification of cells within the labeled broad spermatogenesis cell types.  
557 Using these clusters we assessed the transcriptomic changes throughout the differentiation  
558 trajectory with higher resolution (Fig. S32, S33). Notably, many of the marker genes known  
559 to define stages of spermatogenesis in human were not shared or sometimes not expressed in  
560 macaque and/or mouse scRNA-seq data. For example, human genes *SYCP3*, *YBX2*, *SPACA4*,  
561 *H1FNT*, *PRM1*, and *TNP1* were known to mark human spermatogenesis progression, but  
562 they were absent in the macaque dataset. As only highly variable genes that were expressed  
563 in all three species were considered in the integration process, these genes were not used by  
564 Portal. However, they showed clear expression in the cell clusters where they were expected to  
565 be expressed after integration (Fig. S33), confirming the correctness of Portal's integration  
566 result. The above results show that Portal can provide an accurate integration even for genes  
567 not measured by all three samples. As a comparison, Harmony, Seurat and online iNMF  
568 were also applied. However, Harmony and online iNMF were unable to maintain the gradient  
569 developmental trajectories of spermatogenesis process for at least one species. All of the three  
570 methods showed less satisfactory ability to align cells across the three species (Fig. 6b).

571 Cross-species data integration can be a quick and easy way to generate draft cell atlas  
572 annotations for new species via label transfer from well-annotated species, but moreover, such  
573 integrated data can be used to highlight interesting biological features of shared cell types. In  
574 our Louvain clusters for spermatogenesis, for each species, we selected top 200 highly expressed  
575 genes of every cluster. By taking the intersection of those genes across three species, we then  
576 identified 228 highly variable genes that are shared in the spermatogenesis process across all  
577 three mammalian species (Fig. 6c). For the highly expressed genes that were unique to only one  
578 species, we compared their expressions across all three species (Fig. S34). Such comparisons  
579 could give insight into shared and divergent features of spermatogenesis across different species.

## 580 Discussion

581 Taking advantage of machine learning methodologies, Portal is an efficient and powerful tool for  
582 single-cell data integration that easily scales to handle large datasets with sample sizes in the  
583 millions. As a machine learning-based model, Portal is easy to train, and its training process is  
584 greatly accelerated by using GPUs. Meanwhile, mini-batch optimization allows Portal to be

585 trained with a low memory usage. Besides, it also makes Portal applicable in the situation  
586 where the dataset is not fully observed, but arrives incrementally.

587 The nonlinearity of neural networks makes Portal a flexible approach that can adjust for  
588 complex dataset-specific effects. Nonetheless, according to benchmarking studies, strong ability  
589 for removing dataset-specific effects often comes with the weakness in conserving biological  
590 variation [48, 30], e.g., being prone to overcorrection. Portal overcomes this challenge by its  
591 model and algorithm designs. First, the boundaries of discriminator scores help Portal to  
592 protect dataset-unique cell types from overcorrection. Second, the use of three specifically  
593 designed regularizers not only assists Portal to find correct correspondence across domains, but  
594 also enables Portal to have high-level preservation of subcluster and small cluster identities in  
595 both datasets.

596 Two existing popular methods are Seurat and BBKNN. Seurat often provides integration  
597 results with high accuracy, but also requires high computational cost, preventing its usage on  
598 large-scale datasets; while BBKNN is well-known for its extremely fast speed, its comparatively  
599 less precise results are sometimes a concern for users (Figs. S5 - S17). A major advance of  
600 Portal over these existing state-of-the-art integration approaches is its ability to achieve high  
601 efficiency and accuracy simultaneously. With speed comparable or faster than BBKNN, and  
602 significantly lower memory requirement than BBKNN (Fig. S18), Portal presents similar batch  
603 correction performance as well as superior information preservation performance compared to  
604 that of Seurat (Figs. S5 - S17).

605 Portal also has advantages over several existing deep learning-based methods for single-  
606 cell data integration. Currently, the majority of deep learning-based methods leverages the  
607 variational autoencoders (VAEs) framework [70]. scVI [25], as a prominent representative of  
608 VAE-based methods, is scalable to atlas-level datasets. It utilizes the zero-inflated negative  
609 binomial (ZINB) distribution in its modeling, which may be less efficient in capturing complex  
610 data structures [24]. scANVI [71] is another VAE-based method with similar pros and cons  
611 of scVI, as it is an extension of scVI that incorporates cell type information into its model.  
612 Recently, VIPCCA [24] was proposed to leverage VAE-based networks to perform nonlinear  
613 canonical correlation analysis (CCA) efficiently. However, we empirically found that it favors  
614 the removal of batch effects over the preservation of biological information (Figs. S6, S12 and  
615 S16). scGen [72] utilizes a VAE to find a difference vector in the latent space of each cell type

616 across batches. Similar to scANVI, scGen requires cell type information as its input. There  
617 are also some methods using strategies other than VAEs. One category is deep learning-based  
618 methods utilizing searched MNN pairs as the reference, and then using neural networks to  
619 correct batch effects, such as iMAP [73] and deepMNN [74]. Consequently, the second stage of  
620 correcting batch effects heavily relies on the first stage of constructing MNN pairs. Moreover,  
621 searching MNN pairs is usually performed on CPUs and could be less computationally efficient  
622 for larger datasets. Some deep learning-based methods focus on integrating cross-omics datasets,  
623 including cross-modal autoencoder [75] and scJoint [55]. However, they require additional  
624 information like cell type information or paired data points for data alignment. They may not  
625 be applicable when such information is unavailable. Compared to existing deep learning-based  
626 methods, Portal neither relies on a parametric distribution for single-cell data, nor requires  
627 MNN pairs to serve as anchors for integration. Owing to its unified framework with unique  
628 designs for single-cell datasets, Portal enjoys high flexibility to handle complex datasets and  
629 dataset-specific effects with varying strength, and high scalability to deal with millions of cells  
630 efficiently.

631 By leveraging the adversarial domain translation framework, Portal can build meaningful  
632 alignment between datasets with efficient utilization of information. From single tissue types  
633 to complex cell atlases, Portal showed extraordinary information preservation performance  
634 throughout all integration tasks. This feature of Portal is exemplified by integration of the  
635 spermatogenesis trajectory across three species, where only a limited number of highly variable  
636 genes were shared and utilized by Portal. Improvements can further be made if an effective  
637 way of leveraging the whole transcriptome of all species is developed, which is left for future  
638 work to address. Nonetheless, such cross-species integration allows biologists to easily identify  
639 shared and divergent cellular programs across different species, which is particularly useful  
640 for addressing questions of evolutionary biology. In our example of mouse, macaque, and  
641 human testes tissue integration, identifying genes that are primate-specific can help to generate  
642 hypotheses about the evolution of primates and shed light on the applicability of various animal  
643 models for biological research.

644 It is now clear that using single-cell technologies to assemble comprehensive whole organism  
645 atlases encompassing diverse cell types is accelerating biological discovery, and this demand  
646 will only grow as more datasets are generated. The demand for integration of such datasets,

647 along with the size of these datasets, will expand correspondingly. We expect that Portal, with  
648 its fast, versatile, and robust integration performance, will play a valuable and essential role in  
649 the modern life scientist’s single-cell analysis.

## 650 **Methods**

### 651 **The model of Portal**

652 Expression measurements of cells from two different studies are viewed as datasets originated  
653 from two different domains  $\mathcal{X}$  and  $\mathcal{Y}$ . After standard data preprocessing of the expression  
654 data, Portal performs joint principle component analysis (PCA) across datasets and adopts the  
655 first  $p$  principal components of cells as the low-dimensional representation of cells, namely, cell  
656 embeddings. Portal takes the cell embeddings as the input to achieve data alignment between  
657  $\mathcal{X}$  and  $\mathcal{Y}$ . To learn a harmonized representation of cells, Portal introduces a  $q$ -dimensional  
658 latent space  $\mathcal{Z}$  to connect  $\mathcal{X}$  and  $\mathcal{Y}$ , where the latent codes of cells in  $\mathcal{Z}$  are not affected by  
659 domain-specific effects but capture biological variation.

660 Portal achieves the integration of datasets through training a unified framework of adversarial  
661 domain translation. Let  $\mathbf{x}$  and  $\mathbf{y}$  be the cell embeddings in  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. For domain  
662  $\mathcal{X}$ , Portal first employs encoder  $E_1(\cdot) : \mathcal{X} \rightarrow \mathcal{Z}$  to get a latent code  $E_1(\mathbf{x}) \in \mathcal{Z}$  for all  $\mathbf{x} \in \mathcal{X}$ .  
663 Encoder  $E_1(\cdot)$  is designed to remove domain-specific effects in  $\mathcal{X}$ . To transfer cells from  $\mathcal{X}$  to  
664  $\mathcal{Y}$ , Portal then uses generator  $G_2(\cdot) : \mathcal{Z} \rightarrow \mathcal{Y}$  to model the data generating process in domain  
665  $\mathcal{Y}$ , where domain-specific effects in  $\mathcal{Y}$  are induced.  $E_1(\cdot)$  and  $G_2(\cdot)$  together form a domain  
666 translation network  $G_2(E_1(\cdot))$  that maps cells from  $\mathcal{X}$  to  $\mathcal{Y}$  along  $\mathcal{X} \rightarrow \mathcal{Z} \rightarrow \mathcal{Y}$ . By symmetry,  
667 encoder  $E_2(\cdot) : \mathcal{Y} \rightarrow \mathcal{Z}$  and generator  $G_1(\cdot) : \mathcal{Z} \rightarrow \mathcal{X}$  are utilized to transfer cells from  $\mathcal{Y}$  to  $\mathcal{X}$   
668 along the path  $\mathcal{Y} \rightarrow \mathcal{Z} \rightarrow \mathcal{X}$ .

