

1 Cell BLAST: Searching large-scale scRNA-seq databases via 2 unbiased cell embedding

3 Zhi-Jie Cao¹, Lin Wei^{1,2}, Shen Lu¹, De-Chang Yang¹, Ge Gao^{1,*}

4 ¹ Biomedical Pioneering Innovation Center (BIOPIC), Beijing Advanced Innovation Center for
5 Genomics (ICG), Center for Bioinformatics (CBI), and State Key Laboratory of Protein and Plant Gene
6 Research at School of Life Sciences, Peking University, Beijing, 100871, China

7 ² College of Life Sciences, Beijing Normal University, Beijing, 100875, China

8 * To whom correspondence should be addressed. Tel: +86-010-62755206; Email: gaog@mail.cbi.pku.edu.cn

9 Abstract

10 An effective and efficient cell-querying method is critical for integrating existing scRNA-seq
11 data and annotating new data. Herein, we present Cell BLAST, an accurate and robust cell-
12 querying method. Powered by a well-curated reference database and a user-friendly Web
13 server, Cell BLAST (<http://cblast.gao-lab.org>) provides a one-stop solution for real-world
14 scRNA-seq cell querying and annotation.

15 Main Text

16 Technological advances during the past decade have led to rapid accumulation of large-scale
17 single-cell RNA sequencing (scRNA-seq) data. Analogous to biological sequence analysis¹,
18 identifying expression similarity to well-curated references via a cell-querying algorithm is
19 becoming the first step of annotating newly sequenced cells. Tools have been developed to
20 identify similar cells using approximate cosine distance² or LSH Hamming distance^{3,4}
21 calculated from a subset of carefully selected genes. Such an intuitive approach is efficient,
22 especially for large-scale data, but may suffer from nonbiological variation across datasets
23 (batch effect^{5,6}). Meanwhile, multiple data harmonization methods have been proposed to
24 remove such confounding factors during alignment, for example, via warping canonical
25 correlation vectors⁷ or matching mutual nearest neighbors across batches⁶. While these
26 methods can be applied to align multiple reference datasets, computation-intensive
27 realignment is required to map query cells to the (pre-)aligned reference data space.

28

29 To address these challenges, we introduce a new customized deep generative model together
30 with a novel cell-to-cell similarity metric specifically designed for cell querying (**Fig. 1a**,
31 **Method**). Differing from canonical variational autoencoder (VAE) models⁸⁻¹¹, adversarial
32 batch alignment is applied to correct batch effect during low-dimensional embedding of
33 reference datasets. Such a design also enables a special “online tuning” mode that can handle
34 batch effect between query and reference data when necessary. Moreover, by exploiting the
35 model’s universal approximator posterior to model uncertainty in latent space, we implement
36 a distribution-based metric to measure cell-to-cell similarity. Finally, we also provide a well-
37 curated multispecies single-cell transcriptomics database (ACA) and an easy-to-use Web
38 interface for convenient exploratory analysis.

39

40 To assess our model’s capability of capturing biological similarity in the low-dimensional
41 latent space, we first benchmarked against several popular dimension reduction tools^{8,12,13}
42 using real-world data (**Supplementary Table 1**) and found that our model is overall among
43 the best performing methods (**Supplementary Fig. 1-2**). We further compared batch effect
44 correction performance using combinations of multiple datasets with overlapping cell types
45 profiled (**Supplementary Table 1**). Our model achieves significantly better dataset mixing
46 (**Fig. 1b**) while maintaining comparable cell type resolution (**Fig. 1c**). Latent space

47 visualization also demonstrates that our model can effectively remove batch effect for
48 multiple datasets with a considerable difference in cell type distribution (**Supplementary**
49 **Fig. 3**). Notably, we found that the correction of inter-dataset batch effect does not
50 automatically generalize to that within each dataset, which is most evident in the pancreatic
51 datasets (**Supplementary Fig. 3c-d, Supplementary Fig. 4a-c**). For such complex scenarios,
52 our model is effective in removing multiple levels of batch effect simultaneously
53 (**Supplementary Fig. 4d-h**).

54

55 While the unbiased latent space embedding derived by the nonlinear deep neural network
56 effectively removes confounding factors, the network's random components and nonconvex
57 optimization procedure also lead to serious challenges, especially false-positive hits when
58 cells outside reference types are provided as query. Thus, we propose a novel posterior
59 distribution-based cell-to-cell similarity metric in the latent space, which we term
60 “normalized projection distance” (NPD). Distance metric ROC analysis (**Method**) shows that
61 our posterior NPD metric is more accurate and robust than Euclidean distance which is
62 commonly used in other neural network-based embedding tools (**Fig. 1d, Supplementary**
63 **Fig. 4k**). Additionally, we exploit the stability of query-hit distance across multiple models to
64 improve specificity (**Method, Supplementary Fig. 4l**). An empirical p-value is computed for
65 each query hit as a measure of “confidence” by comparing the posterior distance to the
66 empirical NULL distribution obtained from randomly selected pairs of cells in the queried
67 data.

68

69 The high specificity of Cell BLAST is especially important for discovering novel cell types.
70 Two recent studies (“Montoro”¹⁴ and “Plasschaert”¹⁵) independently reported a rare tracheal
71 cell type named pulmonary ionocyte. We artificially removed ionocytes from the “Montoro”
72 dataset and used it as a reference to annotate query cells from the “Plasschaert” dataset. In
73 addition to accurately annotating 95.9% of query cells, Cell BLAST correctly rejected 12 of
74 19 “Plasschaert” ionocytes (**Fig. 1e**). Moreover, it highlights the existence of a putative novel
75 cell type as a well-defined cluster with large p-values among all 156 rejected cells (**Fig. 1f-g**).
76 Further examination shows that this cluster actually corresponds to ionocytes
77 (**Supplementary Fig. 6a**; also see **Supplementary Fig. 5** for more detailed analysis on the
78 remaining 7 ionocytes). By contrast, scmap-cell² only rejected 7 “Plasschaert” ionocytes
79 despite the higher overall rejection number of 401 (i.e., more false negatives;
80 **Supplementary Fig. 6b-e**).

81

82 We further systematically compared the performance of query-based cell typing with scmap-
83 cell² and CellFishing.jl⁴ (**Method**) using four groups of datasets, each including both positive
84 control and negative control queries (first 4 groups in **Supplementary Table 2**). Of interest,
85 while Cell BLAST shows superior performance than scmap-cell and CellFishing.jl under the
86 default setting (**Supplementary Fig. 7a-c, 8-10**), detailed ROC analysis reveals that the
87 performance of scmap-cell could be further improved to a level comparable to Cell BLAST
88 by employing higher thresholds, while ROC and optimal thresholds of CellFishing.jl show
89 large variation across different datasets (**Supplementary Fig. 7d**). Cell BLAST presents the
90 most robust performance with a default threshold (p-value < 0.05) across different datasets,
91 which will significantly benefit real-world application. Additionally, we assessed their
92 scalability using reference data varying from 1,000 to 1,000,000 cells. Both Cell BLAST and
93 CellFishing.jl scale well with increasing reference size, while scmap-cell's querying time
94 rises dramatically for larger reference datasets with more than 10,000 cells (**Supplementary**
95 **Fig. 7e**).

96

97 Moreover, our deep generative model combined with posterior-based latent-space similarity
98 metric enables Cell BLAST to model the continuous spectrum of cell states accurately. We
99 demonstrate this using a study profiling mouse hematopoietic progenitor cells ("Tusi"¹⁶) in
100 which computationally inferred cell fate distributions are available. For the purpose of
101 evaluation, cell fate distributions inferred by the authors are recognized as ground truth. We
102 selected cells from one sequencing run as query and the other as reference to test whether we
103 can accurately transfer continuous cell fate between experimental batches (**Fig. 2a-b**).
104 Jensen-Shannon divergence between predicted cell fate distributions and ground truth shows
105 that our prediction is again more accurate than scmap (**Fig. 2c**).

106

107 Besides batch effect among different reference datasets, *bona fide* biological similarity could
108 also be confounded by large, undesirable bias between query and reference data. Exploiting
109 the dedicated adversarial batch alignment, we implemented a particular "online tuning" mode
110 to handle such an often-neglected confounding factor. Briefly, the combination of reference
111 and query data is used to fine-tune the existing reference-based model, with the query-
112 reference batch effect added as an additional component to be removed by adversarial batch
113 alignment (**Method**). Using this strategy, we successfully transferred cell fate from the above
114 "Tusi" dataset to an independent human hematopoietic progenitor dataset ("Velten"¹⁷) (**Fig.**

115 **2d**). The expression of known cell lineage markers validates the rationality of transferred cell
116 fates (**Supplementary Fig. 11a-f**). By contrast, scmap-cell incorrectly assigned most cells to
117 monocyte and granulocyte lineages (**Supplementary Fig. 11g**). As another example, we
118 applied “online tuning” to *Tabula Muris*¹⁸ spleen data, which exhibit significant batch effect
119 between 10x- and Smart-seq2-processed cells. The ROC of Cell BLAST improved
120 significantly after “online tuning”, achieving high specificity, sensitivity and Cohen’s κ (a
121 measure of prediction accuracy corrected for chance, see **Methods** for more details)² at the
122 default cutoff (**Supplementary Fig. 11h**, last group in **Supplementary Table 2**).

123
124 A comprehensive and well-curated reference database is crucial for the practical application
125 of Cell BLAST. Based on public scRNA-seq datasets, we curated ACA, a high-quality
126 reference database. With 986,305 cells in total, ACA currently covers 27 distinct organs
127 across 8 species, offering the most comprehensive compendium for diverse species and
128 organs (**Fig. 2e**, **Supplementary Fig. 12a-b**, **Supplementary Table 3**). To ensure a unified
129 and high-resolution cell type description, all records in ACA are collected and annotated
130 using a standard procedure (**Method**), with 98.9% of datasets manually curated with Cell
131 Ontology, a structured controlled vocabulary for cell types. We trained our model on all ACA
132 datasets. Notably, we found that the model works well in most cases with minimal
133 hyperparameter tuning (latent space visualizations, self-projection coverage and accuracy
134 available on our website, **Supplementary Fig. 12e**).

135
136 A user-friendly Web server is publicly accessible at <http://cblast.gao-lab.org>, with all curated
137 datasets and pretrained models available. Based on the wealth of resources, our website
138 provides “off-the-shelf” querying service. Users can obtain querying hits and visualize cell
139 type predictions with minimal effort (**Supplementary Fig. 12c-d**). For advanced users, a
140 well-documented Python package implementing the Cell BLAST toolkit is also available,
141 which enables model training on custom references and diverse downstream analyses.

142
143 By explicitly modeling multilevel batch effect as well as uncertainty in cell-to-cell similarity
144 estimation, Cell BLAST is an accurate and robust querying algorithm for heterogeneous
145 single-cell transcriptome datasets. In combination with a comprehensive, well-annotated
146 database and an easy-to-use Web interface, Cell BLAST provides a one-stop solution for
147 both bench biologists and bioinformaticians.

148 **Software availability**

149 The full package of Cell BLAST is available at <http://cblast.gao-lab.org>. Code necessary to
150 reproduce results in the paper is deposited at https://github.com/gao-lab/Cell_BLAST and
151 https://github.com/gao-lab/Cell_BLAST-notebooks.

152 **Acknowledgments**

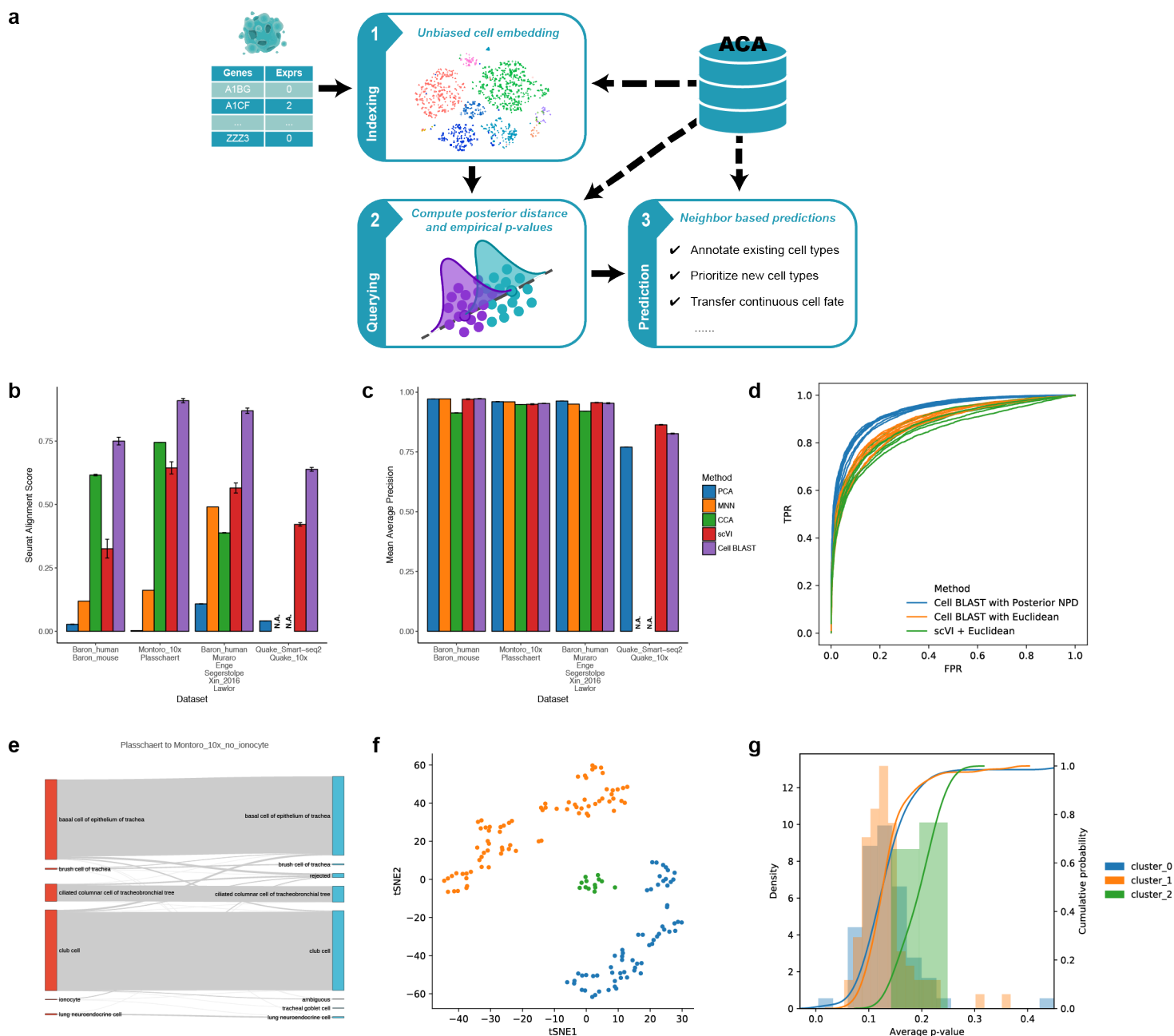
153 The authors thank Drs. Zemin Zhang, Cheng Li, Letian Tao, Jian Lu and Liping Wei at
154 Peking University for their helpful comments and suggestions during the study.

