

Adversarial domain translation networks enable fast and accurate large-scale atlas-level single-cell data integration

Jia Zhao^{1*}, Gefei Wang^{1*}, Jingsi Ming², Zhixiang Lin³, Yang Wang¹
Tabula Microcebus Consortium, Angela Ruohao Wu^{4,5†}, Can Yang^{1†}

¹Department of Mathematics, The Hong Kong University of Science and Technology, Hong Kong SAR, China

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

³Department of Statistics, The Chinese University of Hong Kong, Hong Kong SAR, China

⁴Division of Life Science, The Hong Kong University of Science and Technology, Hong Kong SAR, China

⁵Department 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-sequencing (scRNA-seq) datasets
2 from various sources presents remarkable opportunities for broad and deep biological investiga-
3 tions through integrative analyses. However, harmonizing such datasets requires integration
4 approaches to be not only computationally scalable, but also capable of preserving a wide range
5 of fine-grained cell populations. We created Portal, a unified framework of adversarial domain
6 translation to learn harmonized representations of datasets. With innovation in model and
7 algorithm designs, Portal achieves superior performance in preserving biological variation during
8 integration, while having significantly reduced running time and memory compared to existing
9 approaches, achieving integration of millions of cells in minutes with low memory consumption.
10 We demonstrate the efficiency and accuracy of Portal using diverse datasets ranging from

*These authors contributed to this work equally.

†Correspondence: angelawu@ust.hk, macyang@ust.hk.

11 mouse brain atlas projects, the Tabula Muris project, and the Tabula Microcebus project.
12 Portal has broad applicability and in addition to integrating multiple scRNA-seq datasets, it
13 can also integrate scRNA-seq with single-nucleus RNA-sequencing (snRNA-seq) data. Finally,
14 we demonstrate the utility of Portal by applying it to the integration of cross-species datasets
15 with limited shared-information between them, and are able to elucidate biological insights
16 into the similarities and divergences in the spermatogenesis process between mouse, macaque,
17 and human.

18 **Introduction**

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

31 Although integration methods for single-cell transcriptomics analysis have evolved along
32 with single-cell sequencing technologies, the rapid accumulation of new and diverse single-cell
33 datasets has introduced three major challenges to the integration task. First, as the sample size
34 of each single-cell dataset grows dramatically, numerous extensive datasets with hundreds of
35 thousands or even millions of cells have been produced [8, 9, 12]. The emergence of large-scale
36 datasets requires integration methods to be fast, memory-efficient, and scalable to millions
37 of cells. Second, technology now allows effective, comprehensive characterization of complex
38 organs, containing rare subpopulations of cells that can now be captured, albeit in small
39 numbers, thanks to the scale of profiling that is now possible [7, 13]. Investigation into high-
40 level heterogeneity among cell populations is essential for understanding the mechanism of

41 complex biological systems. Hence, the ideal integration method needs to carefully preserve fine-
42 grained cell populations from each atlas-level dataset. Third, the biological origins of datasets
43 has expanded in diversity, with data now spanning across not only different technological
44 platforms and data types, different individual donors, but even across different species, which
45 can be especially interesting for evolutionary studies [14, 15, 16]. Integrative analysis of such
46 diverse datasets would allow researchers to unify resources to address a wider range of biological
47 questions. Recent single-cell atlasing efforts are a primary example of these challenges – various
48 human tissue atlases [12, 17], mouse multi-tissue atlases [7, 18], and non-human primate atlases
49 [19, 20] have been generated, culminating in data from millions of single cells and single
50 nuclei. Both within and across atlas comparisons are of interest. To perform integrative and
51 comparative analyses based on such diverse data sources, there is an urgent need for methods
52 that can flexibly account for heterogeneous dataset-specific effects, while maintaining a high
53 level of integration accuracy.

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

71 To simultaneously address the above three challenges, we created Portal, a machine learning-

72 based algorithm for aligning atlas-level single-cell datasets with high efficiency, flexibility, and
73 accuracy. Viewing datasets from different studies as distinct domains with domain-specific
74 effects (including technical variation and other sources of unwanted variation), Portal achieves
75 extraordinary data alignment performance through a unified framework of domain translation
76 networks that incorporates an adversarial learning mechanism [27]. To find the correspondence
77 between two domains, our domain translation network utilizes an encoder to embed cells from
78 one domain into a latent space where domain-specific effects are removed, and then uses a
79 generator to map latent codes to another domain. The generator simulates the generation
80 process of domain-specific effects. In each domain, a discriminator is trained to identify where
81 poor alignment between the distributions of original cells and transferred cells occurs. The
82 feedback signal from the discriminator is used to strengthen the domain translation network
83 for better alignment. The nonlinearity of encoders and generators in the adversarial domain
84 translation framework enables Portal to account for complex domain-specific effects. In contrast
85 to existing domain translation methods [28, 29, 30], Portal has the following unique features.
86 First, Portal has a uniquely designed discriminator which can adaptively distinguish domain-
87 shared cell types and domain-unique cell types. Therefore, Portal will not force the alignment
88 of domain-unique cell types, avoiding the risk of overcorrection. Second, without using any cell
89 type label information, three regularizers of Portal can guide domain translation networks to
90 find correct correspondence between domains, account for domain-specific effects, and retain
91 biological variation in the latent space. Third, through a tailored design of light-weight neural
92 networks and mini-batch optimization accelerated by graphics processing units (GPUs), Portal
93 can scale up to datasets containing millions of cells in minutes with nearly constant memory
94 usage. With the above innovations in model and algorithm designs, Portal enables fast and
95 accurate integration of atlas-level datasets across samples, technological platforms, data types,
96 and species.

97 Through integration of heterogeneous collections of atlas-level single-cell RNA sequencing
98 (scRNA-seq) data, Portal shows its superiority over state-of-the-art alignment algorithms
99 in terms of both computational efficiency and accuracy. We then show that Portal can
100 accurately align cells from complex tissues profiled by scRNA-seq and single-nucleus RNA
101 sequencing (snRNA-seq), and also perform cross-species alignment of the gradient of cells in
102 the spermatogenesis process, demonstrating Portal’s versatility and power for a broad range of

103 applications. Comprehensive analyses of real, expert annotated data confirm that integrated
104 cell embeddings provided by Portal can be reliably used for identification of rare cell populations
105 via clustering or label transfer, studies of differentiation trajectories, and transfer learning
106 across data types and across species. Portal is now publicly available as a Python package
107 (<https://github.com/YangLabHKUST/Portal>), serving as an efficient, reliable and flexible
108 tool for integrative analyses.

109 Results

110 Method Overview: Portal learns a harmonized representation of dif- 111 ferent datasets with adversarial domain translation.

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

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

133 However, the well mixing of original cells and transferred cells in each domain does not
 134 imply extraordinary data alignment across domains. First, a domain-unique cell population
 135 should not be mixed with cells from another domain. Second, cell types A and B in domain
 136 \mathcal{X} could be incorrectly aligned with cell types B and A in domain \mathcal{Y} , respectively, although
 137 the distributions of original cells and transferred cells are well mixed. To address these issues,
 138 Portal has the following unique features, which distinguishes it from existing adversarial domain
 139 translation frameworks [28, 29]. On one hand, we deploy the tailored design of discriminators
 140 $D_1(\cdot)$ and $D_2(\cdot)$ such that they can distinguish domain-unique cell types from cell types shared
 141 across different domains. The domain-unique cell types will be treated as outliers and left
 142 in the discriminator’s inactive region (Fig. 1c). In such a way, these cell types will not be
 143 enforced for alignment, avoiding the risk of overcorrection. On the other hand, we design three
 144 regularizers to find correct correspondence across domains and avoid incorrect alignment when
 145 the distributions are well mixed.

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

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

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

160 where q is the dimensionality of \mathcal{Z} , imposes the consistency constraint in the latent space;
161 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
162 the cross-domain correspondence by preserving the cosine similarity between a sample and
163 its transferred version; t_{AE} , t_{LA} and t_{cos} are their corresponding constraint parameters. More
164 detailed explanation can be found in the Method section.

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

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

174 The rapid accumulation of large-scale single-cell datasets requires integration algorithms
175 to efficiently handle datasets containing millions of cells without loss of accuracy. For a
176 comprehensive comparison, we first benchmarked Portal against multiple methods, including
177 Harmony [21], Seurat v3 [22], online iNMF [23], fastMNN [24], Scanorama [25] and BBKNN
178 [26], in terms of integration performance. Using massive scRNA-seq datasets from diverse tissue
179 types with curated cell cluster annotations, including mouse spleen, marrow, and bladder [7], we
180 quantitatively evaluated the integration performance of each method. We evaluated alignment
181 performance, which can sometimes be interpreted as batch effects removal performance, using
182 k-nearest neighbor batch-effect test (kBET) [32]; the higher the kBET score, the higher the
183 degree of mixing across datasets. We also assessed cluster identity preservation performance
184 using the adjust rand index (ARI) and average silhouette width (ASW) metrics. Using the
185 authors' cell type annotations as ground truth, higher ARI and ASW scores denote that correct
186 cell type identities are preserved after integration, while lower scores indicate inappropriate
187 merging of cell types during integration. Based on these metrics, we found that in general,
188 fastMNN, Scanorama, and BBKNN have less satisfactory integration performance compared to
189 the other four methods (Figs. 2, S3 and S4): as indicated by the relatively lower kBET scores

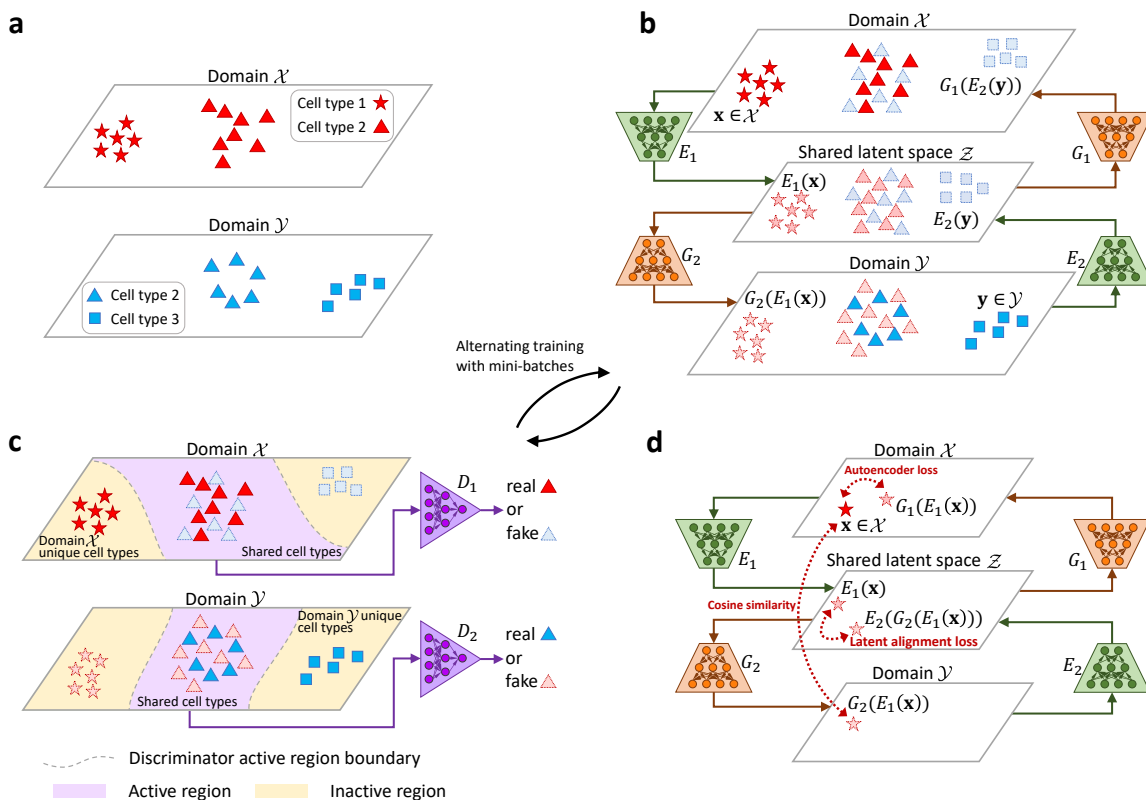


Figure 1: Overview of Portal. **a.** Portal regards different single-cell datasets as different domains. Joint analyses of these datasets are confounded by domain-specific effects, representing the unwanted technical variation. **b.** Portal employs encoders $E_1(\cdot), E_2(\cdot)$ to embed the biological 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.

190 of these three methods, we found that observable batch effects still exist in the integration
 191 results they produced (Fig. 2a); in addition, their ARI and ASW metrics are also lower (Fig.
 192 2b).

193 Among those methods with high user popularity, Harmony, Seurat, and online iNMF

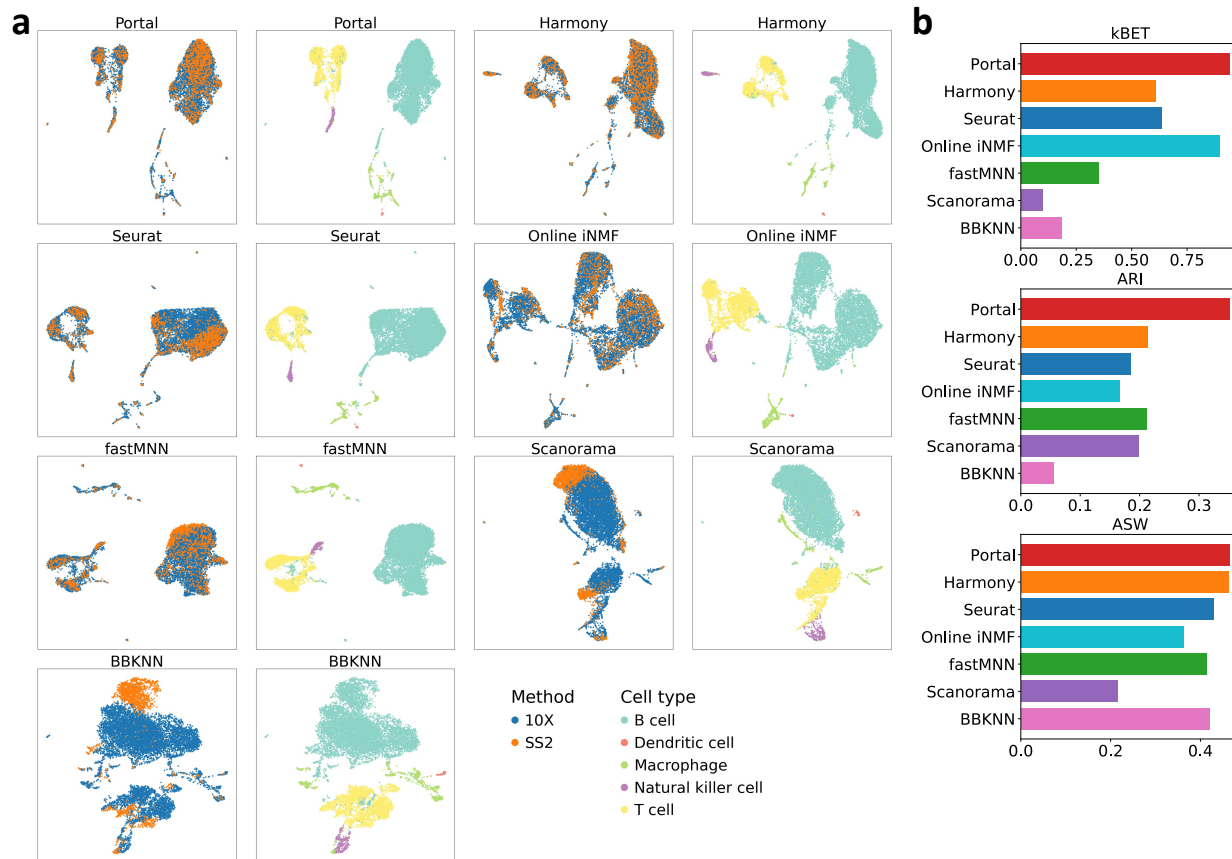


Figure 2: **Comparison of integration methods based on mouse spleen data.** We integrated mouse spleen scRNA-seq datasets profiled by 10X Genomics (10X) and SMART-seq2 (SS2). **a.** UMAP [31] plots colored by profiling methods and cell types. **b.** Alignment (kBET) and cluster preservation performance (ARI and ASW) of compared methods evaluated on the mouse spleen data.

194 also showed the best alignment performance results. To offer precise and robust integration
 195 performance, Seurat [22] utilizes the detection of mutual nearest neighbors (MNN) to build
 196 correspondence between datasets in the shared embedding space obtained by applying canonical
 197 correlation analysis (CCA). Harmony [21] learns a simple linear correction for dataset-specific
 198 effects by running an iterative soft clustering algorithm, enabling fast computation on large
 199 datasets. Online iNMF [23] is a recently developed approach based on widely used integration
 200 method LIGER [33]. It extends LIGER's non-negative matrix factorization to an iterative
 201 and incremental version to improve its scalability, while it has nearly the same performance
 202 as LIGER. For the remainder of this study, we focus our discussion on comparisons between
 203 Portal and these three high-performing and popular methods (Fig. 3) in the main text. The
 204 comparisons with other methods are provided in Supplementary Information (Fig. S5).

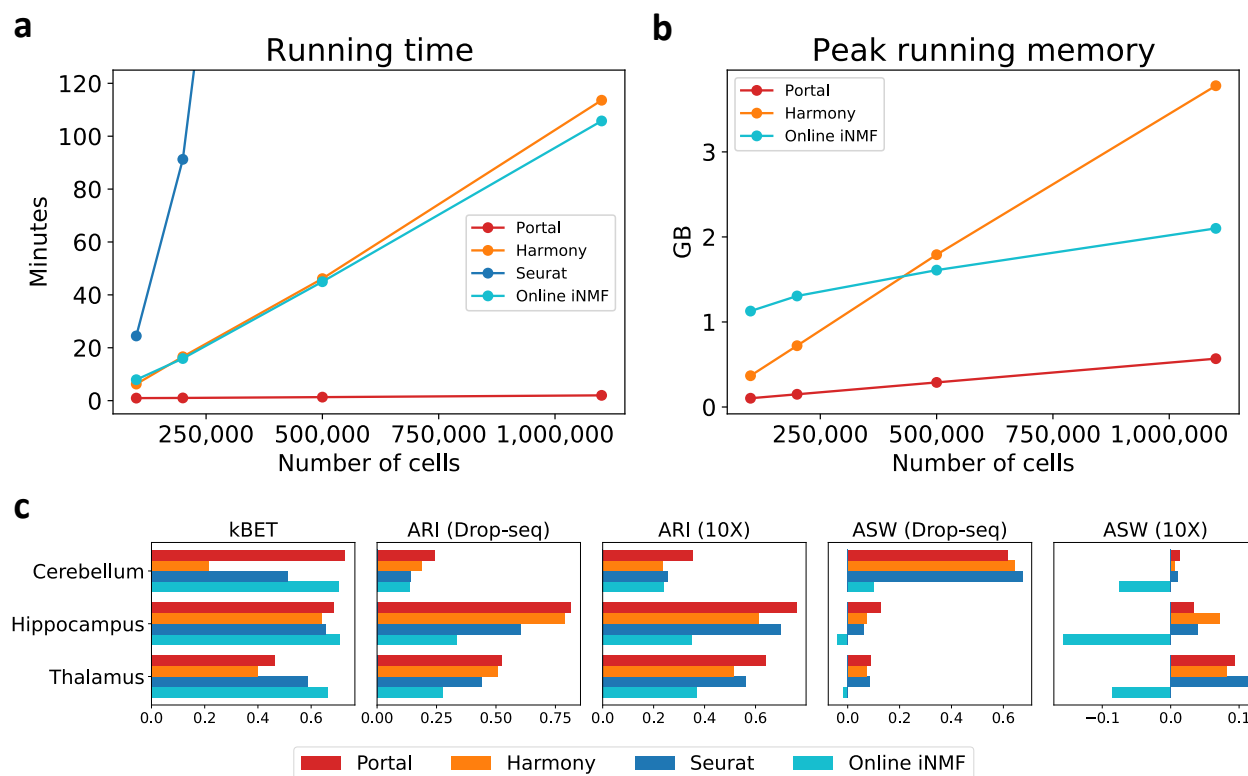


Figure 3: **Benchmark of Portal, Harmony, Seurat and online iNMF.** **a, b.** Running time and peak running memory required by benchmarked methods. The datasets were sampled from two mouse brain atlas datasets ($n = 100,000, 250,000, 500,000,$ and $1,100,167$). Seurat requires 24.52 GB on the dataset with 100,000 cells, which is not comparable to the other three benchmarked methods in terms of peak running memory usage. **c.** Alignment (kBET) and cluster preservation performance (ARI and ASW) evaluated using three shared tissues from two mouse brain atlases (profiled by Drop-seq and 10X), including cerebellum, hippocampus, and thalamus. Cluster preservation performance was assessed based on fine-grained annotations provided by the original publications [8, 9]. A full comparison among all methods is provided in Supplementary Information (Fig. S5).