669 Portal trains domain translation network  $G_2(E_1(\cdot)) : \mathcal{X} \rightarrow \mathcal{Y}$ , such that the distribution  
670 of transferred cells  $G_2(E_1(\mathbf{x}))$  can be mixed with the distribution of cell embeddings  $\mathbf{y}$  in  
671 domain  $\mathcal{Y}$ . Discriminator  $D_2(\cdot)$  is employed in domain  $\mathcal{Y}$  to identify where the poor mixing of  
672 the two distributions occurs. The competition between domain translation network  $G_2(E_1(\cdot))$   
673 and discriminator  $D_2(\cdot)$  is known as adversarial learning [31]. Discriminator  $D_2(\cdot)$  will send a  
674 feedback signal to improve the domain translation network  $G_2(E_1(\cdot))$  until the two distributions  
675 are well mixed. By symmetry, domain translation network  $G_1(E_2(\cdot)) : \mathcal{Y} \rightarrow \mathcal{X}$  and discriminator

676  $D_1(\cdot)$  deployed in domain  $\mathcal{X}$  form another adversarial learning pair. The feedback signal from  
677  $D_1(\cdot)$  improves  $G_1(E_2(\cdot))$  until the well mixing of the transferred cell distribution  $G_1(E_2(\mathbf{y}))$   
678 and the original cell distribution  $\mathbf{x}$  in domain  $\mathcal{X}$ .

679 Notice that the well mixing of the transferred distribution and the original distribution does  
680 not necessarily imply the correct correspondence established between  $\mathcal{X}$  and  $\mathcal{Y}$ . First, cells  
681 from a unique cell population in domain  $\mathcal{X}$  should not be forced to mix with cells in domain  $\mathcal{Y}$ .  
682 Second, cell types  $A$  and  $B$  in domain  $\mathcal{X}$  could be incorrectly aligned with cell types  $B$  and  $A$   
683 in domain  $\mathcal{Y}$ , respectively, even if the two distributions are well mixed. These problems can  
684 occur because we don't have any cell type label information as an anchor for data alignment  
685 across domains. To address these, Portal has the following unique features, distinguishing it  
686 from existing domain translation methods [32, 33]. First, Portal has a tailored discriminator  
687 for the integrative analysis of single-cell data, which can prevent mixing of unique cell types  
688 in one domain with a different type of cell in another domain. Second, Portal deploys three  
689 regularizers to find correct correspondence during adversarial learning; these regularizers also  
690 play a critical role in accounting for domain-specific effects and retaining biological variation in  
691 the shared latent space  $\mathcal{Z}$ .

We propose to train domain translation networks under the following framework:

$$\min_{\{E_1, G_1, E_2, G_2\}} \max_{\{D_1, D_2\}} \mathcal{L}_{\mathcal{X}}(D_1, E_2, G_1) + \mathcal{L}_{\mathcal{Y}}(D_2, E_1, G_2), \quad (2)$$

$$\text{subject to } \mathcal{R}_{\text{AE}}(E_1, G_1, E_2, G_2) \leq t_{\text{AE}}, \quad (3)$$

$$\mathcal{R}_{\text{LA}}(E_1, G_1, E_2, G_2) \leq t_{\text{LA}}, \quad (4)$$

$$\mathcal{R}_{\text{cos}}(E_1, G_1, E_2, G_2) \leq t_{\text{cos}}, \quad (5)$$

692 where component (2) is the objective function of adversarial learning for single-cell data inte-  
693 gration; components (3), (4) and (5) are regularizers for imposing the autoencoder consistency,  
694 the latent alignment consistency and cosine similarity to preserve cross-domain correspondence,  
695 respectively. We have investigated the roles of each component in Portal and provided more  
696 results (Figs. S1 and S2) in the Supplementary Information. We explain each component in  
697 more detail in the next section.

698 **Adversarial learning with discriminator score thresholding.** The adversarial training  
699 between discriminators and domain translation networks is formulated as a min-max opti-  
700 mization problem (2), where  $\mathcal{L}_{\mathcal{X}}(D_1, E_2, G_1) = \mathbb{E}[\log D_1(\mathbf{x})] + \mathbb{E}[\log(1 - D_1(G_1(E_2(\mathbf{y}))))]$  and  
701  $\mathcal{L}_{\mathcal{Y}}(D_2, E_1, G_2) = \mathbb{E}[\log D_2(\mathbf{y})] + \mathbb{E}[\log(1 - D_2(G_2(E_1(\mathbf{x}))))]$  are the objective functions for

702 adversarial learning in domain  $\mathcal{X}$  and domain  $\mathcal{Y}$ , respectively. Given domain translation  
703 network  $G_1(E_2(\cdot))$ , discriminator  $D_1(\cdot) : \mathcal{X} \rightarrow (0, 1)$  is trained to distinguish the transferred  
704 cells  $G_1(E_2(\mathbf{y}))$  from the original cells  $\mathbf{x}$ , where a high score (close to 1) indicates a “real  
705 cell” in domain  $\mathcal{X}$ , and a low score (close to 0) indicates a “transferred cell” from domain  $\mathcal{Y}$ .  
706 This is achieved by maximizing  $\mathcal{L}_{\mathcal{X}}$  with respect to  $D_1(\cdot)$ . Similarly, discriminator  $D_2(\cdot)$  in  
707 domain  $\mathcal{Y}$  is updated by maximizing  $\mathcal{L}_{\mathcal{Y}}$ . Given discriminators  $D_1(\cdot)$  and  $D_2(\cdot)$ , the domain  
708 translation networks are trained by minimizing  $\mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}}$  with respect to  $E_1(\cdot), G_2(\cdot)$  and  
709  $E_2(\cdot), G_1(\cdot)$ , such that the discriminators cannot distinguish transferred cells from real cells.  
710 This is equivalent to  $\min_{\{E_1, G_1, E_2, G_2\}} \mathbb{E}[\log(1 - D_1(G_1(E_2(\mathbf{y}))))] + \mathbb{E}[\log(1 - D_2(G_2(E_1(\mathbf{x}))))]$ .  
711 However, direct optimization of this objective function is known to suffer from severe gradient  
712 vanishing [31, 76]. Therefore, we adopt the “logD-trick” [31] to stabilize the training process.  
713 Denote  $\mathcal{L}_{\mathcal{X}}^{\log D} = -\mathbb{E}[\log D_1(G_1(E_2(\mathbf{y})))]$  and  $\mathcal{L}_{\mathcal{Y}}^{\log D} = -\mathbb{E}[\log D_2(G_2(E_1(\mathbf{x})))]$ . In practice, we  
714 minimize  $\mathcal{L}_{\mathcal{X}}^{\log D} + \mathcal{L}_{\mathcal{Y}}^{\log D} = -\{\mathbb{E}[\log D_1(G_1(E_2(\mathbf{y})))] + \mathbb{E}[\log D_2(G_2(E_1(\mathbf{x})))]\}$  with respect to  
715  $E_1(\cdot), G_2(\cdot)$  and  $E_2(\cdot), G_1(\cdot)$ , instead of minimizing  $\mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}} = \mathbb{E}[\log(1 - D_1(G_1(E_2(\mathbf{y}))))] +$   
716  $\mathbb{E}[\log(1 - D_2(G_2(E_1(\mathbf{x}))))]$ .

717 Although the above adversarial learning can make the transferred cells and real cells well  
718 mixed, it can falsely force cells of a unique cell population in one domain to mix with cells in  
719 another domain, leading to overcorrection. Consider a cell population that is present in  $\mathcal{X}$  but  
720 absent in  $\mathcal{Y}$  as an example. On one hand, discriminator  $D_1(\cdot)$  can easily identify cells from  
721 the unique cell population as real cells in  $\mathcal{X}$ . Cells in the nearby region of this cell population  
722 have extremely high discriminator scores. Some cells in  $\mathcal{Y}$  will be mapped into this region  
723 by the domain translation network  $G_1(E_2(\cdot))$ , leading to incorrect mixing of cell types in  $\mathcal{X}$ .  
724 On the other hand, cells transferred from  $\mathcal{X}$ -unique population will have low  $D_2$  scores in  $\mathcal{Y}$ .  
725 Discriminator  $D_2(\cdot)$  will incorrectly force the domain translation network  $G_2(E_1(\cdot))$  to mix  
726 these cells with real cells in domain  $\mathcal{Y}$ . The cell identity as a domain-unique population in  $\mathcal{X}$   
727 is lost.

728 From the above reasoning, domain-unique cell populations are prone to be assigned with  
729 extreme discriminator scores, either too high in the original domain or too low in the transferred  
730 domain. Such extreme scores can lead to overcorrection. To address this issue in single-cell data  
731 integration tasks, we set boundaries for discriminator scores to make discriminators inactive  
732 on such cells. Specifically, the outputs of standard discriminators are transformed into  $(0, 1)$

733 with the sigmoid function, i.e.,  $D_i(\mathbf{x}) = \text{sigmoid}(d_i(\mathbf{x})) = 1/(1 + \exp(-d_i(\mathbf{x})))$ ,  $i = 1, 2$ , where  
734  $d_i(\mathbf{x}) \in (-\infty, \infty)$  is the logit of the output. We bound the discriminator score by thresholding  
735 its logit to a reasonable range  $[-t, t]$ :

$$\tilde{D}_i(\mathbf{x}) = 1/(1 + \exp(-\text{clamp}(d_i(\mathbf{x})))) \tag{6}$$

736 where  $\text{clamp}(\cdot) = \max(\min(\cdot, t), -t)$ . By clamping the logit  $d_i(\mathbf{x})$ ,  $\tilde{D}_i(\mathbf{x})$  becomes a constant  
737 when  $d_i(\mathbf{x}) < -t$  or  $d_i(\mathbf{x}) > t$ , providing zero gradients for updating the parameters of encoders  
738 and generators. Meanwhile,  $\tilde{D}_i(\mathbf{x})$  remains the same as  $D_i(\mathbf{x})$  when  $d_i(\mathbf{x}) \in [-t, t]$ . By such  
739 design, the adversarial learning mechanism in Portal is only applied to cell populations that  
740 are likely to be common across domains. In Portal, we then use this modified version of  
741 discriminators  $\tilde{D}_i(\cdot)$  to avoid incorrect alignment of domain-unique cell populations. For clarity,  
742 we still use the notation  $D_i(\cdot)$  to represent  $\tilde{D}_i(\cdot)$  hereinafter.