155 This work was supported by funds from the National Key Research and Development
156 Program (2016YFC0901603), the China 863 Program (2015AA020108), as well as the State
157 Key Laboratory of Protein and Plant Gene Research and the Beijing Advanced Innovation
158 Center for Genomics (ICG) at Peking University. The research of G.G. was supported in part
159 by the National Program for Support of Top-notch Young Professionals.

160 Part of the analysis was performed on the Computing Platform of the Center for Life
161 Sciences of Peking University and supported by the High-performance Computing Platform
162 of Peking University.

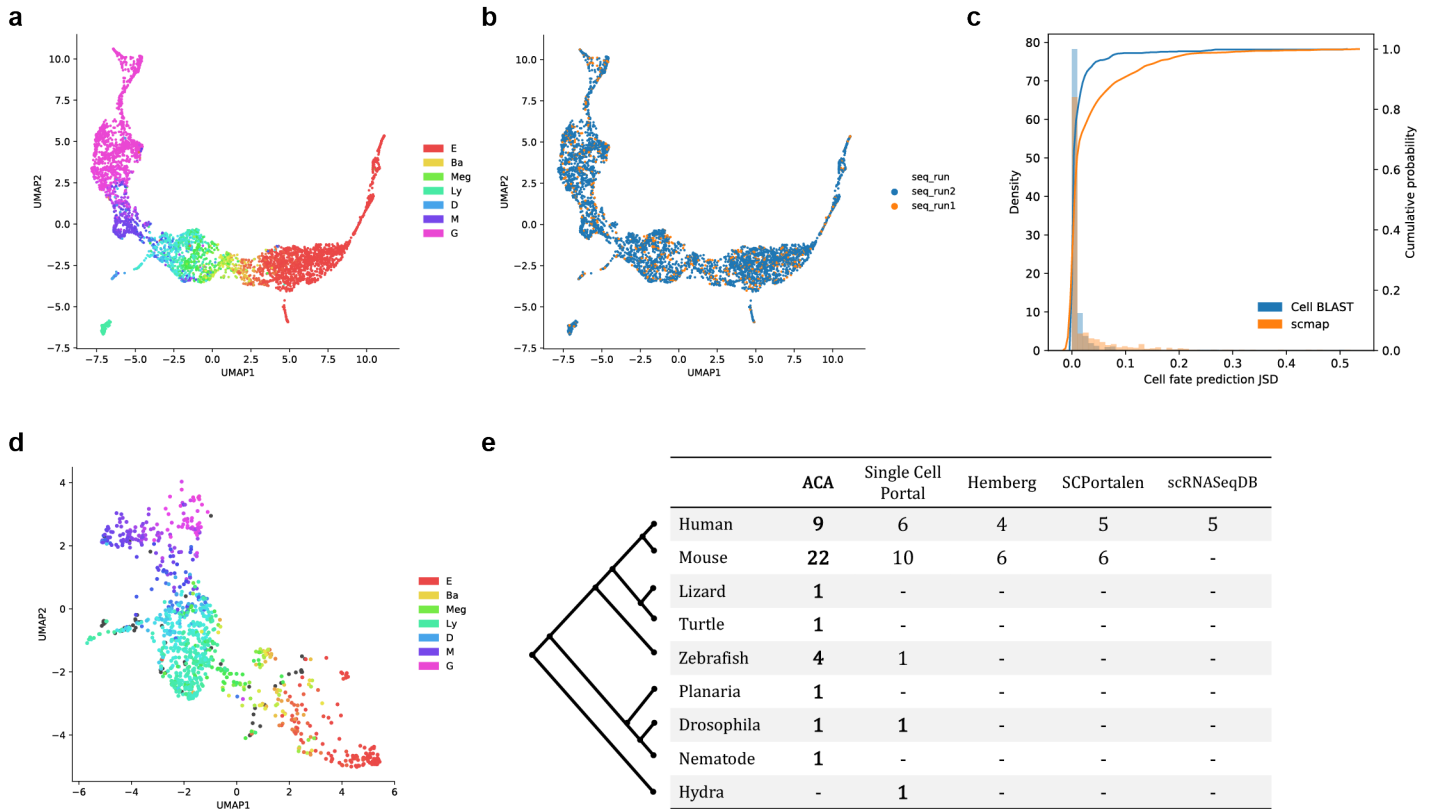
163 **Author contributions**

164 G.G. conceived the study and supervised the research; Z.J.C. and L.W. contributed to the
165 computational framework and data curation; S. L., Z.J.C., and D.C.Y designed, implemented
166 and deployed the website; Z.J.C. and G.G. wrote the manuscript with comments and inputs
167 from all coauthors.



168 **Figure. 1 Cell BLAST benchmarking and application to trachea datasets.**

169 (a) Overall Cell BLAST workflow. (b) Extent of dataset mixing after batch effect correction in four groups of
 170 datasets, quantified by the Seurat alignment score. A high Seurat alignment score indicates that local
 171 neighborhoods consist of cells from different datasets uniformly rather than from the same dataset only. Error
 172 bars indicate mean \pm s.d. Methods that did not finish under the 2-hour time limit are marked as N.A. (c) Cell
 173 type resolution after batch effect correction, quantified by cell type mean average precision (MAP). MAP can be
 174 thought of as a generalization to nearest neighbor accuracy, with larger values indicating higher cell type
 175 resolution, thus more suitable for cell querying. Error bars indicate mean \pm s.d. Methods that did not finish
 176 under the 2-hour time limit are marked as N.A. (d) ROC curve of different distance metrics in discriminating
 177 cell pairs with the same cell type from cell pairs with different cell types. (e) Sankey plot comparing Cell
 178 BLAST predictions and original cell type annotations for the “Plasschaert” dataset. (f) t-SNE visualization of
 179 Cell BLAST-rejected cells, colored by unsupervised clustering. (g) Average p-value distribution of each cluster
 180 in (f).



181 **Figure. 2 Application to hematopoietic progenitor datasets.**

182 (a, b) UMAP visualization of latent space learned on the “Tusi” dataset, colored by sequencing run (a) and cell
 183 cell fate (b). The model is trained solely on cells from run 2 and used to project cells from run 1. Each of the seven
 184 terminal cell fates (E, erythroid; Ba, basophilic or mast; Meg, megakaryocytic; Ly, lymphocytic; D, dendritic;
 185 M, monocytic; G, granulocytic neutrophil) is assigned a distinct color. The color of each single cell is then
 186 determined by the linear combination of these seven colors in hue space, weighed by the cell fate distribution
 187 among these terminal fates. (c) Distribution of Jensen-Shannon divergence between predicted cell fate
 188 distributions and author-provided “ground truth”. (d) UMAP visualization of the “Velten” dataset, colored by
 189 Cell BLAST-predicted cell fates. (e) Number of organs covered in each species for different single-cell
 190 transcriptomics databases, including the Single Cell Portal (https://portals.broadinstitute.org/single_cell),
 191 Hemberg collection², SCPortalen¹⁹, and scRNASeqDB²⁰.

192 Methods

193 The deep generative model

194 The model we used is based on the adversarial autoencoder (AAE)²¹. Below, we denote the
 195 gene expression profile of a cell as $\mathbf{x} \in \mathbb{R}^G$, where G is the number of genes. The data
 196 generative process is modeled by a continuous latent variable $\mathbf{z} \in \mathbb{R}^D$ ($D \ll G$) with standard
 197 Gaussian prior $\mathbf{z} \sim N(\mathbf{0}, I_D)$ which models continuous cell states, as well as a one-hot latent
 198 variable $\mathbf{c} \in \{0,1\}^K$, $\mathbf{c}^T \mathbf{c} = 1$ with categorical prior $\mathbf{c} \sim \text{Cat}(K)$ which aims to model cell
 199 type clusters. A unified latent vector is then determined by $\mathbf{l} = \mathbf{z} + H\mathbf{c}$, where $H \in \mathbb{R}^{D \times K}$. A
 200 neural network (decoder, denoted by Dec below) maps the cell embedding vector \mathbf{l} to two
 201 parameters of the negative binomial distribution $\boldsymbol{\mu}, \boldsymbol{\theta} = \text{Dec}(\mathbf{l})$ that models the distribution of
 202 expression profile \mathbf{x} :

$$p(\mathbf{x}|\mathbf{z}, \mathbf{c}; \text{Dec}, H) = p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\theta}) = \prod_{j=1}^G p(x_j|\mu_j, \theta_j) \quad (1)$$

$$p(x_j|\mu_j, \theta_j) = \frac{\Gamma(x_j + \theta_j)}{\Gamma(\theta_j)\Gamma(x_j + 1)} \left(\frac{\mu_j}{\theta_j + \mu_j}\right)^{x_j} \left(\frac{\theta_j}{\theta_j + \mu_j}\right)^{\theta_j} \quad (2)$$

203 where $\boldsymbol{\mu}$ and $\boldsymbol{\theta}$ are the mean and dispersion of the negative binomial distribution,
 204 respectively. Theoretically, the negative binomial model should be fitted on raw count data⁸,
 205 ^{13, 22}. However, for the purpose of cell querying, datasets have to be normalized to minimize
 206 the influence of capture efficiency and sequencing depth. We empirically found that, using
 207 normalized data, the negative binomial model still produced better results than alternative
 208 distributions like the log-normal distribution. To prevent numerical instability during training
 209 caused by normalization that breaks the mean-variance relationship of the negative binomial
 210 model, we additionally included the variance of the dispersion parameter as a regularization
 211 term.

212 Training objectives for the adversarial autoencoder are:

$$\min_{\text{Dec}, \text{Enc}, H} -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}; \text{Enc}), \mathbf{c} \sim q(\mathbf{c}|\mathbf{x}; \text{Enc})} \log p(\mathbf{x}|\mathbf{z}, \mathbf{c}; \text{Dec}, H) + \lambda_z \right. \\ \left. \cdot \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}; \text{Enc})} \log D_z(\mathbf{z}) + \lambda_c \cdot \mathbb{E}_{\mathbf{c} \sim q(\mathbf{c}|\mathbf{x}; \text{Enc})} \log D_c(\mathbf{c}) \right] \quad (3)$$

$$\max_{D_z} \lambda_z \cdot \left(\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \log D_z(\mathbf{z}) + \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}; \text{Enc})} \log(1 - D_z(\mathbf{z})) \right) \quad (4)$$

$$\max_{D_c} \lambda_c \cdot \left(\mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \log D_c(\mathbf{c}) + \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \mathbb{E}_{\mathbf{c} \sim q(\mathbf{c}|\mathbf{x}; \text{Enc})} \log(1 - D_c(\mathbf{c})) \right) \quad (5)$$

213 $q(\mathbf{z}|\mathbf{x}; \text{Enc})$ and $q(\mathbf{c}|\mathbf{x}; \text{Enc})$ are “universal approximator posteriors” parameterized by
214 another neural network (encoder, denoted by Enc). Expectations with regard to $q(\mathbf{z}|\mathbf{x}; \text{Enc})$
215 and $q(\mathbf{c}|\mathbf{x}; \text{Enc})$ are approximated by sampling $\mathbf{x}' \sim \text{Poisson}(\mathbf{x})$ and feeding to the
216 deterministic encoder network. The choice of Poisson noise is arbitrary as the encoder learns
217 to map this arbitrary noise distribution to an appropriate posterior distribution during training.
218 D_z and D_c are discriminator networks for \mathbf{z} and \mathbf{c} , respectively, which output the probability
219 that a latent sample is from the prior rather than from the posterior. Effectively, adversarial
220 training between the encoder (Enc) and discriminators (D_z and D_c) drives the encoder output
221 to match prior distributions of latent variables $p(\mathbf{z})$ and $p(\mathbf{c})$. λ_z and λ_c are hyperparameters
222 that control prior matching strength. The model is much easier and more stable to train than
223 canonical GANs because of the low dimensionality and simple distribution of \mathbf{z} and \mathbf{c} .
224 At convergence, the encoder learns to map the data distribution to latent variables that follow
225 their respective prior distributions, and the decoder learns to map latent variables from prior
226 distributions back to the data distribution. The key element we use for cell querying is vector
227 \mathbf{l} on the decoding path because it defines a unified latent space in which biological
228 similarities are well captured. The model also works if no categorical latent variable is used,
229 in which case $\mathbf{l} = \mathbf{z}$ directly.
230 Some architectural designs are learned from scVI⁸, including logarithm transformation before
231 encoder input, and softmax output scaled by the library size when computing $\boldsymbol{\mu}$. Stochastic
232 gradient descent with minibatches is applied to optimize the loss functions. Specifically, we
233 use the “RMSProp” optimization algorithm with no momentum term to ensure stability of
234 adversarial training. The model is implemented using the Tensorflow²³ Python library.