205 Next, we evaluated the speed, memory usage, alignment quality, and integration accuracy
 206 using a more challenging integration task. We used two mouse brain atlases [8, 9] as bench-
 207 marking datasets for a more in-depth comparison of Portal and three other methods. One atlas
 208 contains Drop-seq data of 939,489 cells, and another one contains 10X Genomics (10X) data
 209 of 160,678 cells. These two mouse brain atlases have data from three shared brain regions:
 210 cerebellum, hippocampus, and thalamus. There are many small clusters of neuron subtypes
 211 in these datasets, where gene expressions between subclusters could have a relatively small
 212 difference. Thus, these datasets are more challenging to integrate compared to data with clear
 213 clustering patterns.

214 First, Portal has superior integration accuracy even when handling datasets which contain
215 many subclusters with small difference. The ARI and ASW show that Portal outperforms
216 other state-of-the-art methods in cluster identity preservation. In particular, for all three brain
217 regions tested, Portal has the highest ARI score among all the benchmarked methods (Fig. 3c).

218 Second, Portal also outperforms the other three methods on scalability, in terms of time
219 and memory consumption. For this benchmark test, we obtained datasets from the original
220 full-sized datasets by combining the two atlases and subsampling proportionally from each
221 atlas, with each dataset having increasing sample size ranging from 100,000 to 1,100,167 (full
222 dataset). The running time and peak running memory of all methods were recorded using
223 these datasets on the same GPU server. The results show that Portal's running time and peak
224 running memory remained almost constant even when the sample size increased dramatically
225 (Fig. 3a, b). Compared to the other three methods, the running time required by Portal was
226 also substantially less (Fig. 3a). On the dataset containing 500,000 cells, Portal's running
227 time was 80 seconds; when number of cells grew to 1,100,167, Portal's running time only
228 increased to 120 seconds. In comparison, Harmony and online iNMF both needed more than
229 40 minutes to integrate 500,000 cells and more than 100 minutes to complete the integration
230 of the full dataset. The running time of Seurat increased most rapidly among the compared
231 methods. It took as much as 511 minutes (over 8.5 hours) to integrate the 500,000-cell dataset.
232 The computational efficiency of Portal is owing to two important factors in its design: 1) its
233 algorithm takes advantage of GPU-accelerated stochastic optimization, such that Portal reads
234 data in mini-batches from the disk rather than having to load the entire dataset at once, which
235 enables fast integration of large single-cell datasets using small amounts of memory; and 2)
236 lightweight neural networks are adopted in Portal to further improve computational efficiency.
237 As such, Portal is also the most memory-efficient approach among the benchmarked methods
238 (Fig. 3b). Peak running memory required by Portal ranged from 0.10 GB on 100,000-cell
239 dataset to 0.29 GB on the full million-cell dataset. Notably, Portal's lightweight networks and
240 mini-batch stochastic optimization algorithm enable us to control GPU peak running memory
241 usage at a constant level of 0.06 GB. Among compared methods, online iNMF used less memory
242 than Harmony and Seurat when the sample size became larger than 500,000, because it is also
243 trained in mini-batches. However, its peak running memory was 2.10 GB on the million-cell
244 dataset, which is 6 times more than Portal's. Seurat required remarkably more memory usage

245 than the other three methods. For clarity of visualization, we did not display the peak running
246 memory required by Seurat as it ranged from 24.52 GB on the 100,000-cell dataset to 276.41
247 GB on the 500,000-cell dataset.

248 Finally, and importantly, Portal’s high performance in speed and memory consumption does
249 not compromise its ability to align cell type clusters. The kBET shows that Portal’s alignment
250 ability is comparable to, if not better than, the other benchmarked methods, indicating that
251 Portal is capable to effectively remove domain-specific effects.

252 **Portal preserves subcluster and small cluster identities in complex** 253 **tissues thereby facilitating identification of rare subpopulations.**

254 When integrating complex tissues, one problem that can arise is the inadvertent loss of small
255 cell populations and subpopulations. Due to more nuanced differences between clusters, or due
256 to the imbalance in cell numbers for very small cell populations, these “fine-grained” groups of
257 cells may become inappropriately combined with other groups after integration. In the brain,
258 for example, there are many subpopulations of neurons which are distinguished from each other
259 using a few key gene markers while still all bearing the neuron signature; furthermore, some of
260 these neuronal subtypes could be rare compared to other subtypes. To demonstrate that Portal
261 can preserve the nuanced information of such small cell populations and subpopulations, we
262 performed further analysis on the mouse hippocampus tissue integration results. Both mouse
263 brain atlas datasets contain extensive data for this brain region (Fig. 4), and both studies
264 identified a wide range of transcriptionally distinct cell subpopulations, including a variety
265 of neuron subtypes, whose nuanced transcriptional differences should ideally be preserved by
266 integration methods.

267 After applying Portal and the other three benchmarked methods to integrate the data, we
268 used the integrated cell representations to perform clustering. Using the Louvain method [34]
269 with default resolution, we obtained 29 (Portal), 29 (Harmony), 25 (Seurat) and 30 (online
270 iNMF) clusters, respectively (Fig. S6). Particularly, we focused on one region where the cell
271 proportions between two datasets were highly unbalanced, as marked in Fig. 4a. Only a
272 few of cells in this region are from the 10X dataset, making it challenging to build alignment
273 between datasets while preserving subpopulations from the Drop-seq dataset. In the original
274 publication [8], cells from the Drop-seq dataset within the marked region were all annotated

278 cells into three coherent groups in the integrated embedding space. Specifically, clusters 4, 13,
279 26 identified by the Louvain method recovered the *Slc17a6*+ neuron; *Cbln1*+/*Grp88*+ medial
280 entorhinal cortex neuron; and the *Cbln4*+ neuron subpopulations, respectively (Fig. 4b). Each
281 cluster was confirmed by the high expression level of the annotated marker genes (Fig. S7a).
282 Notably, these three groups only accounted for 4.79%, 1.76% and 0.32% of the total sample
283 size, respectively, demonstrating Portal's ability to preserve identities of rare subpopulations.
284 However, the differences among these three subpopulations were not well preserved by the other
285 three methods, making it difficult to detect them each distinctly using the Louvain clustering
286 method (Fig. 4a, b). As shown in Fig. S7c, we also identified eight protein coding genes
287 that were the most significantly differentially expressed among clusters, indicating the different
288 functions of each of the three neuron subtypes. Cluster 4 showed high expression levels of
289 *Camk2n1*, *Map1b*, *Nrgn*, *Syt1*, and no detectable expression of *Camk2d*, *Igfbp5*, *Nr4a2* and
290 *Ntng1*. A different pattern was observed in cluster 13: High expression of *Camk2d*, *Camk2n1*,
291 *Map1b* and *Syt1*, and no detectable expression of the other four genes. Cluster 26, meanwhile,
292 showed moderate levels of expression of all eight genes. In the marked region, cells from the
293 10X dataset were mainly concentrated in cluster 4. The alignment by Portal was confirmed
294 by the consistent gene expression levels seen in cluster 4 between the two datasets (Fig. S7b).
295 Besides the eight differentially expressed genes, we also examined a larger set of genes, and
296 computed the cross correlation of these genes pairwise between cells from all three groups.
297 This analysis showed that cells within each cluster had higher similarity in gene expression
298 than cells from other clusters, further showing the biological difference between these three
299 clusters that should not be mixed after integration. The above results highlight Portal's power
300 to preserve rare cell types (Fig. S7d).

301 The integrative analysis on the hippocampus tissue demonstrates Portal's ability to maintain
302 nuanced transcriptional differences for small subpopulations. This means that Portal can also
303 be used to "call out" rare subpopulations in one dataset based on integration with another
304 dataset via label transfer. To illustrate this feature, we take 10X and SMART-seq2 (SS2) data
305 generated for a mouse lung scRNA-seq atlas [7] as an example: the typically larger sample size
306 of the 10X dataset facilitates powerful clustering analyses for identification of cell types; while
307 the greater sequencing depth and sensitivity of SS2 enables deeper investigation into cell biology
308 [35]. To leverage the different strengths of the two technologies, we used Portal to perform

309 integrated analysis on 1,676 SS2 cells and 5,404 10X cells (Fig. S8a). Specifically, we defined
310 the 10X dataset annotations from the original publication [7] as reference labels (Fig. S8b),
311 then made use of the Portal’s integration results to identify cell types for the SS2 dataset based
312 on these reference labels. After integration, for each SS2 cell, label transfer was performed
313 by detecting its nearest neighbors among 10X cells. From this analysis, we identified four
314 subpopulations of myeloid cells for the SS2 dataset, namely alveolar macrophages, dendritic
315 cell and interstitial macrophages, classical monocytes, and non-classical monocytes (Fig. S8d).
316 Transferred labels of these four subpopulations were validated by known marker gene expression
317 levels [36]. For example, compared to classical monocytes, non-classical monocytes showed
318 lower expression of *Ccr2* and higher expressions of *Trem14* (Fig. S9). Consistent with the gene
319 expression pattern of alveolar macrophages in the 10X dataset, alveolar macrophages annotated
320 by Portal in the SS2 dataset had high expression levels of marker genes *Mrc1* and *Siglec5*.
321 Notably, in the SS2 dataset, the alveolar macrophage subpopulation only accounted for 0.78%
322 of total sample size, and could not be distinguished from the other SS2-profiled macrophages in
323 the original publication [7]. Based on the original labels, alveolar macrophages were unidentified
324 as they were labeled in a more general group named “dendritic cell, alveolar macrophage,
325 and interstitial macrophage” (Fig. S8c). Making good use of the larger 10X dataset, Portal
326 successfully identified extremely rare subpopulations within the SS2 dataset. We then used the
327 mouse lemur bladder scRNA-seq datasets from Tabula Microcebus Consortium [37] as another
328 example to demonstrate Portal’s ability for discovering rare subpopulations via label transfer.
329 In this example, mouse lemur bladder tissue was also profiled by both SS2 and 10X. When we
330 integrated these datasets and transferred labels from the 10X dataset to the SS2 dataset using
331 Portal, we were able to distinguish a very small myofibroblast subpopulation of just 11 cells in
332 the SS2 dataset from the rest of the fibroblasts (Fig. S10a). We verified their myofibroblast
333 identity based on their high expressions of known marker genes *ACTA2*, *MYH11*, *TAGLN* [38]
334 (Fig. S10b).

335 **Integration of comprehensive whole-organism cell atlases.**

336 So far, Portal has shown impressive performance in aligning tissue-level atlases where nuanced
337 transcriptional differences among subpopulations can be maintained after integration. We
338 next assess Portal’s capabilities under another challenging scenario: integrating two atlases

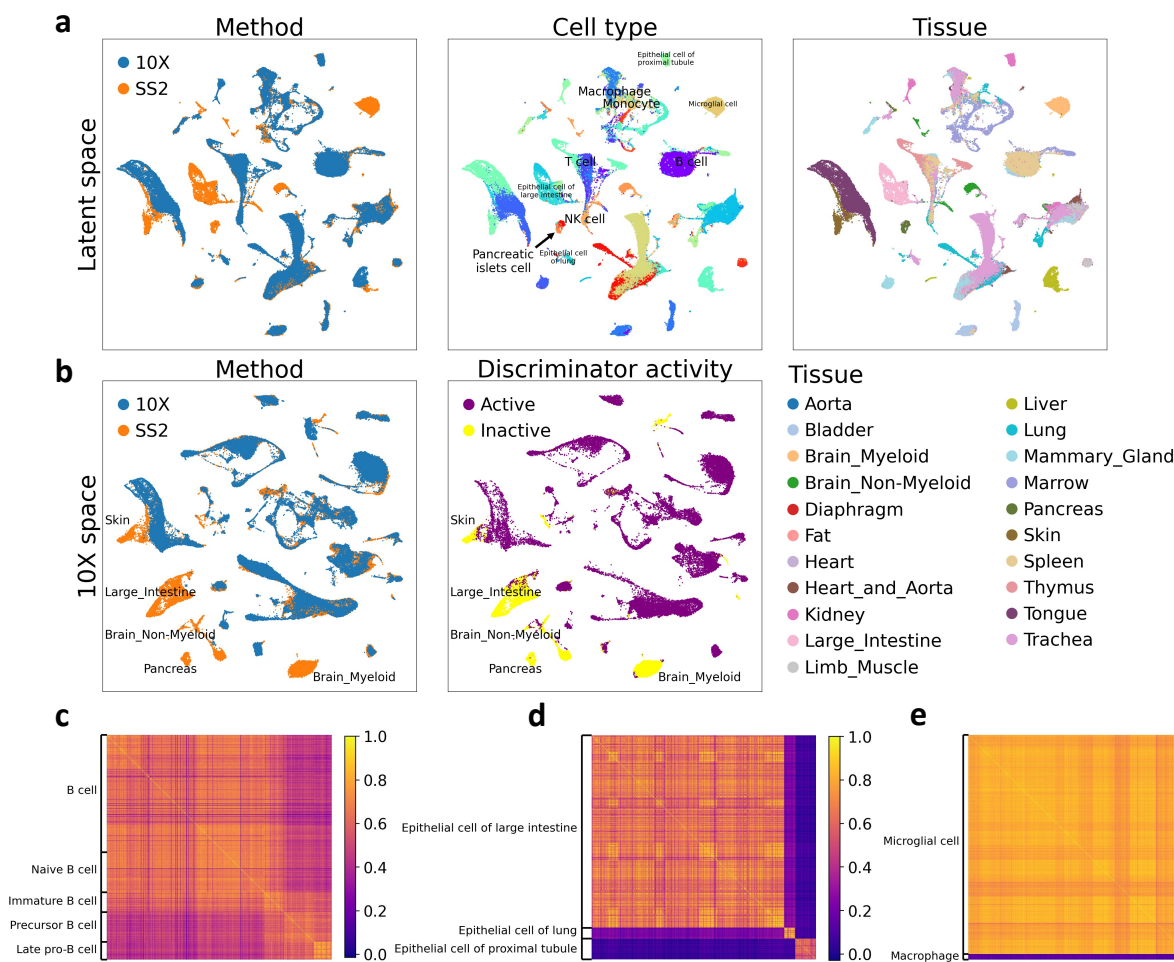


Figure 5: Construction of mouse cell atlas across entire organism by integrating atlas datasets from the Tablula Muris project. We applied Portal to integrate the datasets obtained by 10X and SS2. There were cells from unique tissues presented in the SS2 dataset. **a.** UMAP plots of Portal’s integration 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 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.

339 across an entire organism, where one of the atlases includes many more organs and tissue
 340 types than the other. This is known to be problematic for some integration algorithms due to
 341 having “missing cell types” in one of the datasets [24]. In contrast to these approaches, Portal
 342 uses discriminators with tailored design in the adversarial domain translation framework to
 343 distinguish domain-specific cell types from cell types shared across domains automatically, and

344 is thus robust to non-overlapped tissue samples.

345 To build a foundation for extensive study of cell populations across the whole organism,
346 Tabula Muris Consortium [7] profiled cells from 20 tissues using a combination of SS2 (44,779
347 cells) and 10X (54,865 cells) (Fig. 5). Notably, seven of these 20 tissues were only profiled by
348 SS2 but not 10X: brain (myeloid and non-myeloid), diaphragm, fat, large intestine, pancreas
349 and skin. We used Portal to build a comprehensive integrated mouse atlas that merges all
350 the cells, and we found Portal to show extraordinary accuracy in aligning cells of the same
351 cell type from the two datasets profiled by different platforms, not only in the shared latent
352 space but also in both domains (Figs. 5a, b and S11). After Portal integration, tissue-specific
353 cell types of SS2-only tissues, such as microglial cells in brain (myeloid), cell types in large
354 intestine, and pancreatic islets cells, were all successfully and correctly remained separated
355 from other cell types. The other three benchmarked methods, however, failed to retain many
356 tissue-specific cell types unmixed with other cell types. For instance, they mixed microglial
357 cells together with other macrophage cells, even though the data from these two cell types were
358 clearly transcriptionally different (Figs. 5e, S11).

359 Using this construction of a mouse cell atlas across organs, we also confirmed that the
360 designed boundaries for discriminator active region in Portal (Fig. 1c) indeed helped to
361 maintain the biological variation. By looking into the domain of 10X data (10X space), the
362 discriminator in the 10X domain was found inactive for tissue-specific cell types that were only
363 in the SS2 dataset (Fig. 5b). For these cells, Portal ensured that their identities were preserved
364 by making the adversarial learning objective inactive on them automatically. Portal's ability
365 to conserve information of cell populations indicates its reliability for integrating atlas-level
366 single-cell datasets across entire organisms.

367 Besides the alignment between datasets, Portal's integration result could characterize the
368 similarities and differences among cell types. For example, immune cells such as B cells, T cells,
369 natural killer cells (NK cells), monocytes and macrophages were profiled by both platforms
370 and contained in multiple tissues including brain (myeloid), diaphragm, fat, kidney, limb
371 muscle, liver, lung, mammary gland, marrow, spleen, and thymus. Portal correctly kept the
372 subpopulations belonging to the same type of immune cells close to each other, revealing the
373 resemblance of immune cells across different tissues. For instance, the transcriptional correlation
374 of all types of B cells, containing B cells, naive B cells, immature B cells, precursor B cells, and

375 late pro-B cells confirmed such similarity (Fig. 5c). In addition, the epithelial cells of different
376 tissues were identified by Portal as disjoint clusters, which was consistent with the biological
377 distinction among these cell types (Fig. 5d).