743 **Regularization for autoencoder consistency.** Encoder  $E_1(\cdot) : \mathcal{X} \rightarrow \mathcal{Z}$  and generator  
744  $G_1(\cdot) : \mathcal{Z} \rightarrow \mathcal{X}$  form an autoencoder structure, where  $E_1(\cdot)$  removes domain-specific effects  
745 in  $\mathcal{X}$ , and  $G_1(\cdot)$  recovers them. Similarly,  $E_2(\cdot) : \mathcal{Y} \rightarrow \mathcal{Z}$  and  $G_2(\cdot) : \mathcal{Z} \rightarrow \mathcal{Y}$  form another  
746 autoencoder structure. Therefore, we use the regularizer in (3) for the autoencoder consistency,  
747 where  $\mathcal{R}_{\text{AE}} = \frac{1}{p} \{ \mathbb{E} [\| \mathbf{x} - G_1(E_1(\mathbf{x})) \|_2^2] + \mathbb{E} [\| \mathbf{y} - G_2(E_2(\mathbf{y})) \|_2^2] \}$ ,  $p$  is the dimensionality of  $\mathcal{X}$   
748 and  $\mathcal{Y}$ .

749 **Regularization for cosine similarity correspondence.** Besides the autoencoder consistency,  
750 the cosine similarity regularizer in (5) plays a critical role in data alignment between domains,  
751 where  $\mathcal{R}_{\text{cos}} = \mathbb{E} \left[ 1 - \frac{\langle \mathbf{x}, G_2(E_1(\mathbf{x})) \rangle}{\| \mathbf{x} \|_2 \| G_2(E_1(\mathbf{x})) \|_2} \right] + \mathbb{E} \left[ 1 - \frac{\langle \mathbf{y}, G_1(E_2(\mathbf{y})) \rangle}{\| \mathbf{y} \|_2 \| G_1(E_2(\mathbf{y})) \|_2} \right]$  is the regularizer that imposes  
752 the cross-domain correspondence on domain translation. The key idea is that a cell and its  
753 transferred version should not be largely different from each other in terms of cosine similarity.  
754 This is because cosine similarity is scale invariant and insensitive to domain-specific effects,  
755 including differences in sequencing depth and capture efficiency of protocols used across datasets  
756 [77, 26, 21]. Thus, the cosine similarity regularizer is helpful to uncover robust correspondence  
757 between cells of the same cell type across domains.

758 **Domain-specific effects removal in the shared latent space by latent alignment regu-**  
759 **larization.** Portal decouples domain translation into the encoding process  $\mathcal{X} \rightarrow \mathcal{Z}$  (or  $\mathcal{Y} \rightarrow \mathcal{Z}$ )  
760 and the generating process  $\mathcal{Z} \rightarrow \mathcal{Y}$  (or  $\mathcal{Z} \rightarrow \mathcal{X}$ ). Although adversarial learning enables the do-  
761 main translation networks to effectively transfer cells across domains, it can not remove domain-  
762 specific effects in shared latent space  $\mathcal{Z}$ . To enable encoders  $E_1(\cdot), E_2(\cdot)$  to eliminate domain-

763 specific effects in  $\mathcal{X}$  and  $\mathcal{Y}$ , we propose the latent alignment regularizer in (4) for the consistency  
 764 in latent space  $\mathcal{Z}$ , where  $\mathcal{R}_{\text{LA}} = \frac{1}{q} \{ \mathbb{E} [\|E_1(\mathbf{x}) - E_2(G_2(E_1(\mathbf{x})))\|_2^2] + \mathbb{E} [\|E_2(\mathbf{y}) - E_1(G_1(E_2(\mathbf{y})))\|_2^2] \}$ ,  
 765  $q$  is the dimensionality of  $\mathcal{Z}$ ,  $E_1(\mathbf{x})$  is the latent code of a real cell  $\mathbf{x} \in \mathcal{X}$  and  $E_2(G_2(E_1(\mathbf{x})))$   
 766 is the latent code of its transferred version,  $E_2(\mathbf{y})$  is the latent code of a real cell  $\mathbf{y} \in \mathcal{Y}$  and  
 767  $E_1(G_1(E_2(\mathbf{y})))$  is the latent code of its transferred version. The regularizer (4) encourages the  
 768 latent codes of the same cell to be close to each other. This regularizer helps encoders  $E_1(\cdot)$   
 769 and  $E_2(\cdot)$  to remove domain-specific effects, such that the latent codes in  $\mathcal{Z}$  preserve biological  
 770 variation of cells from different domains.

771 **Algorithm.** Now we develop an alternative updating algorithm to solving the optimization  
 772 problem of adversarial domain translation with the three regularizers. To efficiently solve the  
 773 optimization problem, we replace the constraints (3), (4) and (5) by its Lagrange form. We  
 774 introduce three regularization parameters  $\lambda_{\text{AE}}$ ,  $\lambda_{\text{LA}}$  and  $\lambda_{\text{cos}}$  as coefficients for the regularizers.  
 775 The optimization problem of Portal is rewritten as

$$\min_{\{E_1, G_1, E_2, G_2\}} \max_{\{D_1, D_2\}} \mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}} + \lambda_{\text{AE}} \mathcal{R}_{\text{AE}} + \lambda_{\text{LA}} \mathcal{R}_{\text{LA}} + \lambda_{\text{cos}} \mathcal{R}_{\text{cos}}. \quad (7)$$

776 As we adopt the “logD-trick” for updating domain translation networks formed by  $E_1(\cdot)$ ,  $G_2(\cdot)$   
 777 and  $E_2(\cdot)$ ,  $G_1(\cdot)$ , the optimization problem (7) is modified accordingly as

$$\min_{\{E_1, G_1, E_2, G_2\}} \max_{\{D_1, D_2\}} \mathcal{L}_{\text{adv}} + \lambda_{\text{AE}} \mathcal{R}_{\text{AE}} + \lambda_{\text{LA}} \mathcal{R}_{\text{LA}} + \lambda_{\text{cos}} \mathcal{R}_{\text{cos}},$$

778 where  $\mathcal{L}_{\text{adv}}$  stands for the adversarial learning objective, whose value is  $\mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}}$  when maximizing  
 779 with respect to  $D_1(\cdot)$ ,  $D_2(\cdot)$ , and it is replaced with  $\mathcal{L}_{\mathcal{X}}^{\log D} + \mathcal{L}_{\mathcal{Y}}^{\log D}$  when minimizing with respect  
 780 to  $E_1(\cdot)$ ,  $G_1(\cdot)$ ,  $E_2(\cdot)$ ,  $G_2(\cdot)$ .

781 Let the parameters of the networks  $E_1(\cdot)$ ,  $E_2(\cdot)$ ,  $G_1(\cdot)$ ,  $G_2(\cdot)$ ,  $D_1(\cdot)$  and  $D_2(\cdot)$  be denoted as  
 782  $\theta_{E_1}$ ,  $\theta_{E_2}$ ,  $\theta_{G_1}$ ,  $\theta_{G_2}$ ,  $\theta_{D_1}$  and  $\theta_{D_2}$ . Then we collect the parameter sets as  $\theta_E = \{\theta_{E_1}, \theta_{E_2}\}$ ,  $\theta_G =$   
 783  $\{\theta_{G_1}, \theta_{G_2}\}$  and  $\theta_D = \{\theta_{D_1}, \theta_{D_2}\}$ . We use the Monte Carlo estimators to approximate expectations  
 784 in Portal’s objective. With a mini-batch of  $2m$  samples including  $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$  from  $\mathcal{X}$