235 **Adversarial batch alignment**

236 As a natural extension to the prior matching adversarial training strategy described in the
237 previous section, and following recent work in domain adaptation²⁴⁻²⁶, we propose the
238 adversarial batch alignment strategy to align the latent space distribution of different batches.
239 We denote the batch membership of each cell as $\mathbf{b} \in \{0,1\}^B$, $\mathbf{b}^T \mathbf{b} = 1$. The distribution $p(\mathbf{b})$
240 is categorical:

$$p(b_i = 1) = w_i, \quad \sum_{i=1}^B w_i = 1 \quad (6)$$

241 Adversarial batch alignment introduces an additional loss:

$$\min_{\text{Dec, Enc, } H} \mathbb{E}_{\mathbf{b} \sim p(\mathbf{b}), \mathbf{x} \sim p(\mathbf{x}|\mathbf{b})} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}; \text{Enc}), \mathbf{c} \sim q(\mathbf{c}|\mathbf{x}; \text{Enc})} [\mathcal{L}_{\text{base}} + \lambda_b \cdot \mathbf{b}^T \log D_b(\mathbf{l})] \quad (7)$$

$$\max_{D_b} \mathbb{E}_{\mathbf{b} \sim p(\mathbf{b}), \mathbf{x} \sim p(\mathbf{x}|\mathbf{b})} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}; \text{Enc}), \mathbf{c} \sim q(\mathbf{c}|\mathbf{x}; \text{Enc})} [\lambda_b \cdot \mathbf{b}^T \log D_b(\mathbf{l})] \quad (8)$$

242 \mathcal{L}_{base} denotes the loss function in (3). D_b is a multiclass batch discriminator network that
 243 outputs the probability distribution of batch membership based on the embedding vector \mathbf{l} . λ_b
 244 is a hyperparameter controlling batch alignment strength. Additionally, the generative
 245 distribution is extended to condition on \mathbf{b} as well:

$$\begin{aligned} p(\mathbf{x}|\mathbf{z}, \mathbf{c}, \mathbf{b}; \text{Dec}) &= p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\theta}) \\ \boldsymbol{\mu}, \boldsymbol{\theta} &= \text{Dec}(\mathbf{l}, \mathbf{b}) \\ \mathbf{l} &= \mathbf{z} + H\mathbf{c} \end{aligned} \quad (9)$$

246 Below, we focus on batch alignment and discard the first \mathcal{L}_{base} term and scaling parameter
 247 λ_b . We extend the derivation in the original GAN paper²⁷ to show that adversarial batch
 248 alignment converges when embedding space distributions of different batches are aligned.
 249 To simplify notation, we fuse the data distribution and encoder transformation and replace
 250 the minimization over encoder to minimization over batch-embedding distributions:

$$\min_{p(\mathbf{l})} \sum_{i=1}^B w_i \mathbb{E}_{\mathbf{l} \sim p_i(\mathbf{l})} \log D_{b_i}(\mathbf{l}) \quad (10)$$

$$\max_{D_b} \sum_{i=1}^B w_i \mathbb{E}_{\mathbf{l} \sim p_i(\mathbf{l})} \log D_{b_i}(\mathbf{l}) \quad (11)$$

251 Here $D_{b_i}(\mathbf{l})$ denotes the i^{th} dimension of the discriminator output, i.e., the probability that the
 252 discriminator “thinks” a cell is from the i^{th} batch. D_b is assumed to have sufficient capacity,
 253 which is generally reasonable in the case of neural networks. The global optimum of (11) is
 254 reached when D_b outputs optimal batch membership distribution at every \mathbf{l} :

$$\max_{D_{b_i}(\mathbf{l})} w_i p_i(\mathbf{l}) \log D_{b_i}(\mathbf{l}), \quad s. t. \sum_{i=1}^B D_{b_i}(\mathbf{l}) = 1 \quad (12)$$

255 The solution to the above maximization is given by:

$$D_{b_i}^*(\mathbf{l}) = \frac{w_i p_i(\mathbf{l})}{\sum_{i=1}^B w_i p_i(\mathbf{l})} \quad (13)$$

256 Substituting $D_{b_i}^*(\mathbf{l})$ back into (10), we obtain:

$$\begin{aligned} & \sum_{i=1}^B w_i \mathbb{E}_{\mathbf{l} \sim p_i(\mathbf{l})} \log \frac{w_i p_i(\mathbf{l})}{\sum_{i=1}^B w_i p_i(\mathbf{l})} \\ &= \sum_{i=1}^B w_i \mathbb{E}_{\mathbf{l} \sim p_i(\mathbf{l})} \log \frac{p_i(\mathbf{l})}{\sum_{i=1}^B w_i p_i(\mathbf{l})} + \sum_{i=1}^B w_i \mathbb{E}_{\mathbf{l} \sim p_i(\mathbf{l})} \log w_i \\ &= \sum_{i=1}^B w_i \cdot \text{KL} \left(p_i(\mathbf{l}) \parallel \sum_{i=1}^B w_i p_i(\mathbf{l}) \right) + \sum_{i=1}^B w_i \log w_i \\ &\geq \sum_{i=1}^B w_i \log w_i \end{aligned} \tag{14}$$

257 Thus, $\sum_{i=1}^B w_i \log w_i$ is the global minimum, reached if and only if $p_i(\mathbf{l}) = p_j(\mathbf{l}), \forall i, j$. The
258 minimization of (10) is equivalent to minimizing a form of generalized Jensen-Shannon
259 divergence among multiple batch-embedding distributions.

260 Note that in practice, model training balances between \mathcal{L}_{base} and pure batch alignment.
261 Aligning cells of the same type induces a minimal cost in \mathcal{L}_{base} , while improperly aligning
262 cells of different types could cause \mathcal{L}_{base} to rise dramatically. During training, the gradient
263 from both batch discriminators and decoder provide fine-grain guidance to align different
264 batches, leading to better results than “hand-crafted” alignment strategies like CCA⁷ and
265 MNN⁶. Empirically, given proper values for λ_b , the adversarial approach correctly handles
266 difference in cell type distribution among batches. If multiple levels of batch effect exist, e.g.,
267 within-dataset and cross-dataset, we use an independent batch discriminator for each
268 component, providing extra flexibility.

269 **Data preprocessing for benchmarks**

270 Most informative genes were selected using the Seurat⁷ function “FindVariableGenes”. We
271 set the argument “binning.method” to “equal_frequency” and left other arguments as default.
272 If within-dataset batch effect exists, genes are selected independently for each batch and then
273 pooled together. By default, a gene is retained if it is selected in at least 50% of batches.
274 Downstream benchmarks were all performed using this gene set, except for scmap and
275 CellFishing.jl, which provide their own gene selection method. GNU parallel²⁸ was used to
276 parallelize and manage jobs throughout the benchmarking and data processing pipeline.

277 **Benchmarking dimension reduction**

278 PCA was performed using the R package `irlba`²⁹ (v2.3.2). ZIFA¹² was downloaded from its
279 Github repository, and hard coded random seeds were removed to reveal actual stability.
280 ZINB-WaVE¹³ (v1.0.0) was performed using the R package `zinbwave`. `scVI`⁸ (v0.2.3) was
281 downloaded from its Github repository, and minor changes were made to the original code to
282 address PyTorch³⁰ compatibility issues. Our modified versions of ZIFA and `scVI` are
283 available upon request.

284 For PCA and ZIFA, data were logarithm transformed after normalization and adding a
285 pseudocount of 1. Hyperparameters of all methods above were left as default. For our model,
286 we used the same set of hyperparameters throughout all benchmarks. λ_z and λ_c were both set
287 to 0.001. All neural networks (encoder, decoder and discriminators) used a single layer of
288 128 hidden units. Learning rate of the RMSProp optimizer is set to 0.001, and minibatches of
289 size 128 were used. For comparability, the target dimensionality of each method was set to
290 10. All benchmarked methods were repeated multiple times with different random seeds. 4
291 random seeds were used for PCA, ZIFA and ZINB-WaVE, while 16 random seeds were used
292 for `scVI` and our model, since neural network-based models are typically considered less
293 stable. Run time was limited to 2 hours, after which the jobs were terminated.

294 Cell type nearest neighbor mean average precision (MAP) was computed with K nearest
295 neighbors of each cell based on low-dimensional space Euclidean distance. If we denote the
296 cell type of a cell as y , and the cell types of its ordered nearest neighbors as y_1, y_2, \dots, y_k . The
297 average precision (AP) for that cell is defined as:

$$\text{AP} = \frac{\sum_{k=1}^K \mathbf{1}_{y=y_k} \cdot \frac{\sum_{k'=1}^k \mathbf{1}_{y=y_{k'}}}{k}}{\sum_{k=1}^K \mathbf{1}_{y=y_k}} \quad (15)$$

298 Mean average precision is then given by:

$$\text{MAP} = \frac{1}{N} \sum_{i=1}^N \text{AP}_i \quad (16)$$

299 Note that when $K = 1$, MAP reduces to the nearest neighbor accuracy. We set K to 1% of the
300 total cell number throughout all benchmarks.

301 **Benchmarking batch effect correction**

302 We merged multiple datasets according to shared gene names. If datasets to be merged are
303 from different species, Ensembl ortholog³¹ information was used to map genes to ortholog

304 groups before merging. To obtain informative genes in merged datasets, we take the union of
305 informative genes from each dataset, and then intersect the union with the intersection of
306 detected genes from each dataset.

307 CCA⁷ and MNN⁶ alignments were performed using the R packages Seurat⁷ (v2.3.3) and
308 scran³² (v1.6.9), respectively. Hard-coded random seeds in Seurat were removed to reveal
309 actual stability. The modified version of Seurat is available upon request. For comparability,
310 we evaluated cell type resolution and batch mixing in a 10-dimensional latent space. For
311 MNN alignment, we set the argument “cos.norm.out” to false and left other arguments as
312 default. PCA was applied to reduce the dimensionality to 10 after obtaining the MNN-
313 corrected expression matrix. For CCA alignment, we used the first 10 canonical correlation
314 vectors. Run time was limited to 2 hours, after which the jobs were terminated. Seurat
315 alignment score was computed exactly as described in the CCA alignment paper⁷. For our
316 own model, we consistently used $\lambda_b = 0.01$, and all other hyperparameters remain the same
317 as in dimension reduction benchmarks. 4 random seeds were used for PCA, CCA and MNN,
318 while 16 random seeds were used for scVI and our model, since neural network-based
319 models are typically considered less stable.

320 Cell querying based on posterior distributions

321 We evaluated cell-to-cell similarity based on the posterior distribution distance. Similar to the
322 training phase, we obtained samples from the “universal approximator posterior” by sampling
323 $\mathbf{x}' \sim \text{Poisson}(\mathbf{x})$ and feeding to the encoder network. To obtain a robust estimation of the
324 distribution distance with a small number of posterior samples, we project the posterior
325 samples of two cells onto the line connecting their posterior point estimates in the latent
326 space and use the projected scalar distribution distance to approximate the true distribution
327 distance. Wasserstein distance is computed on normalized projections to account for
328 nonuniform density across the embedding space:

$$NPD(p, q) = \frac{1}{2} \cdot \left(W_1(z_p(p), z_p(q)) + W_1(z_q(p), z_q(q)) \right) \quad (17)$$

329 Where

$$W_1(u, v) = \inf_{\pi \in \Gamma(u, v)} \int |x - y| d\pi(x, y) \quad (18)$$
$$z_u(v) = \frac{v - \mathbb{E}(u)}{\sqrt{\text{var}(u)}}$$

330 We term this distance metric normalized projection distance (NPD). By default, 50 samples
331 from the posterior are used to compute NPD, which produces sufficiently accurate results
332 (**Supplementary Fig. 4i-j**). The definition of posterior NPD does not imply an efficient
333 nearest neighbor searching algorithm. To increase speed, we first use Euclidean distance-
334 based nearest neighbor searching, which is highly efficient in the low-dimensional latent
335 space, and then compute posterior distances only for these nearest neighbors. The empirical
336 distribution of posterior NPD for a dataset is obtained by computing posterior NPD on
337 randomly selected pairs of cells in the reference dataset. Empirical p-values of query hits are
338 computed by comparing the posterior NPD of a query hit to this empirical distribution.
339 We note that even with the querying strategy described above, querying with single models
340 still occasionally leads to many false-positive hits when cell types on which the model has
341 not been trained are provided as query. This is because embeddings of such untrained cell
342 types are mostly random, and they could localize close to reference cells by chance. We
343 reason that embedding randomness of untrained cell types could be utilized to identify and
344 correctly reject them. Practically, we train multiple models with different starting points (as
345 determined by random seeds) and compute query hit significance for each model. A query hit
346 is considered significant only if it is consistently significant across multiple models. To
347 acquire predictions based on significant hits, we use majority voting for discrete variables,
348 e.g., cell type, or averaging for continuous variables, e.g., cell fate distribution.

349 **Distance metric ROC analysis**

350 Our model and scVI⁸ were fitted on reference datasets and applied to positive and negative
351 control query datasets in the pancreas group of **Supplementary Table 2**. We then randomly
352 selected 10,000 query-reference cell pairs. A query-reference pair is defined as “positive” if
353 the query cell and reference cell are of the same cell type, and “negative” otherwise. Each
354 benchmarked similarity metric was then computed on all sampled query-reference pairs and
355 used as predictors for “positive”/“negative” pairs. AUROC values were computed for each
356 benchmarked similarity metric. In addition to the Euclidean distance, we also computed
357 posterior distribution distances for scVI (**Supplementary Fig. 4k**). NPD was computed as
358 described in (17), based on samples from the posterior Gaussian. JSD was computed in the
359 original latent space without projection.

360 **Benchmarking query-based cell typing**

361 Cell ontology annotations in ACA were used as ground truth. Cells without cell ontology
362 annotations were excluded in the analysis. For each querying method, cell type predictions
363 for query cells were obtained based on query hits with a minimal similarity cutoff, i.e., query
364 cells with no significant hits are rejected, while cells not rejected are further assigned cell
365 type predictions. Sensitivity, specificity and Cohen's κ are computed as follows:

$$\text{sensitivity} = 1 - \frac{\# \text{ rejected query cells}}{\# \text{ query cells that match reference cell types}} \quad (19)$$

$$\text{specificity} = \frac{\# \text{ rejected query cells}}{\# \text{ query cells that do not match reference cell types}} \quad (20)$$

$$\text{Cohen's } \kappa = 1 - \frac{1 - \# \text{ correct cell type predictions}}{1 - \# \text{ correct cell type predictions expected by chance}} \quad (21)$$

366 Predictions are considered correct if they exactly match the ground truth, i.e., no flexibility
367 based on cell type similarity. This prevents unnecessary bias introduced in the selection of
368 cell type similarity measure. Cells were inversely weighed by the size of the corresponding
369 dataset when computing average sensitivity, specificity and Cohen's κ . AUROC was
370 computed using linear interpolation. For scmap², we varied the minimal cosine similarity
371 requirement for nearest neighbors. For Cell BLAST, we varied the maximal p-value cutoff
372 used in filtering hits. For CellFishing.jl⁴, the original implementation does not include a
373 dedicated cell type prediction function, so we used the same strategy as that for our own
374 method (majority voting after distance filtering) to acquire final predictions, in which we
375 varied the Hamming distance cutoff used in distance filtering. Finally, 4 random seeds were
376 tested for each cutoff and each method to reflect stability. Several other cell querying tools
377 (CellAtlasSearch³, scQuery³³, scMCA³⁴) were not included in our benchmark because they
378 do not support custom reference datasets.

379 **Benchmarking querying speed**

380 To evaluate the scalability of querying methods, we constructed reference datasets of varying
381 sizes by subsampling from the 1M mouse brain dataset³⁵. For query data, the "Marques"
382 dataset³⁶ was used. Benchmarking was performed on a workstation with 40 CPU cores,
383 100GB RAM and GeForce GTX 1080Ti GPU. For all methods, only the querying time was
384 recorded, not including the time consumed to build reference indices.

