

# 1 **Cell-ID: gene signature extraction and cell identity recognition at individual** 2 **cell level**

3 Cortal Akira<sup>1,2</sup>, Martignetti Loredana<sup>1,2</sup>, Six Emmanuelle<sup>1,3</sup>, Rausell Antonio<sup>1,2,\*</sup>.

4 1. Paris University, Imagine Institute, 75015 Paris, France, EU

5 2. Clinical Bioinformatics Laboratory, INSERM UMR1163, Necker Hospital for Sick Children, 75015 Paris, France,  
6 EU

7 3. Laboratory of Human Lymphohematopoiesis, INSERM UMR 1163, Paris, France;

8 \*Correspondence to: [antonio.rausell@institutimagine.org](mailto:antonio.rausell@institutimagine.org)

9

## 10 **Abstract**

11 The exhaustive exploration of human cell heterogeneity requires the unbiased identification of molecular  
12 signatures that can serve as unique cell identity cards for every cell in the body. However, the stochasticity  
13 associated with high-throughput single-cell sequencing has made it necessary to use clustering-based  
14 computational approaches in which the characterization of cell-type heterogeneity is performed at cell-  
15 subpopulation level rather than at full single-cell resolution. We present here Cell-ID, a clustering-free  
16 multivariate statistical method for the robust extraction of per-cell gene signatures from single-cell sequencing  
17 data. Cell-ID signatures allow unbiased cell identity recognition across different donors, tissues-of-origin, model  
18 organisms and single-cell omics technologies. Cell-ID is distributed as an open-source R software package:  
19 <https://github.com/RausellLab/CelliD>.

20

21 **\body**

22 High-throughput single-cell technologies, such as single-cell RNA-seq, are currently being used for the complete  
23 cellular characterization of human tissues and cell types. The Human Cell Atlas project<sup>1</sup>, the NIH Human  
24 BioMolecular Atlas Program (HuBMAP)<sup>2</sup> and the LifeTime<sup>3</sup> initiative are remarkable examples of current scientific  
25 ambitions in this direction. One of the major goals of these studies is the identification of previously unknown  
26 cell types or cell states with potential physiological roles in health and disease. However, the computational  
27 characterization of cell heterogeneity is rendered more complex by the high dimensionality and high levels of  
28 technical and biological noise associated with single-cell measurements<sup>4</sup>. One common strategy for enhancing  
29 the signal-to-noise ratio involves the use of a low-dimensional representation of cells from which the most salient  
30 relative differences may emerge<sup>5</sup>. The techniques most widely used for this purpose include principal component  
31 analysis (PCA), independent component analysis (ICA), t-stochastic nuclear embedding (t-SNE) and uniform  
32 manifold approximation and projection (UMAP)<sup>6</sup>. Molecular characterization of the observed cell heterogeneity  
33 is then typically performed by analysing differential gene expression between the groups of cells resulting from  
34 clustering in low-dimensional space. However, this reliance on cell clusters results in gene signatures being  
35 assigned at cell subpopulation level rather than at full single-cell resolution. One of the major limitations of  
36 cluster-based approaches is that the gene signature analysis is bound to a certain level of resolution, at which  
37 the cell heterogeneity is partitioned into non-overlapping classes<sup>7</sup>. An exhaustive exploration of transcriptional  
38 heterogeneity requires a statistically robust per-cell gene signature assessment for each cell in a data set. No  
39 such approach has been described to date.

40 We present here Cell-ID, a multivariate approach to extracting a gene signature for each individual cell in a study  
41 (**Fig. 1 and Online Methods**). Cell-ID is based on multiple correspondence analysis (MCA), a statistical technique  
42 based on singular-value decomposition (SVD) that can provide simultaneous projections of individuals (e.g. cells)  
43 and variables (e.g. genes) in the same low-dimensional space<sup>8-10</sup>. This represents a major advantage over  
44 alternative low-dimensional transformations providing only cell projections. Originally MCA applies to binary and  
45 fuzzy-coded data, and this has been so far an obstacle for its use on *omics* data. We describe here a linear scaling  
46 of gene expression values unlocking the use of this technique for quantitative single-cell data analysis, while  
47 keeping the MCA mathematical properties. By representing cells and genes on the same principal axes, analytical  
48 distances can be calculated not only between cells and between genes, but also between each cell and each gene

49 **(Fig. 1A)**. In such a space, the closer a gene  $g$  is to a cell  $c$ , the more specific to such a cell it can be considered.  
50 Gene-to-cell distances can then be ranked for each individual cell, and the top-ranked genes may be regarded as  
51 a unique gene signature representing the identity card of the cell **(Fig. 1A)**. We show here that the unbiased per-  
52 cell gene signatures extracted by Cell-ID reproduce the gene signatures previously established for well-known  
53 cell types. Cell-ID signatures are, in turn, reproducible across independent datasets, despite strong batch effects,  
54 and can be used for cell identity tracing across different donors, model organisms, tissues-of-origin and single-  
55 cell omics protocols.

56 We first evaluated the consistency of MCA-based low-dimensional representations of cells and genes on 100  
57 simulated single-cell RNA-seq datasets **(Supplementary Note 1)**. Consistency was demonstrated at three levels.  
58 The cell representation achieved by MCA dimensionality reduction was largely equivalent to that achieved by  
59 principal component analysis (PCA) on the same dataset: Spearman's correlation coefficients, for the correlation  
60 between MCA and PCA coordinates, with median and standard deviation values across data sets of  $1.00 \pm 0.02$ ,  
61  $1.00 \pm 0.02$ ,  $0.99 \pm 0.02$ ,  $0.99 \pm 0.02$ , and  $0.99 \pm 0.02$  for their first five principal axes, respectively **(Supplementary**  
62 **Fig.1 and Supplementary Note 1)**. Second, the per-cell gene rankings provided by MCA are consistent with the  
63 expression values for the 50 neighbouring cells in MCA space, as reflected by their log-fold change in expression  
64 relative to the other cells **(Supplementary Fig. 2A and 2B)**. Third, MCA-based per-cell gene rankings are robust  
65 to high levels of dropout events. Thus, genes with no expression detected in a cell may nevertheless rank highly  
66 for the cell concerned if more strongly expressed in the surrounding cells than in more distant cells  
67 **(Supplementary Fig. 2A and 2B; Supplementary Note 1)**. Our results highlight the advantages of a multivariate  
68 approach in which per-cell gene assessments are implicitly weighted by their cell neighbourhood in low-  
69 dimensional space. Such patterns would be missed if per-cell gene rankings were naively obtained either (i) from  
70 the log-fold changes in expression in the target cell relative to background cells **(Supplementary Fig. 2C-D)** or (ii)  
71 from highest-to-lowest expression values within a cell, with random ranks for ties, as in another published  
72 approach (AUCell<sup>11</sup>, SCINA<sup>12</sup>) **(Supplementary Fig. 2E-F)**.

73 We also showed that Cell-ID could extract per-cell gene signatures, recovering characteristic marker lists  
74 associated with well-known cell types<sup>13</sup> **(Supplementary Note 2)**. To this end, we used two independent sets of  
75 human blood mononuclear cells for which individual cells were confidently annotated with an actual cell type  
76 through concomitant measurements of single-cell protein marker levels: (i) cord blood mononuclear cells

77 (CBMCs) profiled with a CITE-seq protocol<sup>14</sup>; and (ii) peripheral blood mononuclear cells (PBMCs) profiled with a  
78 REAP-seq protocol<sup>15</sup>. Cell-ID per-cell gene signatures were significantly enriched in the lists of markers associated  
79 with the corresponding cell type (**Figure 2**). This enrichment made it possible to recognize cell types with high  
80 rates of precision (87% and 90%) and recall (84% and 73%) for both datasets (multinomial  $p$ -value  $< 10^{-16}$  for all  
81 figures), outperforming reference methods for cell-type classification on the basis of marker lists (AUCELL<sup>11</sup> and  
82 SCINA<sup>12</sup> ; **Fig. 2B, Supplementary Fig. 3-4**). In more challenging scenarios, Cell-ID was capable of non-disjoint  
83 multi-class cell-type assignments capturing smooth transitions between hematopoietic differentiation states  
84 from the most immature hematopoietic stem cell (HSC) to downstream myeloid (CMP/GMP) and erythroid  
85 progenitors (MEP) (**Fig. 2, C-D**). It was consistently able to identify singleton cells, i.e. rare cell types represented  
86 by only one cell within a dataset (**Supplementary Note 2**). The capacity of Cell-ID to recover well-established cell  
87 types at single-cell resolution supports its use for automated cell-type annotation, even for extremely rare cells,  
88 without the need for clustering.

89 We benchmarked the capacity of Cell-ID to match cells of analogous cell types across independent single-cell  
90 RNA-seq datasets from the same tissue-of-origin, within and across species (**Supplementary Note 3**). Cell-ID  
91 matching across datasets is performed by a per-cell assessment in the *query* dataset evaluating the replication  
92 of gene signatures extracted from the *reference* dataset. Gene signatures from the *reference* dataset can be  
93 automatically derived either from individual cells (Cell-ID(c)), or from previously-established groups of cells (Cell-  
94 ID(g); **Methods**). We thus analysed independent human pancreatic islets<sup>16-18</sup> and human and mouse airway  
95 epithelium datasets<sup>19,20</sup> corresponding to multiple donors and diverse sequencing technologies (**Fig. 3** and  
96 **Supplementary Note 3**). Cell-ID(c) and Cell-ID(g) consistently yielded high precision (>94% and >92%), recall (  
97 >77% and >75%) and F1 values (>88% and >87%, respectively) across all evaluated *reference-to-query* cell-type  
98 assignments (multinomial  $p$ -value  $< 2.2e-16$ ; **Fig. 3A** and **Supplementary Fig. 5**). The overall performance of Cell-  
99 ID was at least as good as that of reference methods for cell matching and label transfer<sup>21-29</sup> (**Fig. 3A,**  
100 **Supplementary Fig. 6, A-B**), and salient scores were obtained for cell types observed at low frequencies (<2%):  
101 epsilon cells, tissue-resident macrophages, mast cells and endothelial cells from pancreatic islet samples<sup>17,18</sup>, and  
102 PNEC, brush cells and ionocytes in mouse and human airway epithelium datasets<sup>19,20</sup> (**Fig. 3B, Supplementary**  
103 **Fig. 6, C-D**). We further probed the capacity of this tool at individual cell resolution, by extracting Cell-ID  
104 signatures from 13 Schwann cells described in a dataset of 8,629 human pancreatic cells<sup>16</sup>. With these signatures,

105 we were able to recognize four and two prototypic Schwann cells that had previously gone unnoticed within two  
106 independent sets of 2126 and 2261 human pancreatic cells, respectively, obtained from different donors and  
107 with different sequencing protocols<sup>17,18</sup> (**Supplementary Figure 7**).

108 We then assessed the ability of Cell-ID to recognize rare cells from the same cell type across different tissues,  
109 and, thus, within diverse cell composition contexts (**Supplementary Note 4**). Thus, based on the unbiased gene  
110 signatures obtained from airway epithelium cells<sup>19</sup>, Cell-ID was able to identify brush/tuft cells, endocrine cells  
111 and goblet cells in the intestinal epithelium<sup>30</sup> with high precision (92%), recall (73%) and F1 scores (81%),  
112 outperforming reference methods for cell matching (multinomial  $p$ -value < 2.2e-16 for all figures; **Fig. 3 C-D**,  
113 **Supplementary Fig. 8**). From a discovery perspective, we used Cell-ID to perform cell-type scanning of two  
114 independent olfactory epithelium datasets<sup>31 32</sup> against brush/tuft signatures from the airway and the intestinal  
115 epithelium, which enabled us to identify putative rare solitary chemosensory cells (SCCs), a type of chemosensory  
116 cells closely related to brush/tuft cells, that had remained unclassified in the original publications. (**Fig. 3 E-F**,  
117 **Supplementary Fig. 9**). Thus, a total of 37 (<0.5%) and 5 (<0.6%) olfactory epithelium cells were found to display  
118 significant enrichment in airway and intestinal brush/tuft signatures (BH corrected  $p$  value < 10e-20), with a  
119 median of 29%±0.5 and 23%±0.4 of genes, respectively, in common. These cells displayed high levels of  
120 expression for the characteristic SCC marker genes *Il25* and *Gnat3*, and their Cell-ID gene signatures were  
121 significantly enriched in cysteinyl leukotriene biosynthesis genes, as reported by Ualiyeva et al<sup>33</sup> for SCCs  
122 (**Supplementary Table 1**). Our findings confirm, at single-cell resolution, the recently reported transcriptional  
123 and functional similarity between rare olfactory SCCs and brush/tuft cells from the airways and intestinal  
124 epithelia<sup>33,34</sup> (**Supplementary Note 4**).

125 Finally, we assessed the reproducibility of Cell-ID gene signatures across datasets profiled with different single-  
126 cell omics technologies: single-cell RNA-seq from the Tabula Muris mouse cell atlas<sup>35</sup> and single-cell ATAC-seq  
127 from the Mouse ATAC Atlas<sup>36</sup> (**Supplementary Note 5**). We benchmarked large-scale cell-type label transfer  
128 between the two expert-annotated atlases, collectively including 50 cell types from the eight tissues common to  
129 both: heart, kidney, liver, lung, bone marrow, spleen, thymus and large intestine (**Fig4. A-B, Supplementary Fig.**  
130 **10**). Both Cell-ID(c) and Cell-ID(g) matched cell types across scRNA-seq and sciATAC-seq datasets with high F1  
131 scores and, together with SingleR<sup>28</sup>, outperformed all the other reference methods evaluated (**Fig4. C-D**,  
132 **Supplementary Fig 11, Supplementary Figure 12**). The capacity of Cell-ID to extract **gene signatures** that are

133 robustly replicated across different single-cell *omics* technologies, in an automated manner, and its  
134 computational scalability (**Supplementary Note 6**) pave the way for the systematic multi-omics scanning of rare  
135 cell types across tissues and whole organisms.

136 Throughout this study, we found that Cell-ID was able to extract unbiased per-cell gene signatures from single-  
137 cell RNA-seq experiments that could then be used as unique cell identity cards without the need for prior  
138 knowledge. Cell-ID signatures were consistently reproducible, across diverse benchmarks collectively involving  
139 14 independent single-cell datasets, 13 organs / tissues, more than 50 cell types, more than 200000 cells, 2 model  
140 organisms, 6 sequencing protocols and 2 single-cell -omics technologies. Such signatures made it possible to  
141 recognize cell identities across datasets from the same or different tissues of origin and model organisms, while  
142 overcoming batch effects arising from the use of different donors and sequencing technologies. In particular, the  
143 automatic cell-type annotations and matching across sets provided by Cell-ID were fully transparent with respect  
144 to the genes driving the hits. Such transparency improves biological interpretation at the individual cell level  
145 (**Supplementary Note 7**), making it possible to discover *bona fide* rare cells, as illustrated by the identification of  
146 Schwann cells and solitary chemosensory cells previously overlooked in published datasets (**Supplementary Note**  
147 **3 and 4**). This contrasts strongly with the capabilities of the other methods currently available, which are based  
148 on assessments of similarity over the entire transcriptome<sup>21</sup>, cell embedding<sup>22,29</sup>, or machine-learning  
149 approach<sup>23,24,26,27</sup>, in which the contributions of individual genes are difficult to interpret.

150 Cell-ID is computationally efficient, can be scaled up for use with large datasets and allows *many-to-many* dataset  
151 comparisons without the need for data integration steps (**Supplementary Note 6**). It can, thus, be used for the  
152 systematic screening of each individual cell in newly sequenced datasets against (a) reference databases of  
153 marker lists associated with well-established cell types (e.g. PanglaoDB<sup>37</sup>, CellMarkers<sup>38</sup>), (b) reference single-cell  
154 atlas databases with manually curated cell-type annotations (Tabula Muris<sup>35</sup>, mouse ATAC atlas<sup>36</sup>, human cell  
155 atlas<sup>1</sup>), and (c) molecular signature collections, functional ontologies and pathway databases (e.g. MSigDB<sup>39</sup>,  
156 GO<sup>40</sup>, Reactome<sup>41</sup>, KEGG<sup>42</sup>, Wikipathways<sup>43</sup>). The automatic single-cell annotations provided by Cell-ID will  
157 alleviate the need for the often-tedious visual inspection of prototypic marker levels based on expert knowledge.  
158 From a discovery perspective, the identification of individual cells presenting distinctive and reproducible gene  
159 signatures consistent with the phenotypes studied would constitute a first step towards the in-depth  
160 experimental characterization of putative novel human cell types or cell states in health and disease. Cell-ID is

161 implemented as an open-source R software package including detailed tutorials and scripts to reproduce all the  
162 analyses and figures presented in this manuscript (**Supplementary Software** and  
163 <https://github.com/RausellLab/CelliD>).

164

#### 165 **Acknowledgements**

166 We thank the Laboratory of Clinical Bioinformatics and the Laboratory of Human Lymphohematopoiesis for  
167 helpful discussions and support. The Laboratory of Clinical Bioinformatics was partly supported by the French  
168 National Research Agency (ANR) “Investissements d’Avenir” Program (Grant ANR-10-IAHU-01). The Laboratory  
169 of Clinical Bioinformatics and the Laboratory of Human Lymphohematopoiesis were partly supported by Christian  
170 Dior Couture, Dior. We also thank Dr. Gloria Fuentes from The Visual Thinker LLP for the creation of the  
171 illustrations in Figures 1-4.