378 **Portal successfully and efficiently aligns single-cell RNA-seq data** 379 **and single-nucleus RNA-seq data.**

380 For frozen samples such as biobanked tissues, and for tissue types that have unique morphology
381 or phenotypes, such as brain, fat, or bone, it can be challenging or sometimes even impossible
382 to extract intact cells for scRNA-seq profiling [39, 40]. To bypass this issue, single-nucleus RNA
383 sequencing (snRNA-seq) has been developed. Although nuclear transcriptomes are shown to be
384 representative of the whole cell [41], distinctions between the whole cell and nucleus in terms
385 of the transcript type and composition make scRNA-seq data and snRNA-seq data intrinsically
386 different [39]. Aligning these two types of data is desirable, as the combined dataset enables
387 joint analysis that can take advantages of both techniques, and help to improve statistical
388 power for the analysis. Especially for comparing multiple complex tissues, with some cell types
389 being shared and others being non-overlapping, researchers could benefit from such integrated
390 joint analysis – one example being the integration of brain snRNA-seq data with scRNA-seq
391 data of blood to examine similarities and differences between immune cells in each tissue
392 milieu. However, due to the inherent difference in these two data types, aligning scRNA-seq
393 and snRNA-seq data is not the same as batch effects correction. Compared to batch effects
394 among scRNA-seq datasets, technical noise and unwanted variation arising from different data
395 types are often more complex and have higher strength [39, 42]. Thus, using standard batch
396 effects correction to integrate across data types may result in loss of alignment accuracy or
397 important biological signals.

398 We evaluated Portal’s ability to integrate snRNA-seq data and scRNA-seq data using three
399 mouse brain atlas datasets, including one snRNA-seq dataset profiled by SPLiT-seq [43], and
400 two scRNA-seq datasets profiled by Drop-seq and 10X [8, 9]. In this task, we applied integration
401 methods to harmonize these three atlases across all brain regions. To test the accuracy of
402 integration results, we only used cells that had annotations provided by the authors in each
403 atlas project. After selecting cells with cell type annotations, 319,359 cells in the Drop-seq
404 dataset, 160,678 cells in the 10X dataset, and 74,159 nuclei in the SPLiT-seq remained for

405 integration.

406 Prior to any integration, the raw datasets were clustered by the experimental method rather
407 than the cell type (Fig. S12a), and shared cell types between the three datasets did not align
408 well, indicating the initial discrepancy between the three large datasets. After integration,
409 UMAP visualizations showed that the different alignment methods gave varying results. Portal
410 (Fig. S12b) and Seurat (Fig. S12d) achieved the best alignment of data across different
411 methods, showing good mixing of cells annotated with the same cell type label, while also
412 preserving subcluster data structure in the integrated results. In particular, the alignment of
413 scRNA-seq (10X, Drop-seq) and snRNA-seq (SPLiT-seq) datasets was comparably good as that
414 of the two scRNA-seq datasets, indicating successful alignment between the two data types
415 without loss of biologically important variations between clusters. Online iNMF (Fig. S12e),
416 although it successfully clustered and aligned the same cell types together, within each cluster
417 the streaky pattern suggested potential numerical artefacts in the integrated data. Furthermore,
418 online iNMF alignment resulted in loss of biological variation, which was most easily observable
419 in the coalescence of the previously distinct neuron subpopulations (Fig. S12a) into one large
420 amorphous cluster (Fig. S12e). Harmony, however, showed poor mixing of the snRNA-seq
421 data in some of the cell types, such as the astrocytes, where the scRNA-seq datasets were
422 well-mixed after alignment, but the snRNA-seq data were not mixed well with the rest (Fig.
423 S12c). Similar to online iNMF, some of the neurons' subcluster structure appeared to be lost
424 after the integration by Harmony. Overall, Portal and Seurat presented the best scRNA-seq
425 and snRNA-seq data alignment performance; however, not including data preprocessing time,
426 Seurat took over 17 hours to complete the task, while Portal only took 87 seconds.

427 **Portal aligns spermatogenesis differentiation process across multiple** 428 **species.**

429 Portal does not need to specify the structure and the strength of unwanted variation when
430 integrating datasets. Instead, it can flexibly account for general difference between datasets,
431 including batch effects, technical noises, and other sources of unwanted variation, by nonlinear
432 encoders and generators in the adversarial domain translation framework. Therefore, Portal is
433 also applicable for merging datasets with intrinsic biological divergence, revealing biologically
434 meaningful connections among these datasets. In this section, we demonstrate that Portal

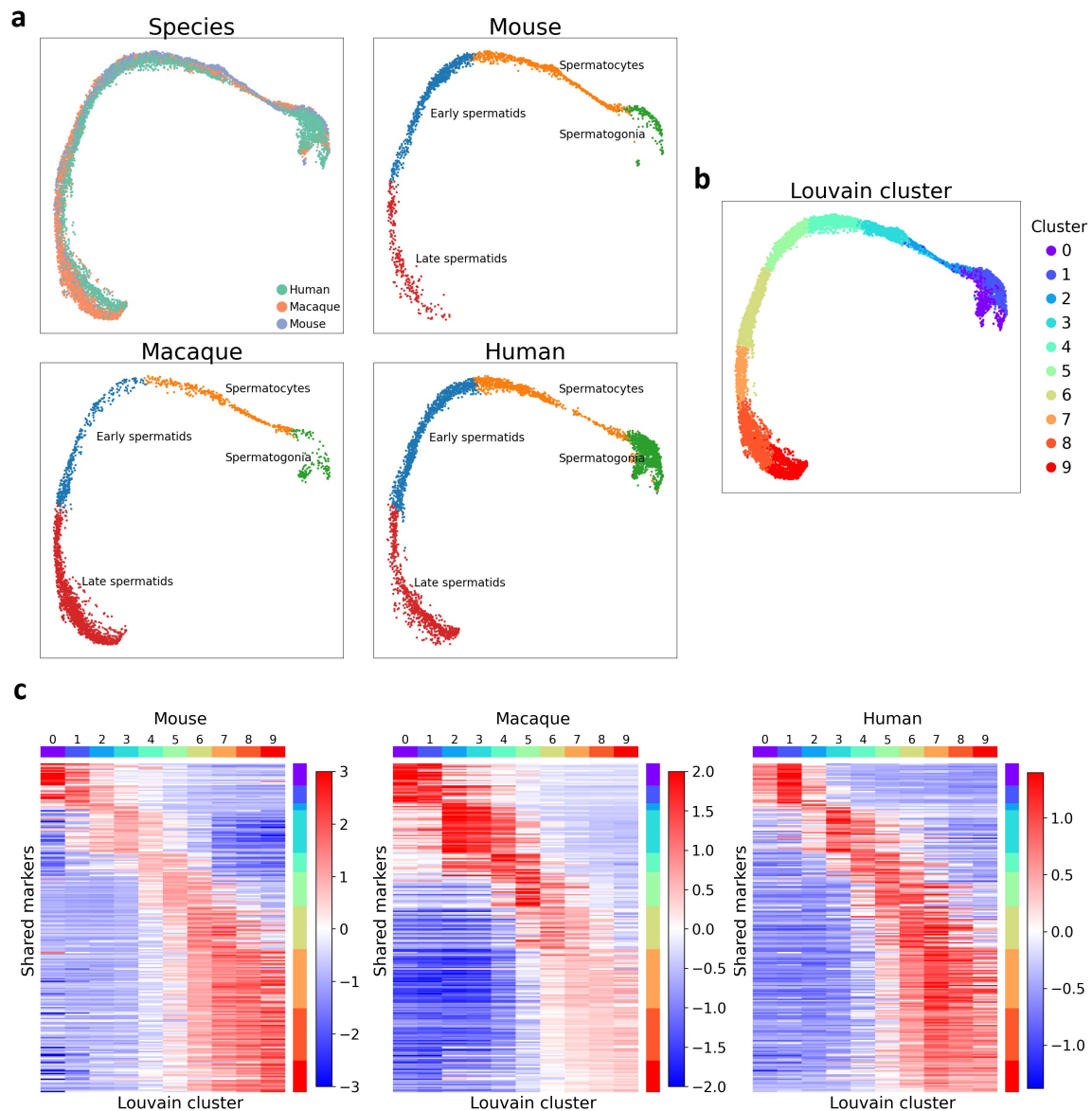


Figure 6: **Integration of spermatogenesis datasets across different species, including mouse, macaque and human.** a. UMAP plot of Portal's result colored by species, as well as UMAP plots of integrated mouse, macaque, human datasets visualized separately. b. Ten clusters were obtained by applying the Louvain clustering algorithm, facilitating detailed comparative analysis across species. c. Portal identified 239 highly variable genes that are shared in the spermatogenesis process across all three mammalian species.

435 can successfully align scRNA-seq datasets of the testes from different species including mouse,
436 macaque and human (Fig. 6).

437 Compared to merging datasets from the same species, cross-species integration poses
438 additional unique challenges. Although the transcriptomes of different species may share
439 expression of homologous or orthologous genes, the number of shared genes varies between

440 different species and is limited. Furthermore, two species may have genes with very similar
441 sequence and be annotated in the transcriptome by the same name, but have altered function,
442 which means that expression of the same gene in different species can denote different cell
443 function [44]. In other words, the amount of information one can utilize for integration becomes
444 limited and fuzzier while the variation across datasets becomes far larger, with limited number
445 of shared genes and even fewer shared highly variable genes across different species. Nonetheless,
446 cross-species integration can be very meaningful despite its challenges, as it can generate quick
447 draft annotations of new or less-studied species' atlases and cell types via label transfer from
448 well-studied species. This saves time in the manual annotation process of single-cell tissue atlas
449 generation for new species. Such integration can also enable detailed comparisons between
450 species, such as comparisons of cell type composition, discovery of cell types unique to a
451 particular species, or cross-species comparisons of the same cell types.

452 Mammalian spermatogenesis is a continuous and irreversible differentiation process from
453 spermatogonial stem cells (SSCs) to sperm cells [45, 46, 47, 48, 16]. Due to the unique
454 degenerate nature of the Y chromosome (Y-chr), Y-chr gene expression is intricately and
455 tightly regulated in the spermatogenesis process through meiotic sex chromosome inactivation
456 (MSCI) [49, 50, 51, 52, 53]. Interestingly, Y-linked genes are highly divergent between different
457 species, including between closely related primates such as the chimpanzee, macaque, and
458 human [49, 54, 55]; yet MSCI as a process is conserved across many species and is required for
459 male fertility [52, 56]. This evidence suggests that while the evolution of genes on the Y-chr
460 generated diverse species-specific genetic combinations, the tight control of gene expression
461 through MSCI is required to ensure genetic stability [49]. Recently, cross-species comparisons of
462 "escape genes" that are able to maintain or re-activate their expression despite MSCI repression
463 during spermatogenesis have generated fascinating insights on evolutionary biology, and on
464 sex chromosome evolution [51, 53, 57, 16]. In this biological context, integrating datasets
465 with continuous and gradient developmental trajectories, such as for spermatogenesis data,
466 requires integration methods to preserve the continuous structure of each dataset, while still
467 providing high accuracy of cell type alignment between datasets. This is more difficult when,
468 like in spermatogenesis data, there are no distinct clusters, making integration of such data a
469 particularly difficult task. Here, we perform cross-species integration of testes datasets from
470 three species, including one mouse [47], one macaque and one human [16], aligning the different

471 stages of spermatogenesis across species thereby highlighting unique features of each. The
472 successful integration of these spermatogenesis trajectories serves as a demonstration of the
473 power of Portal in complex and low-information data alignment, and how it can facilitate the
474 annotation and discovery process for new single-cell tissue atlases.

475 We first annotated the mouse sample according to the pattern of marker genes (Sper-
476 matogonia: *Sycp1*, *Uchl1*, *Crabp1*, *Stra8*; Spermatocytes: *Piwil1*, *Pttg1*, *Insl6*, *Spag6*; Early
477 spermatids: *Tssk1*, *Acrv1*, *Spaca1*, *Tsga8*; Late spermatids: *Prm1*, *Prm2*, *Tnp1*, *Tnp2*) [45, 46].
478 Then we used Portal to harmonize the three samples, where the integration was accomplished
479 in the mouse sample domain: The cells from the mouse sample were used as reference, and
480 cells from the other two species were mapped to the mouse sample domain by Portal. Based on
481 our annotation of the mouse sample, we transferred the broad cell type labels to cells from the
482 macaque and human samples according to the nearest neighbors, using the alignment given by
483 Portal (Fig. 6a). To check whether the alignments were correct for broad cell type identities, we
484 visualized the UMAPs for cells from each species labeled by their original published annotations
485 [16], and we confirmed concordant cell type integration across species (Fig. S13). Then, we used
486 Louvain clustering algorithm to cluster the cells from all three species based on integrated cell
487 representations. Ten clusters were found, and the cluster names were relabeled by their order
488 of progression from the spermatogonia along the developmental trajectory (Fig. 6b). We then
489 visualized the expression of known spermatogenesis markers [45, 46, 16] in each Louvain cluster
490 and found that the Louvain clusters generated by Portal's alignment clearly captured the key
491 transcriptomic features for each stage of spermatogenesis, and correctly identified cells from
492 each stage for all three species (Fig. S14, S15). Furthermore, each Louvain cluster represented
493 a more fine-grained classification of cells within the labeled broad spermatogenesis cell types.
494 Using these clusters we assessed the transcriptomic changes throughout the differentiation
495 trajectory with higher resolution (Fig. S14, S15). Notably, many of the marker genes known
496 to define stages of spermatogenesis in human were not shared or sometimes not expressed in
497 macaque and/or mouse scRNA-seq data. For example, human genes *SYCP3*, *YBX2*, *SPACA4*,
498 *H1FNT*, *PRM1*, and *TNP1* were known to mark human spermatogenesis progression, but they
499 were absent in the macaque dataset. As only highly variable genes that were expressed in all
500 three species were considered in the integration process, these genes were not used by Portal.
501 However, they showed clear expression in the cell clusters where they were expected to be

502 expressed after integration (Fig. S15), confirming the correctness of Portal’s integration result.
503 The above results show that Portal can provide an accurate integration even for genes not
504 measured by all three samples.

505 Cross-species data integration can be a quick and easy way to generate draft cell atlas
506 annotations for new species via label transfer from well-annotated species, but moreover, such
507 integrated data can be used to highlight interesting biological features of shared cell types. In
508 our Louvain clusters for spermatogenesis, for each species, we selected top 200 highly expressed
509 genes of every cluster. By taking the intersection of those genes across three species, we then
510 identified 239 highly variable genes that are shared in the spermatogenesis process across all
511 three mammalian species (Fig. 6c). For the highly expressed genes that were unique to only one
512 species, we compared their expressions across all three species (Fig. S16). Such comparisons
513 could give insight into shared and divergent features of spermatogenesis across different species.

514 Discussion

515 Taking advantage of machine learning methodologies, Portal is an efficient and powerful tool for
516 single-cell data integration that easily scales to handle large datasets with sample sizes in the
517 millions. As a machine learning-based model, Portal is easy to train, and its training process is
518 greatly accelerated by using GPUs. Meanwhile, mini-batch optimization allows Portal to be
519 trained with a low memory usage. Besides, it also makes Portal applicable in the situation
520 where the dataset is not fully observed, but arrives incrementally.

521 The nonlinearity of neural networks makes Portal a flexible approach that can adjust for
522 complex dataset-specific effects. Nonetheless, according to benchmarking studies, strong ability
523 for removing dataset-specific effects often comes with the weakness in conserving biological
524 variation [42, 58], e.g., being prone to overcorrection. Portal overcomes this challenge by its
525 model and algorithm designs. First, the boundaries of discriminator scores help Portal to
526 protect dataset-unique cell types from overcorrection. Second, the use of three specifically
527 designed regularizers not only assists Portal to find correct correspondence across domains, but
528 also enables Portal to have high-level preservation of subcluster and small cluster identities in
529 both datasets.

530 Two existing popular methods are Seurat and BBKNN. Seurat often provides integration
531 results with high accuracy, but also requires high computational cost, preventing its usage on

532 large-scale datasets; while BBKNN is well-known for its extremely fast speed, its comparatively
533 less precise results are sometimes a concern for users (Figs. 2, S3, S4, and S5). A major advance
534 of Portal over these existing state-of-the-art integration approaches is its ability to achieve
535 high efficiency and accuracy simultaneously. With speed comparable or faster than BBKNN,
536 and significantly lower memory requirement than BBKNN (Fig. S5a), Portal presents similar
537 alignment performance as well as superior information preservation performance compared to
538 that of Seurat (Figs. 2, 3a, 4b, S3 and S4).

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

552 Recently, two other generative adversarial networks based approaches have been proposed for
553 single-cell data integration, namely cross-modal autoencoders [59] and iMAP [60]. Cross-modal
554 autoencoders rely on cell type label information or paired data to obtain accurate integration
555 results, and such paired information may not always be available. iMAP software faces scaling
556 challenges when working on datasets of 500,000 cells or more, due to its high GPU memory
557 consumption; further, as a two-stage integration method, iMAP results often rely on the MNN
558 pairs detected in the first stage as anchors. By these reasons, we believe that Portal has made
559 significant progress in the development of single-cell methods, as it is a unified framework
560 which utilizes advanced techniques in domain translation with its tailored designs to achieve
561 efficiency, scalability, flexibility and accuracy simultaneously.

562 It is now clear that using single-cell technologies to assemble comprehensive whole organism

563 atlases encompassing diverse cell types is accelerating biological discovery, and this demand
564 will only grow as more datasets are generated. The demand for integration of such datasets,
565 along with the size of these datasets, will expand correspondingly. We expect that Portal, with
566 its fast, versatile, and robust integration performance, will play a valuable and essential role in
567 the modern life scientist’s single-cell analysis toolkit.

568 **Methods**

569 **The model of Portal**

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

578 Portal achieves the integration of datasets through training a unified framework of adversarial
579 domain translation. Let \mathbf{x} and \mathbf{y} be the cell embeddings in \mathcal{X} and \mathcal{Y} , respectively. For domain
580 \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 $\mathbf{x} \in \mathcal{X}$.
581 Encoder $E_1(\cdot)$ is designed to remove domain-specific effects in \mathcal{X} . To transfer cells from \mathcal{X} to
582 \mathcal{Y} , Portal then uses generator $G_2(\cdot) : \mathcal{Z} \rightarrow \mathcal{Y}$ to model the data generating process in domain
583 \mathcal{Y} , where domain-specific effects in \mathcal{Y} are induced. $E_1(\cdot)$ and $G_2(\cdot)$ together form a domain
584 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,
585 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}
586 along the path $\mathcal{Y} \rightarrow \mathcal{Z} \rightarrow \mathcal{X}$.

587 Portal trains domain translation network $G_2(E_1(\cdot)) : \mathcal{X} \rightarrow \mathcal{Y}$, such that the distribution
588 of transferred cells $G_2(E_1(\mathbf{x}))$ can be mixed with the distribution of cell embeddings \mathbf{y} in
589 domain \mathcal{Y} . Discriminator $D_2(\cdot)$ is employed in domain \mathcal{Y} to identify where the poor mixing of
590 the two distributions occurs. The competition between domain translation network $G_2(E_1(\cdot))$
591 and discriminator $D_2(\cdot)$ is known as adversarial learning [27]. Discriminator $D_2(\cdot)$ will send a

592 feedback signal to improve the domain translation network $G_2(E_1(\cdot))$ until the two distributions
593 are well mixed. By symmetry, domain translation network $G_1(E_2(\cdot)) : \mathcal{Y} \rightarrow \mathcal{X}$ and discriminator
594 $D_1(\cdot)$ deployed in domain \mathcal{X} form another adversarial learning pair. The feedback signal from
595 $D_1(\cdot)$ improves $G_1(E_2(\cdot))$ until the well mixing of the transferred cell distribution $G_1(E_2(\mathbf{y}))$
596 and the original cell distribution \mathbf{x} in domain \mathcal{X} .