385 **Application to trachea datasets**

386 We first removed cells labeled as “ionocytes” in the “Montoro_10x”¹⁴ dataset and used
387 “FindVariableGenes” from Seurat to select informative genes in the remaining cells. Four
388 models with different starting points were trained on the tampered “Montoro_10x” dataset.
389 We used a cutoff of empirical p-value > 0.1 to reject query cells from the “Plasschaert”¹⁵
390 dataset as potential novel cell types. We clustered rejected cells using spectral clustering
391 (Scikit-learn³⁷ v0.20.1) after applying t-SNE³⁸ to latent space coordinates. The average p-
392 value for a query cell was computed as the geometric mean of p-values across all hits.

393 **Online tuning**

394 When significant batch effect exists between reference and query, we support further aligning
395 query data with the reference data in an online-learning manner. All components in the
396 pretrained model, including the encoder, decoder, prior discriminators and batch
397 discriminators, are retained. The reference-query batch effect is added as an extra component
398 to be removed using adversarial batch alignment. Specifically, a new discriminator dedicated
399 to the reference-query batch effect is added, and the decoder is expanded to accept an extra
400 one-hot indicator for reference and query. The expanded model is then fine-tuned using the
401 combination of reference and query data. Two precautions are taken to prevent a decrease in
402 specificity caused by over-alignment. First, adversarial alignment loss is constrained to cells
403 that have mutual nearest neighbors⁶ between reference and query data in each SGD
404 minibatch. Second, we penalize the deviation of tuned model weights from the original
405 weights.

406 **Application to hematopoietic progenitor datasets**

407 For the within- “Tusi”¹⁶ query, we trained four models using only cells from sequencing run
408 2, and cells from sequencing run 1 were used as query cells. PBA inferred cell fate
409 distributions provided by the authors, which are 7-dimensional categorical distributions
410 across 7 terminal cell fates, were used as the ground truth. We took the average cell fate
411 distributions of significant querying hits (p-value < 0.05) as predictions for query cells.
412 Regarding scmap-cell, we filtered nearest neighbors according to a default cosine similarity
413 cutoff of 0.5. Jensen-Shannon divergence (JSD) between true and predicted cell fate
414 distributions was computed as below:

$$JSD(\mathbf{p} \parallel \mathbf{q}) = \frac{1}{2} \cdot \sum_{i \in \{E, Ba, Meg, Ly, D, M, G\}} p_i \log \frac{p_i}{\frac{p_i + q_i}{2}} + q_i \log \frac{q_i}{\frac{p_i + q_i}{2}} \quad (22)$$

415 For cross-species querying between “Tusi” and “Velten”¹⁷, we mapped both mouse and
416 human genes to ortholog groups. Online tuning with 200 epochs was used to increase
417 sensitivity and accuracy. Latent space visualization was performed using UMAP^{39, 40}.

418 **ACA database construction**

419 We searched Gene Expression Omnibus (GEO)⁴¹ using the following search term:

```
420 (
421     "expression profiling by high throughput sequencing"[DataSet Type] OR
422     "expression profiling by high throughput sequencing"[Filter] OR
423     "high throughput sequencing"[Platform Technology Type]
424 ) AND
425 "gse"[Entry Type] AND
426 (
427     "single cell"[Title] OR
428     "single-cell"[Title]
429 ) AND
430 ("2013"[Publication Date] : "3000"[Publication Date]) AND
431 "supplementary"[Filter]
```

432 Datasets in the Hemberg collection (<https://hemberg-lab.github.io/scRNA.seq.datasets/>) were
433 merged into this list. Only animal single-cell transcriptomic datasets profiling samples of
434 normal conditions were selected. We also manually filtered small-scale or low-quality data.
435 Additionally, several other high-quality datasets missing in the previous list were included for
436 comprehensiveness.

437 The expression matrices and metadata of selected datasets were retrieved from GEO,
438 supplementary files of the publication or by directly contacting the authors. Metadata were
439 further manually curated by adding additional descriptions in the paper to acquire the most
440 detailed information of each cell. We unified raw cell type annotation by Cell Ontology⁴², a
441 structured controlled vocabulary for cell types. Closest Cell Ontology terms were manually
442 assigned based on the Cell Ontology description and context of the study.

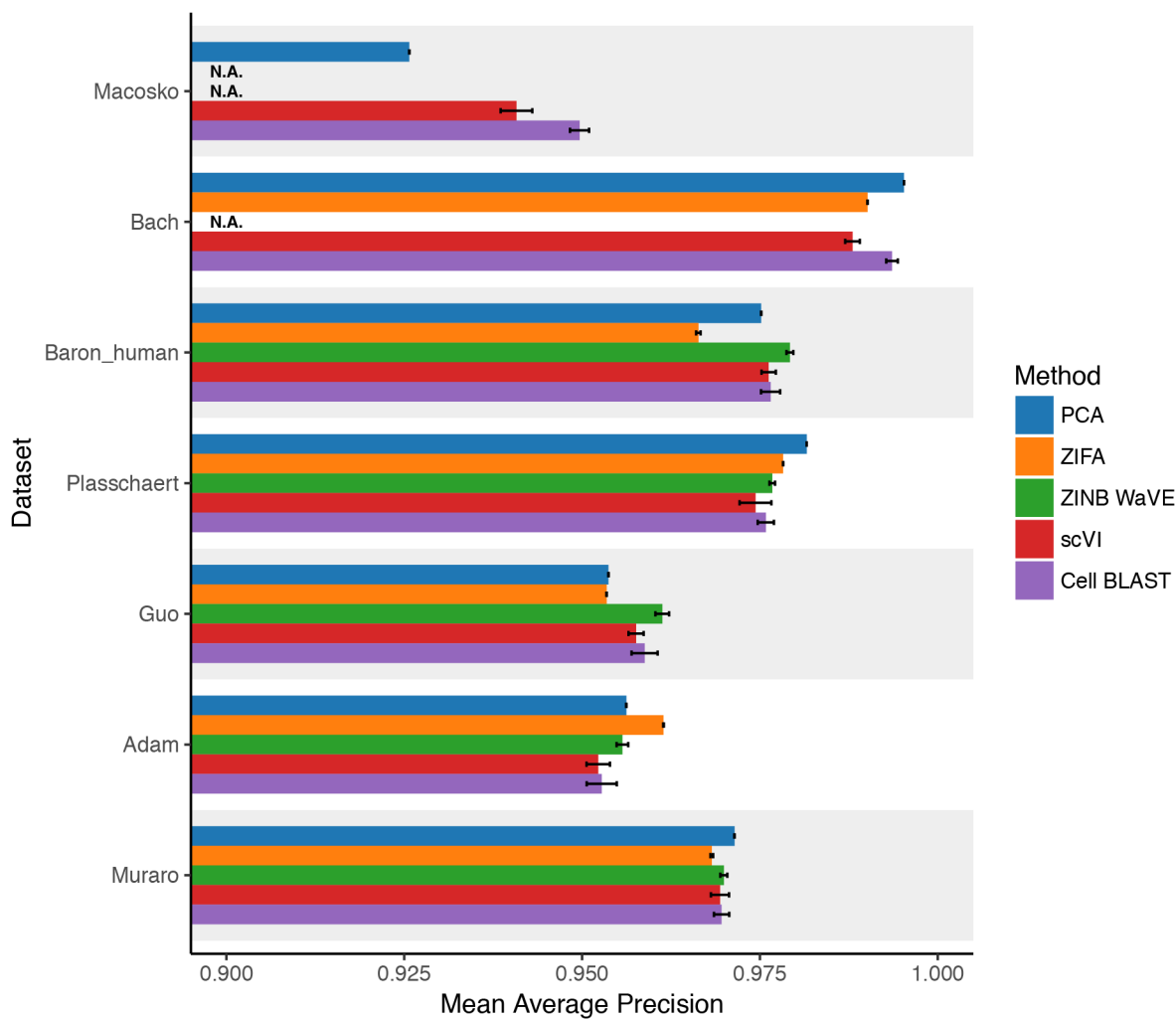
443 **Building reference panels for the ACA database**

444 Two types of searchable reference panels are built for the ACA database. The first consists of
445 individual datasets with dedicated models trained on each, while the second consists of

446 datasets grouped by organ and species, with models trained to align multiple datasets
447 profiling the same species and same organ.
448 Data preprocessing follows the same procedure as in previous benchmarks. Both cross-
449 dataset batch effect and within-dataset batch effect are manually examined and removed
450 when necessary. For the first type of reference panels, datasets too small (typically < 1,000
451 cells sequenced) are excluded because of insufficient training data. These datasets are still
452 included in the second type of panels, where they are trained jointly with other datasets
453 profiling the same organ in the same species. For each reference panel, four models with
454 different starting points are trained.

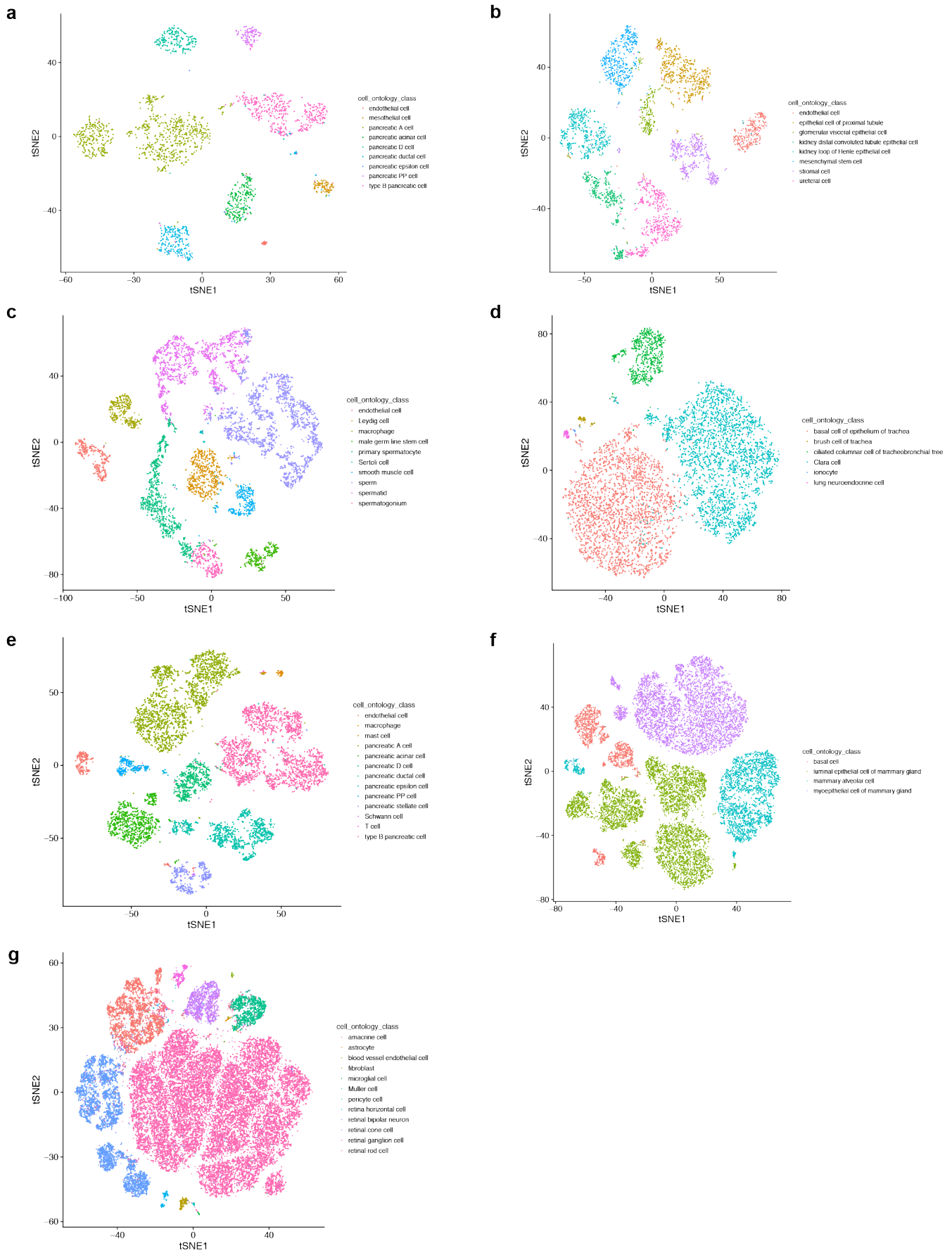
455 **Web interface**

456 For conveniently performing and visualizing Cell BLAST analysis, we built a one-stop Web
457 interface. The client-side was made from Vue.js, a single-page application Javascript
458 framework, and D3.js for cell ontology visualization. We used Koa2, a web framework for
459 Node.js, as the server side. The Cell BLAST Web portal with all accessible curated datasets
460 is deployed on Huawei Cloud.