785 and  $\{\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(m)}\}$  from  $\mathcal{Y}$ , the Monte Carlo estimators are given by

$$\begin{aligned}\widehat{\mathcal{L}}_{\mathcal{X}} &= \frac{1}{m} \sum_{i=1}^m [\log D_1(\mathbf{x}^{(i)}) + \log(1 - D_1(G_1(E_2(\mathbf{y}^{(i)}))))], & \widehat{\mathcal{L}}_{\mathcal{X}}^{\log D} &= -\frac{1}{m} \sum_{i=1}^m \log D_1(G_1(E_2(\mathbf{y}^{(i)}))), \\ \widehat{\mathcal{L}}_{\mathcal{Y}} &= \frac{1}{m} \sum_{i=1}^m [\log D_2(\mathbf{y}^{(i)}) + \log(1 - D_2(G_2(E_1(\mathbf{x}^{(i)}))))], & \widehat{\mathcal{L}}_{\mathcal{Y}}^{\log D} &= -\frac{1}{m} \sum_{i=1}^m \log D_2(G_2(E_1(\mathbf{x}^{(i)}))), \\ \widehat{\mathcal{R}}_{\text{AE}} &= \frac{1}{mp} \sum_{i=1}^m [\|\mathbf{x}^{(i)} - G_1(E_1(\mathbf{x}^{(i)}))\|_2^2 + \|\mathbf{y}^{(i)} - G_2(E_2(\mathbf{y}^{(i)}))\|_2^2], \\ \widehat{\mathcal{R}}_{\text{LA}} &= \frac{1}{mq} \sum_{i=1}^m [\|E_1(\mathbf{x}^{(i)}) - E_2(G_2(E_1(\mathbf{x}^{(i)})))\|_2^2 + \|E_2(\mathbf{y}^{(i)}) - E_1(G_1(E_2(\mathbf{y}^{(i)})))\|_2^2], \\ \widehat{\mathcal{R}}_{\text{cos}} &= \frac{1}{m} \sum_{i=1}^m \left\{ \left[ 1 - \frac{\langle \mathbf{x}^{(i)}, G_2(E_1(\mathbf{x}^{(i)})) \rangle}{\|\mathbf{x}^{(i)}\|_2 \|G_2(E_1(\mathbf{x}^{(i)}))\|_2} \right] + \left[ 1 - \frac{\langle \mathbf{y}^{(i)}, G_1(E_2(\mathbf{y}^{(i)})) \rangle}{\|\mathbf{y}^{(i)}\|_2 \|G_1(E_2(\mathbf{y}^{(i)}))\|_2} \right] \right\}.\end{aligned}$$

786 The implementation of Portal is summarized in Algorithm 1.

---

**Algorithm 1** Stochastic gradient descent training of Portal.

---

**Require:** Batch size  $m$ , coefficients  $\lambda_{\text{AE}}$ ,  $\lambda_{\text{LA}}$  and  $\lambda_{\text{cos}}$

**for** number of training iterations **do**

Sample  $m$  cells  $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$  from  $\mathcal{X}$  and  $m$  cells  $\{\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(m)}\}$  from  $\mathcal{Y}$ .

Calculate  $\widehat{\mathcal{L}}_{\mathcal{X}}$ ,  $\widehat{\mathcal{L}}_{\mathcal{Y}}$ ,  $\widehat{\mathcal{L}}_{\mathcal{X}}^{\log D}$ ,  $\widehat{\mathcal{L}}_{\mathcal{Y}}^{\log D}$ ,  $\widehat{\mathcal{R}}_{\text{AE}}$ ,  $\widehat{\mathcal{R}}_{\text{LA}}$ , and  $\widehat{\mathcal{R}}_{\text{cos}}$ .

Update discriminators by stochastic gradient descent with  $\nabla_{\theta_D}[-(\widehat{\mathcal{L}}_{\mathcal{X}} + \widehat{\mathcal{L}}_{\mathcal{Y}})]$ .

Update encoders and generators simultaneously by stochastic gradient descent with

$$\nabla_{\theta_E, \theta_G}(\widehat{\mathcal{L}}_{\mathcal{X}}^{\log D} + \widehat{\mathcal{L}}_{\mathcal{Y}}^{\log D} + \lambda_{\text{AE}}\widehat{\mathcal{R}}_{\text{AE}} + \lambda_{\text{LA}}\widehat{\mathcal{R}}_{\text{LA}} + \lambda_{\text{cos}}\widehat{\mathcal{R}}_{\text{cos}}).$$

**end for**

---

787 After training, cells from domains  $\mathcal{X}$  and  $\mathcal{Y}$  are encoded into  $\mathcal{Z}$  to construct an integrated  
788 dataset, which can be applied to downstream analysis. In each domain, the original cells and  
789 transferred cells are also well integrated. For integration of multiple datasets, Portal can handle  
790 them incrementally, by transferring all other datasets into the domain formed by one dataset.  
791 **Network structure.** Portal uses lightweight networks which enable computationally efficient  
792 training when dealing with large-scale datasets. The details of Portal's networks, including  
793 network structures, the number of layers and parameters are shown in the Supplementary  
794 Information (Tables S1, S2 and S3).

## 795 Analysis details

796 **Data preprocessing.** We used raw read or unique molecular identifier (UMI) matrices depend-  
797 ing on the data source for all scRNA-seq and snRNA-seq datasets, and gene activity matrices

798 for scATAC-seq datasets. We then performed standard data preprocessing for each count  
799 matrix, including log-normalization, feature selection, scaling and dimensionality reduction. For  
800 each dataset represented by a cell-by-gene count matrix, we first adopted the log-normalization,  
801 following the Seurat and Scanpy pipelines [22, 78]. For each cell, its library size was normalized  
802 to 10,000 reads. Specifically, the counts abundance of each gene was divided by the total counts  
803 for each cell, then multiplied by a scaling factor of 10,000. The normalized dataset was then  
804 transformed to log scale by the function  $\log(1 + x)$ . In order to identify a subset of features  
805 that highlight variability across individual cells, we adopted the feature selection procedure  
806 from the Seurat pipeline. For each dataset, we selected  $K$  top highly variable genes ranked by  
807 dispersion with the control of means. In this paper, we used  $K = 4,000$  throughout all analyses  
808 except for the cross-species analysis. In the cross-species analysis, we used  $K = 3,000$  since the  
809 usage of a larger number of features would result in the situation that correspondence across  
810 species is dominated by the distinction (e.g., altered functions of genes annotated by the same  
811 name). For each selected variable gene, we centered and standardized its expressions across  
812 individual cells to have mean at zero and variance at one. After the above procedures, which  
813 were applied to individual datasets, we continued to preprocess data across datasets. For those  
814 datasets to be integrated, we collected genes that were identified as top highly variable genes in  
815 all of them as features for integration. We extracted the scaled data with these features from  
816 each dataset, and then concatenated them based on features to perform joint PCA. Top  $p = 30$   
817 principle components were kept for all dataset as inputs to Portal. For the shared latent space,  
818 we set its dimensionality to be  $q = 20$  throughout all analyses.

819 **Unifying gene names for cross-species integration.** We retrieved pairwise orthologues  
820 (human vs mouse, human vs macaque) respectively from Ensembl Biomart, and merged them  
821 to obtain one-to-one-to-one orthologues by using human Ensembl gene names as reference.  
822 One-to-one-to-one orthologues across the three species were used to unify gene names. Genes  
823 included in the list were used by Portal. To facilitate the usage of Portal, we have included  
824 the used gene lists (orthologues\_human\_mouse.txt, orthologues\_human\_macaque.txt) as well as  
825 the reproducible code for the cross-species integration, among all details for reproducing the  
826 experiments throughout our paper at <https://github.com/YangLabHKUST/Portal>.

827 **Hyperparameter setting.** Hyperparameters used in Portal are  $m, t, \lambda_{AE}, \lambda_{LA}, \lambda_{\cos}$ , where  $m$   
828 is the batch size used by Portal for mini-batch training;  $t$  is the absolute value of boundaries

829 for the logit of discriminator scores ( $-t < d_i(\mathbf{x}) < t$ ,  $i = 1, 2$ );  $\lambda_{\text{AE}}$ ,  $\lambda_{\text{LA}}$ ,  $\lambda_{\text{cos}}$  are coefficients for  
830 autoencoder consistency regularizer  $\mathcal{R}_{\text{AE}}$ , latent alignment regularizer  $\mathcal{R}_{\text{LA}}$  and cosine similarity  
831 regularizer  $\mathcal{R}_{\text{cos}}$  respectively. Throughout all analyses, we set  $m = 500$ ,  $t = 5.0$ ,  $\lambda_{\text{AE}} = 10.0$ ,  
832  $\lambda_{\text{LA}} = 10.0$ . Hyperparameter  $\lambda_{\text{cos}}$  was tuned within the range  $[10.0, 50.0]$  with interval 5.0  
833 according to the mixing metric, where the mixing metric was designed in Seurat to evaluate  
834 how well the datasets mixed after integration. The insight into tuning  $\lambda_{\text{cos}}$  is as follows: During  
835 domain translations, there is a trade-off between preservation of similarity across domains  
836 and flexibility of modeling domain differences. Since  $\mathcal{R}_{\text{cos}}$  is designed to preserve the cosine  
837 similarity during translations, a higher value of  $\lambda_{\text{cos}}$  can enhance the cosine similarity as the  
838 cross-domain correspondence, and a lower  $\lambda_{\text{cos}}$  allows domain translation networks to deal with  
839 remarkable differences between domains. Following this intuition, we empirically find out that  
840  $\lambda_{\text{cos}} = 10.0$  has a good performance when harmonizing datasets with intrinsic differences, for  
841 example, datasets used in cross-species analysis or cross-modal integration (scRNA-seq and  
842 scATAC-seq). For other integration tasks,  $\lambda_{\text{cos}} = 20.0$  often yields reasonable results, which  
843 is adopted as the default setting in our package. Slightly better alignment results could be  
844 achieved by tuning  $\lambda_{\text{cos}}$ . Through a parameter sensitivity analysis, we have shown that Portal's  
845 performance is insensitive to the choice of hyperparameters (Figs. S3, S4 in the Supplementary  
846 Information).

847 **Label transfer.** Suppose we wish to transfer labels from domain  $\mathcal{X}$  to domain  $\mathcal{Y}$ . As Portal  
848 produces integrated cell representations in each domain and the shared latent space, we can  
849 use any of these representations to perform label transfer. For each cell in domain  $\mathcal{Y}$ , we find  
850 its  $k = 20$ -nearest neighbors among the cells in domain  $\mathcal{X}$  based on the integrated result. The  
851 metric for finding nearest neighbors can be Euclidean distance in shared latent space, or cosine  
852 similarity in domains. The labels in domain  $\mathcal{Y}$  are finally determined by majority voting.