597 Notice that the well mixing of the transferred distribution and the original distribution does
598 not necessarily imply the correct correspondence established between \mathcal{X} and \mathcal{Y} . First, cells
599 from a unique cell population in domain \mathcal{X} should not be forced to mix with cells in domain \mathcal{Y} .
600 Second, cell types A and B in domain \mathcal{X} could be incorrectly aligned with cell types B and A
601 in domain \mathcal{Y} , respectively, even if the two distributions are well mixed. These problems can
602 occur because we don't have any cell type label information as an anchor for data alignment
603 across domains. To address these, Portal has the following unique features, distinguishing it
604 from existing domain translation methods [28, 29]. First, Portal has a tailored discriminator
605 for the integrative analysis of single-cell data, which can prevent mixing of unique cell types
606 in one domain with a different type of cell in another domain. Second, Portal deploys three
607 regularizers to find correct correspondence during adversarial learning; These regularizers also
608 play a critical role in accounting for domain-specific effects and retaining biological variation in
609 the shared latent space \mathcal{Z} .

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

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

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

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

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

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

616 **Adversarial learning with discriminator score thresholding.** The adversarial training
617 between discriminators and domain translation networks is formulated as a min-max opti-

618 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
619 $\mathcal{L}_{\mathcal{Y}}(D_2, E_1, G_2) = \mathbb{E}[\log D_2(\mathbf{y})] + \mathbb{E}[\log(1 - D_2(G_2(E_1(\mathbf{x}))))]$ are the objective functions for
620 adversarial learning in domain \mathcal{X} and domain \mathcal{Y} , respectively. Given domain translation
621 network $G_1(E_2(\cdot))$, discriminator $D_1(\cdot) : \mathcal{X} \rightarrow (0, 1)$ is trained to distinguish the transferred
622 cells $G_1(E_2(\mathbf{y}))$ from the original cells \mathbf{x} , where a high score (close to 1) indicates a “real
623 cell” in domain \mathcal{X} , and a low score (close to 0) indicates a “transferred cell” from domain \mathcal{Y} .
624 This is achieved by maximizing $\mathcal{L}_{\mathcal{X}}$ with respect to $D_1(\cdot)$. Similarly, discriminator $D_2(\cdot)$ in
625 domain \mathcal{Y} is updated by maximizing $\mathcal{L}_{\mathcal{Y}}$. Given discriminators $D_1(\cdot)$ and $D_2(\cdot)$, the domain
626 translation networks are trained by minimizing $\mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}}$ with respect to $E_1(\cdot), G_2(\cdot)$ and
627 $E_2(\cdot), G_1(\cdot)$, such that the discriminators cannot distinguish transferred cells from real cells.
628 This is equivalent to $\min_{\{E_1, G_1, E_2, G_2\}} \mathbb{E}[\log(1 - D_1(G_1(E_2(\mathbf{y}))))] + \mathbb{E}[\log(1 - D_2(G_2(E_1(\mathbf{x}))))]$.
629 However, direct optimization of this objective function is known to suffer from severe gradient
630 vanishing [27, 61]. Therefore, we adopt the “logD-trick” [27] to stabilize the training process.
631 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
632 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
633 $E_1(\cdot), G_2(\cdot)$ and $E_2(\cdot), G_1(\cdot)$, instead of minimizing $\mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}} = \mathbb{E}[\log(1 - D_1(G_1(E_2(\mathbf{y}))))] +$
634 $\mathbb{E}[\log(1 - D_2(G_2(E_1(\mathbf{x}))))]$.

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

646 From the above reasoning, domain-unique cell populations are prone to be assigned with
647 extreme discriminator scores, either too high in the original domain or too low in the transferred
648 domain. Such extreme scores can lead to overcorrection. To address this issue in single-cell data

649 integration tasks, we set boundaries for discriminator scores to make discriminators inactive
650 on such cells. Specifically, the outputs of standard discriminators are transformed into $(0, 1)$
651 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
652 $d_i(\mathbf{x}) \in (-\infty, \infty)$ is the logit of the output. We bound the discriminator score by thresholding
653 its logit to a reasonable range $[-t, t]$:

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

654 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
655 when $d_i(\mathbf{x}) < -t$ or $d_i(\mathbf{x}) > t$, providing zero gradients for updating the parameters of encoders
656 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
657 design, the adversarial learning mechanism in Portal is only applied to cell populations that
658 are likely to be common across domains. In Portal, we then use this modified version of
659 discriminators $\tilde{D}_i(\cdot)$ to avoid incorrect alignment of domain-unique cell populations. For clarity,
660 we still use the notation $D_i(\cdot)$ to represent $\tilde{D}_i(\cdot)$ hereinafter.

661 **Regularization for autoencoder consistency.** Encoder $E_1(\cdot) : \mathcal{X} \rightarrow \mathcal{Z}$ and generator
662 $G_1(\cdot) : \mathcal{Z} \rightarrow \mathcal{X}$ form an autoencoder structure, where $E_1(\cdot)$ removes domain-specific effects
663 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
664 autoencoder structure. Therefore, we use the regularizer in (3) for the autoencoder consistency,
665 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}
666 and \mathcal{Y} .

667 **Regularization for cosine similarity correspondence.** Besides the autoencoder consis-
668 tency, the cosine similarity regularizer in (5) plays a critical role in data alignment between
669 domains, 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
670 imposes the cross-domain correspondence on domain translation. The key idea is that a cell
671 and its transferred version should not be largely different from each other in terms of cosine
672 similarity. This is because cosine similarity is scale invariant and insensitive to domain-specific
673 effects, including differences in sequencing depth and capture efficiency of protocols used
674 across datasets [62, 24, 21]. Thus, the cosine similarity regularizer is helpful to uncover robust
675 correspondence between cells of the same cell type across domains.

676 **Domain-specific effects removal in the shared latent space by latent alignment regu-**
677 **larization.** Portal decouples domain translation into the encoding process $\mathcal{X} \rightarrow \mathcal{Z}$ (or $\mathcal{Y} \rightarrow \mathcal{Z}$)
678 and the generating process $\mathcal{Z} \rightarrow \mathcal{Y}$ (or $\mathcal{Z} \rightarrow \mathcal{X}$). Although adversarial learning enables the do-

679 main translation networks to effectively transfer cells across domains, it can not remove domain-
 680 specific effects in shared latent space \mathcal{Z} . To enable encoders $E_1(\cdot), E_2(\cdot)$ to eliminate domain-
 681 specific effects in \mathcal{X} and \mathcal{Y} , we propose the latent alignment regularizer in (4) for the consistency
 682 in latent space \mathcal{Z} , where $\mathcal{R}_{\text{LA}} = \frac{1}{q} \{ \mathbb{E} [\|E_1(\mathbf{x}) - E_2(G_2(E_1(\mathbf{x})))\|_2^2] + \mathbb{E} [\|E_2(\mathbf{y}) - E_1(G_1(E_2(\mathbf{y})))\|_2^2] \}$,
 683 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 $E_2(G_2(E_1(\mathbf{x})))$
 684 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
 685 $E_1(G_1(E_2(\mathbf{y})))$ is the latent code of its transferred version. The regularizer (4) encourages the
 686 latent codes of the same cell to be close to each other. This regularizer helps encoders $E_1(\cdot)$
 687 and $E_2(\cdot)$ to remove domain-specific effects, such that the latent codes in \mathcal{Z} preserve biological
 688 variation of cells from different domains.

689 **Algorithm.** Now we develop an alternative updating algorithm to solving the optimization
 690 problem of adversarial domain translation with the three regularizers. To efficiently solve the
 691 optimization problem, we replace the constraints (3), (4) and (5) by its Lagrange form. We
 692 introduce three regularization parameters λ_{AE} , λ_{LA} and λ_{cos} as coefficients for the regularizers.
 693 The optimization problem of Portal is rewritten as

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

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

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

696 where \mathcal{L}_{adv} stands for the adversarial learning objective, whose value is $\mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}}$ when maximizing
 697 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
 698 to $E_1(\cdot), G_1(\cdot), E_2(\cdot), G_2(\cdot)$.

699 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
 700 $\theta_{E_1}, \theta_{E_2}, \theta_{G_1}, \theta_{G_2}, \theta_{D_1}$ and θ_{D_2} . Then we collect the parameter sets as $\theta_E = \{\theta_{E_1}, \theta_{E_2}\}$, $\theta_G =$
 701 $\{\theta_{G_1}, \theta_{G_2}\}$ and $\theta_D = \{\theta_{D_1}, \theta_{D_2}\}$. We use the Monte Carlo estimators to approximate expectations
 702 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}

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

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

704 The implementation of Portal is summarized in Algorithm 1.

Algorithm 1 Stochastic gradient descent training of Portal.

Require: Batch size m , coefficients λ_{AE} , λ_{LA} and λ_{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

705 After training, cells from domains \mathcal{X} and \mathcal{Y} are encoded into \mathcal{Z} to construct an integrated
706 dataset, which can be applied to downstream analysis. In each domain, the original cells and
707 transferred cells are also well integrated. For integration of multiple datasets, Portal can handle
708 them incrementally, by transferring all other datasets into the domain formed by one dataset.

709 Analysis details

710 **Data preprocessing.** For all datasets, we used raw read or unique molecular identifier (UMI)
711 matrices depending on the data source. We then performed standard data preprocessing for
712 each count matrix, including log-normalization, feature selection, scaling and dimensionality
713 reduction. For each dataset represented by a cell-by-gene count matrix, we first adopted the
714 log-normalization, following the Seurat and Scanpy pipelines [22, 63]. For each cell, its library
715 size was normalized to 10,000 reads. Specifically, the counts abundance of each gene was

716 divided by the total counts for each cell, then multiplied by a scaling factor of 10,000. The
717 normalized dataset was then transformed to log scale by the function $\log(1 + x)$. In order
718 to identify a subset of features that highlight variability across individual cells, we adopted
719 the feature selection procedure from the Seurat pipeline. For each dataset, we selected K top
720 highly variable genes ranked by dispersion with the control of means. In this paper, we used
721 $K = 4,000$ throughout all analyses except for the cross-species analysis. In the cross-species
722 analysis, we used $K = 3,000$ since the usage of a larger number of features would result in
723 the situation that correspondence across species is dominated by the distinction (e.g., altered
724 functions of genes annotated by the same name). For each selected variable gene, we centered
725 and standardized its expressions across individual cells to have mean at zero and variance
726 at one. After the above procedures, which were applied to individual datasets, we continued
727 to preprocess data across datasets. For those datasets to be integrated, we collected genes
728 that were identified as top highly variable genes in all of them as features for integration. We
729 extracted the scaled data with these features from each dataset, and then concatenated them
730 based on features to perform joint PCA. Top $p = 30$ principle components were kept for all
731 dataset as inputs to Portal. For the shared latent space, we set its dimensionality to be $q = 20$
732 throughout all analyses.

733 **Hyperparameter setting.** Hyperparameters used in Portal are $m, t, \lambda_{\text{AE}}, \lambda_{\text{LA}}, \lambda_{\text{cos}}$, where m
734 is the batch size used by Portal for mini-batch training; t is the absolute value of boundaries
735 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
736 autoencoder consistency regularizer \mathcal{R}_{AE} , latent alignment regularizer \mathcal{R}_{LA} and cosine similarity
737 regularizer \mathcal{R}_{cos} respectively. Throughout all analyses, we set $m = 500, t = 5.0, \lambda_{\text{AE}} = 10.0,$
738 $\lambda_{\text{LA}} = 10.0$. Hyperparameter λ_{cos} was tuned within the range $[10.0, 50.0]$ with interval 5.0
739 according to the mixing metric, where the mixing metric was designed in Seurat to evaluate
740 how well the datasets mixed after integration. The insight into tuning λ_{cos} is as follows: During
741 domain translations, there is a trade-off between preservation of similarity across domains
742 and flexibility of modeling domain differences. Since \mathcal{R}_{cos} is designed to preserve the cosine
743 similarity during translations, a higher value of λ_{cos} can enhance the cosine similarity as the
744 cross-domain correspondence, and a lower λ_{cos} allows domain translation networks to deal with
745 remarkable differences between domains. Following this intuition, we empirically find out that
746 $\lambda_{\text{cos}} = 10.0$ has a good performance when harmonizing datasets with intrinsic differences, for

747 example, datasets used in cross-species analysis. For other integration tasks, $\lambda_{\text{cos}} = 20.0$ often
748 yields reasonable results, which is adopted as the default setting in our package. Slightly better
749 alignment results could be achieved by tuning λ_{cos} .

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

756 **Evaluation metrics.** We assessed all metrics based on Portal's integration results in shared
757 latent space \mathcal{Z} . We used kBET [32] for quantitative evaluation of integration approaches in
758 terms of domain-specific effects removal ability. Firstly, kBET creates a k -nearest neighbour
759 matrix. Then, 10% of the samples are chosen for hypothesis testing, where the null hypothesis
760 is that all batches are well-mixed. For each of selected samples, kBET adopts a Pearson's
761 χ^2 -based test to check whether the batch label distribution in its neighbourhood is similar
762 to the global batch label distribution or not. In our experiments, we ran 100 replicates of
763 kBET with 1,000 random samples, and used the median of the 100 average acceptance rates as
764 the final result. We used the neighbourhood size following the default setting in the official
765 implementation of kBET. To evaluate the cluster preservation performance, we used ARI and
766 ASW. ARI measures the degree to which the two clustering results match. It outputs scores
767 ranging from 0 to 1, where 0 indicates that the two clustering labels are independent to each
768 other, and 1 means that the two clustering labels are the same up to a permutation. We
769 obtained clustering results following the Seurat clustering pipeline with its default setting,
770 and assessed ARI by comparing identified clusters and cell type annotations. ASW is another
771 metric to quantify cluster preservation. The silhouette width for cell \mathbf{x} from cell type C is
772 defined as $(b(\mathbf{x}) - a(\mathbf{x})) / \max(a(\mathbf{x}), b(\mathbf{x}))$, where $a(\mathbf{x})$ is the average distance from cell \mathbf{x} to
773 all cells from cell type C , and $b(\mathbf{x})$ is the minimum value of average distances from cell \mathbf{x} to
774 all cells from each cell type other than C . ASW lies between -1 and 1 , where a higher score
775 means that cells are closer to cells of the same cell type, indicating good cluster preservation.

776 **Benchmarking of running time and memory usage.** Standard data preprocessing such as
777 normalization, feature selection and dimension reduction could be performed incrementally using

778 mini-batches to control memory usage. In Portal’s preprocessing, we adopted the incremental
779 strategy and used a chunk size of 20,000. For example, the preprocessing of Portal took 63.4
780 minutes, requiring 22.0 GB peak running memory on the two mouse brain atlases datasets
781 with 1,100,167 cells. The preprocessing time could be reduced to 37.7 minutes when the chunk
782 size was increased to 200,000, with 36.4 GB peak running memory. Some other methods may
783 not be able to adopt a mini-batch implementation. For the two mouse brain atlases datasets,
784 Harmony took 17.6 minutes to finish preprocessing, but required 127.1 GB memory usage.
785 Online iNMF performed preprocessing with mini-batches. Its default preprocessing procedure
786 on the two mouse brain atlases datasets took 15.9 hours, with 0.6 GB memory usage. For a
787 fair comparison, time and memory usages of data preprocessing procedures were not included
788 in our benchmarking.

789 **Visualization.** We used the UMAP algorithm [31] for visualization of cell representations
790 in a two-dimensional space. In all analyses, the UMAP algorithm was run with 30-nearest
791 neighbors, minimum distance 0.3, and correlation metric.

792 Acknowledgements

793 The authors would like to thank Camille Sophie Ezran (Stanford University), Dr. Angela
794 Oliveira Pisco (CZ Biohub), and Dr. Hosu Sin (Stanford University) for valuable discussions.
795 This work is supported in part by Hong Kong Research Grant Council [16101118, 24301419,
796 14301120, 16307818, 16301419, 16308120], the Hong Kong University of Science and Technology’s
797 startup grant [R9405,R9364], the Hong Kong University of Science and Technology Big Data
798 for Bio Intelligence Laboratory (BDBI), the Lo Ka Chung Foundation through the Hong Kong
799 Epigenomics Project, the Chau Hoi Shuen Foundation, the Chinese University of Hong Kong
800 direct grants [4053360, 4053423, 4053476], the Chinese University of Hong Kong startup grant
801 [4930181], the Chinese University of Hong Kong’s Project Impact Enhancement Fund (PIEF)
802 and Science Faculty’s Collaborative Research Impact Matching Scheme (CRIMS), the East
803 China Normal University startup grant, the Shanghai Sailing Program. The computational
804 task for this work was partially performed using the X-GPU cluster supported by the RGC
805 Collaborative Research Fund: C6021-19EF.

806 **Data availability**

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

- 808 • Mouse Brain cells from Saunders et al [8] (<http://dropviz.org>).
- 809 • Mouse Brain cells from Zeisel et al [9] (<http://mousebrain.org/downloads.html>).
- 810 • Mouse Brain cells from Rosenberg et al [43] (GSE110823).
- 811 • Mouse cell atlas from the Tabula Muris Consortium [7] ([https://figshare.com/projects/
812 Tabula_Muris_Transcriptomic_characterization_of_20_organs_and_tissues_from_
813 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)).
- 814 • Mouse lemur cell atlas from the Tabula Microcebus Consortium ([https://figshare.
815 com/projects/Tabula_Microcebus/112227](https://figshare.com/projects/Tabula_Microcebus/112227)).
- 816 • Mouse spermatogenesis cells from Ernst et al [47] ([https://www.ebi.ac.uk/arrayexpress/
817 experiments/E-MTAB-6946/](https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-6946/)).
- 818 • Human spermatogenesis cells from Shami et al [16] (GSE142585).
- 819 • Macaque spermatogenesis cells from Shami et al [16] (GSE142585).