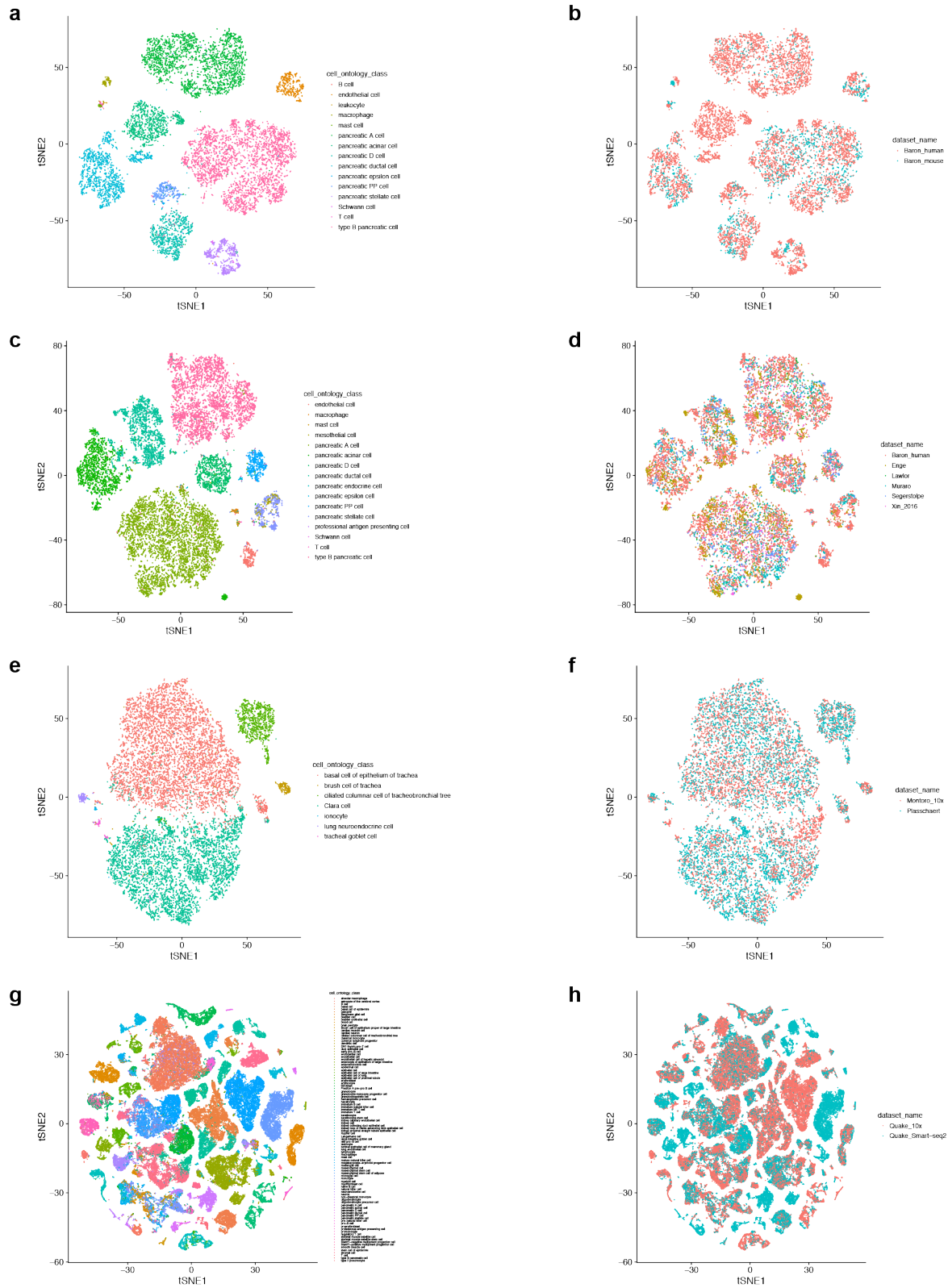
461 **Supplementary Fig. 1 Comparing low-dimensional space cell type resolutions of different dimension**
462 **reduction methods.**

463 Nearest neighbor cell type mean average precision (MAP) is used to evaluate how well biological similarity is
464 captured. MAP can be thought of as a generalization to nearest neighbor accuracy, with larger values indicating
465 higher cell type resolution and, thus, more suitable for cell querying. Error bars indicate mean \pm s.d. Methods
466 that did not finish under the 2-hour time limit are marked as N.A.



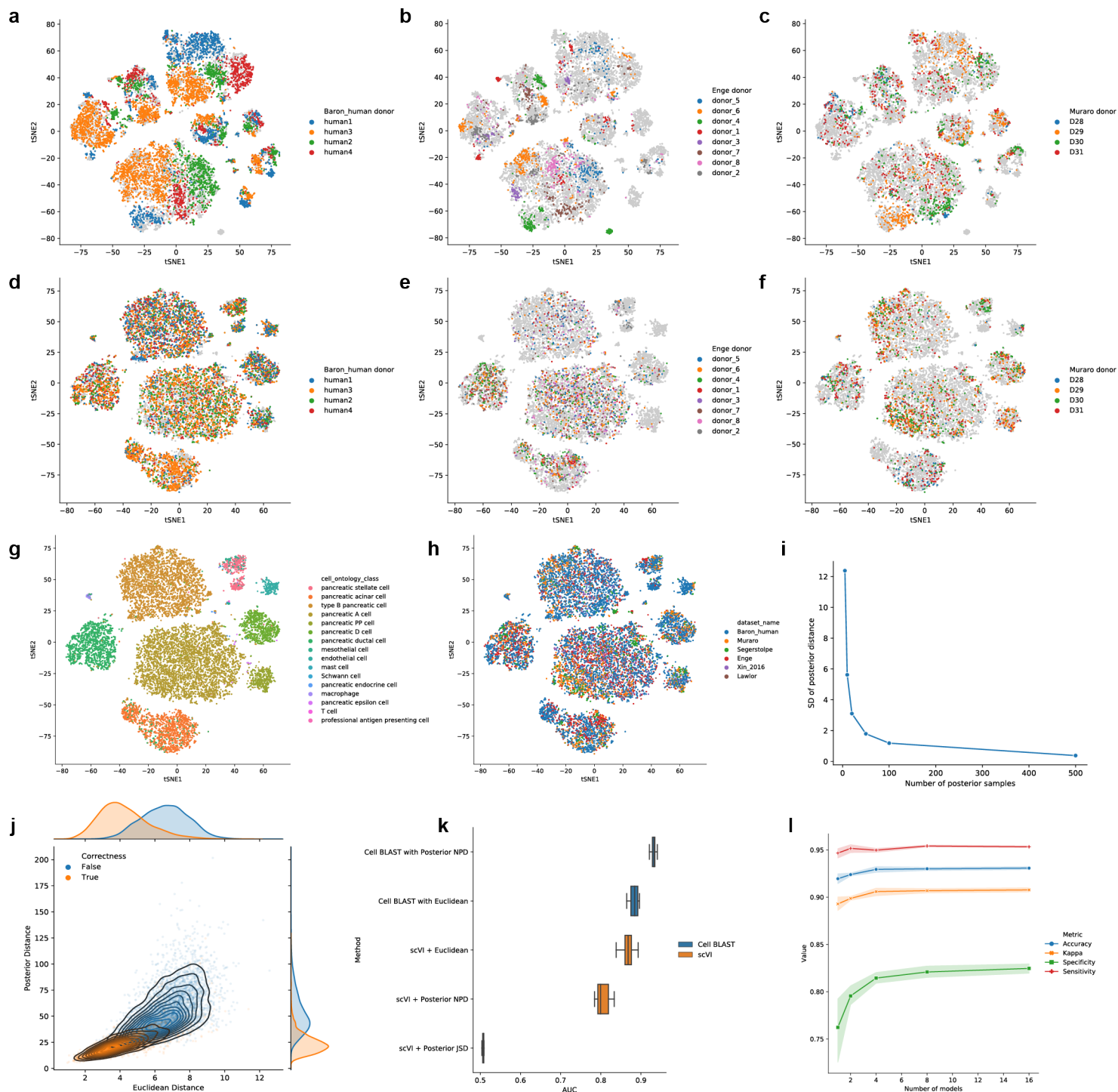
467 **Supplementary Fig. 2 t-SNE visualization of latent spaces learned by our model.**

468 (a) “Muraro”⁴³, (b) “Adam”⁴⁴, (c) “Guo”⁴⁵, (d) “Plasschaert”¹⁵, (e) “Baron_human”⁴⁶, (f) “Bach”⁴⁷, (g)
 469 “Macosko”⁴⁸.



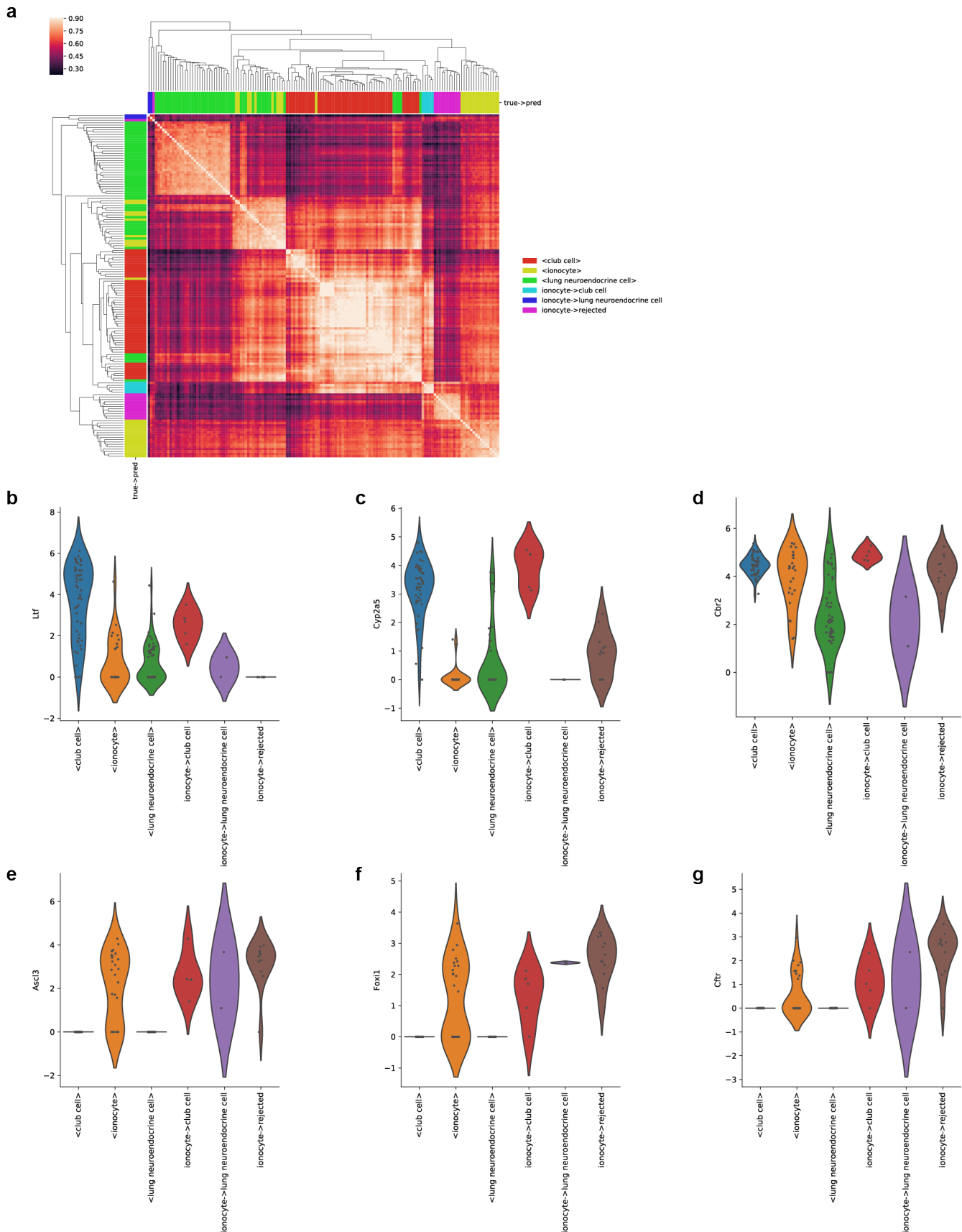
470 **Supplementary Fig. 3 t-SNE visualization of latent spaces learned by our model on combinations of**
 471 **multiple datasets, with batch effect corrected.**

472 Figures in the left column color cells by cell type, while figures in the right column color cells by dataset. (a-b)
 473 “Baron_human”⁴⁶ and “Baron_mouse”⁴⁶; (c-d) “Baron_human”⁴⁶, “Muraro”⁴³, “Enge”⁴⁹, “Segerstolpe”⁵⁰,
 474 “Xin_2016”⁵¹ and “Lawlor”⁵²; (e-f) “Montoro_10x”¹⁴ and “Plasschaert”¹⁵; (g-h) “Quake_Smart-seq2”¹⁸ and
 475 “Quake_10x”¹⁸.



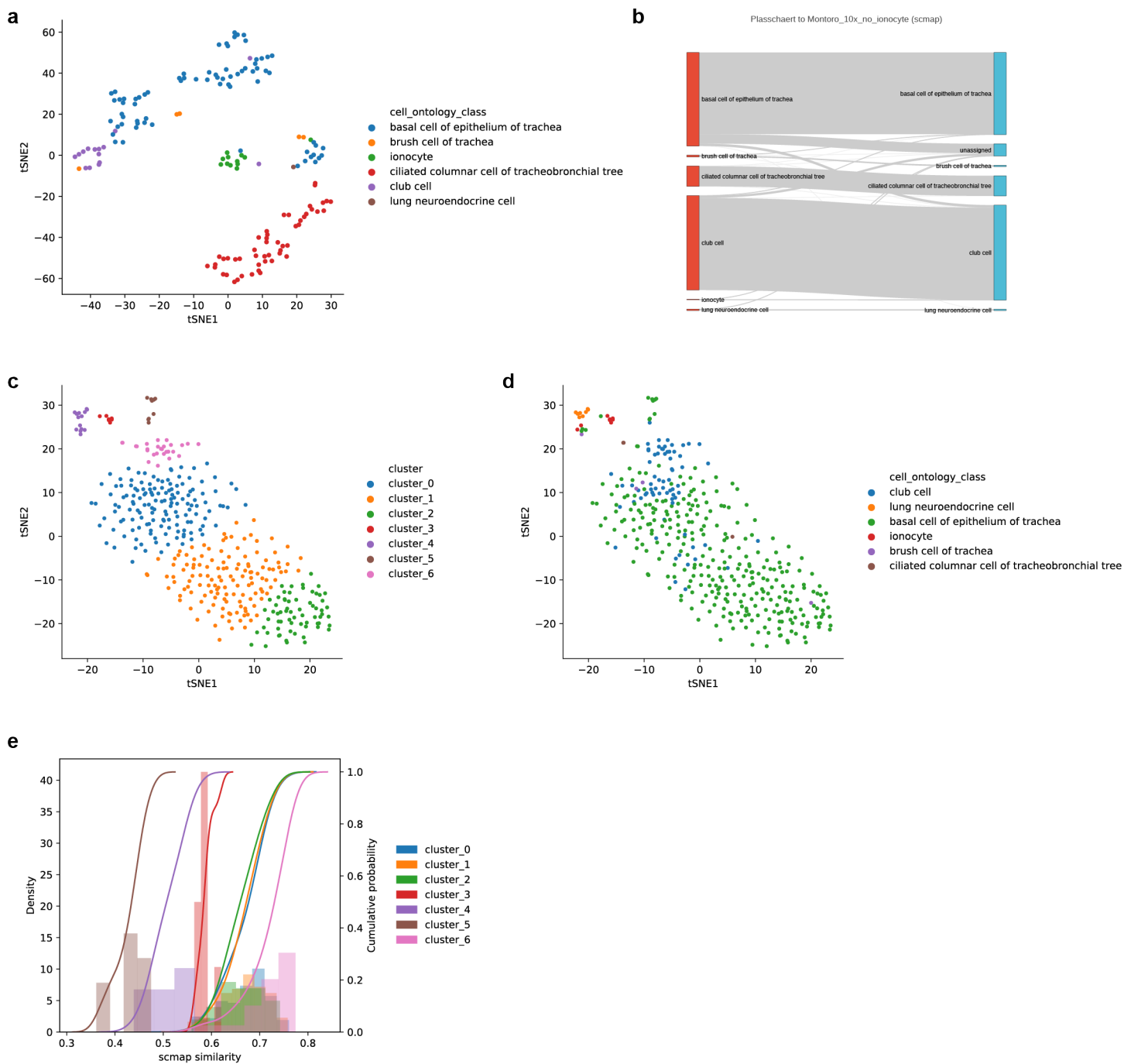
476 **Supplementary Fig. 4 Multilevel batch effect correction and Cell BLAST strategy optimization.**

477 (a-c) Latent space learned with only cross-dataset batch effect correction, colored by (a) donor in
 478 “Baron_human”⁴⁶, (b) donor in “Enge”⁴⁹, (c) donor in “Muraro”⁴³. (d-h) Latent space learned with both cross-
 479 dataset and within-dataset batch effect correction, colored by (d) donor in “Baron_human”⁴⁶, (e) donor in
 480 “Enge”⁴⁹, (f) donor in “Muraro”⁴³, (g) cell type, (h) dataset. (i) Standard deviation decreases as the number of
 481 samples from the posterior increases. (j) Relationship between Euclidean distance and NPD in
 482 “Baron_human”⁴⁶ data. The orange points represent cell pairs that are of the same cell type (“positive pairs”),
 483 while the blue points represent cell pairs of different cell types (“negative pairs”). (k) AUROC of different
 484 distance metrics in discriminating cell pairs with the same cell type from cell pairs with different cell types. Box
 485 plots indicate the median (center lines), interquartile range (hinges), and 1.5 times the interquartile range
 486 (whiskers). Note that the posterior distribution distances for scVI only lead to a decrease in performance,
 487 possibly due to improper Gaussian assumption in the posterior. (l) Accuracy, Cohen’s κ (a measure of
 488 prediction accuracy corrected for chance, see **Methods** for more details)², specificity and sensitivity all increase
 489 as the number of models used for cell querying increases, among which the improvement of specificity is the
 490 most significant. Error bars indicate mean \pm s.d.