853 **Evaluation metrics.** We assessed all metrics based on Portal's integration results in shared  
854 latent space  $\mathcal{Z}$ . We used kBET [35], PCR batch [35], batch ASW [35], graph iLISI [30, 21] and  
855 graph connectivity [30] to assess the ability of batch correction. We used ARI [36], NMI [37],  
856 cell type ASW, graph cLISI [30, 21], isolated label F1 [30], isolated label silhouette [30] and  
857 cell cycle conservation [30] to evaluate the conservation of biological variation. The metrics, if  
858 necessary, were rescaled to  $[0, 1]$  such that a higher value represents a better performance.

859 *kBET.* For each selected cell, kBET adopts a Pearson's  $\chi^2$ -based test to check whether the

860 batch label distribution in its neighbourhood is similar to the global batch label distribution or  
861 not. In our experiments, we ran 100 replicates of kBET with 1,000 random samples, and used  
862 the median of average acceptance rates as the output. The neighbourhood size was chosen  
863 following the default setting in kBET’s official code.

864 *PCR batch.* PCR batch quantifies the removal of batch effects by comparing the variance  
865 contributions of the batch effects to datasets before integration ( $VC_{\text{before}}$ ) and after integration  
866 ( $VC_{\text{after}}$ ), respectively. In our experiments, we concatenated datasets by batches to obtain the  
867 dataset before integration. PCR batch score was calculated as  $\frac{VC_{\text{before}} - VC_{\text{after}}}{VC_{\text{before}}}$ . We clamped  
868 PCR batch score to  $[0, 1]$ , where a higher score means that the impact of batch effects is  
869 eliminated after integration.

870 *Batch ASW.* Batch ASW calculates the silhouette width of cells with respect to batch labels.  
871 If batch effects are corrected in cell embeddings, the evaluated ASW (with respect to batches)  
872 should be close to -1, indicating the good mixing of cells across batches. We rescaled the score  
873 by  $\frac{1 - \text{ASW}(\text{batch})}{2}$ .

874 *Graph iLISI.* The original iLISI is defined as the effective number of datasets in a neigh-  
875 borhood, where 1 means poor mixing, and 2 indicates good mixing of two datasets. Graph  
876 iLISI extends iLISI by enabling the calculation on graphs. The values were rescaled to  $[0, 1]$  by  
877 subtracting 1.

878 *Graph connectivity.* Graph connectivity assesses whether the graph correctly connects  
879 cells of same cell type labels among batches. We used the Scanpy pipeline to derive graph  
880 representation of integrated cell embeddings. The neighborhood size was set to be 15 (default  
881 setting in Scanpy).

882 *ARI.* ARI measures the degree to which the two clustering results match. It ranges from  
883 0 to 1, where 0 indicates that the two clustering labels are independent to each other, and 1  
884 means that the two clustering labels are the same up to a permutation. We obtained clustering  
885 results following the Seurat clustering pipeline with its default setting, and assessed ARI by  
886 comparing identified clusters and cell type annotations.

887 *NMI.* NMI computes normalized mutual information between two clustering results, ranging  
888 from 0 to 1. An NMI value close to 0 means that there is nearly no mutual information, while  
889 a value close to 1 indicates high correlation between the two clustering results. Similar to ARI,  
890 we calculated NMI with clusters identified by the Seurat pipeline and cell type annotations.

891 *Cell type ASW.* Cell type ASW evaluates ASW with respect to cell type labels, where a  
892 higher score means that cells are closer to cells of the same cell type. As ASW lies between -1  
893 and 1, we rescaled the score by cell type  $ASW = \frac{1+ASW(\text{cell type})}{2}$ .

894 *Graph cLISI.* The original cLISI measures the effective number of cell types in a neighbor-  
895 hood, where 1 means that the cell population is well preserved, and larger values indicate the  
896 mixing of different cell populations. Graph cLISI extends cLISI by enabling the calculation on  
897 graphs. The values were rescaled to [0, 1], where higher values indicated good performance of  
898 preserving biological variation.

899 *Isolated label F1.* Isolated label F1 is developed to measure the ability of integration methods  
900 to preserve dataset-specific cell types. We adopted the Seurat pipeline to cluster cells in the  
901 integrated dataset, and evaluated the cluster assignment of dataset-specific cell types based on  
902 the F1 score [79]. Isolated label F1 ranges between 0 and 1, where 1 shows that all cells of  
903 dataset-specific cell types are captured in separate clusters.

904 *Isolated label silhouette.* Isolated label silhouette, which is similar to Isolated label F1,  
905 also measures the conservation of dataset-specific cell types. Instead of using the F1 score, it  
906 evaluates ASW of dataset-specific cell types. In our experiments, we rescaled the score to [0,1].

907 *Cell cycle conservation.* Cell cycle conservation measures how well the cell cycle effect is  
908 preserved by integration approaches. It compares cell cycle scores before integration ( $CC_{\text{before}}$ )  
909 and after integration ( $CC_{\text{after}}$ ) by calculating  $\frac{|CC_{\text{before}} - CC_{\text{after}}|}{CC_{\text{before}}}$ , where score 0 indicates perfect  
910 conservation of cell cycle effects. We used the gene list from the study [80] as reference, and  
911 calculated the cell cycle score based on the Scanpy pipeline. We rescaled the score to [0,1] such  
912 that a higher score indicates a better result.

913 **Benchmarking of the running time and the memory usage.** Standard data preprocessing  
914 such as normalization, feature selection and dimension reduction could be performed incre-  
915 mentally using mini-batches to control memory usage. In Portal's preprocessing, we adopted  
916 the incremental strategy and used a chunk size of 20,000. For example, the preprocessing of  
917 Portal took 63.4 minutes, requiring 22.0 GB peak running memory on the two mouse brain  
918 atlases datasets with 1,100,167 cells. The preprocessing time could be reduced to 37.7 minutes  
919 when the chunk size was increased to 200,000, with 36.4 GB peak running memory. Some  
920 other methods may not be able to adopt a mini-batch implementation. For the two mouse  
921 brain atlases datasets, Harmony took 17.6 minutes to finish preprocessing, but required 127.1

922 GB memory usage. Online iNMF performed preprocessing with mini-batches. Its default  
923 preprocessing procedure on the two mouse brain atlases datasets took 15.9 hours, with 0.6  
924 GB memory usage. For a fair comparison, the time and memory usages of data preprocessing  
925 procedures were not included in our benchmarking.

926 **Integration of multiple datasets.** For multiple datasets, Portal integrates them in an  
927 incremental manner, by transferring all other datasets into the domain constructed by the first  
928 dataset. Here we used the integration of two scRNA-seq datasets (profiled by Drop-seq and  
929 10X) [8, 9] and one snRNA-seq dataset (profiled by SPLiT-seq) [49] to illustrate this procedure.  
930 In this example, Portal ran in two steps:

931 *Step 1.* Portal trained domain translation networks between the 10X dataset (160,678 cells)  
932 and the Drop-seq dataset (319,359 cells), which took 45.48s. Then Portal used the trained  
933 networks to map 10X cells to the Drop-seq dataset domain, which took 0.08s.

934 *Step 2.* Portal trained domain translation networks between the SPLiT-seq dataset (74,159  
935 cells) and the integrated 10X and Drop-seq dataset, which took 41.36s. Then Portal mapped  
936 SPLiT-seq cells to the integrated 10X and Drop-seq dataset domain, which took 0.06s.

937 In total, Portal took 86.98s to integrate all three datasets.

938 The integration of multiple datasets is implemented in one function in Portal package. The  
939 code for reproducing the experiment is available as a Jupyter Notebook at <https://github.com/YangLabHKUST/Portal>, serving as an example for the integration of multiple datasets.

941 **Visualization.** We used the UMAP algorithm [39] for visualization of cell representations  
942 in a two-dimensional space. In all analyses, the UMAP algorithm was run with 30-nearest  
943 neighbors, minimum distance 0.3, and correlation metric.

## 944 Acknowledgements

945 The authors would like to thank Camille Sophie Ezran (Stanford University), Dr. Angela  
946 Oliveira Pisco (CZ Biohub), and Dr. Hosu Sin (Stanford University) for valuable discussions.  
947 This work is supported in part by Hong Kong Research Grant Council [16101118, 24301419,  
948 14301120, 16307818, 16301419, 16308120], the Hong Kong University of Science and Technology's  
949 startup grant [R9405,R9364], the Hong Kong University of Science and Technology Big Data  
950 for Bio Intelligence Laboratory (BDBI), the Lo Ka Chung Foundation through the Hong Kong  
951 Epigenomics Project, the Chau Hoi Shuen Foundation, the Chinese University of Hong Kong

952 direct grants [4053360, 4053423, 4053476], the Chinese University of Hong Kong startup grant  
953 [4930181], the Chinese University of Hong Kong's Project Impact Enhancement Fund (PIEF)  
954 and Science Faculty's Collaborative Research Impact Matching Scheme (CRIMS), the East  
955 China Normal University startup grant, the Shanghai Sailing Program. The computational  
956 task for this work was partially performed using the X-GPU cluster supported by the RGC  
957 Collaborative Research Fund: C6021-19EF.

## 958 **Author contributions**

959 J.Z. and G.W. conceived and developed the method. A.R.W. and C.Y. supervised the project.  
960 J.Z., G.W., Z.L., A.R.W. and C.Y. designed the experiments, performed the analyses and wrote  
961 the paper. J.M., Y.W. and T.M.C. provided critical feedback during the study and helped  
962 revise the manuscript.