820 **Code availability**

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

822 References

- 823 [1] Alexandra-Chloé Villani, Rahul Satija, Gary Reynolds, Siranush Sarkizova, Karthik
824 Shekhar, James Fletcher, Morgane Griesbeck, Andrew Butler, Shiwei Zheng, Suzan Lazo,
825 et al. Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes,
826 and progenitors. *Science*, 356(6335), 2017.
- 827 [2] Maria Brbić, Marinka Zitnik, Sheng Wang, Angela O Pisco, Russ B Altman, Spyros
828 Darmanis, and Jure Leskovec. Mars: discovering novel cell types across heterogeneous
829 single-cell experiments. *Nature methods*, 17(12):1200–1206, 2020.
- 830 [3] Giovanni Iacono, Ramon Massoni-Badosa, and Holger Heyn. Single-cell transcriptomics
831 unveils gene regulatory network plasticity. *Genome biology*, 20(1):1–20, 2019.
- 832 [4] Anna SE Cuomo, Daniel D Seaton, Davis J McCarthy, Iker Martinez, Marc Jan Bonder,
833 Jose Garcia-Bernardo, Shradha Amatya, Pedro Madrigal, Abigail Isaacson, Florian Buet-
834 tner, et al. Single-cell RNA-sequencing of differentiating iPS cells reveals dynamic genetic
835 effects on gene expression. *Nature communications*, 11(1):1–14, 2020.
- 836 [5] Barbara Treutlein, Doug G Brownfield, Angela R Wu, Norma F Neff, Gary L Mantalas,
837 F Hernan Espinoza, Tushar J Desai, Mark A Krasnow, and Stephen R Quake. Reconstruct-
838 ing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. *Nature*,
839 509(7500):371–375, 2014.
- 840 [6] Xiaojie Qiu, Qi Mao, Ying Tang, Li Wang, Raghav Chawla, Hannah A Pliner, and Cole
841 Trapnell. Reversed graph embedding resolves complex single-cell trajectories. *Nature*
842 *methods*, 14(10):979–982, 2017.
- 843 [7] Tabula Muris Consortium et al. Single-cell transcriptomics of 20 mouse organs creates a
844 Tabula Muris. *Nature*, 562(7727):367–372, 2018.
- 845 [8] Arpiar Saunders, Evan Z Macosko, Alec Wysoker, Melissa Goldman, Fenna M Krienen,
846 Heather de Rivera, Elizabeth Bien, Matthew Baum, Laura Bortolin, Shuyu Wang, et al.
847 Molecular diversity and specializations among the cells of the adult mouse brain. *Cell*,
848 174(4):1015–1030, 2018.

- 849 [9] Amit Zeisel, Hannah Hochgerner, Peter Lönnerberg, Anna Johnsson, Fatima Memic, Job
850 Van Der Zwan, Martin Häring, Emelie Braun, Lars E Borm, Gioele La Manno, et al.
851 Molecular architecture of the mouse nervous system. *Cell*, 174(4):999–1014, 2018.
- 852 [10] Hoa Thi Nhu Tran, Kok Siong Ang, Marion Chevrier, Xiaomeng Zhang, Nicole Yee Shin
853 Lee, Michelle Goh, and Jinmiao Chen. A benchmark of batch-effect correction methods
854 for single-cell RNA sequencing data. *Genome biology*, 21(1):1–32, 2020.
- 855 [11] Oliver Stegle, Sarah A Teichmann, and John C Marioni. Computational and analytical
856 challenges in single-cell transcriptomics. *Nature Reviews Genetics*, 16(3):133–145, 2015.
- 857 [12] Monika Litviňuková, Carlos Talavera-López, Henrike Maatz, Daniel Reichart, Catherine L
858 Worth, Eric L Lindberg, Masatoshi Kanda, Krzysztof Polanski, Matthias Heinig, Michael
859 Lee, et al. Cells of the adult human heart. *Nature*, 588(7838):466–472, 2020.
- 860 [13] Kyle J Travaglini, Ahmad N Nabhan, Lolita Penland, Rahul Sinha, Astrid Gillich, Rene V
861 Sit, Stephen Chang, Stephanie D Conley, Yasuo Mori, Jun Seita, et al. A molecular cell
862 atlas of the human lung from single-cell RNA sequencing. *Nature*, 587(7835):619–625,
863 2020.
- 864 [14] Blue B Lake, Simone Codeluppi, Yun C Yung, Derek Gao, Jerold Chun, Peter V
865 Kharchenko, Sten Linnarsson, and Kun Zhang. A comparative strategy for single-nucleus
866 and single-cell transcriptomes confirms accuracy in predicted cell-type expression from
867 nuclear RNA. *Scientific reports*, 7(1):1–8, 2017.
- 868 [15] Haojia Wu, Yuhei Kirita, Erinn L Donnelly, and Benjamin D Humphreys. Advantages of
869 single-nucleus over single-cell RNA sequencing of adult kidney: rare cell types and novel
870 cell states revealed in fibrosis. *Journal of the American Society of Nephrology*, 30(1):23–32,
871 2019.
- 872 [16] Adrienne Niederriter Shami, Xianing Zheng, Sarah K Munyoki, Qianyi Ma, Gabriel L
873 Manske, Christopher D Green, Meena Sukhwani, Kyle E Orwig, Jun Z Li, and Saher Sue
874 Hammoud. Single-cell RNA sequencing of human, macaque, and mouse testes uncovers
875 conserved and divergent features of mammalian spermatogenesis. *Developmental cell*,
876 54(4):529–547, 2020.

- 877 [17] Aviv Regev, Sarah A Teichmann, Eric S Lander, Ido Amit, Christophe Benoist, Ewan
878 Birney, Bernd Bodenmiller, Peter Campbell, Piero Carninci, Menna Clatworthy, et al.
879 Science forum: the human cell atlas. *elife*, 6:e27041, 2017.
- 880 [18] Xiaoping Han, Renying Wang, Yincong Zhou, Lijiang Fei, Huiyu Sun, Shujing Lai, Assieh
881 Saadatpour, Ziming Zhou, Haide Chen, Fang Ye, et al. Mapping the mouse cell atlas by
882 microwell-seq. *Cell*, 172(5):1091–1107, 2018.
- 883 [19] Si Wang, Yuxuan Zheng, Jingyi Li, Yang Yu, Weiqi Zhang, Moshi Song, Zunpeng Liu,
884 Zheyang Min, Huifang Hu, Ying Jing, et al. Single-cell transcriptomic atlas of primate
885 ovarian aging. *Cell*, 180(3):585–600, 2020.
- 886 [20] Shuai Ma, Shuhui Sun, Jiaming Li, Yanling Fan, Jing Qu, Liang Sun, Si Wang, Yiyuan
887 Zhang, Shanshan Yang, Zunpeng Liu, et al. Single-cell transcriptomic atlas of primate
888 cardiopulmonary aging. *Cell research*, 31(4):415–432, 2021.
- 889 [21] Ilya Korsunsky, Nghia Millard, Jean Fan, Kamil Slowikowski, Fan Zhang, Kevin Wei, Yuriy
890 Baglaenko, Michael Brenner, Po-ru Loh, and Soumya Raychaudhuri. Fast, sensitive and
891 accurate integration of single-cell data with Harmony. *Nature methods*, 16(12):1289–1296,
892 2019.
- 893 [22] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi,
894 William M Mauck III, Yuhan Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija.
895 Comprehensive integration of single-cell data. *Cell*, 177(7):1888–1902, 2019.
- 896 [23] Chao Gao, Jialin Liu, April R Kriebel, Sebastian Preissl, Chongyuan Luo, Rosa Castanon,
897 Justin Sandoval, Angeline Rivkin, Joseph R Nery, Margarita M Behrens, et al. Iterative
898 single-cell multi-omic integration using online learning. *Nature Biotechnology*, pages 1–8,
899 2021.
- 900 [24] Laleh Haghverdi, Aaron TL Lun, Michael D Morgan, and John C Marioni. Batch effects
901 in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors.
902 *Nature biotechnology*, 36(5):421–427, 2018.
- 903 [25] Brian Hie, Bryan Bryson, and Bonnie Berger. Efficient integration of heterogeneous
904 single-cell transcriptomes using Scanorama. *Nature biotechnology*, 37(6):685–691, 2019.

- 905 [26] Krzysztof Polański, Matthew D Young, Zhichao Miao, Kerstin B Meyer, Sarah A Teich-
906 mann, and Jong-Eun Park. BBKNN: fast batch alignment of single cell transcriptomes.
907 *Bioinformatics*, 36(3):964–965, 2020.
- 908 [27] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
909 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in*
910 *Neural Information Processing Systems*, pages 2672–2680, 2014.
- 911 [28] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-
912 image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE*
913 *International Conference on Computer Vision*, pages 2223–2232, 2017.
- 914 [29] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation
915 networks. In *Advances in Neural Information Processing Systems*, pages 700–708, 2017.
- 916 [30] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul
917 Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image
918 translation. In *Proceedings of the IEEE conference on Computer Vision and Pattern*
919 *Recognition*, pages 8789–8797, 2018.
- 920 [31] Leland McInnes, John Healy, Nathaniel Saul, and Lukas Grossberger. UMAP: Uniform
921 manifold approximation and projection. *The Journal of Open Source Software*, 3(29):861,
922 2018.
- 923 [32] Maren Büttner, Zhichao Miao, F Alexander Wolf, Sarah A Teichmann, and Fabian J
924 Theis. A test metric for assessing single-cell RNA-seq batch correction. *Nature methods*,
925 16(1):43–49, 2019.
- 926 [33] Joshua D Welch, Velina Kozareva, Ashley Ferreira, Charles Vanderburg, Carly Martin,
927 and Evan Z Macosko. Single-cell multi-omic integration compares and contrasts features
928 of brain cell identity. *Cell*, 177(7):1873–1887, 2019.
- 929 [34] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast
930 unfolding of communities in large networks. *Journal of statistical mechanics: theory and*
931 *experiment*, 2008(10):P10008, 2008.

- 932 [35] Jeanette Baran-Gale, Tamir Chandra, and Kristina Kirschner. Experimental design for
933 single-cell RNA sequencing. *Briefings in functional genomics*, 17(4):233–239, 2018.
- 934 [36] Joey Schyns, Fabrice Bureau, and Thomas Marichal. Lung interstitial macrophages: past,
935 present, and future. *Journal of immunology research*, 2018, 2018.
- 936 [37] Tabula Microcebus Consortium. <https://tabula-microcebus.ds.czbiohub.org>.
- 937 [38] Ting Xie, Yizhou Wang, Nan Deng, Guanling Huang, Forough Taghavifar, Yan Geng,
938 Ningshan Liu, Vrishika Kulur, Changfu Yao, Peter Chen, et al. Single-cell deconvolution
939 of fibroblast heterogeneity in mouse pulmonary fibrosis. *Cell reports*, 22(13):3625–3640,
940 2018.
- 941 [39] Alan Selewa, Ryan Dohn, Heather Eckart, Stephanie Lozano, Bingqing Xie, Eric Gauchat,
942 Reem Elorbany, Katherine Rhodes, Jonathan Burnett, Yoav Gilad, et al. Systematic
943 comparison of high-throughput single-cell and single-nucleus transcriptomes during car-
944 diomyocyte differentiation. *Scientific reports*, 10(1):1–13, 2020.
- 945 [40] Anne-Marie Galow, Markus Wolfien, Paula Müller, Madeleine Bartsch, Ronald M Brunner,
946 Andreas Hoefflich, Olaf Wolkenhauer, Robert David, and Tom Goldammer. Integrative
947 cluster analysis of whole hearts reveals proliferative cardiomyocytes in adult mice. *Cells*,
948 9(5):1144, 2020.
- 949 [41] Blue B Lake, Rizi Ai, Gwendolyn E Kaeser, Neeraj S Salathia, Yun C Yung, Rui Liu, Andre
950 Wildberg, Derek Gao, Ho-Lim Fung, Song Chen, et al. Neuronal subtypes and diversity
951 revealed by single-nucleus RNA sequencing of the human brain. *Science*, 352(6293):1586–
952 1590, 2016.
- 953 [42] Ricard Argelaguet, Anna SE Cuomo, Oliver Stegle, and John C Marioni. Computational
954 principles and challenges in single-cell data integration. *Nature Biotechnology*, pages 1–14,
955 2021.
- 956 [43] Alexander B Rosenberg, Charles M Roco, Richard A Muscat, Anna Kuchina, Paul Sample,
957 Zizhen Yao, Lucas T Graybuck, David J Peeler, Sumit Mukherjee, Wei Chen, et al. Single-
958 cell profiling of the developing mouse brain and spinal cord with split-pool barcoding.
959 *Science*, 360(6385):176–182, 2018.

- 960 [44] Chunmei Cui, Yuan Zhou, and Qinghua Cui. Defining the functional divergence of
961 orthologous genes between human and mouse in the context of miRNA regulation. *Briefings*
962 *in Bioinformatics*, pages 1477–4054, 2021.
- 963 [45] Christopher Daniel Green, Qianyi Ma, Gabriel L Manske, Adrienne Niederriter Shami,
964 Xianing Zheng, Simone Marini, Lindsay Moritz, Caleb Sultan, Stephen J Gurczynski,
965 Bethany B Moore, et al. A comprehensive roadmap of murine spermatogenesis defined by
966 single-cell RNA-seq. *Developmental cell*, 46(5):651–667, 2018.
- 967 [46] Brian P Hermann, Keren Cheng, Anukriti Singh, Lorena Roa-De La Cruz, Kazadi N
968 Mutoji, I-Chung Chen, Heidi Gildersleeve, Jake D Lehle, Max Mayo, Birgit Westernströer,
969 et al. The mammalian spermatogenesis single-cell transcriptome, from spermatogonial
970 stem cells to spermatids. *Cell reports*, 25(6):1650–1667, 2018.
- 971 [47] Christina Ernst, Nils Eling, Celia P Martinez-Jimenez, John C Marioni, and Duncan T
972 Odom. Staged developmental mapping and X chromosome transcriptional dynamics during
973 mouse spermatogenesis. *Nature communications*, 10(1):1–20, 2019.
- 974 [48] Xianzhong Lau, Prabhakaran Munusamy, Mor Jack Ng, and Mahesh Sangrithi. Single-cell
975 RNA sequencing of the *Cynomolgus* macaque testis reveals conserved transcriptional
976 profiles during mammalian spermatogenesis. *Developmental Cell*, 54(4):548–566, 2020.
- 977 [49] Ho-Su Sin and Satoshi H Namekawa. The great escape: Active genes on inactive sex
978 chromosomes and their evolutionary implications. *Epigenetics*, 8(9):887–892, 2013.
- 979 [50] Jennifer F Hughes and David C Page. The biology and evolution of mammalian Y
980 chromosomes. *Annual review of genetics*, 49:507–527, 2015.
- 981 [51] Satoshi H Namekawa and Jeannie T Lee. XY and ZW: is meiotic sex chromosome
982 inactivation the rule in evolution? *PLoS genetics*, 5(5):e1000493, 2009.
- 983 [52] Jeffrey M Cloutier and James MA Turner. Meiotic sex chromosome inactivation. *Current*
984 *Biology*, 20(22):R962–R963, 2010.
- 985 [53] Erica L Larson, Emily EK Kopania, and Jeffrey M Good. Spermatogenesis and the
986 evolution of mammalian sex chromosomes. *Trends in Genetics*, 34(9):722–732, 2018.

- 987 [54] Jennifer F Hughes, Helen Skaletsky, Laura G Brown, Tatyana Pyntikova, Tina Graves,
988 Robert S Fulton, Shannon Dugan, Yan Ding, Christian J Buhay, Colin Kremitzki, et al.
989 Strict evolutionary conservation followed rapid gene loss on human and rhesus Y chromo-
990 somes. *Nature*, 483(7387):82–86, 2012.
- 991 [55] Jennifer F Hughes, Helen Skaletsky, Tatyana Pyntikova, Tina A Graves, Saskia KM
992 Van Daalen, Patrick J Minx, Robert S Fulton, Sean D McGrath, Devin P Locke, Cynthia
993 Friedman, et al. Chimpanzee and human Y chromosomes are remarkably divergent in
994 structure and gene content. *Nature*, 463(7280):536–539, 2010.
- 995 [56] Polly Campbell, Jeffrey M Good, and Michael W Nachman. Meiotic sex chromosome
996 inactivation is disrupted in sterile hybrid male house mice. *Genetics*, 193(3):819–828, 2013.
- 997 [57] Ho-Su Sin, Yosuke Ichijima, Eitetsu Koh, Mikio Namiki, and Satoshi H Namekawa. Human
998 postmeiotic sex chromatin and its impact on sex chromosome evolution. *Genome research*,
999 22(5):827–836, 2012.
- 1000 [58] Malte D Luecken, Maren Büttner, Kridsakorn Chaichoompu, Anna Danese, Marta
1001 Interlandi, Michaela Fee Mueller, Daniel C Strobl, Luke Zappia, Martin Dugas, Maria
1002 Colomé-Tatché, et al. Benchmarking atlas-level data integration in single-cell genomics.
1003 *BioRxiv*, 2020.
- 1004 [59] Karren Dai Yang, Anastasiya Belyaeva, Saradha Venkatachalapathy, Karthik Damodaran,
1005 Abigail Katcoff, Adityanarayanan Radhakrishnan, GV Shivashankar, and Caroline Uh-
1006 ler. Multi-domain translation between single-cell imaging and sequencing data using
1007 autoencoders. *Nature Communications*, 12(1):1–10, 2021.
- 1008 [60] Dongfang Wang, Siyu Hou, Lei Zhang, Xiliang Wang, Baolin Liu, and Zemin Zhang. imap:
1009 integration of multiple single-cell datasets by adversarial paired transfer networks. *Genome*
1010 *biology*, 22(1):1–24, 2021.
- 1011 [61] Martin Arjovsky and Léon Bottou. Towards principled methods for training generative
1012 adversarial networks. In *ICLR*, 2017.
- 1013 [62] Sean C Bendall, Kara L Davis, El-ad David Amir, Michelle D Tadmor, Erin F Simonds,
1014 Tiffany J Chen, Daniel K Shenfeld, Garry P Nolan, and Dana Pe’er. Single-cell trajectory

- 1015 detection uncovers progression and regulatory coordination in human B cell development.
1016 *Cell*, 157(3):714–725, 2014.
- 1017 [63] F Alexander Wolf, Philipp Angerer, and Fabian J Theis. SCANPY: large-scale single-cell
1018 gene expression data analysis. *Genome biology*, 19(1):1–5, 2018.

1019 Supplementary Information

1020 Investigation of the role of each component in Portal

1021 In this section, we investigate the role of each component in Portal. The optimization problem
1022 solved by Portal is

$$\min_{\{E_1, G_1, E_2, G_2\}} \max_{\{D_1, D_2\}} \mathcal{L}_{\text{adv}} + \lambda_{\text{cos}} \mathcal{R}_{\text{cos}} + \lambda_{\text{LA}} \mathcal{R}_{\text{LA}} + \lambda_{\text{AE}} \mathcal{R}_{\text{AE}}, \quad (\text{S1})$$

1023 where $E_1(\cdot)$ and $E_2(\cdot)$ are encoder networks; $G_1(\cdot)$ and $G_2(\cdot)$ are generator networks; $D_1(\cdot)$
1024 and $D_2(\cdot)$ are discriminator networks; \mathcal{L}_{adv} stands for the adversarial learning objective, whose
1025 value is $\mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{Y}}$ when maximizing with respect to $D_1(\cdot), D_2(\cdot)$, and it is replaced with
1026 $\mathcal{L}_{\mathcal{X}}^{\log D} + \mathcal{L}_{\mathcal{Y}}^{\log D}$ when minimizing with respect to $E_1(\cdot), G_1(\cdot), E_2(\cdot), G_2(\cdot)$ according to the “logD-
1027 trick”; $\mathcal{R}_{\text{cos}} = \mathcal{R}_{\text{cos}}(E_1, G_1, E_2, G_2)$ is a regularizer for the cosine similarity correspondence
1028 cross domains; $\mathcal{R}_{\text{LA}} = \mathcal{R}_{\text{LA}}(E_1, G_1, E_2, G_2)$ is a regularizer for the alignment consistency in
1029 the latent space; $\mathcal{R}_{\text{AE}} = \mathcal{R}_{\text{AE}}(E_1, G_1, E_2, G_2)$ is a regularizer for the autoencoder consistency;
1030 $\lambda_{\text{cos}}, \lambda_{\text{LA}}, \lambda_{\text{AE}}$ are coefficients for the three regularizers respectively. To demonstrate the roles
1031 of the objective function \mathcal{L}_{adv} and three regularizers ($\mathcal{R}_{\text{cos}}, \mathcal{R}_{\text{LA}}$, and \mathcal{R}_{AE}), we rewrite the
1032 optimization problem (S1) as

$$\min_{\{E_1, G_1, E_2, G_2\}} \max_{\{D_1, D_2\}} \lambda_{\text{adv}} \mathcal{L}_{\text{adv}} + \lambda_{\text{cos}} \mathcal{R}_{\text{cos}} + \lambda_{\text{LA}} \mathcal{R}_{\text{LA}} + \lambda_{\text{AE}} \mathcal{R}_{\text{AE}}, \quad (\text{S2})$$

1033 with λ_{adv} set to 1.0 in Portal’s algorithm. Based on (S2), we are able to study on the impact
1034 of each component of Portal by manually setting the corresponding coefficient to zero, and
1035 then compare its performance with that of the standard algorithm empirically. Recall that
1036 the discriminators are designed to deal with domain-unique cell types by discriminator score
1037 thresholding. In this section, we also experimentally verify the effectiveness of such design.
1038 Here we took mouse mammary gland scRNA-seq atlas from the Tabula Muris consortium as an
1039 example. In the mouse mammary gland data, 4,481 cells were profiled by 10X Genomics (10X),
1040 and 2,405 cells were profiled by SMART-seq2 (SS2). With these two datasets, we investigate
1041 the role of each component in Portal.