172

#### 173 **Author contributions**

174 A.C. and A.R. conceived and designed research; A.C. performed research; A.C and L.M. contributed with  
175 materials/analysis tools; A.C and A.R. analyzed data; A.C, E.S and A.R. interpreted results; and A.C., E.S. and A.R.  
176 wrote the paper. All authors read and approved the final manuscript.

177

#### 178 **Competing interests**

179 Authors declare that they have no competing financial and/or non-financial interests, or other interests that  
180 might be perceived to influence the results and/or discussion reported in this paper.

181

## 182 References

- 183 1. Teichmann, S. *et al.* The Human Cell Atlas. *eLife* **6**, (2017).
- 184 2. The Human BioMolecular Atlas Program - HuBMAP | NIH Common Fund. <https://commonfund.nih.gov/HuBMAP>.
- 185 3. The LifeTime Initiative. *LifeTime FET Flagship* <https://lifetime-fetflagship.eu/>.
- 186 4. Lähnemann, D. *et al.* Eleven grand challenges in single-cell data science. *Genome Biol.* **21**, 31 (2020).
- 187 5. Sun, S., Zhu, J., Ma, Y. & Zhou, X. Accuracy, robustness and scalability of dimensionality reduction methods for single-cell
- 188 RNA-seq analysis. *Genome Biol.* **20**, 269 (2019).
- 189 6. Becht, E. *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.* **37**, 38–44 (2019).
- 190 7. Kiselev, V. Y., Andrews, T. S. & Hemberg, M. Challenges in unsupervised clustering of single-cell RNA-seq data. *Nat. Rev.*
- 191 *Genet.* **20**, 273–282 (2019).
- 192 8. Greenacre, M. J. *Theory and applications of correspondence analysis.* (1984).
- 193 9. *Multiple correspondence analysis and related methods. Selected papers based on the presentations at the international*
- 194 *conference (CARME 2003), Barcelona, Spain, 29 June to 2 July 2003. Chapman & Hall/CRC Statistics in the Social and*
- 195 *Behavioral Sciences Series* (Chapman & Hall/CRC, 2006).
- 196 10. Aşan, Z. & Greenacre, M. Biplots of fuzzy coded data. *Fuzzy Sets Syst.* **183**, 57–71 (2011).
- 197 11. Aibar, S. *et al.* SCENIC: Single-cell regulatory network inference and clustering. *Nat. Methods* **14**, 1083–1086 (2017).
- 198 12. Zhang *et al.* SCINA: Semi-Supervised Analysis of Single Cells in Silico. *Genes* **10**, 531–531 (2019).
- 199 13. Aran, D., Hu, Z. & Butte, A. J. xCell: digitally portraying the tissue cellular heterogeneity landscape. *Genome Biol.* **18**, 220
- 200 (2017).
- 201 14. Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells. *Nat. Methods* **14**, 865–868
- 202 (2017).
- 203 15. Peterson, V. M. *et al.* Multiplexed quantification of proteins and transcripts in single cells. *Nat. Biotechnol.* **35**, 936–939
- 204 (2017).
- 205 16. Baron, M. *et al.* A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and Intra-cell
- 206 Population Structure. *Cell Syst.* **3**, 346–360 (2016).
- 207 17. Segerstolpe, Å. *et al.* Single-Cell Transcriptome Profiling of Human Pancreatic Islets in Health and Type 2 Diabetes. *Cell*
- 208 *Metab.* **24**, 593–607 (2016).
- 209 18. Muraro, M. J. *et al.* A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Syst.* **3**, 385–394.e3 (2016).
- 210 19. Plasschaert, L. W. *et al.* A single-cell atlas of the airway epithelium reveals the CFTR-rich pulmonary ionocyte. *Nature*
- 211 **560**, 377–381 (2018).

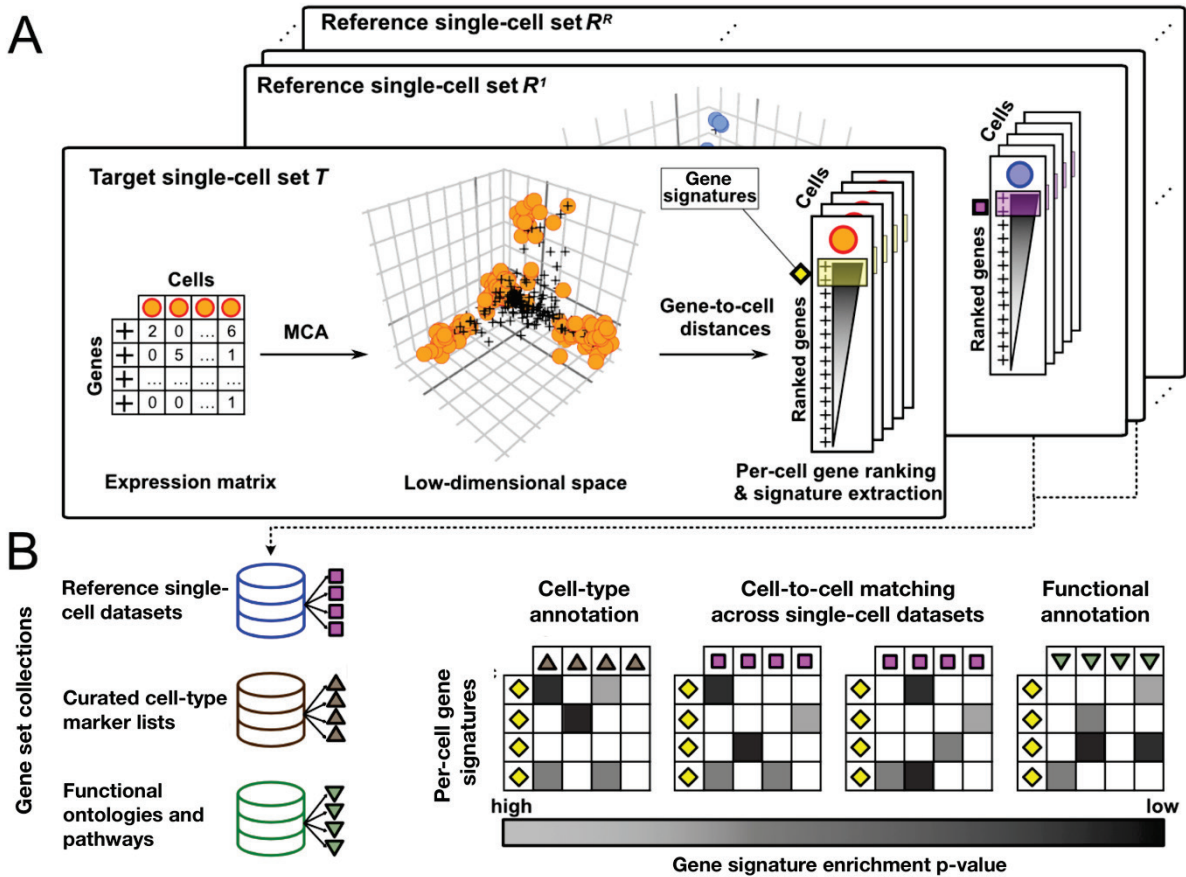


- 212 20. Montoro, D. T. *et al.* A revised airway epithelial hierarchy includes CFTR-expressing ionocytes. *Orit Rozenblatt-Rosen* **4**,  
213 19–19 (2018).
- 214 21. Kiselev, V. Y., Yiu, A. & Hemberg, M. scmap: projection of single-cell rna-seq data across data sets. **15**, 359–359 (2018).
- 215 22. Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-sequencing data are corrected  
216 by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–427 (2018).
- 217 23. De Kanter, J. K., Lijnzaad, P., Candelli, T., Margaritis, T. & Holstege, F. C. P. CHETAH: a selective, hierarchical cell type  
218 identification method for single-cell RNA sequencing. *Nucleic Acids Res.* **47**, (2019).
- 219 24. Lieberman, Y., Rokach, L. & Shay, T. CaSTLe – Classification of single cells by transfer learning: Harnessing the power of  
220 publicly available single cell RNA sequencing experiments to annotate new experiments. *PLOS ONE* **13**, e0205499–  
221 e0205499 (2018).
- 222 25. Boufeaa, K., Seth, S. & Batada, N. N. scID Uses Discriminant Analysis to Identify Transcriptionally Equivalent Cell Types  
223 across Single-Cell RNA-Seq Data with Batch Effect. *iScience* **23**, 100914 (2020).
- 224 26. Tan, Y. & Cahan, P. SingleCellNet: A Computational Tool to Classify Single Cell RNA-Seq Data Across Platforms and Across  
225 Species. *Cell Syst.* **9**, 207-213.e2 (2019).
- 226 27. Alquicira-Hernandez, J., Sathe, A., Ji, H. P., Nguyen, Q. & Powell, J. E. ScPred: Accurate supervised method for cell-type  
227 classification from single-cell RNA-seq data. *Genome Biol.* **20**, 264–264 (2019).
- 228 28. Aran, D. *et al.* Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage. *Nat.*  
229 *Immunol.* **20**,
- 230 29. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21 (2019).
- 231 30. Haber, A. L. *et al.* A single-cell survey of the small intestinal epithelium. (2017) doi:10.1038/nature24489.
- 232 31. Wu, Y. *et al.* A Population of Navigator Neurons Is Essential for Olfactory Map Formation during the Critical Period  
233 Article A Population of Navigator Neurons Is Essential for Olfactory Map Formation during the Critical Period. *Neuron*  
234 **100**, 1066-1082.e6 (2018).
- 235 32. Fletcher, R. B. *et al.* Deconstructing Olfactory Stem Cell Trajectories at Single-Cell Resolution. *Cell Stem Cell* **20**, 817-  
236 830.e8 (2017).
- 237 33. Ualiyeva, S. *et al.* Airway brush cells generate cysteinyl leukotrienes through the ATP sensor P2Y2. *Sci. Immunol.* **5**,  
238 eaax7224–eaax7224 (2020).
- 239 34. Bankova, L. G. *et al.* The cysteinyl leukotriene 3 receptor regulates expansion of IL-25–producing airway brush cells  
240 leading to type 2 inflammation. *Sci. Immunol.* **3**, (2018).
- 241 35. Schaum, N. *et al.* Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris the tabula Muris consortium\*.  
242 (2018) doi:10.1038/s41586-018-0590-4.

- 243 36. Cusanovich, D. A. *et al.* A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility. *Cell* **174**, 1309-1324.e18  
244 (2018).
- 245 37. Franzén, O., Gan, L.-M. & Björkegren, J. L. M. PanglaoDB: a web server for exploration of mouse and human single-cell  
246 RNA sequencing data. *Database* **2019**, (2019).
- 247 38. Zhang, X. *et al.* CellMarker: a manually curated resource of cell markers in human and mouse. *Nucleic Acids Res.* **47**,  
248 D721–D728 (2019).
- 249 39. Liberzon, A. *et al.* The Molecular Signatures Database Hallmark Gene Set Collection. *Cell Syst.* **1**, 417–425 (2015).
- 250 40. The Gene Ontology (GO) database and informatics resource. *Nucleic Acids Res.* **32**, D258–D261 (2004).
- 251 41. Jassal, B. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498–D503 (2020).
- 252 42. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for gene and protein  
253 annotation. *Nucleic Acids Res.* **44**, 457–462 (2015).
- 254 43. Slenter, D. N. *et al.* WikiPathways: a multifaceted pathway database bridging metabolomics to other omics research.  
255 *Nucleic Acids Res.* **46**, D661–D667 (2018).
- 256

257 **Figures**

258



259

260

261

262 **Figure 1. Overview of the Cell-ID approach. (A)** Cell-ID performs a dimensionality reduction of the gene

263 expression matrix through multiple correspondence analysis (MCA), where both cells and genes are projected in

264 a common orthogonal space. In such space, the closer a gene is to a cell the more specific it is to it. Thus, a gene

265 ranking is obtained for each cell in a dataset based on their distance in the MCA space. The top-ranked genes for

266 a given cell define its gene signature, which can be regarded as a unique cell identity card. Per-cell gene signatures

267 can be independently extracted for a collection of single-cell datasets for downstream analyses. **(B)** Per-cell gene

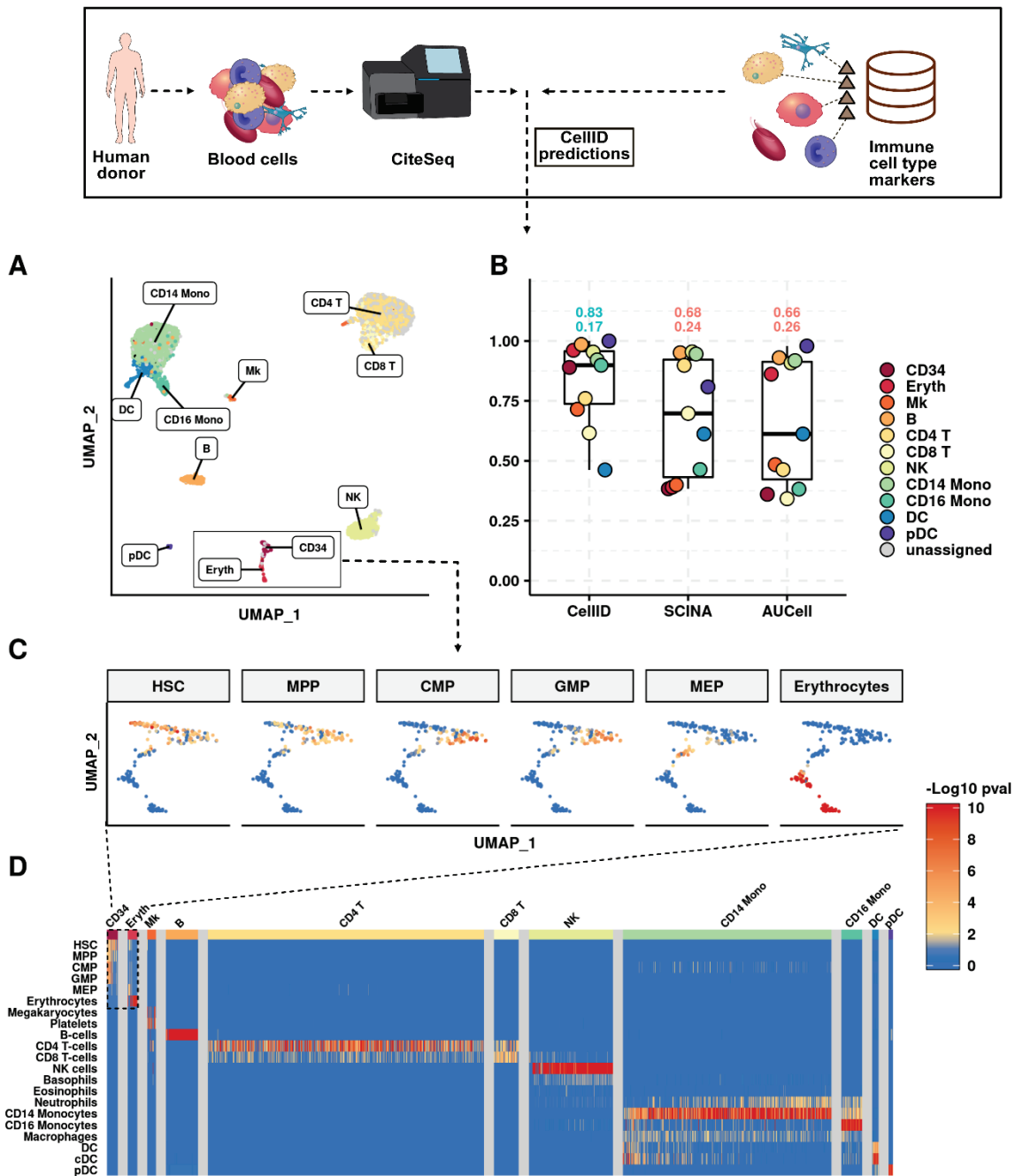
268 signatures from a dataset can be evaluated for their enrichment against (i) collections of pre-established cell type

269 markers, in order to perform automatic cell type annotation, (ii) per-cell gene signatures from independent

270 single-cell datasets, allowing cell matching and label transferring, and (iii) gene sets representing biological

271 functions or molecular pathways, allowing functional annotation and interpretation of cell states.