491 **Supplementary Fig. 5 Ionocytes predicted to be club cells are potentially doublets or of an intermediate**
 492 **cell state.**

493 **(a)** Cell-cell correlation heatmap for several cell types of interest. Cells labeled as “<X>” are reference cells in
 494 the “Montoro”¹⁴ dataset. Cells labeled as “X->Y” are cells annotated as “X” in the original “Plasschaert”¹⁵
 495 dataset but predicted to be “Y”. **(b-d)** Expression levels of several club cell markers in the cell groups of
 496 interest. **(e-g)** Expression levels of several ionocyte markers in the cell groups of interest.



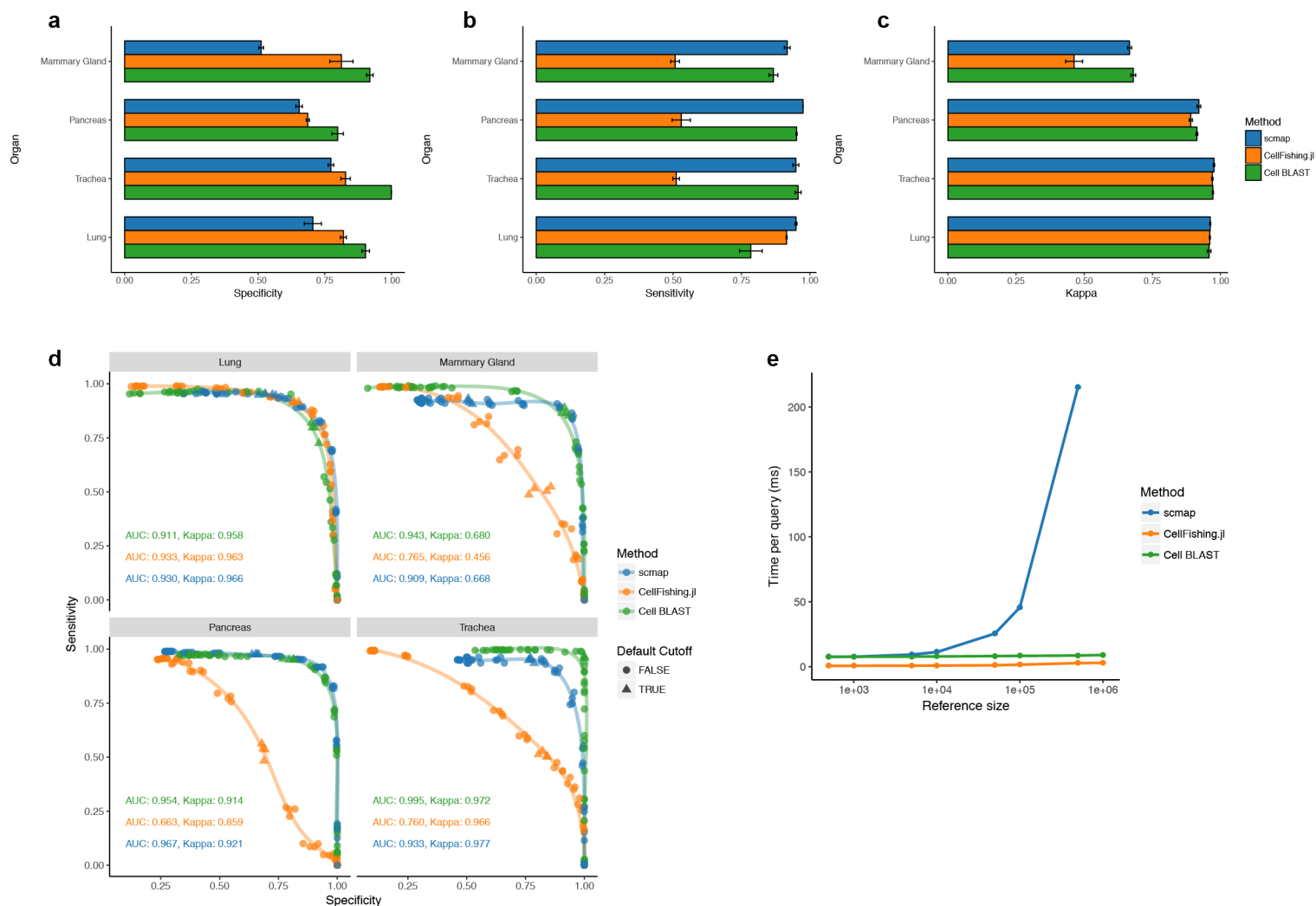
497 **Supplementary Fig. 6 Rejected cells in the “Montoro” - “Plasschaert” query.**

498 (a) t-SNE visualization of Cell BLAST-rejected cells, colored by cell type. (b) Sankey plot of scmap prediction.

499 (c, d) t-SNE visualization of scmap-rejected cells, colored by unsupervised clustering (c) and cell type (d). (e)

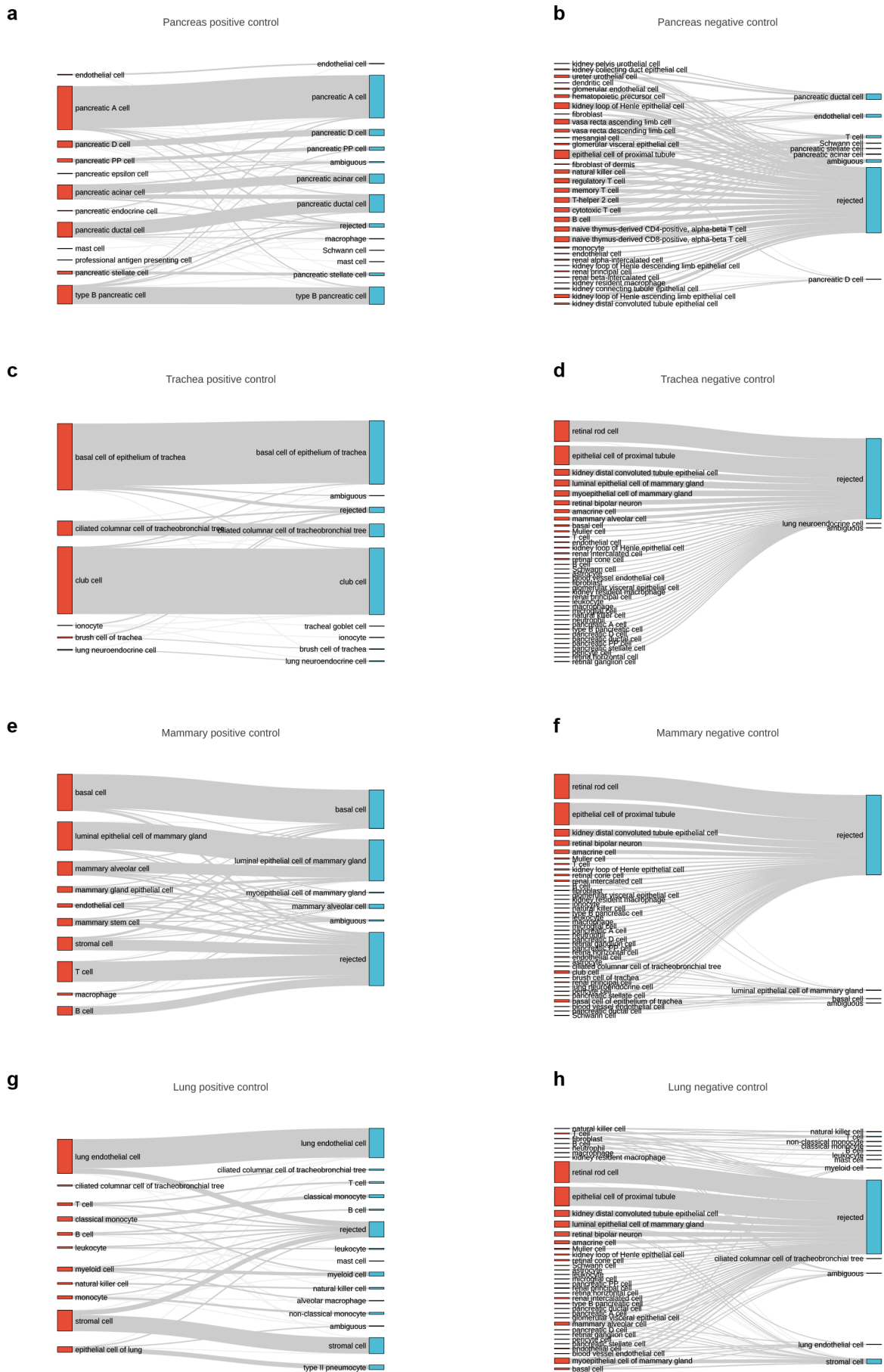
500 scmap similarity distribution in each cluster of scmap-rejected cells. The rejected ionocytes do not have the

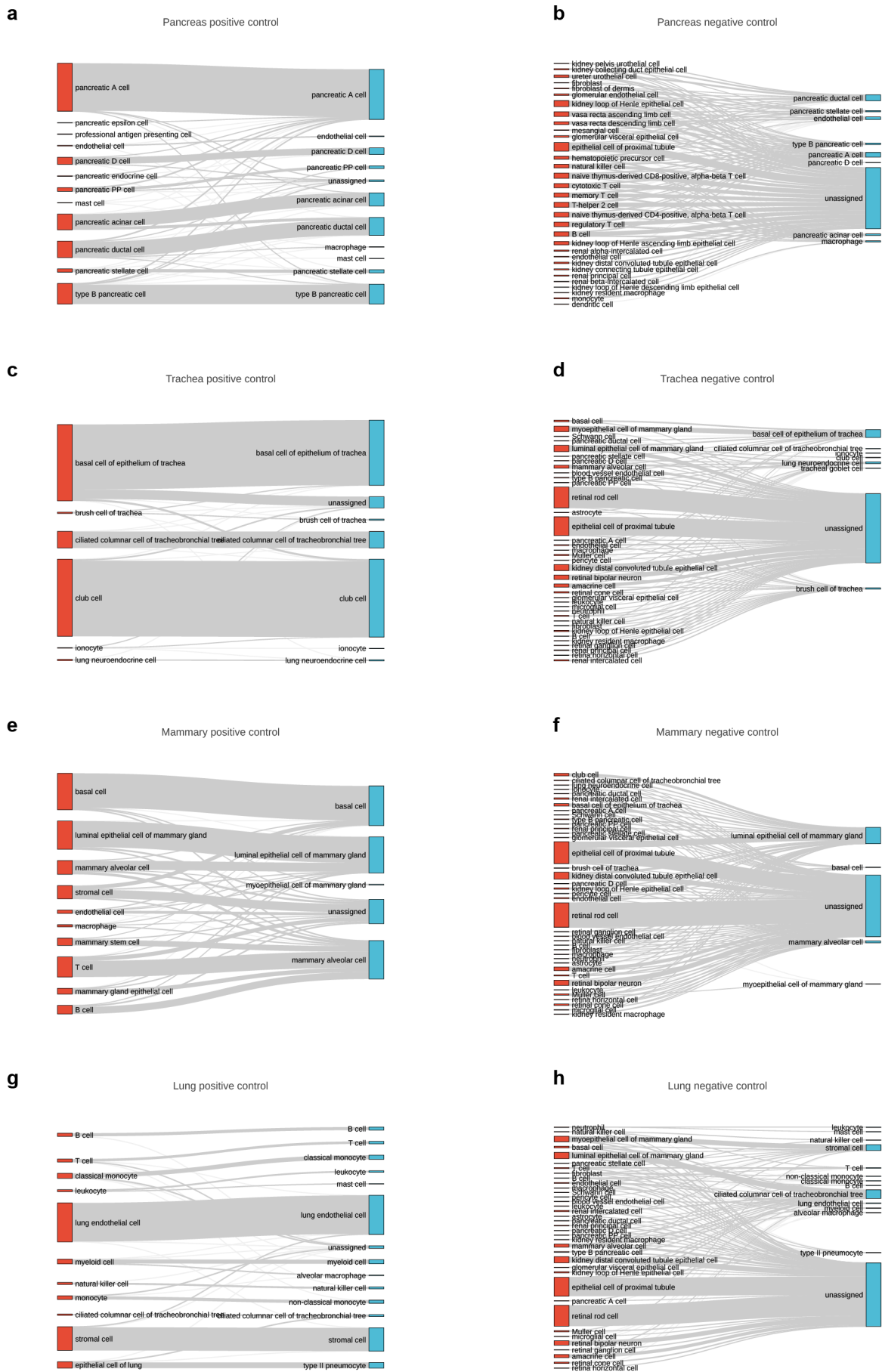
501 lowest cosine similarity scores to draw sufficient attention.