## 963 **Data availability**

964 All data used in this work are publicly available through online sources.

- 965 • Mouse brain cells from Saunders et al. [8] (<http://dropviz.org>).
- 966 • Mouse brain cells from Zeisel et al. [9] (<http://mousebrain.org/downloads.html>).
- 967 • Mouse brain cells from Rosenberg et al. [49] (GSE110823).
- 968 • Mouse cell atlas from the Tabula Muris Consortium [7] ([https://figshare.com/projects/  
969 Tabula\\_Muris\\_Transcriptomic\\_characterization\\_of\\_20\\_organs\\_and\\_tissues\\_from\\_  
970 Mus\\_musculus\\_at\\_single\\_cell\\_resolution/27733](https://figshare.com/projects/Tabula_Muris_Transcriptomic_characterization_of_20_organs_and_tissues_from_Mus_musculus_at_single_cell_resolution/27733)).
- 971 • Mouse lemur cell atlas from Tabula Microcebus Consortium ([https://figshare.com/  
972 projects/Tabula\\_Microcebus/112227](https://figshare.com/projects/Tabula_Microcebus/112227)).
- 973 • Human peripheral blood mononuclear cells from Mimitou et al. [54] (GSE156478).
- 974 • Human peripheral blood mononuclear cells from Ding et al. [81] (GSE132044).
- 975 • Human peripheral blood mononuclear cells from 10X Genomics. [81] ([https://support.  
976 10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/pbmc3k](https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/pbmc3k)).

- 977 • Mouse spermatogenesis cells from Ernst et al. [59] ([https://www.ebi.ac.uk/arrayexpress/](https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-6946/)  
978 [experiments/E-MTAB-6946/](https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-6946/)).
- 979 • Human spermatogenesis cells from Shami et al. [16] (GSE142585).
- 980 • Macaque spermatogenesis cells from Shami et al. [16] (GSE142585).
- 981 • Hematopoietic stem cells from Paul et al. [82] (GSE72857).
- 982 • Hematopoietic stem cells from Nestorowa et al. [83] (GSE81682).
- 983 • Reprogramming of induced pluripotent stem cells from Schiebinger et al. [84] (GSE122662).
- 984 • Human brain cells from Fullard et al. [51] (GSE164485).
- 985 • Human brain cells from Tran et al. [52] ([https://github.com/LieberInstitute/](https://github.com/LieberInstitute/10xPilot_snRNAseq-human#work-with-the-data)  
986 [10xPilot\\_snRNAseq-human#work-with-the-data](https://github.com/LieberInstitute/10xPilot_snRNAseq-human#work-with-the-data)).

## 987 Code availability

988 Portal software is available at <https://github.com/YangLabHKUST/Portal>.



## 989 References

- 990 [1] Alexandra-Chloé Villani, Rahul Satija, Gary Reynolds, Siranush Sarkizova, Karthik  
991 Shekhar, James Fletcher, Morgane Griesbeck, Andrew Butler, Shiwei Zheng, Suzan Lazo,  
992 et al. Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes,  
993 and progenitors. *Science*, 356(6335), 2017.
- 994 [2] Maria Brbić, Marinka Zitnik, Sheng Wang, Angela O Pisco, Russ B Altman, Spyros  
995 Darmanis, and Jure Leskovec. Mars: discovering novel cell types across heterogeneous  
996 single-cell experiments. *Nature methods*, 17(12):1200–1206, 2020.
- 997 [3] Giovanni Iacono, Ramon Massoni-Badosa, and Holger Heyn. Single-cell transcriptomics  
998 unveils gene regulatory network plasticity. *Genome biology*, 20(1):1–20, 2019.
- 999 [4] Anna SE Cuomo, Daniel D Seaton, Davis J McCarthy, Iker Martinez, Marc Jan Bonder,  
1000 Jose Garcia-Bernardo, Shradha Amatya, Pedro Madrigal, Abigail Isaacson, Florian Buet-  
1001 tner, et al. Single-cell RNA-sequencing of differentiating iPS cells reveals dynamic genetic  
1002 effects on gene expression. *Nature communications*, 11(1):1–14, 2020.
- 1003 [5] Barbara Treutlein, Doug G Brownfield, Angela R Wu, Norma F Neff, Gary L Mantalas,  
1004 F Hernan Espinoza, Tushar J Desai, Mark A Krasnow, and Stephen R Quake. Reconstruct-  
1005 ing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. *Nature*,  
1006 509(7500):371–375, 2014.
- 1007 [6] Xiaojie Qiu, Qi Mao, Ying Tang, Li Wang, Raghav Chawla, Hannah A Pliner, and Cole  
1008 Trapnell. Reversed graph embedding resolves complex single-cell trajectories. *Nature*  
1009 *methods*, 14(10):979–982, 2017.
- 1010 [7] Tabula Muris Consortium et al. Single-cell transcriptomics of 20 mouse organs creates a  
1011 Tabula Muris. *Nature*, 562(7727):367–372, 2018.
- 1012 [8] Arpiar Saunders, Evan Z Macosko, Alec Wysoker, Melissa Goldman, Fenna M Krienen,  
1013 Heather de Rivera, Elizabeth Bien, Matthew Baum, Laura Bortolin, Shuyu Wang, et al.  
1014 Molecular diversity and specializations among the cells of the adult mouse brain. *Cell*,  
1015 174(4):1015–1030, 2018.

- 1016 [9] Amit Zeisel, Hannah Hochgerner, Peter Lönnerberg, Anna Johnsson, Fatima Memic, Job  
1017 Van Der Zwan, Martin Häring, Emelie Braun, Lars E Borm, Gioele La Manno, et al.  
1018 Molecular architecture of the mouse nervous system. *Cell*, 174(4):999–1014, 2018.
- 1019 [10] Hoa Thi Nhu Tran, Kok Siong Ang, Marion Chevrier, Xiaomeng Zhang, Nicole Yee Shin  
1020 Lee, Michelle Goh, and Jinmiao Chen. A benchmark of batch-effect correction methods  
1021 for single-cell RNA sequencing data. *Genome biology*, 21(1):1–32, 2020.
- 1022 [11] Oliver Stegle, Sarah A Teichmann, and John C Marioni. Computational and analytical  
1023 challenges in single-cell transcriptomics. *Nature Reviews Genetics*, 16(3):133–145, 2015.
- 1024 [12] Monika Litviňuková, Carlos Talavera-López, Henrike Maatz, Daniel Reichart, Catherine L  
1025 Worth, Eric L Lindberg, Masatoshi Kanda, Krzysztof Polanski, Matthias Heinig, Michael  
1026 Lee, et al. Cells of the adult human heart. *Nature*, 588(7838):466–472, 2020.
- 1027 [13] Kyle J Travaglini, Ahmad N Nabhan, Lolita Penland, Rahul Sinha, Astrid Gillich, Rene V  
1028 Sit, Stephen Chang, Stephanie D Conley, Yasuo Mori, Jun Seita, et al. A molecular cell  
1029 atlas of the human lung from single-cell RNA sequencing. *Nature*, 587(7835):619–625,  
1030 2020.
- 1031 [14] Blue B Lake, Simone Codeluppi, Yun C Yung, Derek Gao, Jerold Chun, Peter V  
1032 Kharchenko, Sten Linnarsson, and Kun Zhang. A comparative strategy for single-nucleus  
1033 and single-cell transcriptomes confirms accuracy in predicted cell-type expression from  
1034 nuclear RNA. *Scientific reports*, 7(1):1–8, 2017.
- 1035 [15] Haojia Wu, Yuhei Kirita, Erinn L Donnelly, and Benjamin D Humphreys. Advantages of  
1036 single-nucleus over single-cell RNA sequencing of adult kidney: rare cell types and novel  
1037 cell states revealed in fibrosis. *Journal of the American Society of Nephrology*, 30(1):23–32,  
1038 2019.
- 1039 [16] Adrienne Niederriter Shami, Xianing Zheng, Sarah K Munyoki, Qianyi Ma, Gabriel L  
1040 Manske, Christopher D Green, Meena Sukhwani, Kyle E Orwig, Jun Z Li, and Saher Sue  
1041 Hammoud. Single-cell RNA sequencing of human, macaque, and mouse testes uncovers  
1042 conserved and divergent features of mammalian spermatogenesis. *Developmental cell*,  
1043 54(4):529–547, 2020.