272

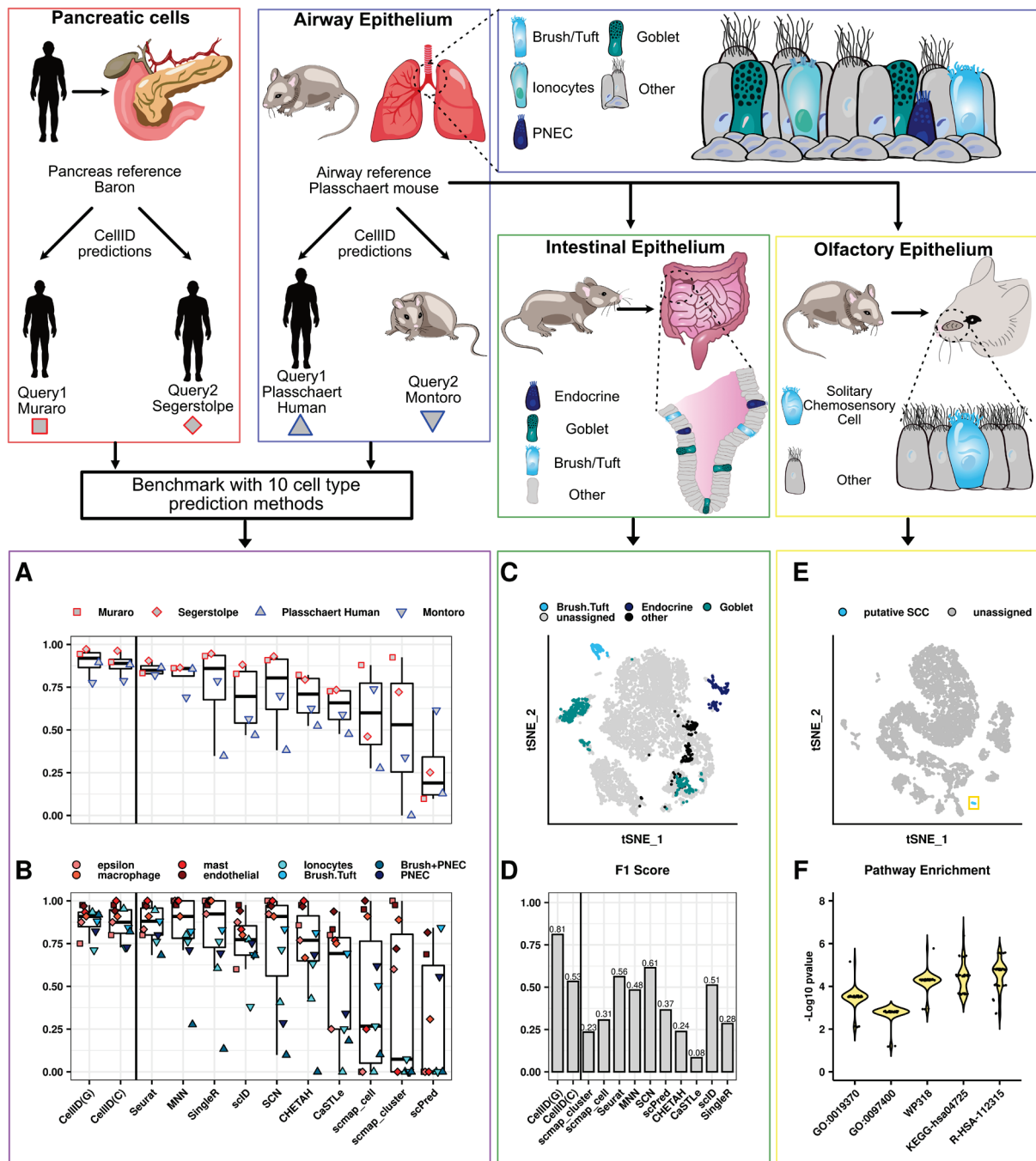


273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288

**Figure 2. Cell-ID cell type prediction of human cord blood mononuclear cells using pre-established marker lists.**

(A) UMAP representation of 8005 cord blood mononuclear cells profiled by CITE-seq<sup>3</sup>. Dots representing cells are colored according to Cell-ID cell type predictions using pre-established immune cell signatures, as indicated by the labels in the figure. (B) Performance measured through the F1 score achieved by Cell-ID, AUCell and SCINA cell type predictions for each of the blood cell types reported in the original publication<sup>3</sup>. Boxplots summarize the F1 scores for each method. The numbers above boxplots denote the global performance (macro F1 score, upper digits) and its standard deviation (lower digits), where the maximum and minimum values across methods are colored in black and grey, respectively. (C) Zoomed UMAP representation on Erythrocytes and CD34+ cells showing that the Cell-ID multi-class cell assignments capture transient cell states consistent with the cell-type hierarchy associated with immature hematopoietic stem cell differentiation. Cells are color-coded according to the -log<sub>10</sub> enrichment p-value obtained by Cell-ID in tests of the association of their gene signature with the cell-type signatures associated to their pre-cursor cell types: HSC, MPP, CMP, GMP, MEP and erythrocytes. The color scale for cells extends from white (*p* value=1) to dark red (*p* value = 1e-10), with *p*-values < 1e-10 fixed at this value). (D) Heatmap representing, for each individual cell (displayed in columns), the -log<sub>10</sub> transformed *p*-value

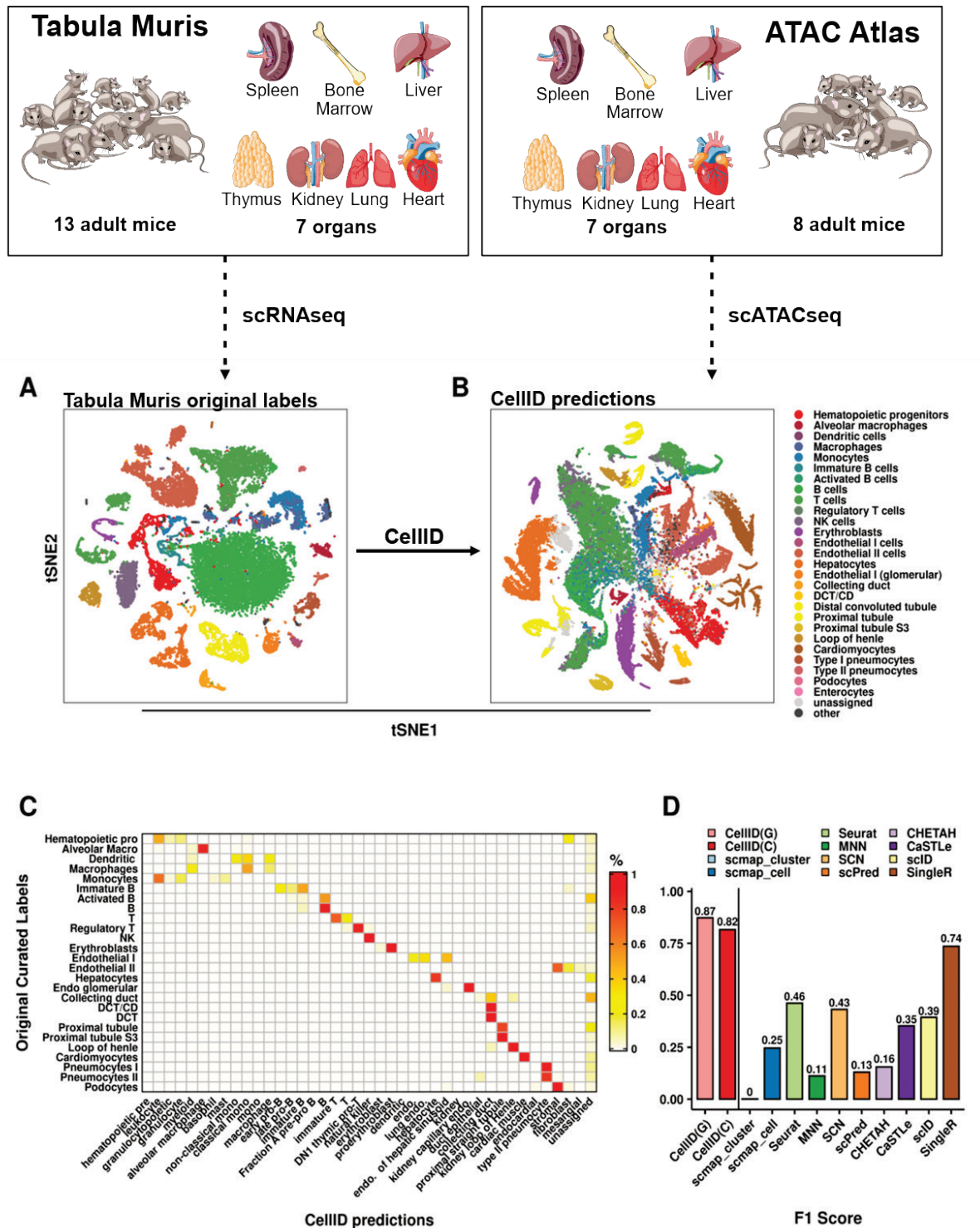
289 obtained by Cell-ID in tests of the association of the gene signature with each of the evaluated pre-established  
290 marker lists, representing a total of 21 blood cell types (displayed in rows). The color code of the heatmap extends  
291 from dark blue ( $p$  value=1) to yellow ( $p$  value =  $10^{-2}$ ) to dark red ( $p$  value =  $10^{-10}$ ), with  $p$  values  $<10^{-10}$  fixed at this  
292 value). Non-significant associations ( $p$  value  $>10^{-2}$  after Benjamini Hochberg correction for the number of gene  
293 signatures tested) are shown in blue. The columns in the heatmaps were grouped by the reference cell-type  
294 label, as indicated by the colored bands at the top and in the associated legend. CD34: CD34+ hematopoietic  
295 stem cells; Eryth: Erythrocytes; Mk: Megakaryocytes, B: B cells, CD4 T: CD4+ T cells, CD8: CD8+ T cells, CD14:  
296 CD14+ monocytes, CD16 Mono: CD16+ monocytes, NK: natural killer cells, DC: dendritic cells, pDC: plasmacytoid  
297 dendritic cell, HSC: hematopoietic stem cells, MPP: multi-potent-progenitor, CMP: common-myeloid-  
298 progenitors, GMP: granulocyte-monocyte-progenitor, MEP: megakaryocyte-erythrocyte-progenitor.  
299



300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315

**Figure 3. Performance of Cell-ID cell-to-cell matching across independent scRNA-seq datasets from the same or different tissue of origin, within and across species. (A and B)** Performance measured through the F1 score (y-axis) achieved by Cell-ID(g), Cell-ID(c) and 10 alternative state-of-the-art methods (x-axis), covering the major approaches for cell-matching or label transfer across scRNA-seq datasets (**Supplementary Table 9**). **(A)** The performance for each method is represented for each of the label transferring evaluated (as schematically represented in the top left panels), corresponding to cell-to-cell matching across datasets from pancreatic islets cells (red squares and red diamonds), and across datasets from airway epithelium (blue triangles). Boxplots summarize the global performance (macro F1 scores) for each method. **(B)** The performance for each method is represented for each of the rare cell types reported in the original publications associated to the pancreatic islets cells (squares and diamonds) and airway epithelium datasets (blue triangles) evaluated in **(A)**. Rare cell types are represented following the color palette indicated in the legend and representing epsilon, macrophages, mast, and endothelial cells (pancreatic cells), and ionocytes, brush/tuft, and PNEC (airway epithelium). A cell type label gathering together brush and PNEC cells was used for consistency with the labelling provided in the human sample from Plasschaert et. al. **(C)** t-SNE representation of 7216 cells from mouse small intestinal epithelium<sup>29</sup>.

316 Dots representing cells are colored according to Cell-ID(g) cell type predictions, using as a reference the mouse  
317 airway epithelial gene signatures extracted from Plasschaert et al, as schematically represented in the panels  
318 above. Intestinal epithelium cell types with significant enrichment p-values are colored in green (goblet),  
319 brush/tuft (light blue) and dark blue (endocrine). Cells with significant enrichments in airway epithelial signatures  
320 without an analogous cell type in intestinal epithelium are represented in black. Intestinal epithelium cells with  
321 no significant enrichments were left unassigned and are displayed in grey. For comparison, the associated manual  
322 cell type annotations provided in <sup>29</sup> are represented in **Supplementary Figure 10 A. (D)** F1 score (y-axis) achieved  
323 by Cell-ID(g), Cell-ID(c) and 10 alternative state-of-the-art methods (x-axis), for the label transferring depicted in  
324 **h. (E)** UMAP representation of 9126 mouse olfactory epithelium cells from <sup>30</sup>. Dots representing cells are colored  
325 according to Cell-ID(g) cell type predictions using as a reference the mouse brush/tuft gene signatures extracted  
326 from (i) mouse airway epithelium, and (ii) mouse small intestinal epithelium, as schematically represented in the  
327 panels above. The 37 cells significantly enriched with airway Brush/Tuft gene signatures are highlighted in blue  
328 and were interpreted as putative SCCs. Identical results were obtained when using intestinal Brush/Tuft gene  
329 signatures. **(F)** Violin plots corresponding to the distribution of -log<sub>10</sub> enrichment p-values (y-axis) across the 37  
330 cells identified in **(E)** (black dots) are represented for 5 significant functional terms (x-axis). GO-0019370: Gene  
331 Ontology term “leukotriene biosynthetic process”; GO-0097400: Gene Ontology term “interleukin-17-mediated  
332 signalling pathway”; WP318: WikiPathways term “Eicosanoid Synthesis”; KEGG-hsa04725: KEGG term  
333 “Cholinergic synapse”; R-HSA-112315: Reactome term “Transmission across Chemical Synapses”.  
334



335  
336

337 **Figure 4. Performance of Cell-ID cell-to-cell matching across independent datasets from different single-cell**  
338 **omics technologies: scRNAseq and scATACseq**

339 **(A)** t-SNE representation of 25332 cells from 7 organs profiled with scRNAseq 10X genomics from the Tabula  
340 Muris mouse cell atlas. Cells are colored according to the manually annotated cell types provided by the atlas,  
341 regrouped by the categories represented in the legend in **(B)** and described in **Supplementary Table 10**. **(B)** t-  
342 SNE representation of 50284 cells from 7 organs profiled with scATAC-seq from the Mouse ATAC atlas. Cells are  
343 colored according to the Cell-ID(g) predictions using the group gene signatures extracted from **(A)**, as  
344 represented in the legend. **(C)** Heatmap representing the confusion matrix between the manually curated cell  
345 type annotations from the Mouse ATAC atlas (displayed per rows) and the Cell-ID(g) cell type predictions  
346 (displayed per columns) using the gene-signatures extracted from the manually annotated cell types provided by



347 the Tabula Muris mouse cell atlas. The color code in the heatmap represents the ratio  $r$  of the cell types displayed  
348 per rows that are allocated in the cell types represented per columns, ranging from white ( $r = 0$ ) to red ( $r = 1$ ). **(D)**  
349 Global performance measured through the macro F1 score (y-axis) achieved by Cell-ID(g), Cell-ID(c) and 10  
350 alternative state-of-the-art methods (x-axis) on the label transferring from the scRNA-seq to the scATAC-seq  
351 mouse cell atlas. DCT/CD: Distal convoluted tubule/ collecting duct, endo: endothelial, mono: monocytes.  
352

## 353 **Online methods**

### 354 **Data availability**

355 All single-cell data sets used in this paper are publicly available (**Supplementary Table 2**). scRNA-seq datasets for  
356 human blood cells profiled by Cite Seq<sup>1</sup> and Reap Seq<sup>2</sup> were downloaded from the gene expression omnibus  
357 (GEO; accession numbers GSE100866 and GSE100501, respectively). Cell-type labels for these two datasets were  
358 obtained following the Multimodal Analysis vignette of the *Seurat*<sup>3</sup> R package  
359 ([https://satijalab.org/seurat/multimodal\\_vignette.html](https://satijalab.org/seurat/multimodal_vignette.html)). Pancreas scRNA-seq datasets from Baron<sup>4</sup>, Muraro<sup>5</sup>,  
360 and Segerstolpe<sup>6</sup>, as well as their associated cell-type annotations were downloaded via the scRNAseq<sup>7</sup> R package  
361 as a SingleCellExperiment format R object. Plasschaert<sup>8</sup> mouse and human and Montoro<sup>9</sup> mouse airway  
362 epithelium scRNA-seq datasets, and their annotations were downloaded from GEO (GSE102580, GSE103354).  
363 Haber<sup>10</sup> intestinal epithelium scRNA-seq dataset was downloaded from GEO accession code GSE92332. Olfactory  
364 epithelium scRNA-seq datasets from Fletcher<sup>11</sup> and Wu<sup>12</sup> were downloaded from GEO (GSE95601, GSE120199),  
365 and their cell type annotations were obtained from the associated Github repositories:  
366 <https://github.com/rufletch/p63-HBC-diff> and <https://www.stowers.org/research/publications/odr> for  
367 Fletcher<sup>11</sup> and Wu<sup>12</sup>, respectively. Tabula Muris<sup>13</sup> 10X and Smart-seq mouse scRNAseq datasets were  
368 downloaded from <https://tabula-muris.ds.czbiohub.org/>. Gene activity score matrices from the Mouse sci-ATAC-  
369 seq atlas datasets from Cusanovich<sup>14</sup> were obtained from <http://atlas.gs.washington.edu/mouse-atac/data/>, as  
370 provided by the authors and resulting from the aggregation of information across all differentially accessible  
371 chromatin sites linked to a target gene.

### 372 **Preprocessing and normalization of single-cell datasets**

373 All single-cell RNA-seq datasets analyzed in the study took as input the raw count gene expression matrices  
374 provided by the original sources. Library size normalization was carried out by rescaling counts to a common  
375 library size of 10000. Log transformation was performed after adding a pseudo-count of 1. All analyses  
376 throughout the manuscript were restricted to a background set of 19308 and 21914 protein-coding genes from  
377 human and mouse, respectively, obtained from BioMart Ensembl release 100, version April 2020 (GrCH38.p13  
378 for human, and GRCm38.p6 for mouse<sup>15,16</sup>). Genes expressing at least one count in less than 5 cells were  
379 removed. No filtering of cells was done unless the original sources provided “doublet” or contamination

380 annotations, which were in such cases filtered out. In the case of sci-ATAC gene activity score matrices, the same  
381 preprocessing as for scRNA-seq data was applied.