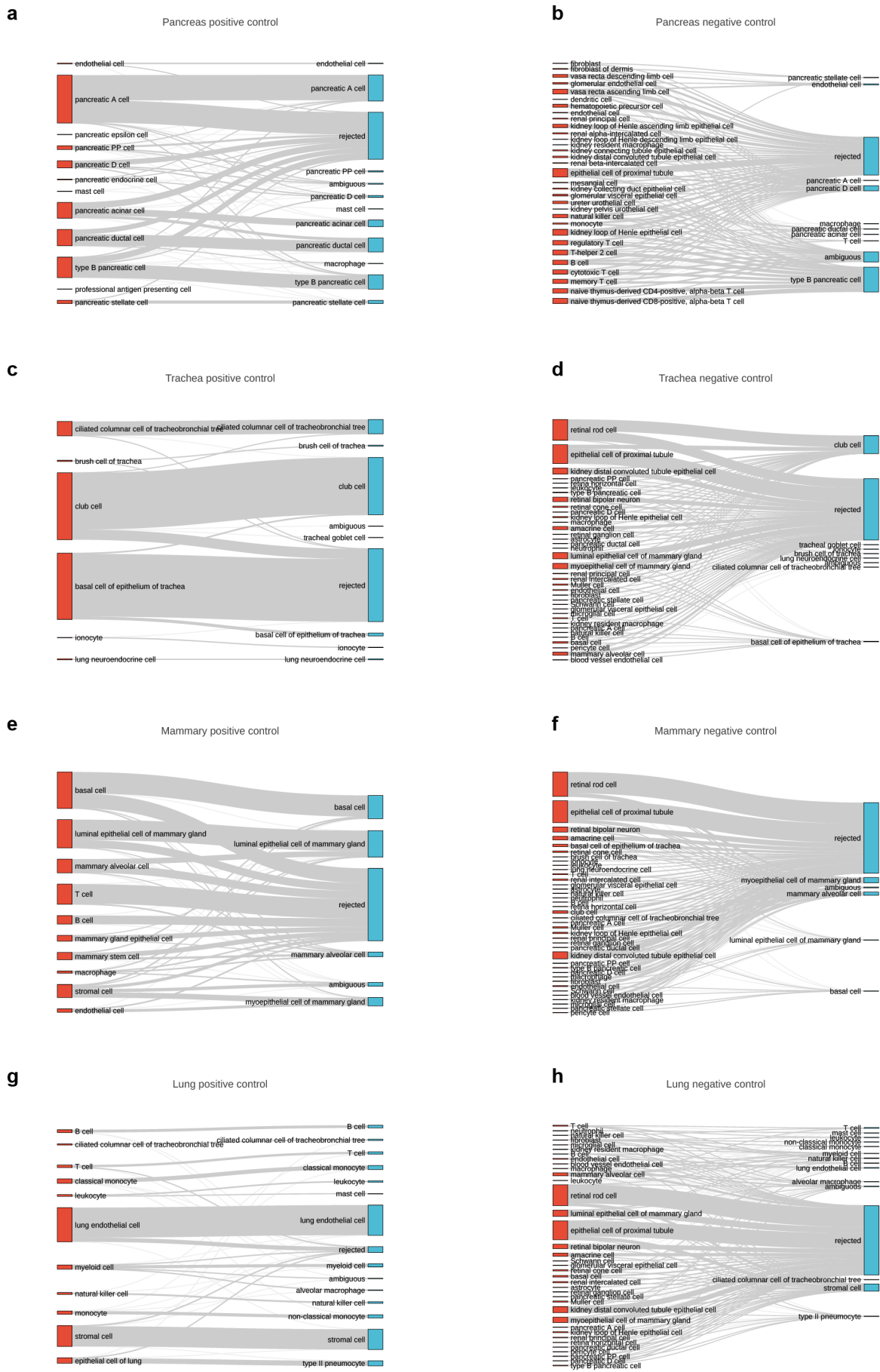


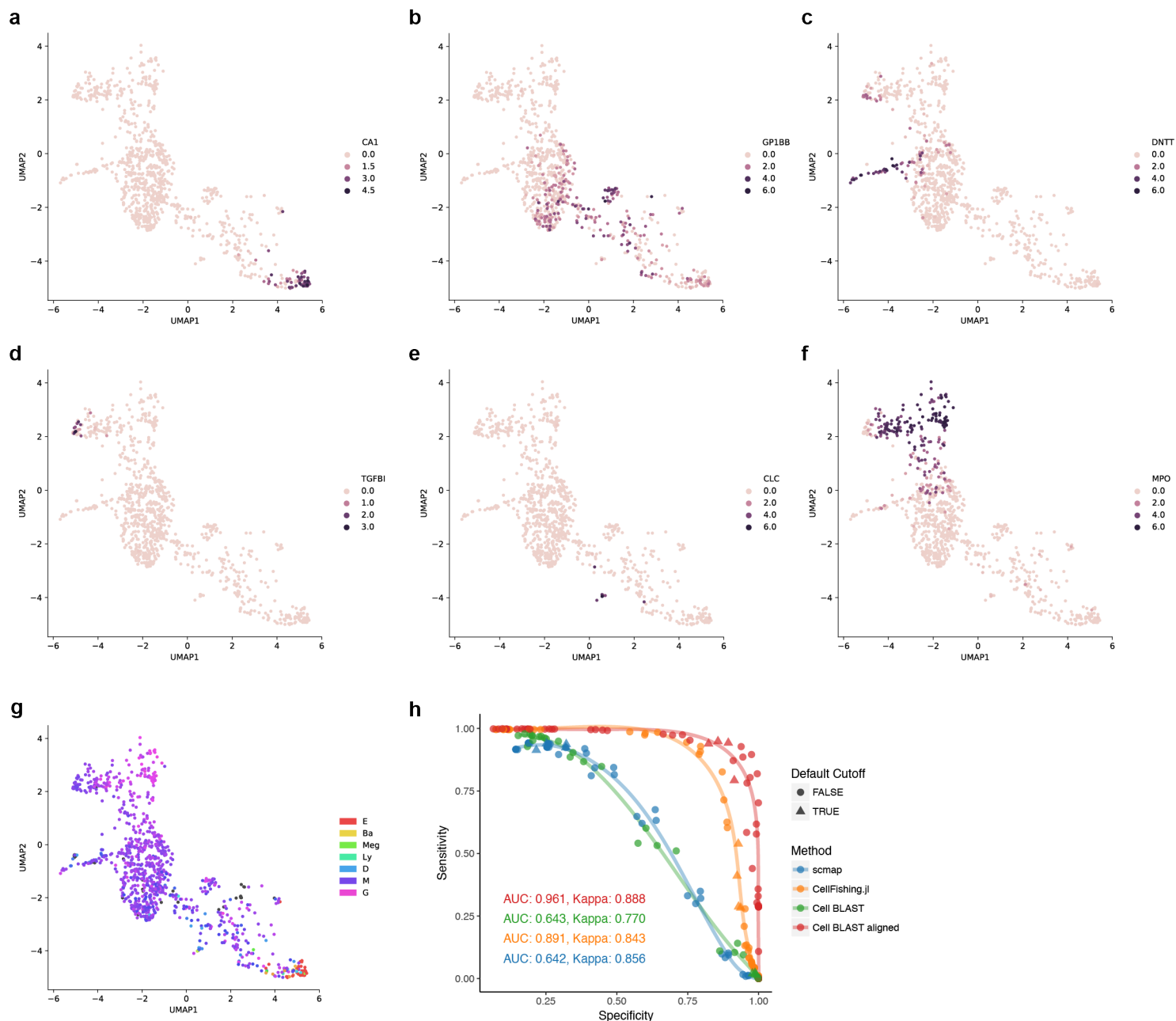
502 **Supplementary Fig. 7 Benchmarking query-based cell typing.**

503 (a-c) Querying specificity (a), sensitivity (b) and Cohen's κ (c) for different methods under the default setting.
 504 Error bars indicate mean \pm s.d. (d) ROC curve of cell querying in four different groups of test datasets. Cohen's
 505 κ values in the bottom left of each subpanel correspond to the optimal point on the ROC curve. Points
 506 corresponding to each method's default cutoff (scmap: cosine distance = 0.5, CellFishing.jl: Hamming distance
 507 = 110, Cell BLAST: p-value = 0.05) are marked as triangles. Note that CellFishing.jl does not provide a default
 508 cutoff, so we chose a Hamming distance of 110, which is the closest to balancing sensitivity and specificity, but
 509 it is still far from being stable across different datasets. (e) Querying speed on reference datasets of different
 510 sizes subsampled from the 1M mouse brain dataset³⁵. Error bars indicate mean \pm s.d.



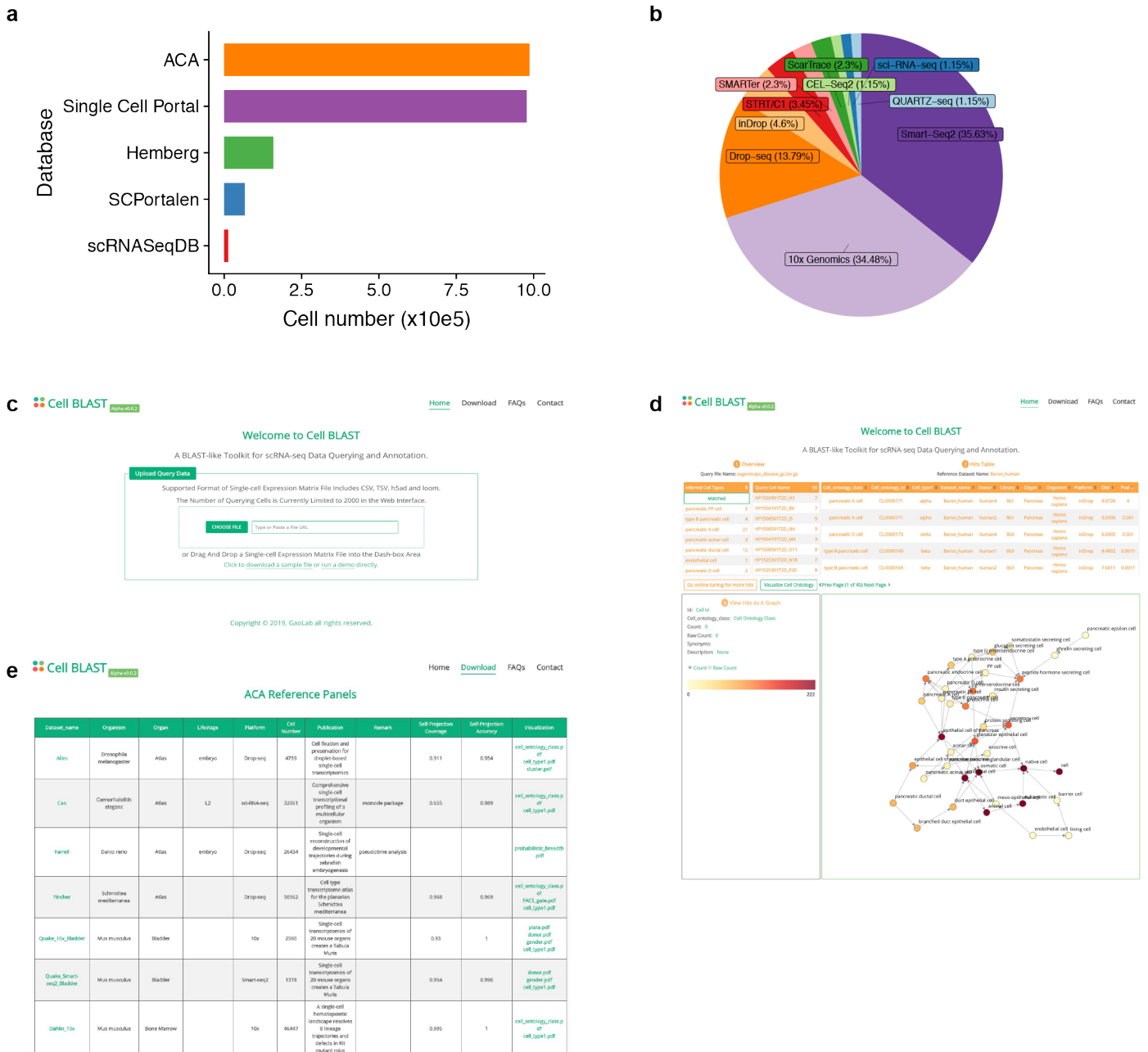






514 **Supplementary Fig. 11 Using “online tuning” in hematopoietic progenitor and *Tabula Muris*¹⁸ spleen**
 515 **data.**

516 UMAP visualization of the “Velten”¹⁷ dataset, colored by the expression of known lineage markers, including
 517 CA1 for the erythrocyte lineage (a), GP1BB for the megakaryocyte lineage (b), DNNT for the B-cell lineage (c),
 518 TGFB1 for monocyte and dendritic cell lineages (d), CLC for eosinophil, basophil, and mast cell lineages (e),
 519 MPO for the neutrophil lineage (f), and scmap predicted cell fate distribution (g). (h) ROC curve of cell
 520 querying in *Tabula Muris*¹⁸ spleen data. Cohen’s κ values in the bottom left of each subpanel correspond to the
 521 optimal point on the ROC curve. Points corresponding to each method’s default cutoff (scmap: cosine distance
 522 = 0.5, CellFishing.jl: Hamming distance = 110, Cell BLAST: p-value = 0.05) are marked as triangles.



523 **Supplementary Fig. 12 ACA database and Cell BLAST Web portal.**

524 **(a)** Comparison of cell numbers in different single-cell transcriptomics databases. **(b)** Composition of different
 525 single-cell sequencing platforms in ACA. **(c)** Home page of the Cell BLAST Web interface. **(d)** Web interface
 526 showing the results of a sample query. **(e)** A full list of ACA reference panels available in our Web interface.