- 1044 [17] Aviv Regev, Sarah A Teichmann, Eric S Lander, Ido Amit, Christophe Benoist, Ewan  
1045 Birney, Bernd Bodenmiller, Peter Campbell, Piero Carninci, Menna Clatworthy, et al.  
1046 Science forum: the human cell atlas. *elife*, 6:e27041, 2017.
- 1047 [18] Xiaoping Han, Renying Wang, Yincong Zhou, Lijiang Fei, Huiyu Sun, Shujing Lai, Assieh  
1048 Saadatpour, Ziming Zhou, Haide Chen, Fang Ye, et al. Mapping the mouse cell atlas by  
1049 microwell-seq. *Cell*, 172(5):1091–1107, 2018.
- 1050 [19] Si Wang, Yuxuan Zheng, Jingyi Li, Yang Yu, Weiqi Zhang, Moshi Song, Zunpeng Liu,  
1051 Zheyang Min, Huifang Hu, Ying Jing, et al. Single-cell transcriptomic atlas of primate  
1052 ovarian aging. *Cell*, 180(3):585–600, 2020.
- 1053 [20] Shuai Ma, Shuhui Sun, Jiaming Li, Yanling Fan, Jing Qu, Liang Sun, Si Wang, Yiyuan  
1054 Zhang, Shanshan Yang, Zunpeng Liu, et al. Single-cell transcriptomic atlas of primate  
1055 cardiopulmonary aging. *Cell research*, 31(4):415–432, 2021.
- 1056 [21] Ilya Korsunsky, Nghia Millard, Jean Fan, Kamil Slowikowski, Fan Zhang, Kevin Wei, Yuriy  
1057 Baglaenko, Michael Brenner, Po-ru Loh, and Soumya Raychaudhuri. Fast, sensitive and  
1058 accurate integration of single-cell data with Harmony. *Nature methods*, 16(12):1289–1296,  
1059 2019.
- 1060 [22] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi,  
1061 William M Mauck III, Yuhan Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija.  
1062 Comprehensive integration of single-cell data. *Cell*, 177(7):1888–1902, 2019.
- 1063 [23] Chao Gao, Jialin Liu, April R Kriebel, Sebastian Preissl, Chongyuan Luo, Rosa Castanon,  
1064 Justin Sandoval, Angeline Rivkin, Joseph R Nery, Margarita M Behrens, et al. Iterative  
1065 single-cell multi-omic integration using online learning. *Nature Biotechnology*, pages 1–8,  
1066 2021.
- 1067 [24] Jialu Hu, Mengjie Chen, and Xiang Zhou. Effective and scalable single-cell data alignment  
1068 with non-linear canonical correlation analysis. *Nucleic acids research*, 50(4):e21–e21, 2022.
- 1069 [25] Romain Lopez, Jeffrey Regier, Michael B Cole, Michael I Jordan, and Nir Yosef. Deep  
1070 generative modeling for single-cell transcriptomics. *Nature methods*, 15(12):1053–1058,  
1071 2018.

- 1072 [26] Laleh Haghverdi, Aaron TL Lun, Michael D Morgan, and John C Marioni. Batch effects  
1073 in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors.  
1074 *Nature biotechnology*, 36(5):421–427, 2018.
- 1075 [27] Brian Hie, Bryan Bryson, and Bonnie Berger. Efficient integration of heterogeneous  
1076 single-cell transcriptomes using Scanorama. *Nature biotechnology*, 37(6):685–691, 2019.
- 1077 [28] Krzysztof Polański, Matthew D Young, Zhichao Miao, Kerstin B Meyer, Sarah A Teich-  
1078 mann, and Jong-Eun Park. BBKNN: fast batch alignment of single cell transcriptomes.  
1079 *Bioinformatics*, 36(3):964–965, 2020.
- 1080 [29] Ruben Chazarra-Gil, Stijn van Dongen, Vladimir Yu Kiselev, and Martin Hemberg.  
1081 Flexible comparison of batch correction methods for single-cell RNA-seq using batchbench.  
1082 *Nucleic acids research*, 49(7):e42–e42, 2021.
- 1083 [30] Malte D Luecken, Maren Büttner, Kridsakorn Chaichoompu, Anna Danese, Marta  
1084 Interlandi, Michaela F Müller, Daniel C Strobl, Luke Zappia, Martin Dugas, Maria Colomé-  
1085 Tatché, et al. Benchmarking atlas-level data integration in single-cell genomics. *Nature*  
1086 *methods*, 19(1):41–50, 2022.
- 1087 [31] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
1088 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in*  
1089 *Neural Information Processing Systems*, pages 2672–2680, 2014.
- 1090 [32] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-  
1091 image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE*  
1092 *International Conference on Computer Vision*, pages 2223–2232, 2017.
- 1093 [33] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation  
1094 networks. In *Advances in Neural Information Processing Systems*, pages 700–708, 2017.
- 1095 [34] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul  
1096 Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image  
1097 translation. In *Proceedings of the IEEE conference on Computer Vision and Pattern*  
1098 *Recognition*, pages 8789–8797, 2018.

- 1099 [35] Maren Büttner, Zhichao Miao, F Alexander Wolf, Sarah A Teichmann, and Fabian J  
1100 Theis. A test metric for assessing single-cell RNA-seq batch correction. *Nature methods*,  
1101 16(1):43–49, 2019.
- 1102 [36] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*,  
1103 2(1):193–218, 1985.
- 1104 [37] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,  
1105 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al.  
1106 Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*,  
1107 12:2825–2830, 2011.
- 1108 [38] Joshua D Welch, Velina Kozareva, Ashley Ferreira, Charles Vanderburg, Carly Martin,  
1109 and Evan Z Macosko. Single-cell multi-omic integration compares and contrasts features  
1110 of brain cell identity. *Cell*, 177(7):1873–1887, 2019.
- 1111 [39] Leland McInnes, John Healy, Nathaniel Saul, and Lukas Grossberger. UMAP: Uniform  
1112 manifold approximation and projection. *The Journal of Open Source Software*, 3(29):861,  
1113 2018.
- 1114 [40] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast  
1115 unfolding of communities in large networks. *Journal of statistical mechanics: theory and  
1116 experiment*, 2008(10):P10008, 2008.
- 1117 [41] Jeanette Baran-Gale, Tamir Chandra, and Kristina Kirschner. Experimental design for  
1118 single-cell RNA sequencing. *Briefings in functional genomics*, 17(4):233–239, 2018.
- 1119 [42] Joey Schyns, Fabrice Bureau, and Thomas Marichal. Lung interstitial macrophages: past,  
1120 present, and future. *Journal of immunology research*, 2018, 2018.
- 1121 [43] Camille Ezran, Shixuan Liu, Stephen Chang, Jingsi Ming, Olga Botvinnik, Lolita Penland,  
1122 Alexander Tarashansky, Antoine de Morree, Kyle J Travaglini, Kazuteru Hasegawa, et al.  
1123 Tabula microcebus: A transcriptomic cell atlas of mouse lemur, an emerging primate  
1124 model organism. *bioRxiv*, 2021.
- 1125 [44] Ting Xie, Yizhou Wang, Nan Deng, Guanling Huang, Forough Taghavifar, Yan Geng,  
1126 Ningshan Liu, Vrishika Kulur, Changfu Yao, Peter Chen, et al. Single-cell deconvolution

- 1127 of fibroblast heterogeneity in mouse pulmonary fibrosis. *Cell reports*, 22(13):3625–3640,  
1128 2018.
- 1129 [45] Alan Selewa, Ryan Dohn, Heather Eckart, Stephanie Lozano, Bingqing Xie, Eric Gauchat,  
1130 Reem Elorbany, Katherine Rhodes, Jonathan Burnett, Yoav Gilad, et al. Systematic  
1131 comparison of high-throughput single-cell and single-nucleus transcriptomes during car-  
1132 diomyocyte differentiation. *Scientific reports*, 10(1):1–13, 2020.
- 1133 [46] Anne-Marie Galow, Markus Wolfien, Paula Müller, Madeleine Bartsch, Ronald M Brunner,  
1134 Andreas Hoeflich, Olaf Wolkenhauer, Robert David, and Tom Goldammer. Integrative  
1135 cluster analysis of whole hearts reveals proliferative cardiomyocytes in adult mice. *Cells*,  
1136 9(5):1144, 2020.
- 1137 [47] Blue B Lake, Rizi Ai, Gwendolyn E Kaeser, Neeraj S Salathia, Yun C Yung, Rui Liu, Andre  
1138 Wildberg, Derek Gao, Ho-Lim Fung, Song Chen, et al. Neuronal subtypes and diversity  
1139 revealed by single-nucleus RNA sequencing of the human brain. *Science*, 352(6293):1586–  
1140 1590, 2016.
- 1141 [48] Ricard Argelaguet, Anna SE Cuomo, Oliver Stegle, and John C Marioni. Computational  
1142 principles and challenges in single-cell data integration. *Nature Biotechnology*, pages 1–14,  
1143 2021.
- 1144 [49] Alexander B Rosenberg, Charles M Roco, Richard A Muscat, Anna Kuchina, Paul Sample,  
1145 Zizhen Yao, Lucas T Graybuck, David J Peeler, Sumit Mukherjee, Wei Chen, et al. Single-  
1146 cell profiling of the developing mouse brain and spinal cord with split-pool barcoding.  
1147 *Science*, 360(6385):176–182, 2018.
- 1148 [50] 3k peripheral blood mononuclear cells (PBMCs) from a healthy donor from 10X Genomics.  
1149 <https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/pbmc3k>.
- 1150 [51] John F Fullard, Hao-Chih Lee, Georgios Voloudakis, Shengbao Suo, Behnam Javidfar,  
1151 Zhiping Shao, Cyril Peter, Wen Zhang, Shan Jiang, André Corvelo, et al. Single-nucleus  
1152 transcriptome analysis of human brain immune response in patients with severe covid-19.  
1153 *Genome medicine*, 13(1):1–13, 2021.
- 1154 [52] Matthew N Tran, Kristen R Maynard, Abby Spangler, Louise A Huuki, Kelsey D Mont-  
1155 gomery, Vijay Sadashivaiah, Madhavi Tippani, Brianna K Barry, Dana B Hancock,