## 382 **Overview of the Cell-ID approach.**

383 The main steps of the Cell-ID workflow are schematized in **Figure 1**. First, the library size normalized and log-  
384 transformed count matrix is transformed into a fuzzy coded indicator matrix where expression values are  
385 represented in a continuous scale between 0 and 1. Second, Cell-ID performs a dimensionality reduction of the  
386 indicator matrix using Multiple Correspondence Analysis where both cells and genes are represented into the  
387 same vector space<sup>17</sup>. Third, per-cell gene rankings are calculated from the gene-to-cell distances in MCA space,  
388 where the top closest genes to a cell will define its gene signature. If a grouping of cells is provided, per-group  
389 gene rankings may be obtained in an analogous way by using the geometric centroid in MCA space of the cells  
390 belonging to a given group. The enrichment of per-cell and/or per-group gene signatures is then evaluated  
391 through hypergeometric tests against (i) reference marker gene lists, and/or (ii) per-cell and/or per-group gene  
392 signatures extracted through Cell-ID from reference single-cell datasets. Per-cell and per-group gene signatures  
393 represent thus identity cards allowing automatic cell type and functional annotation, as well as cell matching  
394 across datasets. Each of these steps is described in detail in the following sections:

## 395 **Multiple Correspondence Analysis of the gene expression matrix**

396 Multiple Correspondence Analysis (MCA) is a multivariate descriptive statistical technique conceptually  
397 equivalent to Principal Component Analysis for qualitative/binary data<sup>18,19</sup>. MCA can be applied to quantitative  
398 data through an intermediate step of a so-called *fuzzy coding*. Here, each continuous variable  $p$  is coded through  
399 user-defined functions into a number of disjoint categories  $Q_p$  where membership  $x$  to each category  $q$  is  
400 represented in a continuous scale between 0 and 1, and  $\sum_{j=1}^{Q_p} x_j = 1$ . Following (Aşan and Greenacre, 2011),  
401 fuzzy-coding of a cases-by-variables matrix of continuous data can be performed in its simplest form by doubling  
402 each variable into  $Q_p = 2$  categories as follows: Let  $M_{N,P}$  be the gene expression matrix of  $N$  cells (i.e. cases) and  
403  $P$  genes (i.e. variables), with general term  $m_{np}$  gathering the expression level of gene  $p$  in cell  $n$ . For each column  
404 vector  $Mp$  in  $M$ , two membership functions  $x^+ = f^+(m): \mathfrak{R} \rightarrow \mathfrak{R}: [0,1]$  and  $x^- = f^-(m): \mathfrak{R} \rightarrow \mathfrak{R}: [0,1]$ , where  
405  $f^-(m) = 1 - f^+(m)$ , can be defined by linearly scaling between 0 and 1 the expression values for each gene  
406 across all cells as follows:

407 
$$x^+_{np} = \frac{m_{np} - \min(M_p)}{\max(M_p) - \min(M_p)} ; \quad x^-_{np} = 1 - x^+_{np} ;$$

408

409 From such functions, a fuzzy-coded indicator matrix  $X_{N,K}$  can be built, representing a total of  $K=2P$  categories :

410 
$$X_{N,K} = \begin{bmatrix} x^+_{11} & x^-_{11} & \dots & x^+_{1P} & x^-_{1P} \\ x^+_{21} & x^-_{21} & \dots & x^+_{2P} & x^-_{2P} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x^+_{N1} & x^-_{N1} & \dots & x^+_{NP} & x^-_{NP} \end{bmatrix}$$

411 The grand total of  $X_{N,K}$  is thus  $N*P$ , since each of the  $N$  cells has  $P$  sets of fuzzy coded-values, each adding up to

412 1. The Multiple Correspondence Analysis (MCA) of a fuzzy-coded indicator matrix follows that of a regular MCA<sup>20</sup>.

413 From the matrix  $X_{N,K}$ , a matrix of relative frequencies  $F_{N,K}$  is defined as :

414 
$$F_{N,K} = \frac{1}{NP} X$$

415 From the row sums and column sums of  $F$ , two diagonal matrices  $D_r$  and  $D_c$  are build, respectively, with general

416 terms

417 
$$d_{r,n,n} = \sum_{k=1}^K f_{nk} ; \quad d_{c,k,k} = \sum_{n=1}^N f_{nk}.$$

418 Let  $S_{NK}$  be the matrix of standardized relative frequencies resulting from

419 
$$S = D_r^{-1/2} F D_c^{-1/2}$$

420 The singular-value decomposition (SVD) of the matrix  $S_{NK}$  leads to

421 
$$S = U D_\alpha V^T$$

422 Where  $U$  and  $V$  contain by columns the singular vectors of norm 1 ( $U^T U = 1 ; V^T V = 1$ ), and  $D_\alpha$  is a diagonal matrix

423 of singular values  $\alpha_i$ , which are positive and displayed in descending order:  $\alpha_1 \geq \alpha_2 \geq \dots > 0$ . Alternatively,  $U$

424 and  $V$  can be obtained as the matrices displaying by columns the eigen vectors of norm 1 of the product  $SS^T$  and

425  $S^T S$ , respectively, with eigen values  $\lambda_i$ , where  $\alpha_i = \lambda_i^{1/2}$ . Thus,

426 
$$SS^T U = U D_\lambda ; \quad S^T S V = V D_\lambda ; \quad U^T U = 1 ; \quad V^T V = 1$$

427 Alternatively,  $V$  can be calculated from  $U$  with the transition formula:

428 
$$V = S^T U D_\alpha^{-1} = S^T U D_\lambda^{-1/2}$$

429 In the previous expressions, the first vectors  $\vec{u}_1$  and  $\vec{v}_1$  are associated with the trivial solution  $\alpha_1 = \lambda_1 = 1$ ; and  
430 are thus removed from the analysis at this stage. After eliminating the trivial solution, the sum of all the  
431 eigenvalues  $\lambda_i$  from  $SS'$  (so-called *total Inertia* in MCA terminology) equals the Chi-squared statistic  $\chi^2$  of the  
432 indicator matrix  $X_{N,K}$  divided by  $N$ . Thus, the orthogonal vectoral space generated by the eigenvectors of  $SS'$  can  
433 be viewed as a decomposition of the Chi-squared statistic  $\chi^2$  in its independent sources of variation, each  
434 accounting for a fraction given by  $\lambda_i / \sum_{i=1}^I \lambda_i$ . At this stage, further dimensionality reduction can be performed  
435 by retaining the first  $J$  eigen vectors as the most informative components, while disregarding the rest of  
436 dimensions from downstream analysis. Here we established  $J = 50$  as the default parameter throughout all the  
437 analyses performed.

438 The orthogonal dimensions given by the eigenvectors  $U$  and  $V$  allow to simultaneously represent both rows (i.e.  
439 cells) and columns (i.e. gene categories) in  $S$  into the same orthogonal vectoral space, where coordinates are  
440 obtained as follows<sup>21</sup>:

441 Row coordinates:  $\Phi = D_r^{-1/2} U = D_r^{-1/2} S V D_\alpha^{-1}$

442 Column coordinates:  $G = D_c^{-1/2} S^T U = D_c^{-1/2} V D_\alpha$

443

444 The previous expressions correspond to so-called *standard* and *principal* coordinates for rows and columns,  
445 respectively, in MCA terminology<sup>21</sup>. The simultaneous representation of cells and genes into the same vector  
446 space is a main advantage of MCA as compared to alternative dimensionality reduction techniques such as PCA,  
447 where only cells are projected.

448

#### 449 **Gene-to-cell distances in MCA space and per-cell and per-group gene signature extraction**

450 In the vectoral space provided by MCA, a barycentric relationship is fulfilled between the rows and the column  
451 coordinates: the general term  $g_{kj}$  of  $G$  representing the coordinate of a column  $k$  in the dimension represented  
452 by the eigenvector  $\vec{u}_j$  of  $U$ , corresponds to the weighted average (centroid) of the  $N$  row coordinates  $\phi_{nj}$  from  
453  $\Phi$ , where weights are given by  $x_{nk} / \sum_{n=1}^N x_{nk}$ , i.e. the frequency conditioned by columns of the corresponding  
454 values in the fuzzy indicator matrix  $X_{N,K}$ :

455 
$$g_{kj} = \frac{1}{\sum_{n=1}^N x_{nk}} \sum_{n=1}^N x_{nk} * \phi_{nj}$$

456 Thus, in MCA space the closer a column (i.e. gene category) is to a row (i.e. gene), the more specific it is to it. In  
 457 addition, each set of  $Q_p = 2$  categories for each gene  $p$  is centered at the origin: i.e.  $\vec{g}_p^+ + \vec{g}_p^- = \vec{0}$ . At this stage,  
 458 only gene category coordinates  $\vec{g}_p^+$ , conveying presence of gene expression relative to the maximum per gene,  
 459 are retained for downstream analysis. From the previous expressions, the Euclidean distances  $d_{np}(\vec{\phi}_n, \vec{g}_p^+)$  can  
 460 be computed for each cell  $n$  and each gene  $p$  in the dataset. The genes  $g_p$  constituting the signature  $\Gamma_n$  associated  
 461 to a cell  $n$  are obtained from its top  $\gamma$  closest genes in MCA space:

462 
$$\Gamma_n = \left\{ g_p \mid \forall p: \text{rank}_p(d_{np}(\vec{\phi}_n, \vec{g}_p^+) \leq \gamma) \right\}$$

463  
 464 A default value of  $\gamma = 200$  was established throughout this work and ties resolved with random ranks. More  
 465 generally, the genes  $g_p$  constituting the signature  $\Gamma_\Theta$  associated to a group of cells  $\Theta$  can be obtained from the  
 466 Euclidean distances  $d_p(\vec{\phi}_\Theta, \vec{g}_p^+)$  between each gene  $p$  and the group centroid  $\vec{\phi}_\Theta$ , obtained from the geometric  
 467 center of the  $\vec{\phi}_n$  vectors associated to the cells  $n \in \Theta$ .

468 
$$\Gamma_\Theta = \left\{ g_p \mid \forall p: \text{rank}_p(d_{np}(\vec{\phi}_\Theta, \vec{g}_p^+) \leq \gamma) \right\}$$

469 We note however that Cell-ID does not perform or relies in any clustering step whatsoever. Notwithstanding,  
 470 cell grouping information may optionally be used here as input, as provided by a external reference source (e.g.  
 471 database or publication).

472 **Per-cell gene signature enrichment analyses against reference gene sets**

473 The gene signatures  $\Gamma_n$  extracted for each cell  $n$  in a dataset can be assessed through their enrichment against  
 474 reference gene set (e.g. a marker gene list) associated to well characterized cell types and/or functional terms.  
 475 Cell-ID evaluates such enrichment through a hypergeometric test as follows: Let  $P$  be the set of genes retained  
 476 in the gene expression matrix  $M_{N,P}$  previously defined, after the initial steps of cell and gene filtering described  
 477 above. Let  $W$  be the set of genes within a reference gene set which are contained on  $P$  ( $W \subset P$ ). Let  $w$  be the  
 478 number of genes overlapping between the signature  $\Gamma_n$  of size  $\gamma$  and the gene set  $W$ :

479 
$$w = |\Gamma_n \cap W|$$

480 The observed overlap  $w$  can be modelled as a random variable  $X$  distributed hypergeometrically, with probability  
481 mass function given by:

482 
$$\mathbf{Prob}_{nW}(X = w) = \frac{\binom{W}{w} \binom{P-W}{\gamma-w}}{\binom{P}{\gamma}}$$

483 Only reference gene sets of size  $W \geq 10$  were considered throughout this work. When the gene signature  $\Gamma_n$  of  
484 a cell  $n$  in a dataset  $D$  is evaluated against a collection of reference gene sets  $W_1, W_2, \dots, W_\Omega$ , (e.g. a repository of  
485 cell-type marker lists or a pathway database), the above hypergeometric test  $p$ -values are corrected by multiple  
486 testing for the number of gene sets  $\Omega$  evaluated. Thus, a cell  $n$  is considered as enriched in those gene sets for  
487 which the hypergeometric test  $p$ -value is  $< 1e-02$ , after Benjamini Hochberg multiple<sup>22</sup> correction. In addition,  
488 when a disjointed classification is required, a cell  $n$  may be assigned to the gene set  $W_\omega$  with the lowest significant  
489 corrected  $p$ -value. On the contrary, if no significant hits are found, a cell  $n$  will remain unassigned.

490 **Per-cell gene signature enrichment analyses against per-cell and per-group gene-signatures extracted from**  
491 **reference single-cell datasets.**

492 The gene signatures  $\Gamma_n$  extracted for each cell  $n$  in a dataset  $D$  can be assessed through their enrichment against  
493 the gene signatures  $\Gamma'_{n'}$  extracted for each cell  $n'$  in a reference dataset  $D'$ , an approach called here Cell-ID(c).  
494 Analogous to the previous section, Cell-ID(c) evaluates such enrichment through a hypergeometric test as  
495 follows: Let  $P$  be the set of genes retained in the gene expression matrix  $M_{N,P}$  associated to dataset  $D$  as  
496 previously defined. Let  $\Gamma'_{n'|P}$  be the set of genes of size  $W'$  within a per-cell gene signature  $\Gamma'_{n'}$  extracted for a  
497 cell  $n'$  in the dataset  $D'$ , which are contained on  $P$ , i.e.:  $\Gamma'_{n'|P} = \Gamma'_{n'} \cap P$ .

498 Let  $w$  be the number of genes overlapping between the signature  $\Gamma_n$  of size  $\gamma$  and the gene set  $P$ :

499 
$$w' = |\Gamma_n \cap \Gamma'_{n'|P}|$$

500 The observed overlap  $w'$  between two per-cell gene signatures can be modelled as a random variable  $X$   
501 distributed hypergeometrically, with probability mass function given by:

502 
$$\mathbf{Prob}_{n'}(X = w') = \frac{\binom{W'}{w'} \binom{P-W'}{\gamma-w'}}{\binom{P}{\gamma}}$$

503 For each cell  $n$  in a dataset  $D$ , the above hypergeometric test  $p$ -values are corrected by multiple testing for the  
504 number of cells  $N'$  in the reference dataset  $D'$  against which it is evaluated. Thus, a cell  $n$  in  $D$  is considered as  
505 enriched in those signatures  $n'$  in  $D$  for which the hypergeometric test  $p$ -value is  $<1e-02$ , after Benjamini  
506 Hochberg correction on the number  $n'$  of tested gene signatures. In addition, when a disjointed classification is  
507 required, a cell  $n$  may be assigned to the cell  $n'$  in  $D'$  with the lowest significant corrected  $p$ -value. Best hits can  
508 be used for cell-to-cell matching and label transferring across datasets. On the contrary, if no significant hits are  
509 found, a cell  $n$  will remain unassigned.

510 Alternatively, if a grouping  $\Theta'_1, \Theta'_2, \dots, \Theta'_\theta$  of the  $N'$  cells in  $D'$  is provided, the gene signatures  $\Gamma_n$  for each cell  $n$  in  
511 a dataset  $D$  can be assessed through their enrichment against the corresponding per-group gene signatures  
512  $\Gamma'_{\Theta'_1}, \Gamma'_{\Theta'_2}, \dots, \Gamma'_{\Theta'_\theta}$  extracted from  $D'$  as described above. We call this approach Cell-ID(g). Here, a cell  $n$  in  $D$  is  
513 considered as enriched in those cell groups  $\Theta'_g$  from  $D'$  for which the hypergeometric test  $p$ -value is  $<1e-02$ , after  
514 Benjamini Hochberg correction for the number of groups evaluated. In addition, when a disjointed classification  
515 is required, a cell  $n$  may be assigned to the group  $\Theta'_g$  in  $D'$  with the lowest significant corrected  $p$ -value. Best hits  
516 can be used for cell-to-group matching and group-based label transferring across datasets. On the contrary, if no  
517 significant hits are found, a cell  $n$  will remain unassigned. Cell-ID(g) can handle both disjoint and non-disjoint cell  
518 groupings (i.e. overlapping groups), as well as complete or non-complete groupings (i.e. when not all cells in  $D'$   
519 have been assigned to a group).

## 520 **Simulated datasets**

521 Simulated scRNA-seq datasets were obtained with the Splatter<sup>23</sup> Bioconductor package (version 1.4.1;  
522 <https://bioconductor.org/packages/release/bioc/html/splatter.html>). For the generation of structured datasets,  
523 each simulation was set to originate from five underlying subpopulations with relative sizes of 30%, 25%, 20%,  
524 15% and 10% cells, respectively. No clustering or cluster labels were used in any way for the purpose of the  
525 analysis. Splatter's logistic function, modeling the probability of a gene having zero counts, was defined by a  
526 midpoint parameter  $x_0 = 3$  counts, to obtain a simulated dataset with about 60% ~ 70% of the count matrix  
527 content equal to zero after default normalization and filtering, a dropout rate consistent with those observed on  
528 other datasets used in this manuscript. Default values were used for all the other parameters. Centered and  
529 scaled principal component analysis (PCA) was performed with the base R `prcomp` function.