1042 **Role of objective function \mathcal{L}_{adv} .** The objective function \mathcal{L}_{adv} plays an essential role in
1043 learning effective domain translation across different datasets. To demonstrate the importance
1044 of adversarial training using \mathcal{L}_{adv} , we removed it from Portal by setting λ_{adv} to zero, then
1045 applied this version of Portal (Portal ($\lambda_{\text{adv}} = 0$)) to integrate mouse mammary gland datasets.

1046 Comparison between integration results obtained by Portal (Fig. S2a) and Portal ($\lambda_{\text{adv}} = 0$)
1047 (Fig. S2b) confirmed that cells from different datasets could not be well mixed without the
1048 objective function.

1049 **Role of regularizer \mathcal{R}_{cos} .** Regularizer \mathcal{R}_{cos} helps to establish reliable alignment between
1050 different domains. It guides domain translation networks to find correspondence of the same cell
1051 type across domains. To confirm this, we fixed λ_{cos} in (S2) as zero to remove \mathcal{R}_{cos} from Portal,
1052 and we denoted this version of Portal as Portal ($\lambda_{\text{cos}} = 0$). After applying Portal ($\lambda_{\text{cos}} = 0$),
1053 cells from the two datasets were well mixed, however, the obtained alignment between these
1054 datasets was problematic. For example, basal cells from SS2 dataset were incorrectly aligned
1055 with B cells and T cells from 10X dataset (Fig. S1c). In contrast, the standard version of
1056 Portal built the alignment correctly (Fig. S1a). The difference between results obtained by
1057 Portal and Portal ($\lambda_{\text{cos}} = 0$) verified the usefulness of \mathcal{R}_{cos} to establish robust correspondence
1058 between datasets.

1059 **Role of regularizer \mathcal{R}_{LA} .** Regularizer \mathcal{R}_{LA} is introduced to impose the consistency constraint
1060 for latent representations of cells. It is helpful to remove domain-specific effects in the latent
1061 space. To demonstrate the effectiveness of \mathcal{R}_{LA} , we set $\lambda_{\text{LA}} = 0$. For Portal ($\lambda_{\text{LA}} = 0$), the
1062 learned representation in the latent space showed a poor alignment of two datasets (Fig. S1d).
1063 This result indicated that the learned representation in the latent space would not be a valid
1064 integration result without adopting regularizer \mathcal{R}_{LA} .

1065 **Role of regularizer \mathcal{R}_{AE} .** $\{E_1(\cdot), G_1(\cdot)\}$ and $\{E_2(\cdot), G_2(\cdot)\}$ form two autoencoder structures
1066 in Portal's framework, \mathcal{R}_{AE} is hence introduced for regularizing autoencoder consistency. Here
1067 we set $\lambda_{\text{AE}} = 0$ to evaluate the role of \mathcal{R}_{AE} with Portal ($\lambda_{\text{AE}} = 0$). Comparison between results
1068 obtained by Portal (Fig. S1a) and Portal ($\lambda_{\text{AE}} = 0$) (Fig. S1e) indicated that \mathcal{R}_{AE} was useful
1069 to improve the accuracy of Portal's results by imposing the consistency between encoders and
1070 generators.

1071 **Role of discriminator score thresholding.** The discriminator score thresholding in Portal
1072 is a tailored design for single-cell integration tasks. With such design, Portal does not force the
1073 alignment of domain-unique cell types, preventing overcorrection of domain-specific effects. To
1074 illustrate the role of the discriminators, we used the same mouse mammary gland data. We
1075 manually removed all basal cells from the 10X dataset and thereby basal cell type became a
1076 domain-unique cell type in the SS2 dataset. We applied standard Portal and Portal without

1077 discriminator score thresholding (denoted as “Portal w/o D score thresholding”) for integration.
 1078 The results in Fig. S2 indicated that, without discriminator score thresholding, Portal could
 1079 not retain the identity of basal cells in the SS2 dataset, and incorrectly aligned them with T
 1080 cells and B cells in the 10X dataset.

1081 Supplementary Information: Figures

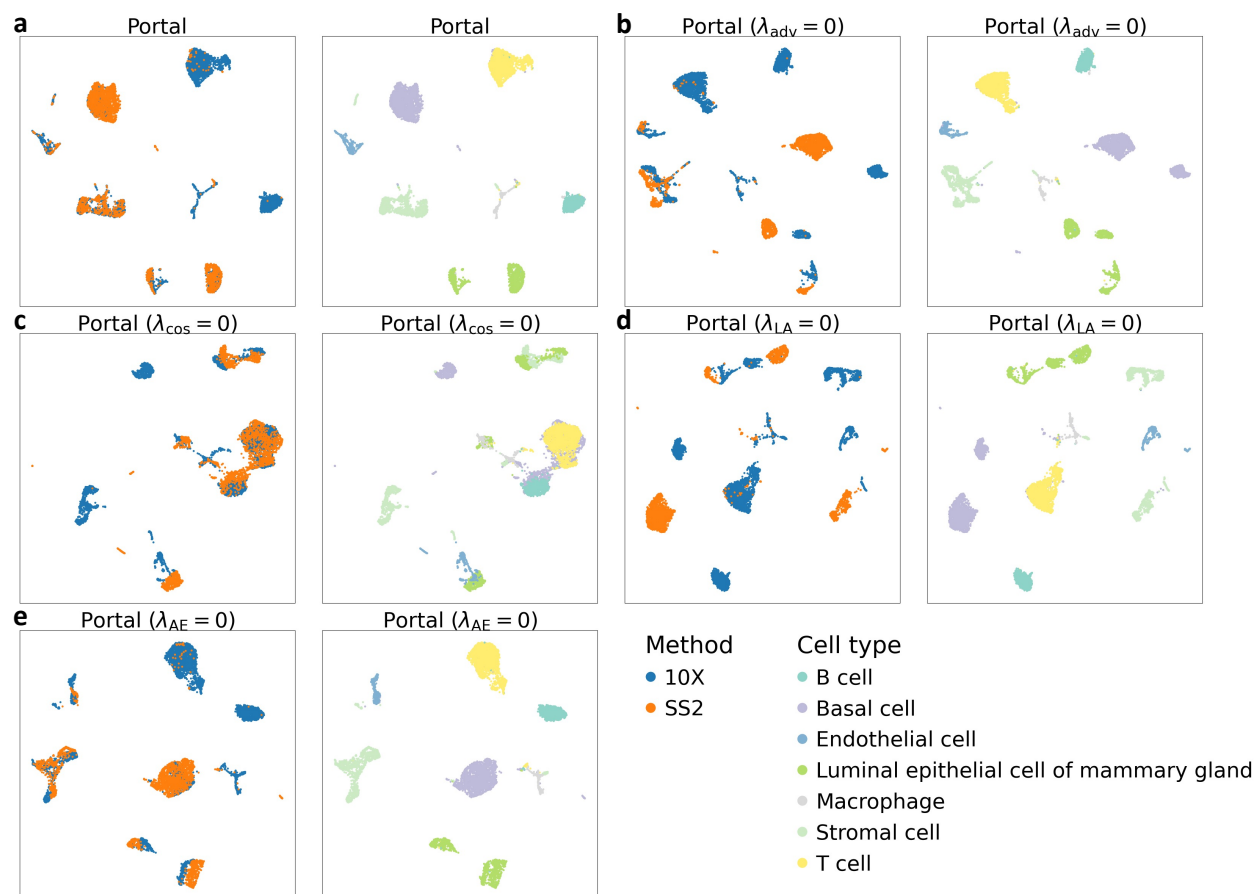


Figure S1: **Investigation of the roles of objective function \mathcal{L}_{adv} and three regularizers \mathcal{R}_{cos} , \mathcal{R}_{LA} and \mathcal{R}_{AE} in Portal.** We used the mouse mammary gland scRNA-seq datasets from Tabluma Muris Consortium in this study. **a.** We applied Portal to integrate the two datasets as a baseline. **b-e.** Then we fixed λ_{adv} , λ_{cos} , λ_{LA} , λ_{AE} in (S2) at zero to evaluate the effectiveness of \mathcal{L}_{adv} , \mathcal{R}_{cos} , \mathcal{R}_{LA} and \mathcal{R}_{AE} , respectively. Clearly, each component of Poral plays its important role in data integration.

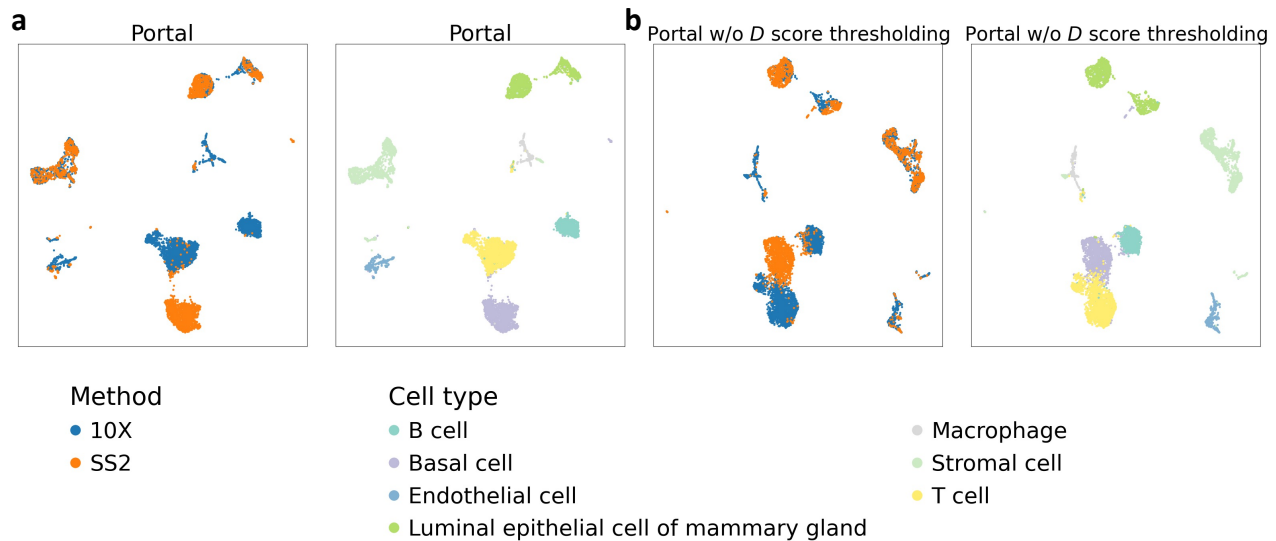


Figure S2: **Investigation of the role of discriminator score thresholding in Portal.** We used the same mouse mammary gland data from Tabluma Muris Consortium, and removed all basal cells from the 10X dataset to make basal cell a domain-unique cell type in the SS2 dataset. **a.** We applied Portal to integrating the two datasets as a baseline. **b.** We removed discriminator score thresholding in Portal to integrate the two datasets.



Figure S3: **Comparison of integration methods based on mouse bladder data.** **a.** UMAP plots colored by profiling methods and cell types. **b.** Alignment (kBET) and cluster preservation performance (ARI and ASW) evaluated using the mouse bladder data.

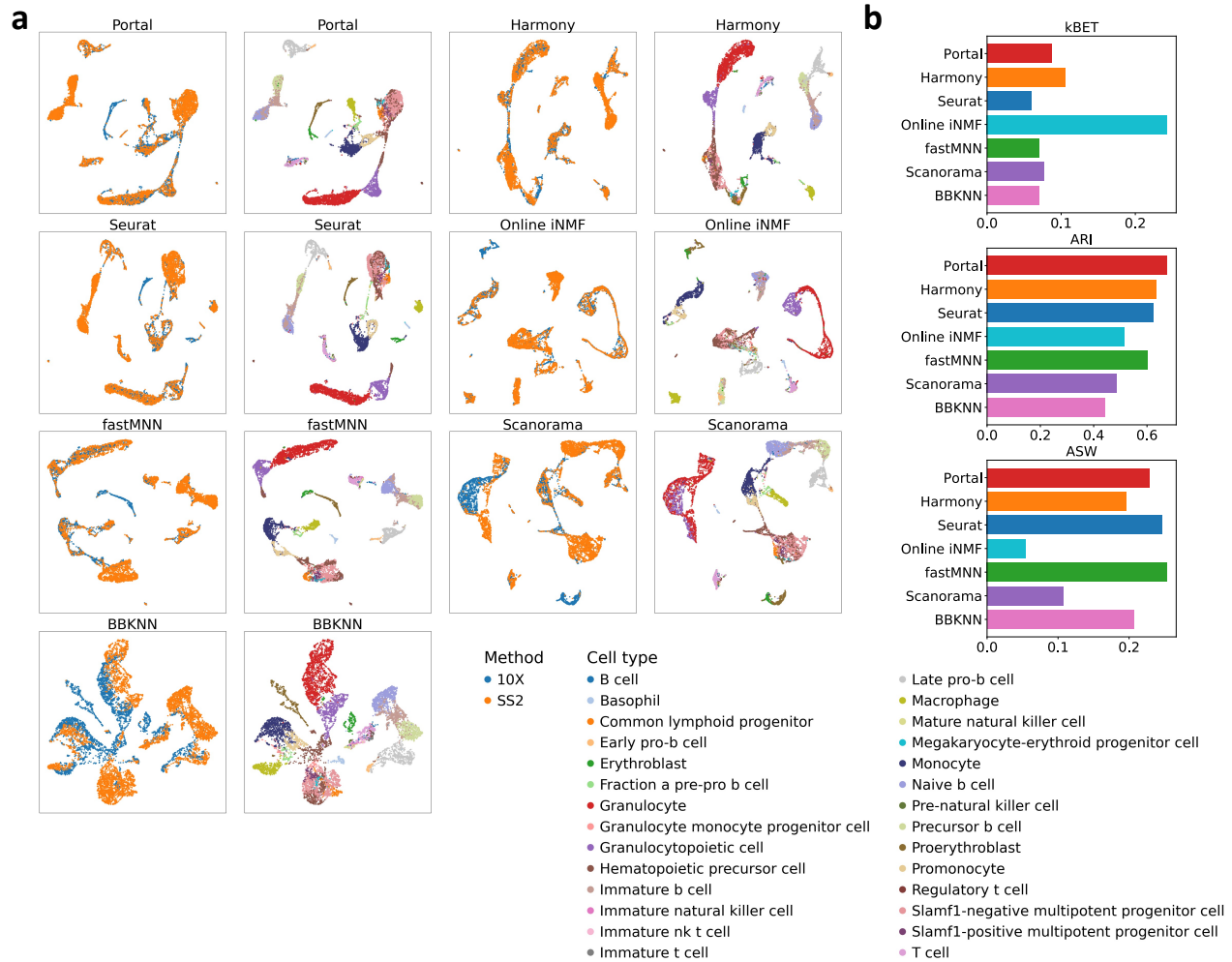


Figure S4: **Comparison of integration methods based on mouse marrow data.** **a.** UMAP plots colored by methods and cell types. **b.** Alignment (kBET) and cluster preservation performance (ARI and ASW) evaluated using the mouse marrow data.

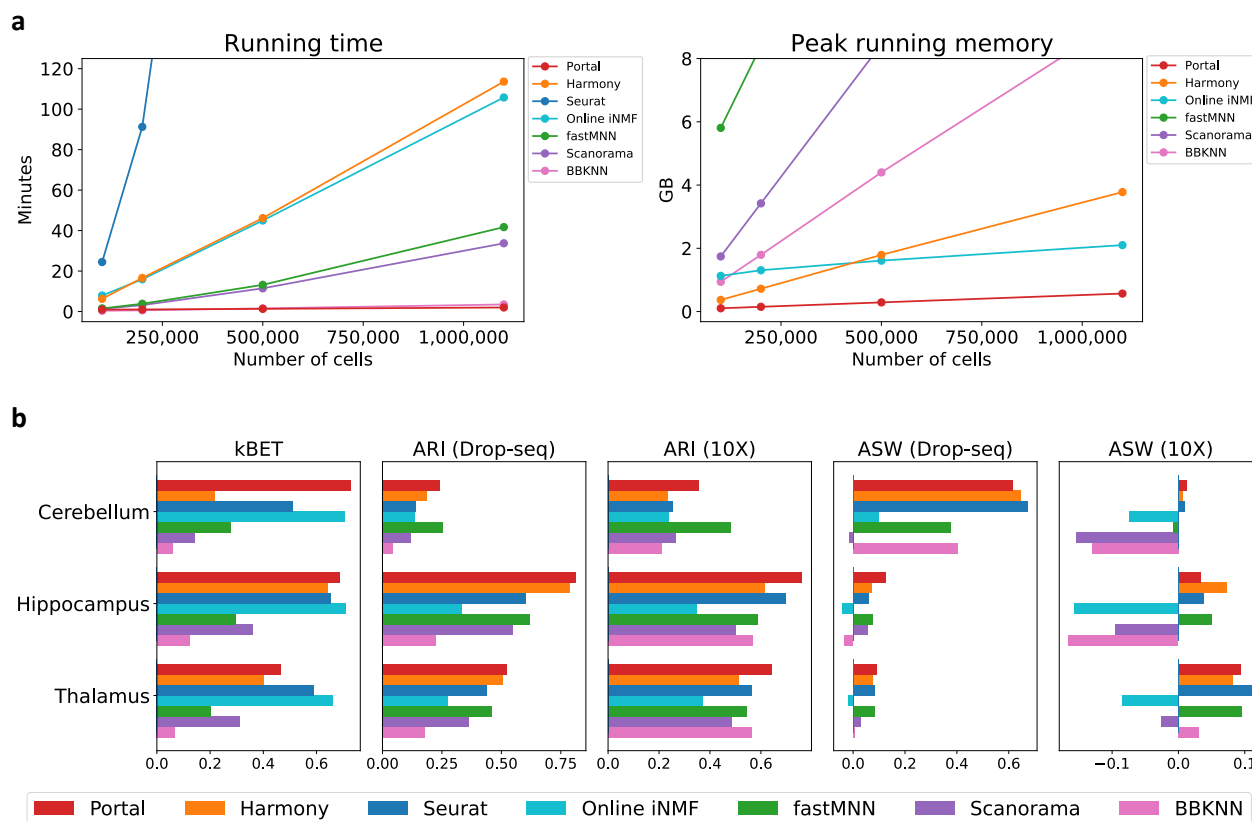


Figure S5: **Benchmark of Portal, Harmony, Seurat, online iNMF, fastMNN, Scanorama and BBKNN.** **a.** We evaluated running time and memory required by all compared methods. Datasets with a total sample size $n = 100,000, 250,000, 50,000,$ and $1,100,167$ were sampled from two mouse brain atlas datasets. Considering running time and peak running memory usage, Portal was the most efficient method. Since comparison among Portal, Harmony, Seurat and online iNMF have been discussed in the main text, here we focus on investigating the performance of fastMNN, Scanorama and BBKNN. Among all compared methods, Portal and BBKNN were remarkably faster than other methods. However, BBKNN required much more memory usage than Portal as sample size increased. More importantly, BBKNN often provided less satisfactory integration performance as indicated by UMAP plots and quantitative metrics in Figs. 2, S3, S4, and (b). Similar to BBKNN, the two methods Scanorama and fastMNN also showed their comparatively limited performance compared to that of Portal, Harmony, Seurat and online iNMF. These two methods showed similar pattern of time and memory usage. Specifically, running Scanorama, fastMNN, BBKNN on full datasets with $1,100,167$ cells required $33.7, 41.7, 3.5$ minutes, and $15.6, 57.8, 9.3$ GB respectively. As a comparison, running Portal used 2.0 minutes and 0.57 GB in the same experiment. Seurat required 24.5 GB on the datasets with $100,000$ cells, so we did not include it in the comparison of peak running memory for clarity. **b.** Alignment (kBET) and cluster preservation performance (ARI and ASW) evaluated using datasets of three shared tissues, including cerebellum, hippocampus and thalamus, in two mouse brain atlas projects. Consistent with previous benchmarking results, fastMNN, Scanorama and BBKNN presented less accurate alignment results, indicated by low kBET, ARI and ASW scores.