Dataset Name	Organism	Organ Profiled	Experimental Platform	Used In
Guo ⁴⁵	Human	Testis	10x ⁵³	DR
Muraro ⁴³		Pancreas	CEL-Seq2 ⁵⁴	DR, BC
Xin_2016 ⁵¹			SMARTer ⁵⁵	BC
Lawlor ⁵²			SMARTer ⁵⁵	BC
Segerstolpe ⁵⁰			Smart-seq2 ⁵⁶	BC
Enge ⁴⁹			Smart-seq2 ⁵⁶	BC
Baron_human ⁴⁶			inDrop ⁵⁷	DR, BC
Baron_mouse ⁴⁶		inDrop ⁵⁷	BC	
Adam ⁴⁴	Mouse	Kidney	Drop-seq ⁴⁸	DR
Plasschaert ¹⁵		Trachea	inDrop ⁵⁷	DR, BC
Montoro_10x ¹⁴			10x ⁵³	BC
Macosko ⁴⁸		Retina	Drop-seq ⁴⁸	DR
Bach ⁴⁷		Mammary Gland	10x ⁵³	DR
Quake_Smart-seq2 ¹⁸		20 Organs	Smart-seq2 ⁵⁶	BC
Quake_10x ¹⁸		12 Organs	10x ⁵³	BC

527

528 **Supplementary Table 1. Datasets used in dimensionality reduction and batch effect**
 529 **correction benchmarking.**

530 DR, dimension reduction benchmarking; BC, batch effect correction benchmarking.

Group	Role	Dataset Name	Organism	Organ Profiled	Experimental Platform
Pancreas	Reference	Baron_human ⁴⁶	Human	Pancreas	inDrop ⁵⁷
		Xin_2016 ⁵¹			SMARTer ⁵⁵
		Lawlor ⁵²			SMARTer ⁵⁵
	Positive control query	Muraro ⁴³			CEL-Seq2 ⁵⁴
		Segerstolpe ⁵⁰			Smart-seq2 ⁵⁶
		Enge ⁴⁹			Smart-seq2 ⁵⁶
	Negative control query	Wu_human ⁵⁸		Kidney	10x ⁵³
		Zheng ⁵³		PBMC	10x ⁵³
		Philippeos ⁵⁹		Skin	Smart-seq2 ⁵⁶
Trachea	Reference	Montoro_10x ¹⁴	Mouse	Trachea	10x ⁵³
	Positive control query	Plasschaert ¹⁵			inDrop ⁵⁷
	Negative control query	Baron_mouse ⁴⁶		Pancreas	inDrop ⁵⁷
		Park ⁶⁰		Kidney	10x ⁵³
		Bach ⁴⁷		Mammary Gland	10x ⁵³
		Macosko ⁴⁸		Retina	Drop-seq ⁴⁸
Mammary Gland	Reference	Bach ⁴⁷	Mouse	Mammary Gland	10x ⁵³
	Positive control query	Giraddi_10x ⁶¹			10x ⁵³
		Quake_Smart-seq2_Mammary_Gland ¹⁸			Smart-seq2 ⁵⁶
		Quake_10x_Mammary_Gland ¹⁸			10x ⁵³
	Negative control query	Baron_mouse ⁴⁶		Pancreas	inDrop ⁵⁷
		Park ⁶⁰		Kidney	10x ⁵³
		Plasschaert ¹⁵		Trachea	inDrop ⁵⁷
		Macosko ⁴⁸		Retina	Drop-seq ⁴⁸

Lung	Reference	Quake_10x_ Lung ¹⁸	Mouse	Lung	10x ⁵³
	Positive control query	Quake-Smart-seq2_ Lung ¹⁸			Smart-seq2 ⁵⁶
	Negative control query	Baron_mouse ⁴⁶		Pancreas	inDrop ⁵⁷
		Park ⁶⁰		Kidney	10x ⁵³
		Bach ⁴⁷		Mammary Gland	10x ⁵³
Plasschaert ¹⁵	Trachea	inDrop ⁵⁷			
Spleen	Reference	Quake_10x_ Spleen ¹⁸	Mouse	Spleen	10x ⁵³
	Positive control query	Quake_Smart-seq2_ Spleen ¹⁸			Smart-seq2 ⁵⁶
	Negative control query	Baron_mouse ⁴⁶		Pancreas	inDrop ⁵⁷
		Park ⁶⁰		Kidney	10x ⁵³
		Macosko ⁴⁸		Retina	Drop-seq ⁴⁸
Bach ⁴⁷	Mammary Gland	10x ⁵³			

531

532 **Supplementary Table 2. Datasets used in cell query benchmarking.**

533

534 **Supplementary Table 3. Raw data of benchmarking results.**

535

536 **Supplementary Table 4. Datasets in ACA.**

537 Main text references

- 538 1. Altschul, S.F., Gish, W., Miller, W., Myers, E.W. & Lipman, D.J. Basic local
539 alignment search tool. *J Mol Biol* **215**, 403-410 (1990).
- 540 2. Kiselev, V.Y., Yiu, A. & Hemberg, M. scmap: projection of single-cell RNA-seq data
541 across data sets. *Nat Methods* **15**, 359-362 (2018).
- 542 3. Srivastava, D., Iyer, A., Kumar, V. & Sengupta, D. CellAtlasSearch: a scalable search
543 engine for single cells. *Nucleic Acids Res* **46**, W141-W147 (2018).
- 544 4. Sato, K., Tsuyuzaki, K., Shimizu, K. & Nikaido, I. CellFishing.jl: an ultrafast and
545 scalable cell search method for single-cell RNA sequencing. *Genome Biol* **20**, 31
546 (2019).
- 547 5. Tung, P.Y. et al. Batch effects and the effective design of single-cell gene expression
548 studies. *Sci Rep* **7**, 39921 (2017).
- 549 6. Haghverdi, L., Lun, A.T.L., Morgan, M.D. & Marioni, J.C. Batch effects in single-
550 cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat*
551 *Biotechnol* **36**, 421-427 (2018).
- 552 7. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell
553 transcriptomic data across different conditions, technologies, and species. *Nat*
554 *Biotechnol* **36**, 411-420 (2018).
- 555 8. Lopez, R., Regier, J., Cole, M.B., Jordan, M.I. & Yosef, N. Deep generative modeling
556 for single-cell transcriptomics. *Nat Methods* **15**, 1053-1058 (2018).
- 557 9. Ding, J., Condon, A. & Shah, S.P. Interpretable dimensionality reduction of single
558 cell transcriptome data with deep generative models. *Nat Commun* **9**, 2002 (2018).
- 559 10. Grønbech, C.H. et al. scVAE: Variational auto-encoders for single-cell gene
560 expression data. *bioRxiv preprint*, 318295 (2019).
- 561 11. Wang, D. & Gu, J. VASC: Dimension Reduction and Visualization of Single-cell
562 RNA-seq Data by Deep Variational Autoencoder. *Genomics, proteomics*
563 *bioinformatics* (2018).
- 564 12. Pierson, E. & Yau, C. ZIFA: Dimensionality reduction for zero-inflated single-cell
565 gene expression analysis. *Genome Biol* **16**, 241 (2015).
- 566 13. Risso, D., Perraudeau, F., Gribkova, S., Dudoit, S. & Vert, J.P. A general and flexible
567 method for signal extraction from single-cell RNA-seq data. *Nat Commun* **9**, 284
568 (2018).
- 569 14. Montoro, D.T. et al. A revised airway epithelial hierarchy includes CFTR-expressing
570 ionocytes. *Nature* **560**, 319-324 (2018).
- 571 15. Plasschaert, L.W. et al. A single-cell atlas of the airway epithelium reveals the CFTR-
572 rich pulmonary ionocyte. *Nature* **560**, 377-381 (2018).
- 573 16. Tusi, B.K. et al. Population snapshots predict early haematopoietic and erythroid
574 hierarchies. *Nature* **555**, 54-60 (2018).
- 575 17. Velten, L. et al. Human haematopoietic stem cell lineage commitment is a continuous
576 process. *Nat Cell Biol* **19**, 271-281 (2017).
- 577 18. Tabula Muris, C. et al. Single-cell transcriptomics of 20 mouse organs creates a
578 Tabula Muris. *Nature* **562**, 367-372 (2018).
- 579 19. Abugessaisa, I. et al. SCPortalen: human and mouse single-cell centric database.
580 *Nucleic Acids Res* **46**, D781-D787 (2018).
- 581 20. Cao, Y., Zhu, J., Jia, P. & Zhao, Z. scRNASeqDB: A Database for RNA-Seq Based
582 Gene Expression Profiles in Human Single Cells. *Genes (Basel)* **8** (2017).

583 Supplementary references

- 584 21. Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I. & Frey, B. Adversarial
585 autoencoders. *arXiv preprint* (2015).
- 586 22. Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change and
587 dispersion for RNA-seq data with DESeq2. *Genome Biol* **15**, 550 (2014).
- 588 23. Abadi, M. et al. in 12th USENIX Symposium on Operating Systems Design and
589 Implementation 265-283 (2016).
- 590 24. Ganin, Y. & Lempitsky, V. Unsupervised domain adaptation by backpropagation.
591 *arXiv preprint* (2014).
- 592 25. Xie, Q., Dai, Z., Du, Y., Hovy, E. & Neubig, G. in Advances in Neural Information
593 Processing Systems 585-596 (2017).
- 594 26. Tzeng, E., Hoffman, J., Saenko, K. & Darrell, T. in Proceedings of the IEEE
595 Conference on Computer Vision and Pattern Recognition 7167-7176 (2017).
- 596 27. Goodfellow, I. et al. in Advances in neural information processing systems 2672-2680
597 (2014).
- 598 28. Tange, O. Gnu parallel 2018. (2018).
- 599 29. Baglama, J., Reichel, L. & Lewis, B.J.R.p.v. irlba: Fast truncated singular value
600 decomposition and principal components analysis for large dense and sparse matrices.
601 **2** (2017).
- 602 30. Paszke, A. et al. Automatic differentiation in pytorch. (2017).
- 603 31. Herrero, J. et al. Ensembl comparative genomics resources. *Database (Oxford)* **2016**
604 (2016).
- 605 32. Lun, A.T., McCarthy, D.J. & Marioni, J.C. A step-by-step workflow for low-level
606 analysis of single-cell RNA-seq data with Bioconductor. *F1000Res* **5**, 2122 (2016).
- 607 33. Alavi, A., Ruffalo, M., Parvangada, A., Huang, Z. & Bar-Joseph, Z. A web server for
608 comparative analysis of single-cell RNA-seq data. *Nat Commun* **9**, 4768 (2018).
- 609 34. Han, X. et al. Mapping the Mouse Cell Atlas by Microwell-Seq. *Cell* **172**, 1091-1107
610 e1017 (2018).
- 611 35. 10x Genomics in 1.3 Million Brain Cells from E18 Mice (2017).
- 612 36. Marques, S. et al. Oligodendrocyte heterogeneity in the mouse juvenile and adult
613 central nervous system. *Science* **352**, 1326-1329 (2016).
- 614 37. Pedregosa, F. et al. Scikit-learn: Machine learning in Python. **12**, 2825-2830 (2011).
- 615 38. Maaten, L.v.d. & Hinton, G. Visualizing data using t-SNE. *Journal of machine*
616 *learning research* **9**, 2579-2605 (2008).
- 617 39. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and
618 projection for dimension reduction. *arXiv preprint* (2018).
- 619 40. Becht, E. et al. Dimensionality reduction for visualizing single-cell data using UMAP.
620 *Nat Biotechnol* (2018).
- 621 41. Barrett, T. et al. NCBI GEO: archive for functional genomics data sets--update.
622 *Nucleic Acids Res* **41**, D991-995 (2013).
- 623 42. Diehl, A.D. et al. The Cell Ontology 2016: enhanced content, modularization, and
624 ontology interoperability. *J Biomed Semantics* **7**, 44 (2016).
- 625 43. Muraro, M.J. et al. A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell*
626 *Syst* **3**, 385-394 e383 (2016).
- 627 44. Adam, M., Potter, A.S. & Potter, S.S. Psychrophilic proteases dramatically reduce
628 single-cell RNA-seq artifacts: a molecular atlas of kidney development. *Development*
629 **144**, 3625-3632 (2017).
- 630 45. Guo, J. et al. The adult human testis transcriptional cell atlas. *Cell Res* **28**, 1141-1157
631 (2018).

- 632 46. Baron, M. et al. A Single-Cell Transcriptomic Map of the Human and Mouse
633 Pancreas Reveals Inter- and Intra-cell Population Structure. *Cell Syst* **3**, 346-360 e344
634 (2016).
- 635 47. Bach, K. et al. Differentiation dynamics of mammary epithelial cells revealed by
636 single-cell RNA sequencing. *Nat Commun* **8**, 2128 (2017).
- 637 48. Macosko, E.Z. et al. Highly Parallel Genome-wide Expression Profiling of Individual
638 Cells Using Nanoliter Droplets. *Cell* **161**, 1202-1214 (2015).
- 639 49. Enge, M. et al. Single-Cell Analysis of Human Pancreas Reveals Transcriptional
640 Signatures of Aging and Somatic Mutation Patterns. *Cell* **171**, 321-330 e314 (2017).
- 641 50. Segerstolpe, A. et al. Single-Cell Transcriptome Profiling of Human Pancreatic Islets
642 in Health and Type 2 Diabetes. *Cell Metab* **24**, 593-607 (2016).
- 643 51. Xin, Y. et al. RNA Sequencing of Single Human Islet Cells Reveals Type 2 Diabetes
644 Genes. *Cell Metab* **24**, 608-615 (2016).
- 645 52. Lawlor, N. et al. Single-cell transcriptomes identify human islet cell signatures and
646 reveal cell-type-specific expression changes in type 2 diabetes. *Genome Res* **27**, 208-
647 222 (2017).
- 648 53. Zheng, G.X. et al. Massively parallel digital transcriptional profiling of single cells.
649 *Nat Commun* **8**, 14049 (2017).
- 650 54. Hashimshony, T. et al. CEL-Seq2: sensitive highly-multiplexed single-cell RNA-Seq.
651 *Genome Biol* **17**, 77 (2016).
- 652 55. Verboom, K. et al. SMARTer single cell total RNA sequencing. *bioRxiv preprint*,
653 430090 (2018).
- 654 56. Picelli, S. et al. Full-length RNA-seq from single cells using Smart-seq2. *Nat Protoc*
655 **9**, 171-181 (2014).
- 656 57. Klein, A.M. et al. Droplet barcoding for single-cell transcriptomics applied to
657 embryonic stem cells. *Cell* **161**, 1187-1201 (2015).
- 658 58. Wu, H. et al. Comparative Analysis and Refinement of Human PSC-Derived Kidney
659 Organoid Differentiation with Single-Cell Transcriptomics. *Cell Stem Cell* **23**, 869-
660 881 e868 (2018).
- 661 59. Philippeos, C. et al. Spatial and Single-Cell Transcriptional Profiling Identifies
662 Functionally Distinct Human Dermal Fibroblast Subpopulations. *J Invest Dermatol*
663 **138**, 811-825 (2018).
- 664 60. Park, J. et al. Single-cell transcriptomics of the mouse kidney reveals potential
665 cellular targets of kidney disease. *Science* **360**, 758-763 (2018).
- 666 61. Girardi, R.R. et al. Single-Cell Transcriptomes Distinguish Stem Cell State Changes
667 and Lineage Specification Programs in Early Mammary Gland Development. *Cell*
668 *Rep* **24**, 1653-1666 e1657 (2018).
- 669