530



531 **Comparative benchmark of approaches performing marker-based cell type annotation**

532 In the comparative benchmark assessing the method's capacity to classify cells using pre-established gene  
533 signatures, two alternative semi-supervised classifiers were used: SCINA<sup>24</sup> and AUCell<sup>25</sup> (**Supplementary Table**  
534 **3**). Cell-ID predictions were based on a reference collection of blood cell markers from the XCell repository <sup>26</sup>  
535 (**Supplementary Table 4**). Only cell types of the hematopoietic lineage supported by more than three bulk  
536 RNAseq samples in XCell were considered: hematopoietic stem cells (HSC), multipotent progenitors (MPP), B  
537 cells, CD4+ T cells, CD8+ T cells, natural killer (NK) cells, plasmacytoid dendritic cells (pDC), common myeloid  
538 progenitors (CMP), granulocyte myeloid progenitors (GMP), megakaryocyte and erythrocyte progenitors (MEP),  
539 erythrocytes, megakaryocytes, platelets, basophils, eosinophils, neutrophils, CD14+ monocytes, CD16+  
540 monocytes, macrophages, dendritic c(DC), conventional dendritic cells (cDC). For each cell type, the marker list  
541 used included genes replicated in at least 20% of the reported sources. Raw count matrices were used as input  
542 for AUCell and log-transformed normalized matrices were used for SCINA, following their associated vignettes.  
543 SCINA was run with default parameters except for i) the maximum number of iterations and the convergence  
544 rate, which were increased to 20 and 0.999 respectively to ensure a stable results, and ii) the rejection parameter,  
545 which was set to true to enable cells to be labelled as unassigned when there is a low confidence on the cell type  
546 prediction. For AUCell, the gene set with the highest AUC score was used to classify cells unless the maximum  
547 AUC score was <0.1, what left a cell as 'unassigned'.

548 **Comparative benchmark of approaches performing cell-type label-transferring across datasets**

549 In addition to Cell-ID, we evaluated 10 alternative approaches for cell-type label-transferring across scRNA-seq  
550 (**Supplementary Table 3**). All methods were run using default parameters unless otherwise stated. When default  
551 settings were not explicitly defined, setting used in the associated vignettes were followed. Methods used as  
552 input either the raw or the normalized count data (after gene filtering as described above), following each  
553 method's documentation. For those methods that stipulate it, gene expression matrices were further restricted  
554 to genes in common between the reference and the query datasets. In the case of MNN<sup>27</sup>, we transferred labels  
555 from the reference to the query datasets between closest mutual nearest neighbor cells, and cells were left  
556 unassigned when no mutual nearest neighbor cells were found. We modified the default parameter of k nearest  
557 neighbor of MNN to k = 50 as the default k = 20 failed to find mutual nearest neighbor matches for a large fraction  
558 of cells, which negatively affected the benchmark metrics evaluated. Furthermore, for the selection of

559 hypervariable genes in MNN, we used the default Seurat function for highly variable gene detection (2000 genes),  
560 and we took the intersection between the reference and the query highly variable genes to perform the query  
561 and reference dataset integration following the package's vignette. In the case of SCN<sup>28</sup>, all cells that were  
562 classified as *rand* were considered as unassigned, as well as all cells classified as "nodes" in CHETAH<sup>29</sup>. For Seurat  
563 cells were labelled as unassigned when the projection score was below 0.5.

#### 564 **Classification performance assessment**

565 Cell annotation performance was assessed through 3 complementary metrics i.e.: precision, recall and F1 score,  
566 which is the harmonic mean between the precision and recall. Each metric was first calculated for each cell type  
567 in the query dataset. Second, each metric (i.e. recall, precision and F1 score) was calculated for the global set as  
568 the arithmetic mean of the corresponding metric across the evaluated cell types. In such a way, an overweighed  
569 contribution of largely populated cell types is avoided, allowing thus a larger contribution of more rare cell types  
570 to the metrics evaluated. Mapping across datasets of cell type nomenclature was performed through manual  
571 curation and is reported in **Supplementary Table 5**. For label-transferring performance assessment, cells left  
572 unassigned in a query dataset were considered as a false assignment when their actual cell type was indeed  
573 represented in the reference dataset, or considered as a true assignment otherwise. Therefore, for intestinal  
574 dataset label transferring, any cell types apart from endocrine, brush/tuft and goblet were considered as negative  
575 cell type since there are not present in the reference dataset and hence should be labeled unassigned. Thus,  
576 unassigned cells were considered as true positive if labelled against negative cells and the negative cells were  
577 evaluated in the calculation of performance metrics in the same manner as the three other cell types present in  
578 the reference and was considered also in the global performance metrics. For interspecies label-transferring  
579 across human and mouse datasets, the initial raw matrix of the query dataset was restricted to genes with one-  
580 to-one orthologs. Ortholog relations were obtained from BioMart (release 100, version April 2020, GrCH38.p13  
581 for human, and GRCm38.p6 for mouse<sup>15,16</sup>) using gene symbols.

#### 582 **Functional Enrichment**

583 Functional enrichment analyses were performed using 6 sources of functional annotations: Reactome<sup>30</sup>, KEGG<sup>31</sup>,  
584 WikiPathways<sup>32</sup>, GO biological process<sup>33</sup>, GO molecular function and GO cellular component, collectively  
585 gathering a total of 8709 terms. Gene sets associated to functional pathways and ontology terms were obtained  
586 as provided in enrichr<sup>34</sup> website <http://amp.pharm.mssm.edu/Enrichr/#stats>. (**Supplementary Table 6**)

## 587 **Visualization**

588 UMAP<sup>35</sup> or tSNE<sup>36</sup> representations were alternatively used for visualization purposes throughout the manuscript.  
589 However, no biological conclusions whatsoever were drawn from visual inspection of such representations.  
590 UMAP and tSNE representations were obtained with Seurat default parameters following this vignette  
591 ([https://satijalab.org/seurat/v3.1/pbmc3k\\_tutorial.html](https://satijalab.org/seurat/v3.1/pbmc3k_tutorial.html)).

## 592 **Computational resources**

593 All analyses presented in the manuscript were run in a workstation with 64-GB RAM memory and a AMD Ryzen  
594 2700X processor with 8 3.6-GHz physical cores, with the exception of the scalability benchmark presented in  
595 **Supplementary Note 6** where an Intel Xeon Gold 6140 with 36 2.3 GHz cores processor and 640Gb of RAM was  
596 used.

## 597 **Code availability**

598 Cell-ID is implemented as an R package and is available on GitHub (<https://github.com/RausellLab/CellID>) under  
599 the GPL-3 open source license. Complete documentation is provided with step-by-step procedures for MCA  
600 dimensionality reduction, per-cell gene signature extraction, cell type prediction, label-transferring across  
601 datasets and functional enrichment analysis. All R scripts and intermediate data representations required to  
602 reproduce all figures in the manuscript are provided as a supplementary file.

## 603 **Methods references**

- 604 1. Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells. *Nat. Methods* **14**,  
605 865–868 (2017).
- 606 2. Peterson, V. M. *et al.* Multiplexed quantification of proteins and transcripts in single cells. *Nat. Biotechnol.* **35**,  
607 936–939 (2017).
- 608 3. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21 (2019).
- 609 4. Baron, M. *et al.* A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and Intra-  
610 cell Population Structure. *Cell Syst.* **3**, 346–360 (2016).
- 611 5. Muraro, M. J. *et al.* A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Syst.* **3**, 385-394.e3 (2016).
- 612 6. Segerstolpe, Å. *et al.* Single-Cell Transcriptome Profiling of Human Pancreatic Islets in Health and Type 2  
613 Diabetes. *Cell Metab.* **24**, 593–607 (2016).

- 614 7. scRNAseq. *Bioconductor* <http://bioconductor.org/packages/scRNAseq/>.
- 615 8. Plasschaert, L. W. *et al.* A single-cell atlas of the airway epithelium reveals the CFTR-rich pulmonary ionocyte.  
616 *Nature* **560**, 377–381 (2018).
- 617 9. Montoro, D. T. *et al.* A revised airway epithelial hierarchy includes CFTR-expressing ionocytes. *Orit Rozenblatt-*  
618 *Rosen* **4**, 19–19 (2018).
- 619 10. Haber, A. L. *et al.* A single-cell survey of the small intestinal epithelium. (2017) doi:10.1038/nature24489.
- 620 11. Fletcher, R. B. *et al.* Deconstructing Olfactory Stem Cell Trajectories at Single-Cell Resolution. *Cell Stem Cell*  
621 **20**, 817-830.e8 (2017).
- 622 12. Wu, Y. *et al.* A Population of Navigator Neurons Is Essential for Olfactory Map Formation during the Critical  
623 Period Article A Population of Navigator Neurons Is Essential for Olfactory Map Formation during the Critical  
624 Period. *Neuron* **100**, 1066-1082.e6 (2018).
- 625 13. Schaum, N. *et al.* Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris the tabula Muris  
626 consortium\*. (2018) doi:10.1038/s41586-018-0590-4.
- 627 14. Cusanovich, D. A. *et al.* A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility. *Cell* **174**, 1309-  
628 1324.e18 (2018).
- 629 15. Zerbino, D. R. *et al.* Ensembl 2018. *Nucleic Acids Res.* **46**, D754–D761 (2018).
- 630 16. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping Identifiers for the Integration of Genomic Datasets  
631 with the R/Bioconductor package biomaRt. *Nat. Protoc.* **4**, 1184–1191 (2009).
- 632 17. *Multiple correspondence analysis and related methods. Selected papers based on the presentations at the*  
633 *international conference (CARME 2003), Barcelona, Spain, 29 June to 2 July 2003. Chapman & Hall/CRC*  
634 *Statistics in the Social and Behavioral Sciences Series* (Chapman & Hall/CRC, 2006).
- 635 18. Greenacre, M. J. *Theory and applications of correspondence analysis.* (1984).
- 636 19. Lebart, L., Morineau, A. & Warwick, K. M. *Multivariate descriptive statistical analysis. Correspondence analysis*  
637 *and related techniques for large matrices. Transl. from French by Elisabeth Morailon Berry. With a foreword*  
638 *by Herman P. Friedman. Wiley Series in Probability and Mathematical Statistics* (John Wiley & Sons, Hoboken,  
639 NJ, 1984).
- 640 20. Aşan, Z. & Greenacre, M. Biplots of fuzzy coded data. *Fuzzy Sets Syst.* **183**, 57–71 (2011).
- 641 21. Greenacre. Chapter 8. in *Biplots in practice* 79–88 (2010).

- 642 22. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to  
643 Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* **57**, 289–300 (1995).
- 644 23. Zappia, L., Phipson, B. & Oshlack, A. Splatter: Simulation of single-cell RNA sequencing data. *Genome Biol.* **18**,  
645 174–174 (2017).
- 646 24. Zhang *et al.* SCINA: Semi-Supervised Analysis of Single Cells in Silico. *Genes* **10**, 531–531 (2019).
- 647 25. Aibar, S. *et al.* SCENIC: Single-cell regulatory network inference and clustering. *Nat. Methods* **14**, 1083–1086  
648 (2017).
- 649 26. Aran, D., Hu, Z. & Butte, A. J. xCell: digitally portraying the tissue cellular heterogeneity landscape. *Genome*  
650 *Biol.* **18**, 220 (2017).
- 651 27. Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-sequencing data are  
652 corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–427 (2018).
- 653 28. Tan, Y. & Cahan, P. SingleCellNet: A Computational Tool to Classify Single Cell RNA-Seq Data Across Platforms  
654 and Across Species. *Cell Syst.* **9**, 207-213.e2 (2019).
- 655 29. De Kanter, J. K., Lijnzaad, P., Candelli, T., Margaritis, T. & Holstege, F. C. P. CHETAH: a selective, hierarchical  
656 cell type identification method for single-cell RNA sequencing. *Nucleic Acids Res.* **47**, (2019).
- 657 30. Jassal, B. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498–D503 (2020).
- 658 31. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for gene  
659 and protein annotation. *Nucleic Acids Res.* **44**, 457–462 (2015).
- 660 32. Slenter, D. N. *et al.* WikiPathways: a multifaceted pathway database bridging metabolomics to other omics  
661 research. *Nucleic Acids Res.* **46**, D661–D667 (2018).
- 662 33. The Gene Ontology (GO) database and informatics resource. *Nucleic Acids Res.* **32**, D258–D261 (2004).
- 663 34. Chen, E. Y. *et al.* Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC*  
664 *Bioinformatics* **14**, 128 (2013).
- 665 35. Becht, E. *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.* **37**, 38–  
666 44 (2019).
- 667 36. Kobak, D. & Berens, P. The art of using t-SNE for single-cell transcriptomics. *Nat. Commun.* **10**, 5416 (2019).
- 668

669 **Cell-ID: gene signature extraction and cell identity recognition at individual**  
670 **cell level**

671

672

## Supplementary Notes 1 - 7

673

### 674 **Supplementary Note 1. Consistency of MCA low-dimensional representation of cells and genes.**

675 The consistency of MCA-based low-dimensional representation of cells and genes was first evaluated with 100  
676 simulated scRNA-seq datasets, each containing 1000 cells and 5000 genes (**Online Methods**). We first compared  
677 the correspondence between MCA and PCA low-dimensional representations of cells by performing Spearman's  
678 rank correlation analysis on their principal axes coordinates. MCA and PCA cell representations were largely  
679 equivalent (**Supplementary Fig. 1**). We then determined whether the per-cell gene rankings obtained with MCA  
680 were consistent with the gene expression values for neighboring cells in the MCA space. As expected, genes  
681 specific to a given cell had higher log-fold changes in expression in the 5% of cells closest to target cell ( $n = 50$ )  
682 than in the other cells (**Supplementary Fig. 2A**). We then investigated how the ranking of genes with zero-counts  
683 in a cell related to the specificity of the genes concerned in neighboring cells (**Supplementary Fig. 2A**). We found  
684 that genes not detected in a given cell were nevertheless attributed a high ranking in this cell if the surrounding  
685 cells displayed high levels of expression for these genes (**Supplementary Fig. 2A**). This is an important result,  
686 highlighting the capacity of multivariate approaches to consider a gene to be specific to a cell in which it was not  
687 detected, provided that the gene concerned is specific to very similar cells. The MCA approach is thus robust to  
688 zero-count values that probably correspond to technical dropouts. These results could be generalized to all  
689 individual cells in a given dataset: Spearman's rank correlation coefficient between rank and log-fold change =  
690 0.72,  $p$ -value  $< 10e-16$  (**Supplementary Fig. 2B**). By contrast, the correlation was weaker for per-cell gene rankings  
691 obtained via a naïve approach based on either (i) the log-fold changes in gene expression observed in one cell  
692 relative to all other cells (Spearman's rank correlation coefficient = 0.38,  $p$ -value  $< 10e-16$ , **Supplementary Fig. 2,**  
693 **C-D**) or (ii) highest-to-lowest rankings of expression values within a cell, with random ranks for ties, as previously  
694 described (AUCell<sup>1</sup>, Spearman's rank correlation coefficient = 0.08,  $p$ -value =  $3.18e-04$ , **Supplementary Fig. 2, E-**  
695 **F**, respectively).

696

697 **Supplementary Note 2. Capacity of Cell-ID to identify well-established cell types on the basis of reference**  
698 **marker lists.**

699 We evaluated the ability of Cell-ID per-cell gene signatures to reproduce reference marker lists associated with  
700 well-established cell types. We searched for scRNA-seq datasets in which concomitant measurements of single-  
701 cell protein marker levels had been performed. Protein marker levels provide additional evidence for cell-type  
702 annotations. Cell-type labels in such settings can be taken as a gold standard reference for this study, thereby  
703 avoiding the potential circularity associated with the use of transcription-based annotations. Two independent  
704 scRNA-seq datasets for human blood cells profiled by CITE-seq<sup>2</sup> and REAP-seq<sup>3</sup> protocols met this criterion  
705 (**Online Methods**). The Cite-Seq dataset contained data for a total of 8005 human cord blood cells whereas the  
706 REAP-seq dataset contained data for 7488 PBMCs. Cell-ID predictions were based on a reference collection of  
707 blood cell markers from the XCell<sup>4</sup> repository (**Supplementary Table 4, Online Methods**).

708 The gene signature of each cell, extracted from the CITE-seq and REAP-seq datasets, was first evaluated with  
709 Cell-ID against each of the 21 marker lists associated with previously described cell types (**Figure 2** and  
710 **Supplementary Figure 3, Supplementary Table 4; Online Methods**). Each cell was assigned to the cell type for  
711 which it displayed the strongest enrichment (**Online Methods**). We assessed the accuracy of Cell-ID by calculating  
712 its precision (positive predictive value), recall (true positive rate), and overall agreement (F1) with the previously  
713 assigned reference cell-type labels. In both datasets, per-cell gene signatures displayed significant enrichment in  
714 the marker genes for at least one reference cell type in 83% and 73% of cells for CITE-seq and REAP-seq datasets,  
715 respectively. The best match to the gene signature was used for automatic cell-type prediction with high  
716 precision (0.87 and 0.9), recall (0.84 and 0.73), and F1 (0.83 and 0.78) values (multinomial  $p$ -value  $< 10^{-16}$  for all  
717 figures; **Figure 2, Supplementary Figure 4** and **Supplementary Table 7**). Cell-ID outperformed reference methods  
718 for cell-type classification based on pre-established signatures, such as SCINA<sup>5</sup> and AUCell<sup>1</sup>, for most of the cell  
719 types considered, achieving remarkably high levels of precision. CellAssign<sup>6</sup> could not be used due to execution  
720 errors of the original software, which we attribute to the length of the signatures used here, ranging from 13 to  
721 174 genes in total, which is much longer than the gene sets of 2 to 20 genes generally used by CellAssign. Overall  
722 performance was better for the CITE-seq dataset than for the REAP-seq dataset, for all methods evaluated, due

723 to a significant difference in the median number of genes detected per cell between the two datasets (2346 and  
724 1260, respectively, two-sided Wilcoxon test  $p$  value  $< 10^{-16}$ ).