- 1156 Stephanie C Hicks, et al. Single-nucleus transcriptome analysis reveals cell-type-specific  
1157 molecular signatures across reward circuitry in the human brain. *Neuron*, 109(19):3088–  
1158 3103, 2021.
- 1159 [53] Tim Stuart and Rahul Satija. Integrative single-cell analysis. *Nature reviews genetics*,  
1160 20(5):257–272, 2019.
- 1161 [54] Eleni P Mimitou, Caleb A Lareau, Kelvin Y Chen, Andre L Zorzetto-Fernandes, Yuhan  
1162 Hao, Yusuke Takeshima, Wendy Luo, Tse-Shun Huang, Bertrand Z Yeung, Efthymia  
1163 Papalexi, et al. Scalable, multimodal profiling of chromatin accessibility, gene expression  
1164 and protein levels in single cells. *Nature Biotechnology*, pages 1–13, 2021.
- 1165 [55] Yingxin Lin, Tung-Yu Wu, Sheng Wan, Jean YH Yang, Wing H Wong, and YX Wang.  
1166 scJoint integrates atlas-scale single-cell RNA-seq and ATAC-seq data with transfer learning.  
1167 *Nature Biotechnology*, pages 1–8, 2022.
- 1168 [56] Chunmei Cui, Yuan Zhou, and Qinghua Cui. Defining the functional divergence of  
1169 orthologous genes between human and mouse in the context of miRNA regulation. *Briefings  
1170 in Bioinformatics*, pages 1477–4054, 2021.
- 1171 [57] Christopher Daniel Green, Qianyi Ma, Gabriel L Manske, Adrienne Niederriter Shami,  
1172 Xianing Zheng, Simone Marini, Lindsay Moritz, Caleb Sultan, Stephen J Gurczynski,  
1173 Bethany B Moore, et al. A comprehensive roadmap of murine spermatogenesis defined by  
1174 single-cell RNA-seq. *Developmental cell*, 46(5):651–667, 2018.
- 1175 [58] Brian P Hermann, Keren Cheng, Anukriti Singh, Lorena Roa-De La Cruz, Kazadi N  
1176 Mutoji, I-Chung Chen, Heidi Gildersleeve, Jake D Lehle, Max Mayo, Birgit Westernströer,  
1177 et al. The mammalian spermatogenesis single-cell transcriptome, from spermatogonial  
1178 stem cells to spermatids. *Cell reports*, 25(6):1650–1667, 2018.
- 1179 [59] Christina Ernst, Nils Eling, Celia P Martinez-Jimenez, John C Marioni, and Duncan T  
1180 Odom. Staged developmental mapping and X chromosome transcriptional dynamics during  
1181 mouse spermatogenesis. *Nature communications*, 10(1):1–20, 2019.
- 1182 [60] Xianzhong Lau, Prabhakaran Munusamy, Mor Jack Ng, and Mahesh Sangrithi. Single-cell  
1183 RNA sequencing of the *Cynomolgus* macaque testis reveals conserved transcriptional  
1184 profiles during mammalian spermatogenesis. *Developmental Cell*, 54(4):548–566, 2020.

- 1185 [61] Ho-Su Sin and Satoshi H Namekawa. The great escape: Active genes on inactive sex  
1186 chromosomes and their evolutionary implications. *Epigenetics*, 8(9):887–892, 2013.
- 1187 [62] Jennifer F Hughes and David C Page. The biology and evolution of mammalian Y  
1188 chromosomes. *Annual review of genetics*, 49:507–527, 2015.
- 1189 [63] Satoshi H Namekawa and Jeannie T Lee. XY and ZW: is meiotic sex chromosome  
1190 inactivation the rule in evolution? *PLoS genetics*, 5(5):e1000493, 2009.
- 1191 [64] Jeffrey M Cloutier and James MA Turner. Meiotic sex chromosome inactivation. *Current*  
1192 *Biology*, 20(22):R962–R963, 2010.
- 1193 [65] Erica L Larson, Emily EK Kopania, and Jeffrey M Good. Spermatogenesis and the  
1194 evolution of mammalian sex chromosomes. *Trends in Genetics*, 34(9):722–732, 2018.
- 1195 [66] Jennifer F Hughes, Helen Skaletsky, Laura G Brown, Tatyana Pyntikova, Tina Graves,  
1196 Robert S Fulton, Shannon Dugan, Yan Ding, Christian J Buhay, Colin Kremitzki, et al.  
1197 Strict evolutionary conservation followed rapid gene loss on human and rhesus Y chromo-  
1198 somes. *Nature*, 483(7387):82–86, 2012.
- 1199 [67] Jennifer F Hughes, Helen Skaletsky, Tatyana Pyntikova, Tina A Graves, Saskia KM  
1200 Van Daalen, Patrick J Minx, Robert S Fulton, Sean D McGrath, Devin P Locke, Cynthia  
1201 Friedman, et al. Chimpanzee and human Y chromosomes are remarkably divergent in  
1202 structure and gene content. *Nature*, 463(7280):536–539, 2010.
- 1203 [68] Polly Campbell, Jeffrey M Good, and Michael W Nachman. Meiotic sex chromosome  
1204 inactivation is disrupted in sterile hybrid male house mice. *Genetics*, 193(3):819–828, 2013.
- 1205 [69] Ho-Su Sin, Yosuke Ichijima, Eitetsu Koh, Mikio Namiki, and Satoshi H Namekawa. Human  
1206 postmeiotic sex chromatin and its impact on sex chromosome evolution. *Genome research*,  
1207 22(5):827–836, 2012.
- 1208 [70] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International*  
1209 *Conference on Learning Representations*, 2014.
- 1210 [71] Chenling Xu, Romain Lopez, Edouard Mehlman, Jeffrey Regier, Michael I Jordan, and  
1211 Nir Yosef. Probabilistic harmonization and annotation of single-cell transcriptomics data  
1212 with deep generative models. *Molecular systems biology*, 17(1):e9620, 2021.



- 1213 [72] Mohammad Lotfollahi, F Alexander Wolf, and Fabian J Theis. scGen predicts single-cell  
1214 perturbation responses. *Nature methods*, 16(8):715–721, 2019.
- 1215 [73] Dongfang Wang, Siyu Hou, Lei Zhang, Xiliang Wang, Baolin Liu, and Zemin Zhang.  
1216 iMAP: integration of multiple single-cell datasets by adversarial paired transfer networks.  
1217 *Genome biology*, 22(1):1–24, 2021.
- 1218 [74] Bin Zou, Tongda Zhang, Ruilong Zhou, Xiaosen Jiang, Huanming Yang, Xin Jin, and  
1219 Yong Bai. deepMNN: Deep learning-based single-cell rna sequencing data batch correction  
1220 using mutual nearest neighbors. *Frontiers in Genetics*, page 1441, 2021.
- 1221 [75] Karren Dai Yang, Anastasiya Belyaeva, Saradha Venkatachalapathy, Karthik Damodaran,  
1222 Abigail Katcoff, Adityanarayanan Radhakrishnan, GV Shivashankar, and Caroline Uh-  
1223 ler. Multi-domain translation between single-cell imaging and sequencing data using  
1224 autoencoders. *Nature Communications*, 12(1):1–10, 2021.
- 1225 [76] Martin Arjovsky and Léon Bottou. Towards principled methods for training generative  
1226 adversarial networks. In *ICLR*, 2017.
- 1227 [77] Sean C Bendall, Kara L Davis, El-ad David Amir, Michelle D Tadmor, Erin F Simonds,  
1228 Tiffany J Chen, Daniel K Shenfeld, Garry P Nolan, and Dana Pe’er. Single-cell trajectory  
1229 detection uncovers progression and regulatory coordination in human B cell development.  
1230 *Cell*, 157(3):714–725, 2014.
- 1231 [78] F Alexander Wolf, Philipp Angerer, and Fabian J Theis. SCANPY: large-scale single-cell  
1232 gene expression data analysis. *Genome biology*, 19(1):1–5, 2018.
- 1233 [79] David Powers. Evaluation: From precision, recall and F-measure to ROC, informedness,  
1234 markedness & correlation. *Journal of Machine Learning Technologies*, 2(1):37–63, 2011.
- 1235 [80] Itay Tirosh, Benjamin Izar, Sanjay M Prakadan, Marc H Wadsworth, Daniel Treacy, John J  
1236 Trombetta, Asaf Rotem, Christopher Rodman, Christine Lian, George Murphy, et al.  
1237 Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq.  
1238 *Science*, 352(6282):189–196, 2016.
- 1239 [81] Jiarui Ding, Xian Adiconis, Sean K Simmons, Monika S Kowalczyk, Cynthia C Hession,  
1240 Nemanja D Marjanovic, Travis K Hughes, Marc H Wadsworth, Tyler Burks, Lan T Nguyen,

- 1241 et al. Systematic comparison of single-cell and single-nucleus RNA-sequencing methods.  
1242 *Nature biotechnology*, 38(6):737–746, 2020.
- 1243 [82] Franziska Paul, Ya’ara Arkin, Amir Giladi, Diego Adhemar Jaitin, Ephraim Kenigsberg,  
1244 Hadas Keren-Shaul, Deborah Winter, David Lara-Astiaso, Meital Gury, Assaf Weiner,  
1245 et al. Transcriptional heterogeneity and lineage commitment in myeloid progenitors. *Cell*,  
1246 163(7):1663–1677, 2015.
- 1247 [83] Sonia Nestorowa, Fiona K Hamey, Blanca Pijuan Sala, Evangelia Diamanti, Mairi Shepherd,  
1248 Elisa Laurenti, Nicola K Wilson, David G Kent, and Berthold Göttgens. A single-cell  
1249 resolution map of mouse hematopoietic stem and progenitor cell differentiation. *Blood*,  
1250 *The Journal of the American Society of Hematology*, 128(8):e20–e31, 2016.
- 1251 [84] Geoffrey Schiebinger, Jian Shu, Marcin Tabaka, Brian Cleary, Vidya Subramanian, Aryeh  
1252 Solomon, Joshua Gould, Siyan Liu, Stacie Lin, Peter Berube, et al. Optimal-transport anal-  
1253 ysis of single-cell gene expression identifies developmental trajectories in reprogramming.  
1254 *Cell*, 176(4):928–943, 2019.