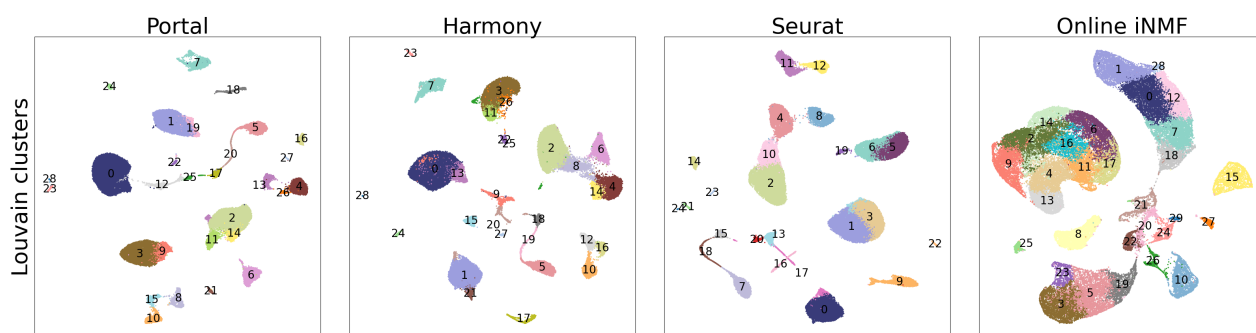


Figure S6: Clusters identified by applying the Louvain algorithm to cell embeddings obtained by Portal, Harmony, Seurat and online iNMF after integration. With default resolution setting, Louvain algorithm detected 29 (Portal), 29 (Harmony), 25 (Seurat), 30 (online iNMF) clusters as shown in UMAP plots. The UMAP plots were drawn separately and colored by clusters identified in the cell embedding space of each method, respectively.

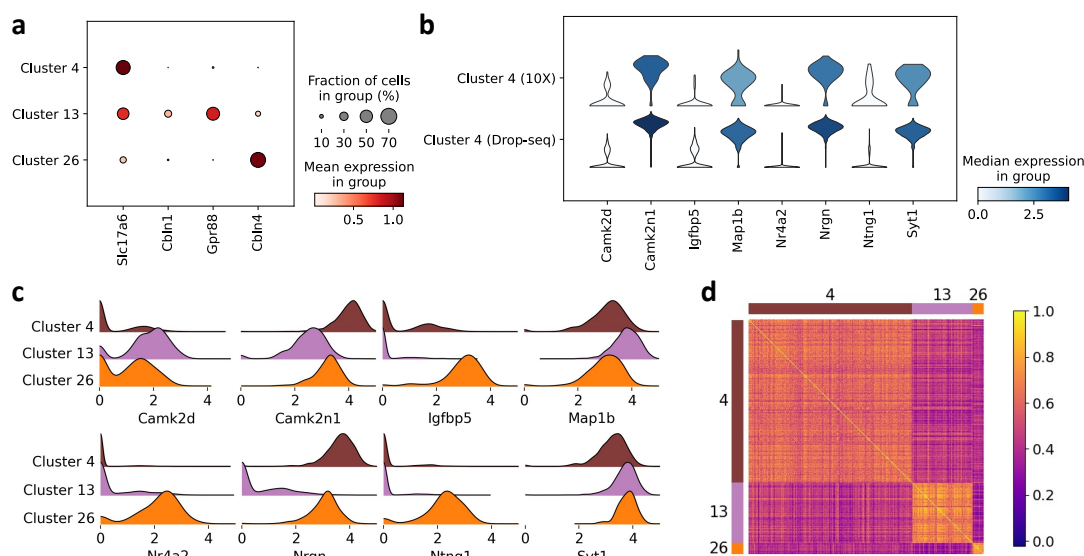


Figure S7: Detailed verification of Portal's integration result on hippocampus datasets at transcriptome level. **a.** We confirmed Portal's alignments of cluster 4, 13, 26 and the three neuron subpopulations by investigating the pattern of marker genes. **b.** The integration result from Portal was validated by the consistent pattern of differentially expressed genes across distinct clusters. Here we only investigated into cluster 4 as cells from the 10X dataset only concentrated in cluster 4 in the marked region. **c.** We identified eight genes that showed distinct expression patterns across the three clusters. **d.** Transcriptional difference among the three clusters was further reflected by examining the correlation between cells using more genes.

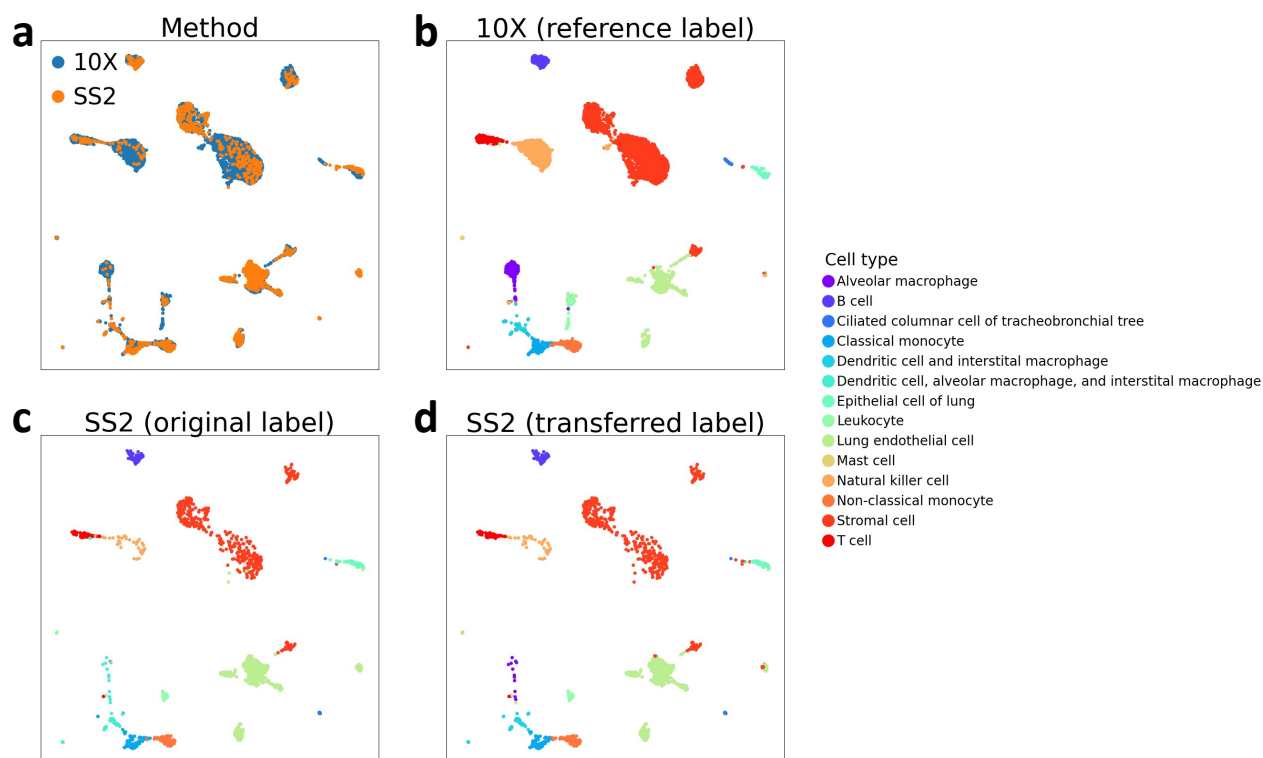


Figure S8: Identification of rare subpopulations in mouse lung scRNA-seq data via label transfer. Utilizing Portal's integration result (**a**), we transferred annotations from the 10X dataset (**b**) to the SS2 dataset (**d**). Portal's integration helped to identify fine-grained subpopulation alveolar macrophage (**d**), which was not identified in its original labels (**c**). **a**. UMAP plot of Portal's integration result colored by profiling methods. **b**, **c**. UMAP plots of integrated 10X, SS2 data colored by cell types obtained from their original publication [7]. **d**. UMAP plot of integrated SS2 data colored by transferred labels provided by Portal.

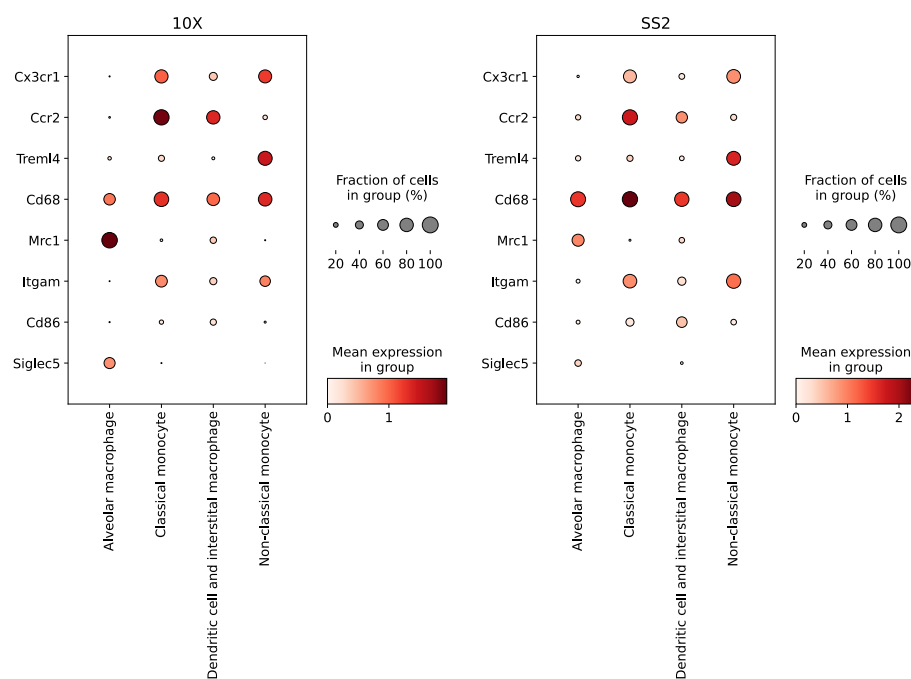
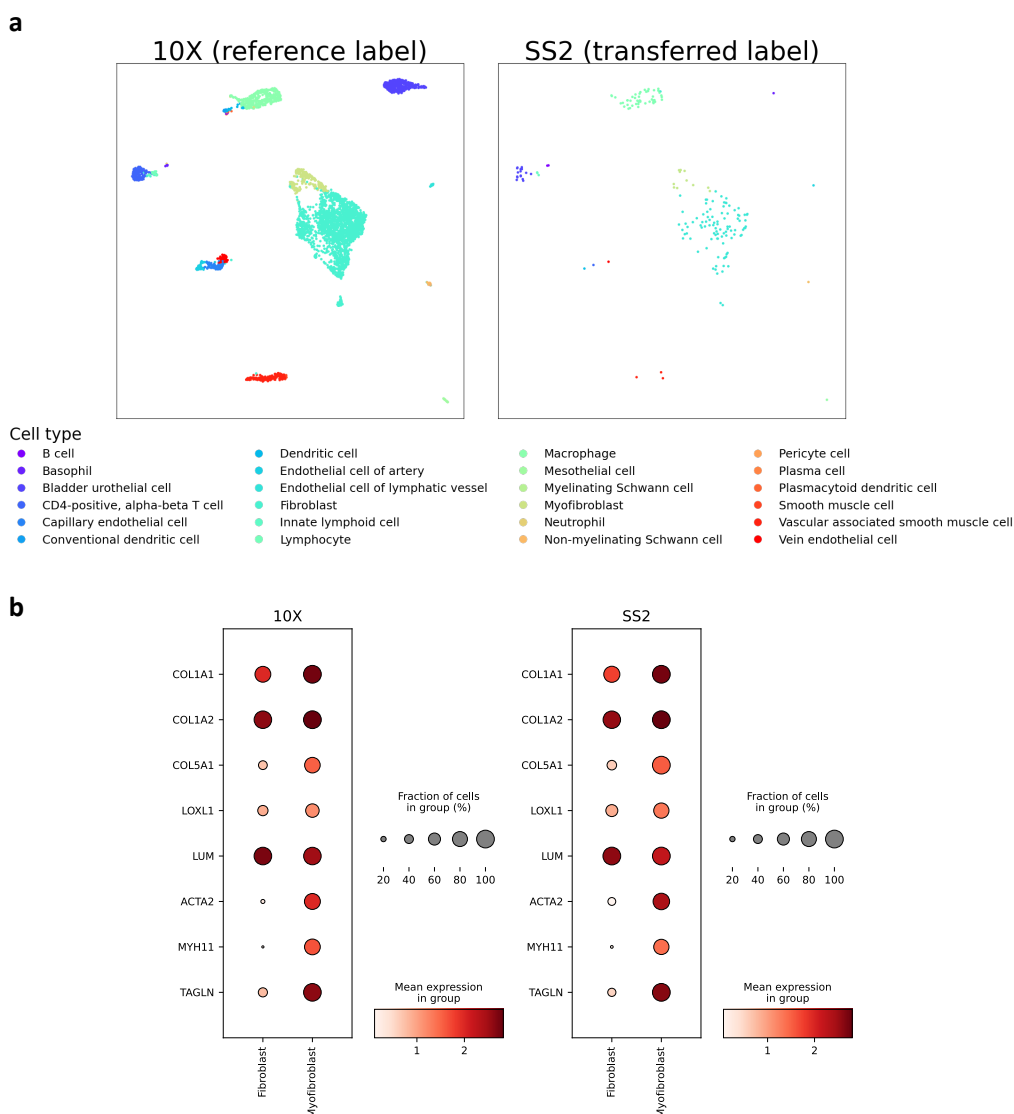


Figure S9: **Marker gene pattern of identified rare subpopulations in mouse lung scRNA-seq data.** By transferring labels from the 10X dataset to SS2 dataset, Portal identified four subpopulations of myeloid cells in SS2 dataset, including alveolar macrophage, dendritic cell and interstitial macrophage, classical monocyte, and non-classical monocyte. To validate the result, we examined four subpopulations' expression levels of marker genes: *Cd68* is a marker of macrophages and monocytes. Between classical monocytes and non-classical monocytes, *Ccr2* is a marker of classical monocytes, *Cx3cr1*, *Trem14* are markers of non-classical monocytes. Between alveolar macrophages and interstitial macrophages, *Mrc1*, *Siglec5* are markers of alveolar macrophages, *Itgam*, *Cd86* are markers of interstitial macrophages.



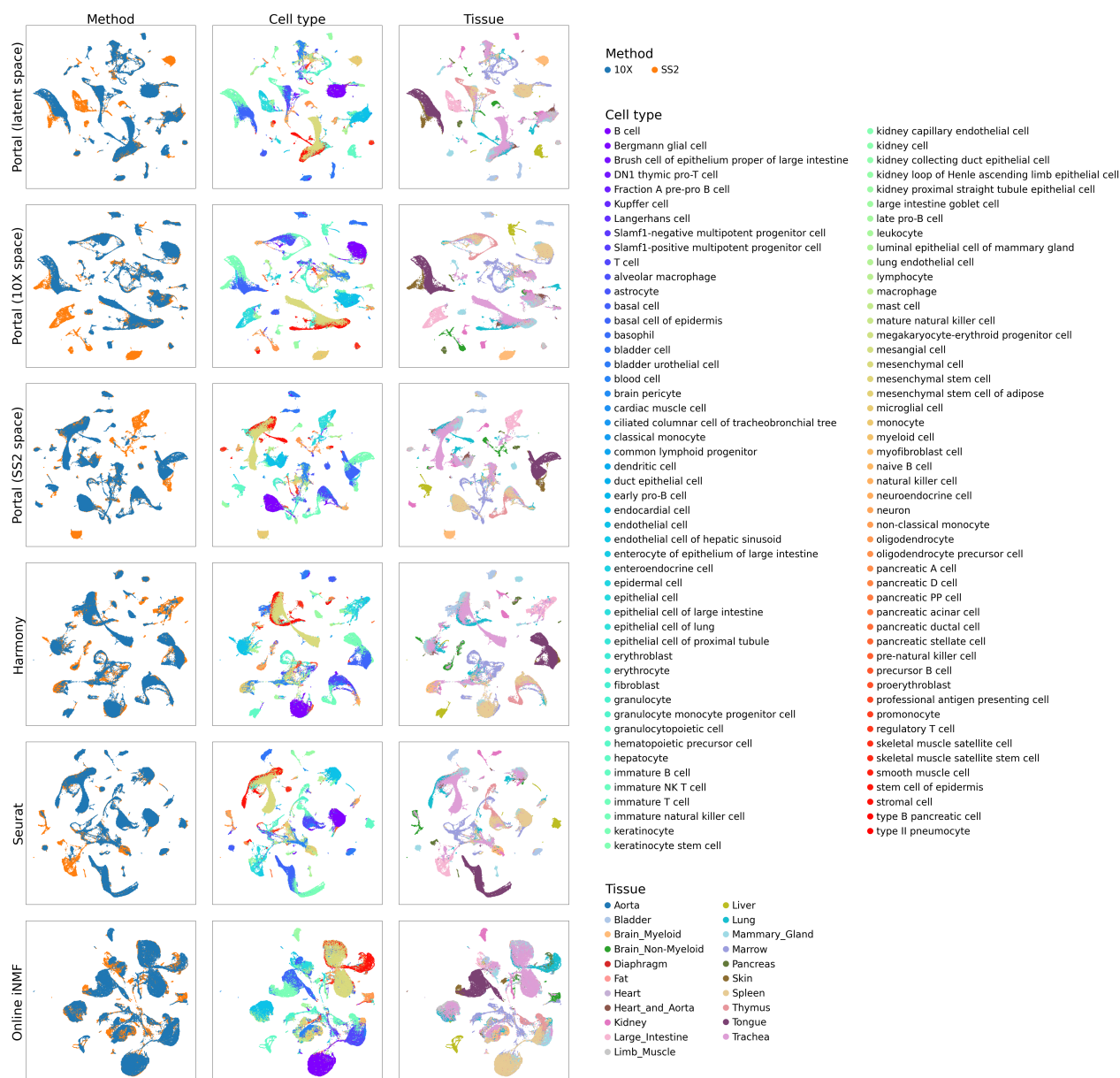


Figure S11: Comparison of the capability of Portal, Harmony, Seurat and online iNMF to construct a comprehensive cell atlas across entire organism. We applied the four integration approaches to harmonize the SS2 dataset and the 10X dataset from the Tabula Muris project, where mouse cells from 20 tissues were profiled. For a comprehensive investigation into Portal's performance, we visualized integration results of Portal in three spaces, namely shared latent space, 10X data space, and SS2 data space. Notably, among the 20 tissues, only 13 of them were included in the 10X data, while all of them were included in the SS2 data. In such a integration task, Portal preserved unique cell types contained in SS2 data, e.g., microglial cells, cell types in large intestine, and pancreatic islets cells. In contrast, Harmony, Seurat and online iNMF provided less accurate results, e.g., they incorrectly mixed microglial cells in brain myeloid with macrophage cells.