725 In addition to identifying the best match between gene lists, Cell-ID was able to identify cells displaying significant  
726 enrichment in the genes of more than one of the pre-established cell-type marker lists. This contrasts with  
727 clustering-based approaches, which partition cells into disjointed groups. This situation is illustrated by the CITE-  
728 seq CBMC dataset, for which Cell-ID identified fine-grained transitions within the Hematopoietic Stem and  
729 Progenitor cells subset (CD34<sup>+</sup>) according to the hematopoietic hierarchy (**Figure 2C**). No immature CD34<sup>+</sup> cells  
730 in the REAP-seq dataset were reported in the original publication<sup>3</sup>. The gradient of Cell-ID enrichment scores  
731 was consistent with the Hematopoietic Stem Cells (HSC) lineage differentiation process. Thus, Cell-ID multi-class  
732 scores reflect smooth transitions at individual cell level from the most immature HSC to multipotent progenitors  
733 (MPP), branching between myeloid (CMP/GMP) and erythroid (MEP) progenitors, which ultimately differentiate  
734 into erythrocytes. Such fine-grained transitions would have been missed by both clustering-based approaches  
735 and alternative methods. Thus, SCINA and AUCELL provided a coarse-grained classification between HSCs and  
736 erythrocytes, with no assignment to intermediate states between these two classes.

737 We then analyzed more challenging scenarios based on real datasets in which we evaluated the capacity of Cell-  
738 ID to identify the only cell of a given rare cell type ( $n=1$ ). We focused on cell types observed at low frequencies  
739 ( $<2\%$ ) in the previous CBMC dataset: plasmacytoid dendritic cells (pDC, 0.6%,  $n=49$ ), erythrocytes (1.3%,  $n=105$ )  
740 and hematopoietic stem and progenitor cells (HSPC, CD34<sup>+</sup> subset, 1.8%,  $n=134$ ) from the CITE-Seq dataset. For  
741 each of these rare cell types containing a given number  $n$  of cells, we generated  $n$  corresponding datasets, each  
742 retaining only one cell at a time, while the rest of the subpopulations remained unchanged. Using the pre-  
743 established immune cells marker lists, as described above, Cell-ID identified the singleton cell with a high mean  
744 recall (92%, 94%, and 86%), but modest precision (32%, 8% and 76%) due to the experimental design with very  
745 high number of negatives compared to positive (1 positive and 999 negatives), which results in mid-range F1  
746 score (48%, 0.15%, 0.81%) (**Supplementary Table 8**). Singleton cells pose a problem in clustering-based  
747 approaches, in which they are often merged with larger sub-populations or treated as outliers and are thus  
748 filtered out of downstream analyses. Conversely, alternative reference methods able to work at individual cell  
749 level without a clustering step, such as SCINA<sup>5</sup> and AUCell<sup>1</sup>, could not be applied in this setting either, because  
750 they are not suitable for independent evaluations one cell type at a time.



751

752 **Supplementary Note 3. Capacity of Cell-ID to match cells of analogous cell types across independent scRNA-**  
753 **seq datasets from the same tissue of origin, within and across species.**

754 Cell-ID provides two options for cell type matching across datasets: (i) Cell-ID(c), in which individual cells in a  
755 *query* dataset are matched to transcriptionally analogous cells in a *reference* dataset; and (ii) Cell-ID(g), in which  
756 individual cells in a *query* dataset are matched to cell group labels previously established in a *reference* dataset  
757 through clustering and/or expert annotation. Thus Cell-ID(c) provides *cell-to-cell* matching, whereas Cell-ID(g)  
758 performs *group-to-cell* matching from *reference-to-query* datasets. It should be noted, however, that Cell-ID(g)  
759 performs no clustering whatsoever, as the cell groups in the *reference* dataset are provided as an input by the  
760 *reference* source. For the evaluation of Cell-ID(c) and Cell-ID(g), we searched for tissues profiled in several  
761 independent scRNA-seq datasets and meeting the following three conditions: (i) cell type labels curated through  
762 expert annotation in the original publications, (ii) containing at least one rare cell type (cell types accounting for  
763 less than 2% of the cells in a sample), and (iii) different sequencing protocols and/or different model organisms  
764 (i.e. humans and mice) used. Two tissues meeting these requirements were identified (**Figure 3**):

765 (A) Pancreatic islets: we considered three scRNA-seq datasets for pancreatic cells from human donors: (i) 8659  
766 cells from four deceased donors healthy at the time of death profiled with the inDrop protocol (Baron et al.<sup>7</sup>), (ii)  
767 2126 cells from four deceased organ donors sequenced with CEL-seq2 (Muraro et al.<sup>8</sup>) and (iii) 2168 cells from  
768 six healthy donors and four donors with type-2 diabetes sequenced with Smart-seq2 (Segerstolpe et al.<sup>9</sup>). The  
769 Baron dataset reported the greatest diversity of cell types and was, thus, chosen as the *reference* for this analysis.  
770 The cell types identified in Baron's human dataset included two exocrine cell types (acinar and duct cells), five  
771 endocrine cell types (alpha, beta, delta, gamma and epsilon cells), three immune cell types (tissue-resident  
772 macrophages, mast cells, cytotoxic T cells), pancreatic stellate cells, endothelial cells and Schwann cells of neural  
773 crest origin. All of these cell types, with the exception of Schwann cells and cytotoxic T cells were reported by  
774 Muraro et al. and Segerstolpe et al.

775 (B) Airway epithelium: three independent airway epithelial cell datasets from mice and human donors were  
776 obtained, including (i) 7662 tracheal epithelial cells from four mice profiled with 10X Genomics technology  
777 (Plasschaert et al.<sup>10</sup>); (ii) 7586 airway tracheal cells from six mice sequenced with inDrop (Montoro et al.<sup>11</sup>); (iii)  
778 2970 primary bronchial epithelial cells from three human donors profiled with 10X Genomics technology

779 (Plasschaert et al<sup>10</sup>) (**Online Methods**). All three datasets reported the six major airway epithelial cell types,  
780 including basal, secretory cells (also known as club cells), ciliated cells, rare pulmonary neuroendocrine cells  
781 (PNEC), brush cells (also known as tuft cells) and ionocytes. Plasschaert's mouse dataset was used as the  
782 *reference* for this study, as it allowed a within-species analysis (mouse-to-mouse) as well as a between-species  
783 analysis (mouse-to-human) based on datasets originating from the same laboratory.

784 Cell-ID gene signatures at either individual cell level or group level were obtained independently for each of the  
785 datasets described above (**Online Methods**). Four independent *reference-to-query* assignments were evaluated:  
786 Baron's human pancreatic cells against Muraro's and against Segerstolpe's human pancreatic datasets, and  
787 Plasschaert's mouse airway epithelial cells against Montoro's mouse airway epithelial cells and against  
788 Plasschaert's human airway epithelial cells. For each of the *reference-to-query* assessments, both *cell-to-cell*  
789 matching, and *group-to-cell* matching were performed, with Cell-ID(c) and Cell-ID(g), respectively (**Methods**).  
790 Thus, each cell in the query dataset was assigned to the cell type from the *reference* cell or *reference* group for  
791 which it presented the lowest significant gene signature enrichment *p* value (or was left unassigned if no  
792 significant hits were found; **Supplementary Figure 5**). Cell-ID *cell-to-cell* matching and *group-to-cell* matching  
793 were thus evaluated by assessing the cell-type agreement between the transferred labels and the original labels  
794 in the *query* dataset. Comparative benchmarking was performed against a representative set of alternative state-  
795 of-the-art methods covering major approaches for cell-matching or label transfer across scRNA-seq datasets  
796 (**Supplementary Table 3**): (i) integration-based methods for cell-to-cell matching across datasets based on  
797 reciprocal *k*-nearest neighbor analysis and projection in a common low-dimensional space (Seurat<sup>12</sup>, MNN<sup>13</sup>); (ii)  
798 transcriptome similarity assessment without an integration step (scmap\_cluster and scmap cell<sup>14</sup>, scID<sup>15</sup> and  
799 singleR<sup>16</sup>); and (iii) machine learning-based methods training models through cross-validation in the reference  
800 dataset followed by prediction in the query dataset (CaSTle<sup>17</sup>, scPred, SCN<sup>18</sup> and CHETAH<sup>19</sup>).

801 Both Cell-ID(c) and Cell-ID(g) consistently reached high precision (>82% and >83%), recall (>82% and >=76%) and  
802 F1 values (>74% and >74%, respectively) across all *reference-to-query* assignments evaluated (multinomial *p*  
803 value < **2.2e-16** for all figures; **Figure 3A, Supplementary Figure 6 A-B, Supplementary Table 9**). Cell-ID  
804 performance was at least as good as that of all alternative state-of-the-art methods. Notably, Cell-ID obtained  
805 high performance scores even in for cell-type matching across species, which proved challenging for most of the  
806 methods evaluated. We assessed the robustness of Cell-ID gene signatures further by focusing on the cell types

807 present at low frequencies (<2%) in the previous query datasets: epsilon cells, tissue-resident macrophages, mast  
808 cells and endothelial cells from pancreatic islet samples (Muraro et al. and Segerstolpe et al), and PNEC, brush  
809 cells and ionocytes in the mouse and human airway epithelium datasets (Montoro et al and Plasschaert et al.,  
810 respectively). As expected, *reference-to-query* assignments for such rare cell types were generally more error-  
811 prone for all the methods evaluated, including Cell-ID (**Figure 3B, Supplementary Figure 6 C-D**). Nevertheless,  
812 Cell-ID performed better than the alternative methods, with salient scores for Cell-ID(g) for most of the rare cell  
813 types evaluated (median F1 values greater than 88%, respectively, with  $p$  value <  $1e-5$  for all evaluated rare cell  
814 populations relative to expectations from a random binomial distribution).

815 From a discovery perspective, the datasets used here provided us with an opportunity to evaluate the capacity  
816 of Cell-ID to detect ultra-rare cell types potentially missed in the original publications. Indeed, Baron et al.  
817 reported the presence of Schwann cells ( $n=13$ ; 0.17 %) in pancreatic islet datasets, an ultra-rare cell type of neural  
818 crest origin. However, Schwann cells were not described in the samples of Muraro et al. and Segerstolpe et al.  
819 We thus explored possible enrichment in the Cell-ID gene signatures extracted for each of the 13 Schwann cells  
820 in the Baron et al. dataset for the 2126 and 2168 cells, respectively of the query datasets. Both Cell-ID(c) and  
821 Cell-ID(g) identified  $n=4$  cells in the Muraro et al. dataset and  $n=2$  cells in the Segerstolpe et al. dataset that had  
822 been initially labeled as ductal (Muraro et al.<sup>8</sup>) or unclassified (Segerstolpe et al.<sup>9</sup>) in the associated publications.  
823 These cells presented high expression levels of Schwann cell markers (SOX10, S100B, CRYAB, NGFR, CDH19, and  
824 PMP22) and of markers of response to nerve injury (SOX2, ID4, and FOXD3), as described by Baron et al.<sup>7</sup>.  
825 Accordingly, SOX10, S100B, CRYAB, NGFR, CDH19, and PMP22 ranked among the top 100 gene signatures  
826 associated with the four and two cells annotated as Schwann cells by Cell-ID (**Supplementary Figure 7**), and these  
827 six genes were significantly upregulated in these cells relative to the other cells (Wilcoxon test  $p$  value <  $10^{-3}$ , in  
828 both datasets). The identification of putative Schwann cells was replicated by Seurat, MNN, scID, SCN, SingleR  
829 and scmap cluster, but missed by CaSTLe, scmap cell, scPred and CHETAH, which labeled them as “other” or  
830 “unassigned”. Moreover, functional enrichment analysis of the gene signatures of each of the four and two  
831 putative Schwann cells revealed significant enrichment in 62 Gene Ontology (GO) biological processes (median  
832  $p$ -value across cells < 0.01) for the Muraro et al. dataset and 26 GO biological processes for the Segerstolpe et al.  
833 dataset, with 10 terms common to the two sets (**Supplementary Table 1**). Nervous system development  
834 (GO:0007399), including glial cell differentiation (GO:0010001) and axonogenesis (GO:0007409), were among

835 the top 10 terms displaying enrichment in both sets (p-values less than or equal to 2.18e-05 and 2.98e-03,  
836 respectively), consistent with the neural crest origin of Schwann cells<sup>20</sup>. The prominent role played by Schwann  
837 cells in the myelination process<sup>21</sup> was, in turned, reflected in several functional terms and pathways for which  
838 significant enrichment was detected, including collagen fibril organization (GO:0030199) in the four cells from  
839 the Muraro dataset and myelination (GO:0042552) and collagen binding (GO:0005518) in the two cells from  
840 the Segerstolpe dataset. These signals were consistently reproduced with alternative pathway annotation  
841 sources such as KEGG<sup>22</sup>, Reactome<sup>23</sup> and WikiPathways<sup>24</sup> (**Supplementary Table 1, Supplementary Table 6**).

842 Overall, these results show that Cell-ID gene signatures are highly reproducible at both the cell and group levels,  
843 across independent datasets from the same tissue of origin, despite the use of different sequencing protocols,  
844 donors, and species. This reproducibility was robust even for rare cell types present at extremely low frequencies,  
845 demonstrating the ability of Cell-ID to extract biologically relevant gene signatures at individual cell resolution.  
846 In particular, matching across datasets by Cell-ID is fully transparent in terms of the set of genes driving the hits,  
847 and is therefore fully interpretable in biological terms. This contrasts sharply with the situation for label transfer  
848 methods based on assessments of similarity over the entire transcriptome (e.g. Seurat, MNN, scmap) or machine-  
849 learning approaches, for which individual gene contributions are difficult to interpret (e.g. scPred, SCN).

850

#### 851 **Supplementary Note 4. Capacity of Cell-ID to match cell types across independent scRNA-seq datasets from** 852 **different tissues of origin**

853 In a more challenging scenario, we then evaluated the capacity of Cell-ID to recognize gene signatures of rare  
854 cells across independent sets from different tissues of origin. We thus searched for rare cell types meeting the  
855 following conditions: (i) that they were common to at least two different tissues, and (ii) for which single-cell  
856 RNA-seq datasets were available, in which cell-type labels had been curated by expert annotation in the original  
857 publications. Chemosensory epithelial cells appeared to be an ideal rare cell type present in various epithelium  
858 types, presenting all the characteristics for the purpose of the analysis. Recent studies have shown that  
859 chemosensory epithelial cells, referred to as tuft cells in the intestinal mucosa, and solitary chemosensory cells  
860 (SCCs) in the nasal respiratory mucosa, belong to the brush/tuft cell family, together with the brush cells in the  
861 tracheal airway epithelium<sup>25 26</sup>. This chemosensory epithelial cell family has common transcriptional programs  
862 and functions in all three tissues<sup>27</sup>. In particular, these cells are the primary source of interleukin-25<sup>28</sup>, a

863 proinflammatory protein mediating type 2 inflammation induced by diverse pathogens in various mucosal  
864 tissues.

865 We evaluated the capacity of Cell-ID to identify rare cell types across tissues, by focusing on cell matching  
866 between the mouse tracheal airway epithelium (Plasschaert <sup>10</sup> and Montoro <sup>11</sup> datasets, as described in  
867 **Supplementary Note 3**) and the mouse small intestinal epithelium (single-cell dataset from Haber et al.<sup>29</sup>, for  
868 7216 cells from four mice sequenced with the inDrop protocol). The mouse airway and intestinal epithelia display  
869 major differences in terms of their cell-type composition. Both Plasschaert et al. and Montoro et al. reported the  
870 presence of basal, secretory (including goblet in Montoro), ciliated cells, PNEC, ionocytes and brush/tuft cells in  
871 the airway epithelium (**Supplementary Note 3**). Haber et al. performed single-cell analyses on the intestinal  
872 epithelium in which they identified enterocytes (45%), transit-amplifying cells (21%), stem cells (18%), goblet  
873 cells (7%), Paneth cells (4%) enteroendocrine cells (4%) and brush/tuft cells (2%) (**Supplementary Figure 8A**).  
874 Cell-ID(g) and Cell-ID(c) were applied with default parameters (**Online Methods**) for *reference-to-query*  
875 assignments from Plasschaert and Montoro's mouse airway epithelial cells to Haber's mouse intestinal epithelial  
876 cells. Each cell in the mouse intestinal epithelial dataset was assigned to the cell type from the mouse airway  
877 epithelial for which it presented the lowest significant *p* value for gene signature enrichment (or was left  
878 unassigned if no significant hits were found). Both Cell-ID(g) and Cell-ID(c) identified brush/tuft, but also  
879 endocrine cells and goblet cells in the intestinal epithelium, based on brush/tuft, PNEC and secretory/goblet gene  
880 signatures, respectively, extracted from the airway epithelium (**Figure 3C** and **Supplementary Figure 8B**). Global  
881 matching accuracy was evaluated by analyzing the concordance of cell-type labeling with the original cell-type  
882 annotations from the corresponding publications (**Supplementary Table 5**). Comparative benchmarking was  
883 performed against a representative set of alternative state-of-the-art methods covering major approaches for  
884 label transfer or cell-matching across scRNA-seq datasets (**Supplementary Table 3**), as previously described in  
885 **Supplementary Note 3**. Cell-ID(g) yielded high precision, recall and F1 scores, clearly outperforming all the other  
886 methods evaluated (**Figure 3D**, **Supplementary Figure 8C** and **Supplementary Table 10**). Similar results were  
887 obtained when the Plasschaert or Montoro mouse airway epithelial cells were used as the *reference set*.  
888 Alternative methods, including Cell-ID(c), performed less well, by contrast to the results obtained for label  
889 transfer between samples from the same tissue of origin (**Supplementary Note 3**). This lower performance was  
890 driven by many false positives (i.e. incorrect label transfer from the reference to the query dataset), resulting in

891 low precision and low F1 scores. By contrast, correct label assignment rates remained high for Cell-ID(g), which  
892 did not assign labels from the reference to the query set in the absence of significant statistical support. The  
893 higher performance of Cell-ID(g) than of Cell-ID(c) suggests that group-based gene signatures are more robust  
894 than individual cell-based signatures for cell-type matching across samples from different tissues of origin.