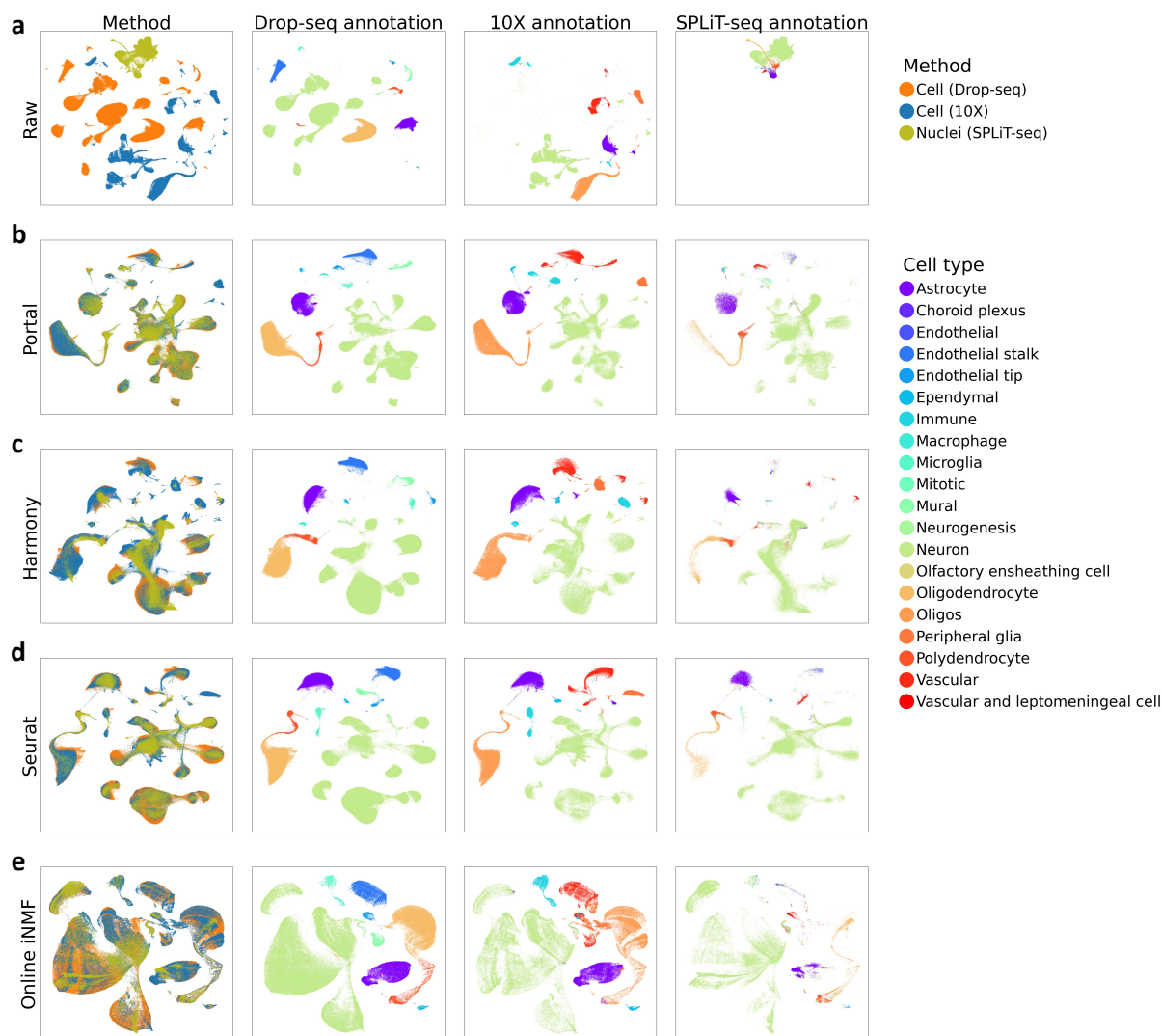


Figure S12: **Comparison of the ability of Portal, Harmony, Seurat and online iNMF to build alignment across one snRNA-seq dataset and two scRNA-seq datasets.** We applied the four integration approaches to align one snRNA-seq dataset profiled by SPLiT-seq [43], and two scRNA-seq datasets profiled by Drop-seq and 10X [8, 9]. We combined the cell type annotations provided by the three datasets together, although they contained slightly different annotations for non-neuron cells, e.g. immune cells and endothelial cells. **a-e**, UMAP visualizations of combined raw data (**a**), integration results of Portal (**b**), Harmony (**c**), Seurat (**d**) and online iNMF (**e**).

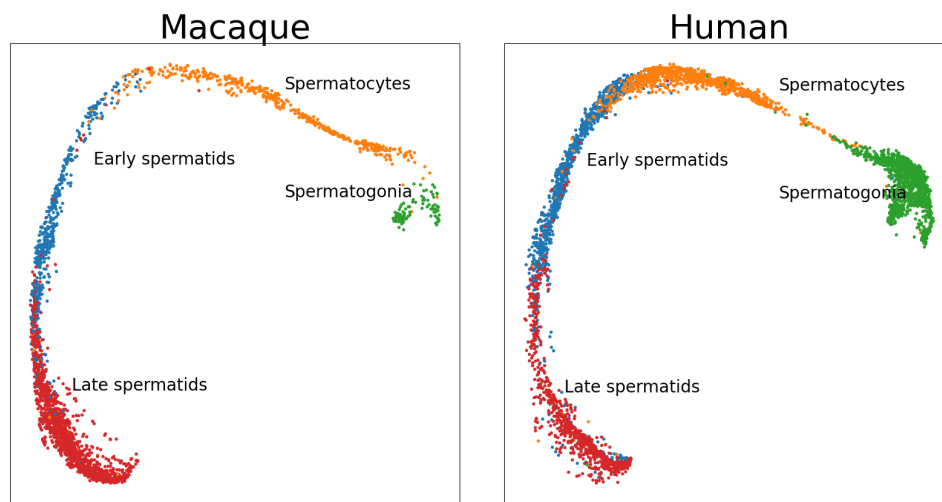


Figure S13: UMAP visualization of the original published annotations for spermatogenesis data of macaque and human [16].

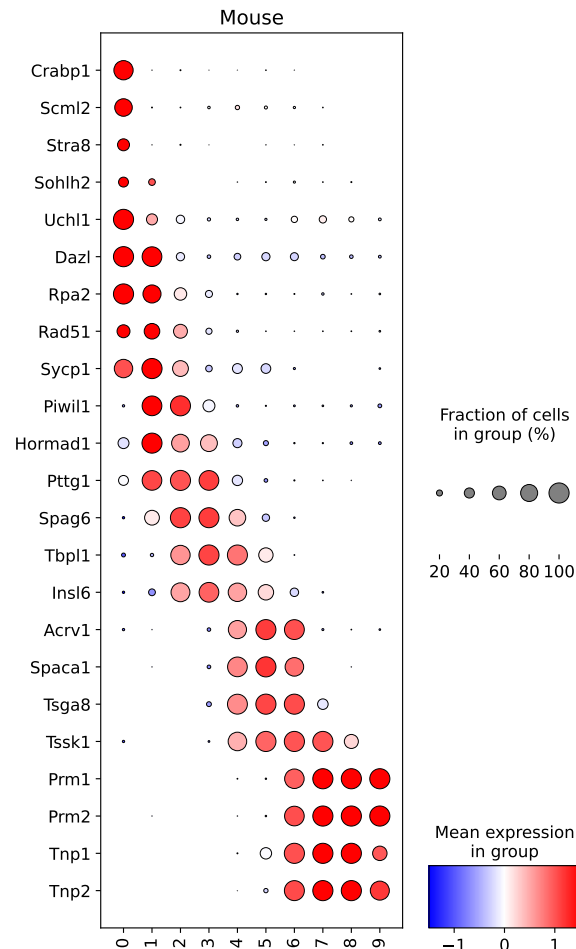


Figure S14: **Marker gene patterns of mouse in Louvain clusters in cross-species integration of spermatogenesis differentiation process.** For mouse, *Crabp1*, *Scml2*, *Stra8*, *Sohlh2*, *Uchl1*, *Dazl*, *Rpa2*, *Rad51*, *Sycp1* are markers of spermatogonia, *Piwil1*, *Hormad1*, *Pttg1*, *Spag6*, *Tbp11*, *Insl6* are markers of spermatocytes, *Acrv1*, *Spaca1*, *Tsga8*, *Tssk1* are markers of early spermatids, and *Prm1*, *Prm2*, *Tnp1*, *Tnp2* are markers of late spermatids.

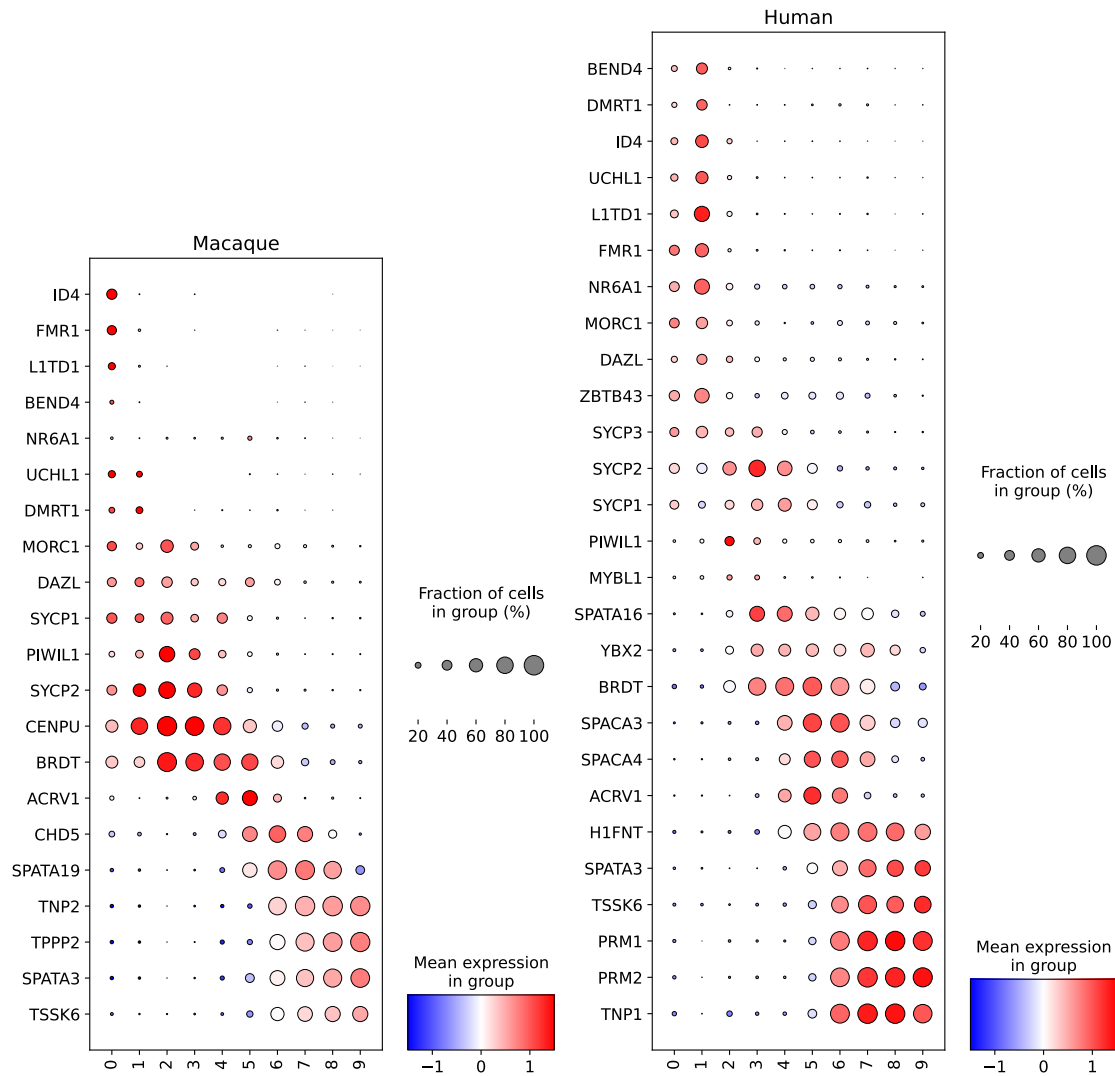


Figure S15: **Marker gene patterns of macaque and human in Louvain clusters in cross-species integration of spermatogenesis differentiation process.** For macaque, *ID4*, *FMR1*, *L1TD1*, *BEND4*, *NR6A1*, *UCHL1*, *DMRT1*, *MORC1*, *DAZL* are markers of spermatogonia, *SYCP1*, *PIWIL1*, *SYCP2*, *CENPU* are markers of spermatocytes, *BRDT*, *ACRV1*, *CHD5*, *SPATA19* are markers of early spermatids, and *TNP2*, *TPPP2*, *SPATA3*, *TSSK6* are markers of late spermatids. For human, *BEND4*, *DMRT1*, *ID4*, *UCHL1*, *L1TD1*, *FMR1*, *NR6A1*, *MORC1*, *DAZL*, *ZBTB43*, *SYCP3* are markers of spermatogonia, *SYCP2*, *SYCP1*, *PIWIL1*, *MYBL1*, *SPATA16*, *YBX2* are markers of spermatocytes, *BRDT*, *SPACA3*, *SPACA4*, *ACRV1*, *H1FNT* are markers of early spermatids, and *TSSK6*, *PRM1*, *PRM2*, *TNP1*, *SPATA3* are markers of late spermatids. The marker gene patterns validated the label transfer results given by Portal.

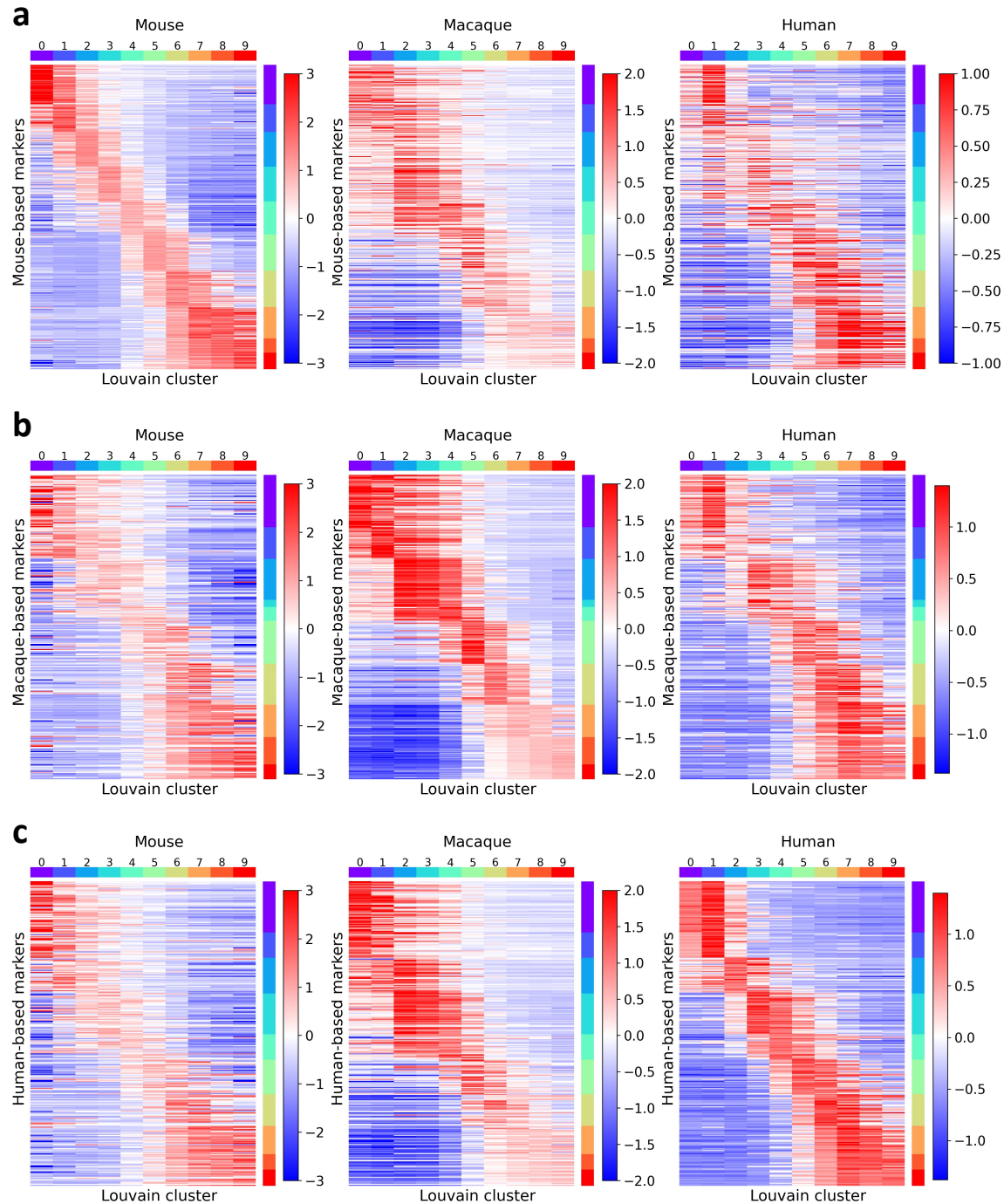


Figure S16: **Gene expression heatmaps in Louvain clusters in cross-species integration of spermatogenesis differentiation process.** For each species, we selected highly-expressed genes for each cluster and combined them together. Gene expression patterns on genes selected based on mouse (a), macaque (b) and human (c) showed connection and distinction among spermatogenesis differentiation processes of different species.