895 We then searched for olfactory epithelium single-cell datasets of use for the evaluation of the capacity of Cell-ID  
896 to identify solitary chemosensory cells (SCCs) by projecting brush/tuft signatures extracted from the airway and  
897 intestinal epithelia. Two datasets from published works were identified: (i) Wu et al<sup>30</sup>, in which a total of 9126  
898 olfactory epithelium cells from mice aged 0, 3, 7 and 21 days were sequenced with the 10X Genomics protocol,  
899 and (ii) Fletcher et al.<sup>31</sup>, in which 849 cells from mice aged 21-28 days were sequenced with the Smart-Seq2  
900 protocol. The mouse olfactory epithelium has a cell-type composition different from those of the airway and  
901 intestinal epithelia. Both Wu et al and Fletcher et al reported the presence of transitional horizontal basal cells,  
902 resting horizontal basal cells, microvillous cells, mature sustentacular cells, mature olfactory sensory neurons,  
903 immediate neuronal precursors, immature sustentacular cells, immature olfactory sensory neurons, globose  
904 basal cells and unknown/unlabeled cells. However, neither of these studies reported the presence of SCCs among  
905 the cells sequenced. Nevertheless, these two datasets provided us with an opportunity to illustrate the capacity  
906 of Cell-ID for exploratory cell-type scanning on single-cell datasets of a *query* tissue, using gene signatures  
907 extracted from single-cell datasets for a set of *reference* tissues. Cell-ID(g) was used to evaluate the presence in  
908 the olfactory epithelium of cells with the tuft/brush cell gene signatures extracted from the **airway** and **intestinal**  
909 **epithelium**, as described above. Cell-ID(g) identified a total of 37 cells in the Wu dataset and 5 cells in the Fletcher  
910 dataset displaying gene signatures with significant enrichment in the genes of the tuft/brush cell signatures  
911 extracted from the **airway epithelium** (dataset from Plasschaert et al<sup>10</sup>). Cell-ID(g) also identified the same 37  
912 and 5 cells as having signatures significantly enriched in the genes of the tuft/brush cell signatures extracted from  
913 the intestinal **epithelium** (dataset from Haber et al). Interestingly, these cells were labeled “unidentified” in the  
914 original publications<sup>30 31</sup>. They displayed higher levels of expression for IL25 and the chemosensory cell marker  
915 Gnat3 than any of the other cells in the corresponding datasets (two-tailed Wilcoxon test,  $p < 10^{-16}$  for both  
916 datasets), further supporting the classification of these cells as SCCs (Ualiyeva et al. 2020; **Figure 3E-F**,  
917 **Supplementary Figure 9**). Moreover, functional enrichment analysis performed on the gene signatures of each  
918 of the 37 and 5 putative SCCs identified 20 and 14 Gene Ontology (GO) biological processes (median  $p$  value

919 across cells <0.01) displaying significant enrichment in the datasets of Wu and Fletcher, respectively, with nine  
920 terms common to the two datasets (**Supplementary Table 1**). Unsaturated fatty acid biosynthetic process  
921 (GO:0006636), including eicosanoid (GO:0046456) and leukotriene (GO:0019370) biosynthetic processes, were  
922 among the top five terms displaying the highest level of enrichment in both datasets (5.76E-04 and 1,77E-05,  
923 respectively), driven by the *Alox5*, *Alox5ap*, *Ltc4s*, *Pla2g4a* genes in both sets of putative SCCs. This finding is  
924 consistent with recent studies showing that SCCs are a primary source of cysteinyl leukotriene in the olfactory  
925 epithelium<sup>25 26</sup>), and that these cells have immune effector functions in common with airway brush/tuft cells  
926<sup>25</sup>. These signals were consistently reproduced with alternative pathway annotation sources (KEGG<sup>22</sup>, Reactome<sup>23</sup>  
927 and WikiPathways<sup>24</sup>; **Supplementary Table 1**). Taken together, the gene signatures and functional hallmarks  
928 revealed by Cell-ID(g) confirm, at single-cell transcriptional level, the presence of rare cells in the olfactory  
929 epithelium similar to the airway and intestinal brush/tuft cells. In accordance with recent findings<sup>25 26</sup>, the 37  
930 cells from Wu et al and the five cells from Fletcher et al described above should, thus, be annotated as SCCs.

931

932 **Supplementary Note 5. Capacity of Cell-ID to match cell types across independent datasets from different**  
933 **single-cell *omics* technologies: scRNAseq and scATACseq**

934 We investigated whether the gene signatures extracted with Cell-ID from scRNA-seq datasets were replicated in  
935 independent single-cell ATAC-seq experiments on the same or different tissues of origin. We made use of two  
936 recent large-scale single-cell projects annotating cell-type heterogeneity through in-depth expert curation in a  
937 comprehensive set of organs and tissues for the model organism *Mus musculus*: the Tabula Muris<sup>32</sup> mouse cell  
938 atlas, based on scRNA-seq, and the Mouse ATAC Atlas<sup>33</sup>, based on scATAC-seq. The Tabula Muris project profiled  
939 20 mouse organs and tissues in a total of eight adult mice, with two alternative sequencing technologies: SMART-  
940 Seq2 (53760 cells) and 10X genomics (55656 cells). The Mouse ATAC Atlas reported genome-wide chromatin  
941 accessibility in ~100,000 single cells from 17 samples spanning 13 different tissues in 13 adult mice. Only cells  
942 from the eight organs/tissues common to the Tabula Muris and ATAC Atlas datasets were retained for  
943 downstream analyses: heart, kidney, liver, lung, bone marrow, spleen, thymus and large intestine (large  
944 intestine available only for SmartSeq scRNA-seq and scATACseq). Collectively, the tissues retained contained 50,  
945 43 and 37 different cell types for the 10X Genomics, SMART-Seq2 and scATAC-seq datasets, respectively. Cell-  
946 type nomenclature equivalences between the Tabula Muris and the ATAC Atlas datasets are presented in

947 **Supplementary Table 5.** Comparisons were made across datasets based on gene expression matrices (scRNA-  
948 seq) and gene activity scores (scATACseq) (**Online Methods**).

949 Cell-ID(g) and Cell-ID(c) were applied with default parameters (**Online Methods**) for *reference-to-query*  
950 assignments from (i) SMART-Seq2 and (ii) 10X genomics Tabula Muris scRNA-seq data to the Mouse ATAC Atlas.  
951 Thus, each cell in the Mouse ATAC Atlas was assigned to the cell type from the Tabula Muris scRNA-seq dataset  
952 for which it presented the lowest significant  $p$ -value for gene signature enrichment (or was left unassigned if no  
953 significant hits were found). We performed this procedure independently for the SMART-Seq2 and 10X genomics  
954 datasets. Global matching accuracy was evaluated by assessing the concordance of cell-type labeling between  
955 the manually curated cell type annotations provided by the corresponding databases (**Supplementary Table 5**).  
956 Comparative benchmarking was performed against a representative set of alternative state-of-the-art methods,  
957 as described in **Supplementary Note 3 and 4**, except for scID<sup>15</sup>, which could not be included here due to  
958 computational time limitations (>48 h in our computing infrastructure; **Methods**). By detecting shared gene  
959 signatures, both Cell-ID(g) and Cell-ID(c) correctly matched equivalent cell types across scRNA-seq and scATAC-  
960 seq datasets (**Figure 4, Supplementary Figure 10-12**). Thus, both Cell-ID(g) and Cell-ID(c) maintained a good  
961 balance between precision and recall rates, resulting in high F1 scores (>0.74 for 10X and >0.6 for SmartSeq2),  
962 clearly outperforming all other methods evaluated, with the exception of SingleR, which yielded slightly lower  
963 performances (**Supplementary Figure 12, Supplementary Table 11**). Overall, these results confirm that Cell-ID  
964 can extract, in an unbiased manner, gene signatures that are robustly reproduced across diverse single-cell *omics*  
965 technologies applied across highly heterogeneous cell types from multiple tissues and organs.

966

#### 967 **Supplementary Note 6. Computational details, timing and memory consumption**

968 To evaluate the Cell-ID scalability to massive single-cell RNA-seq dataset analysis, we evaluated its time and  
969 memory consumption for different input sizes, and benchmarked them against the state-of-the-art methods  
970 previously considered in our study (**Supplementary Notes 3, 4 and 5**). To that aim, we performed large-scale cell  
971 mapping and label transfer tasks between different *reference-to-query* datasets using single-cell RNA-seq data  
972 from the Tabula Muris atlas (see **Supplementary Note 5**). All 20 tissues and organs were considered here. For  
973 the purpose of the analysis, we first randomly subset an increasing number of cells (200, 500, 1000, 2000, 5000,  
974 10000, 20000 and 50000 cells) from the 10X and -independently- from the Smart-seq datasets. Sampling of the



975 reference dataset was restricted to the 10 most abundant cell types from 10X. We fixed the dataset to have 10  
976 different subpopulations since the number subpopulations in the reference dataset impacts significantly  
977 computation time for some methods such as scID or SingleR. We then evaluated the computation time and the  
978 total memory allocation of a label transfer task as a function of the number of cells in the *query* dataset. Thus,  
979 label transfer from a fixed *reference* subset of 5000 cells from the SmartSeq dataset was performed against 8  
980 different *query* datasets of increasing size, corresponding to the 8 subsets from the 10X Genomics dataset  
981 previously described (**Supplementary Figure 13 A-B**). In an analogous manner, we evaluated the computation  
982 time and the total memory allocation of a label transfer tasks as a function of the number of cells in the *reference*  
983 dataset. Here, label transfer from 8 different *reference* datasets of increasing size (corresponding to the 8 subsets  
984 from the Smart-seq dataset previously described) against a fixed *query* subset of 5000 cells from the 10X  
985 Genomics dataset (**Supplementary Figure 15 C-D**). Overall, Cell-ID computational time behaved in a comparable  
986 way to other state-of-the-art methods, yet with slightly higher memory consumption. The previous benchmark  
987 is based on the cell matching and label transfer between a pair of datasets. In the case of Cell-ID, such process  
988 involves (i) the MCA low dimensionality reduction and per-cell gene signature extraction for each dataset, and  
989 (ii) cell-to-cell matching between datasets, based on hypergeometric tests on the pre-extracted per-cell gene  
990 signatures. From a computational point of view, each of such cell-to-cell evaluations constitutes a fully  
991 independent job, and thus they can be readily parallelized in a multi-core and multi-node computing cluster. Cell-  
992 ID thus allows the creation of reference libraries of individual cell's gene signatures from massive collections of  
993 single-cell RNA-seq datasets, enabling efficient large-scale cell-to-cell matching and label transfer across sets.

994

#### 995 **Supplementary Note 7. Novel visualization options for enhanced biological interpretation of cell heterogeneity**

996 Cell-ID provides two novel visualization options for the explorative analysis of single-cell RNA-seq data. First, the  
997 dimensionality reduction performed through Multiple Correspondence Analysis provides a simultaneous  
998 representation of cells and genes on the same principal axes. Thus, Cell-ID allows to visualize cells in the MCA  
999 principal components (analogous to a PCA representation), but also to map key gene markers together with the  
1000 cell representation. In such biplots, multiple gene markers can be displayed at once in a way that, the closer a  
1001 marker is represented to a given cell, the more specific to them it is. This is illustrated for two independent  
1002 datasets corresponding to human pancreatic cells <sup>7</sup> and mouse airway epithelial cells <sup>10</sup>, where simultaneous

1003 projection of prototypical markers allows to rapidly identify the cell identity of the corresponding regions  
1004 **(Supplementary Figure 14 A-D).**

1005 Second, Cell-ID provides functional enrichment analysis of the gene signatures obtained for each cell in a dataset,  
1006 using Gene Ontology terms and pathway annotation databases such KEGG, Reactome and WikiPathways.  
1007 Functional enrichment analysis may help on the functional interpretation of cell heterogeneity and assist on cell  
1008 type identification, as illustrated for Schwann cells **(Supplementary Note 3)** and Solitary Chemosensory Cells  
1009 (SCCs, **Supplementary Note 4**) in the exploratory analysis of pancreas and airway epithelium single-cell RNA-seq  
1010 datasets, respectively. Thus, functional enrichments can be visualized in a low-dimensionality representation of  
1011 cells by coloring each cell with an intensity proportional to its  $-\log_{10}$  p-value for a query functional term or  
1012 biological pathway **(Supplementary Figure 14 E-F)**. Cell-ID R package provides per-cell functional enrichment  
1013 analysis as a built-in function, which stores the enrichment outputs as additional cell attributes in standard R  
1014 single-cell data structures, such as SingleCellExperiment and Seurat objects. Such seamless integration allows  
1015 other single-cell RNA-seq tools to use enrichment scores in alternative cell visualizations, e.g. UMAP<sup>34</sup> or cell  
1016 diffusion maps<sup>35</sup>.

## 1017 **Supplementary references**

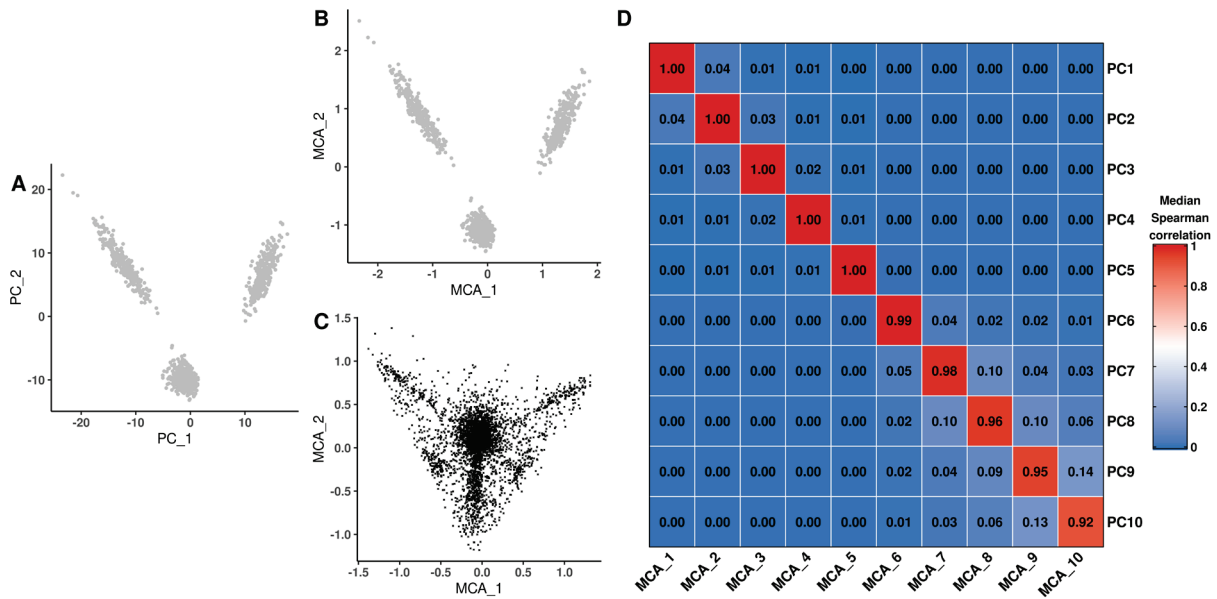
- 1018 1. Aibar, S. *et al.* SCENIC: Single-cell regulatory network inference and clustering. *Nat. Methods* **14**, 1083–1086  
1019 (2017).
- 1020 2. Stoekius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells. *Nat. Methods* **14**,  
1021 865–868 (2017).
- 1022 3. Peterson, V. M. *et al.* Multiplexed quantification of proteins and transcripts in single cells. *Nat. Biotechnol.*  
1023 **35**, 936–939 (2017).
- 1024 4. Aran, D., Hu, Z. & Butte, A. J. xCell: digitally portraying the tissue cellular heterogeneity landscape. *Genome*  
1025 *Biol.* **18**, 220 (2017).
- 1026 5. Zhang *et al.* SCINA: Semi-Supervised Analysis of Single Cells in Silico. *Genes* **10**, 531–531 (2019).
- 1027 6. Zhang, A. W. *et al.* Probabilistic cell-type assignment of single-cell RNA-seq for tumor microenvironment  
1028 profiling. *Nat. Methods* **16**, 1007–1015 (2019).

- 1029 7. Baron, M. *et al.* A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and  
1030 Intra-cell Population Structure. *Cell Syst.* **3**, 346–360 (2016).
- 1031 8. Muraro, M. J. *et al.* A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Syst.* **3**, 385-394.e3 (2016).
- 1032 9. Segerstolpe, Å. *et al.* Single-Cell Transcriptome Profiling of Human Pancreatic Islets in Health and Type 2  
1033 Diabetes. *Cell Metab.* **24**, 593–607 (2016).
- 1034 10. Plasschaert, L. W. *et al.* A single-cell atlas of the airway epithelium reveals the CFTR-rich pulmonary  
1035 ionocyte. *Nature* **560**, 377–381 (2018).
- 1036 11. Montoro, D. T. *et al.* A revised airway epithelial hierarchy includes CFTR-expressing ionocytes. *Orit*  
1037 *Rozenblatt-Rosen* **4**, 19–19 (2018).
- 1038 12. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21 (2019).
- 1039 13. Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-sequencing data  
1040 are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–427 (2018).
- 1041 14. Kiselev, V. Y., Yiu, A. & Hemberg, M. scmap: projection of single-cell rna-seq data across data sets. **15**, 359–  
1042 359 (2018).
- 1043 15. Boufeua, K., Seth, S. & Batada, N. N. scID Uses Discriminant Analysis to Identify Transcriptionally Equivalent  
1044 Cell Types across Single-Cell RNA-Seq Data with Batch Effect. *iScience* **23**, 100914 (2020).
- 1045 16. Aran, D. *et al.* Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic  
1046 macrophage. *Nat. Immunol.* **20**,
- 1047 17. Lieberman, Y., Rokach, L. & Shay, T. CaSTLe – Classification of single cells by transfer learning: Harnessing  
1048 the power of publicly available single cell RNA sequencing experiments to annotate new experiments. *PLOS*  
1049 *ONE* **13**, e0205499–e0205499 (2018).
- 1050 18. Tan, Y. & Cahan, P. SingleCellNet: A Computational Tool to Classify Single Cell RNA-Seq Data Across  
1051 Platforms and Across Species. *Cell Syst.* **9**, 207-213.e2 (2019).
- 1052 19. De Kanter, J. K., Lijnzaad, P., Candelli, T., Margaritis, T. & Holstege, F. C. P. CHETAH: a selective, hierarchical  
1053 cell type identification method for single-cell RNA sequencing. *Nucleic Acids Res.* **47**, (2019).
- 1054 20. Bhatheja, K. & Field, J. Schwann cells: Origins and role in axonal maintenance and regeneration. *Int. J.*  
1055 *Biochem. Cell Biol.* **38**, 1995–1999 (2006).
- 1056 21. Chernousov, M. A. *et al.* Development/Plasticity/Repair Glypican-1 and 4(V) Collagen Are Required for  
1057 Schwann Cell Myelination. (2006) doi:10.1523/JNEUROSCI.2544-05.2006.

- 1058 22.Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for gene  
1059 and protein annotation. *Nucleic Acids Res.* **44**, 457–462 (2015).
- 1060 23.Jassal, B. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498–D503 (2020).
- 1061 24.Slenter, D. N. *et al.* WikiPathways: a multifaceted pathway database bridging metabolomics to other omics  
1062 research. *Nucleic Acids Res.* **46**, D661–D667 (2018).
- 1063 25.Bankova, L. G. *et al.* The cysteinyl leukotriene 3 receptor regulates expansion of IL-25–producing airway  
1064 brush cells leading to type 2 inflammation. *Sci. Immunol.* **3**, (2018).
- 1065 26.Ualiyeva, S. *et al.* Airway brush cells generate cysteinyl leukotrienes through the ATP sensor P2Y2. *Sci.*  
1066 *Immunol.* **5**, eaax7224–eaax7224 (2020).
- 1067 27.Bouchery, T. & Marsland, B. J. Airway brush cells: Not as “tuft” as you might think. *Sci. Immunol.* **3**, (2018).
- 1068 28.Kohanski, M. A. *et al.* Solitary chemosensory cells are a primary epithelial source of IL-25 in patients with  
1069 chronic rhinosinusitis with nasal polyps. *J. Allergy Clin. Immunol.* **142**, 460-469.e7 (2018).
- 1070 29.Haber, A. L. *et al.* A single-cell survey of the small intestinal epithelium. (2017) doi:10.1038/nature24489.
- 1071 30.Wu, Y. *et al.* A Population of Navigator Neurons Is Essential for Olfactory Map Formation during the Critical  
1072 Period Article A Population of Navigator Neurons Is Essential for Olfactory Map Formation during the Critical  
1073 Period. *Neuron* **100**, 1066-1082.e6 (2018).
- 1074 31.Fletcher, R. B. *et al.* Deconstructing Olfactory Stem Cell Trajectories at Single-Cell Resolution. *Cell Stem Cell*  
1075 **20**, 817-830.e8 (2017).
- 1076 32.Schaum, N. *et al.* Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris the tabula Muris  
1077 consortium\*. (2018) doi:10.1038/s41586-018-0590-4.
- 1078 33.Cusanovich, D. A. *et al.* A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility. *Cell* **174**, 1309-  
1079 1324.e18 (2018).
- 1080 34.Becht, E. *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.* **37**, 38–  
1081 44 (2019).
- 1082 35.Haghverdi, L., Buettner, F. & Theis, F. J. Diffusion maps for high-dimensional single-cell analysis of  
1083 differentiation data. *Bioinformatics* **31**, 2989–2998 (2015).
- 1084
- 1085

1086 **Supplementary Figures**

1087



1088

1089

1090

1091 **Supplementary Figure 1. Correspondence between PCA and MCA low-dimensional representations of cells.**

1092 Low-dimensional representation of a simulated scRNA-seq dataset on the two first principal axes obtained

1093 through PCA (**A**, cell projection) or MCA (**B**, cell projection, and **C**, gene projection; **Supplementary Note 1**).

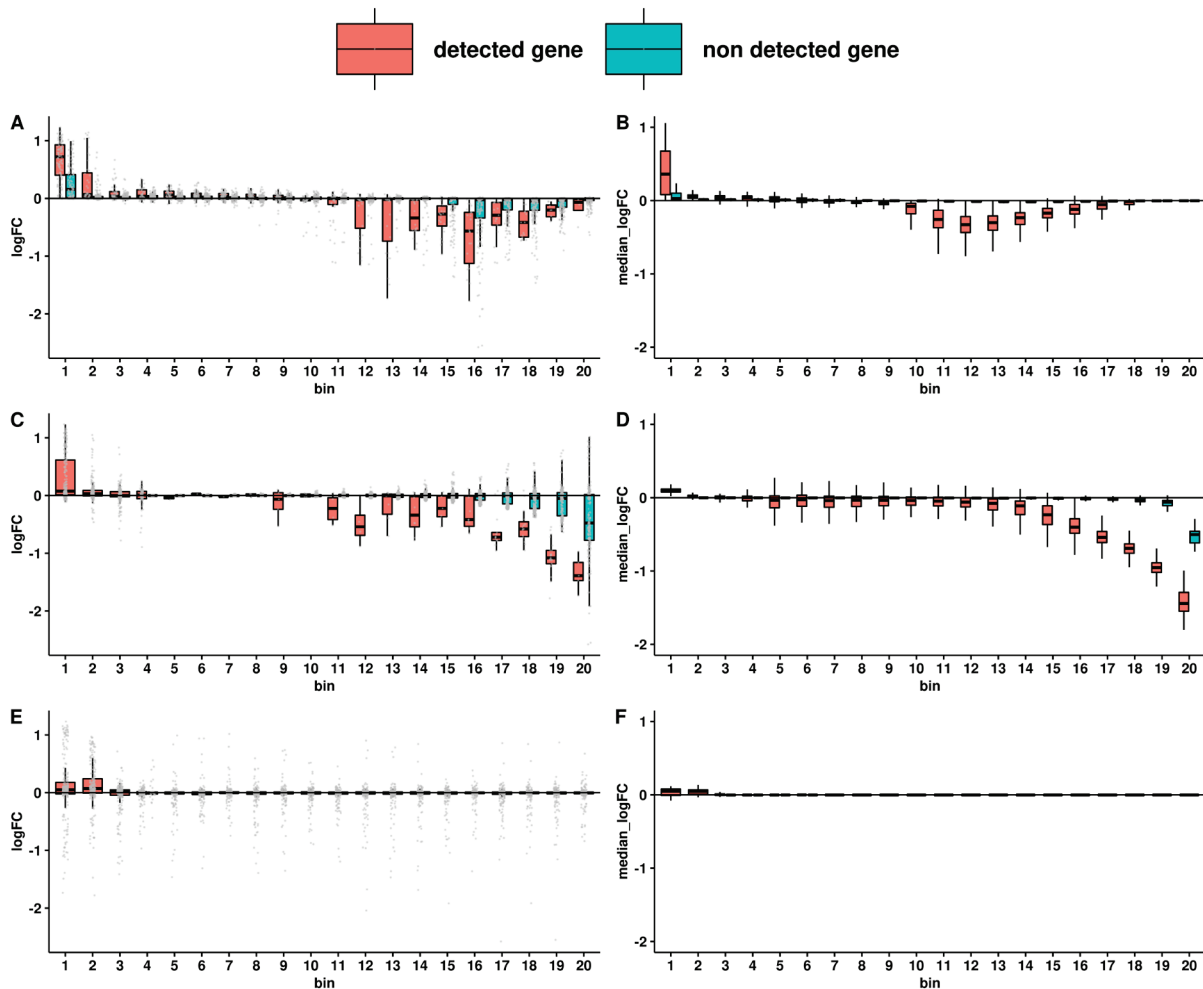
1094 (**D**) Median values of Spearman's correlation coefficient for the relationship between cell coordinates on the first 10

1095 principal axes from MCA ( $x$ -axis) and the first 10 principal axes from PCA ( $y$ -axis) over 100 simulated scRNA-seq

1096 datasets. Median absolute Spearman's correlation coefficient values are represented on a color scale ranging

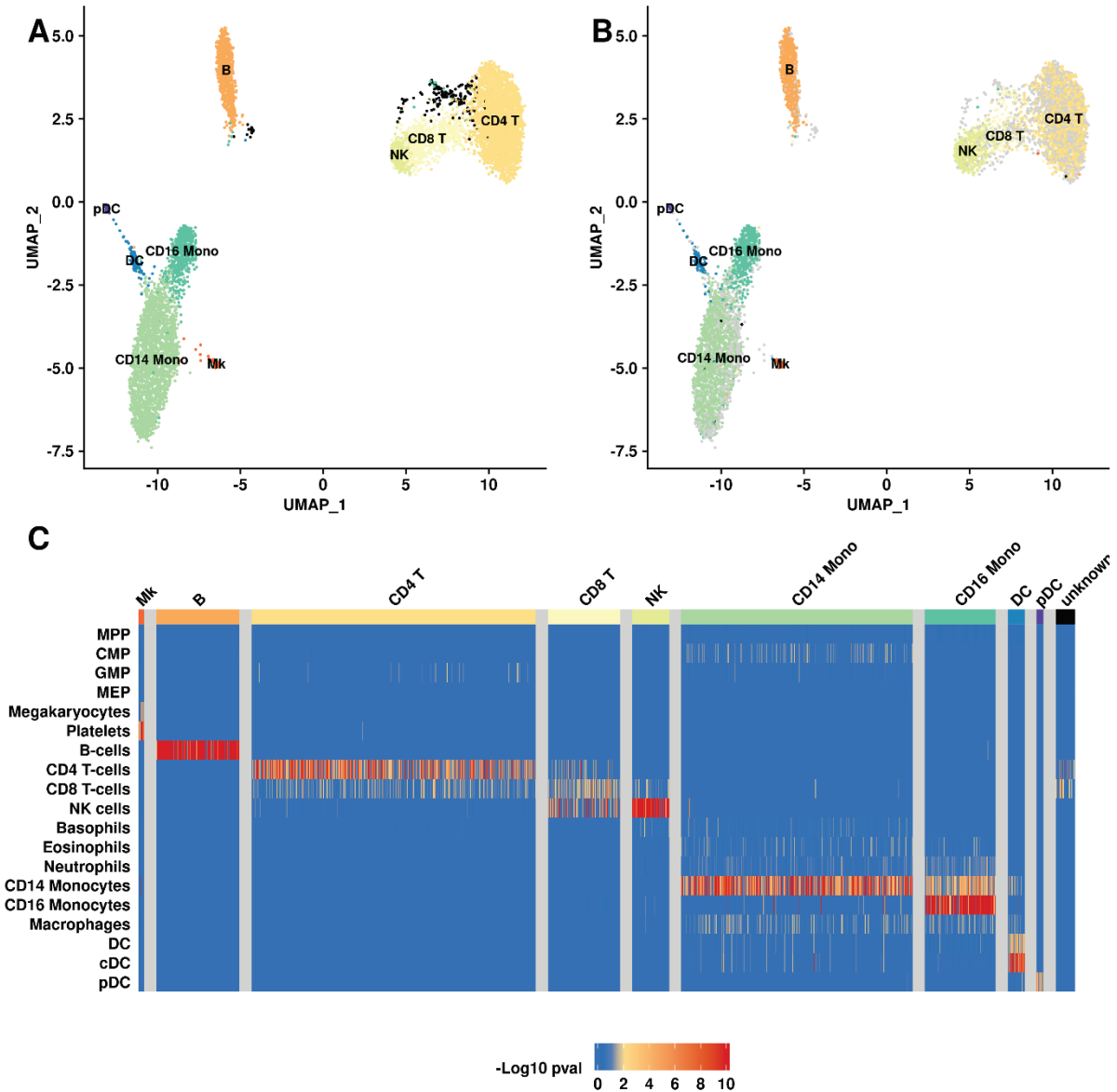
1097 from 0 (dark blue) to 1 (dark red), and the precise value is indicated inside.

1098



1099  
1100

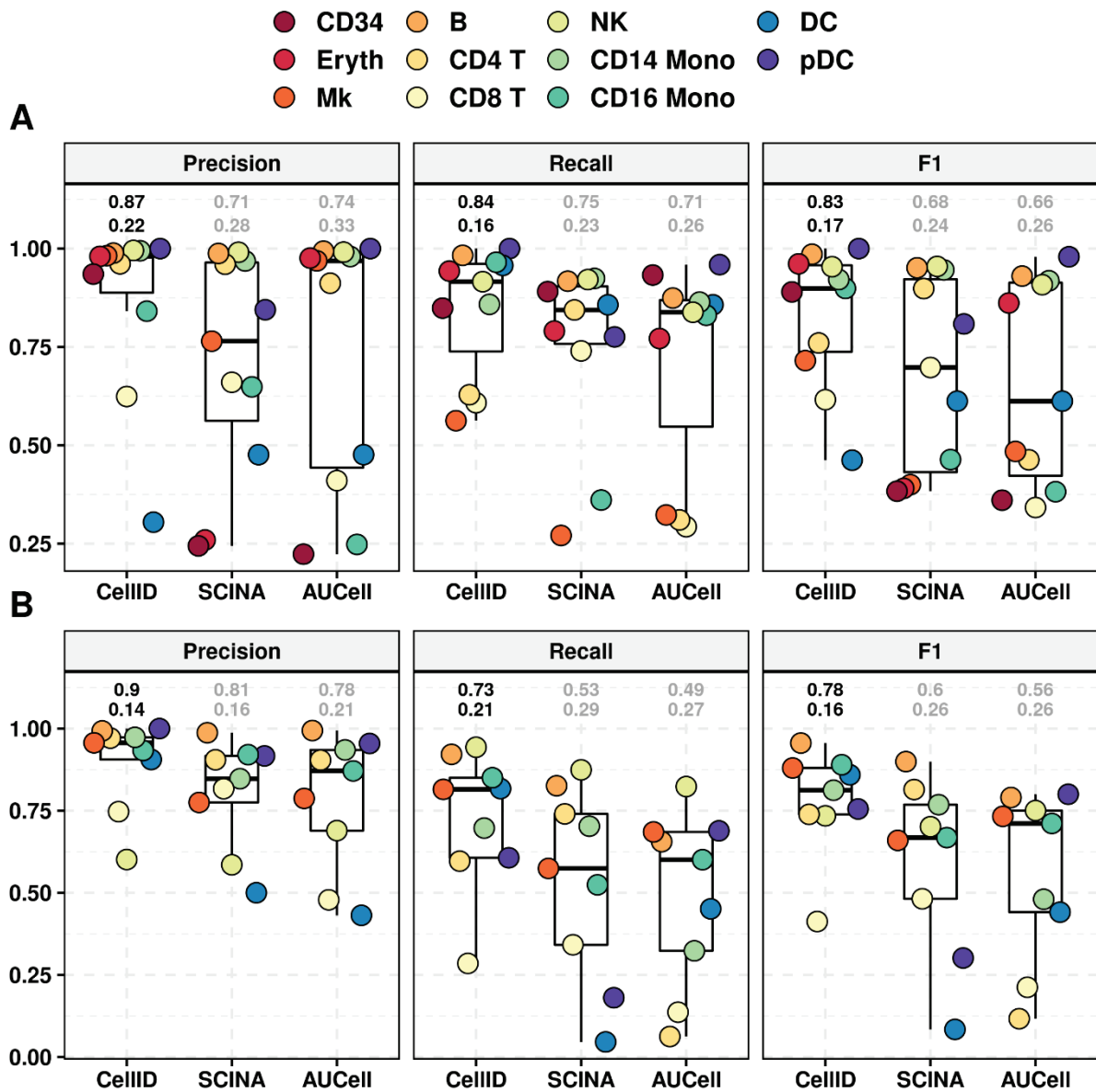
1101 **Supplementary Figure 2. Consistency of the per-cell gene rankings with the gene expression levels of**  
 1102 **neighboring cells in the MCA space. (A)** An individual cell was randomly selected from a dataset randomly  
 1103 selected from 100 simulated scRNA-seq sets, each containing 1000 cells and 5000 genes (**Supplementary Note**  
 1104 **1**). For each gene (grey dot), this figure shows the log-fold change in expression relative to (i) its mean expression  
 1105 value in the 5% of cells closest to the target cell ( $n=50$ ), and (ii) its mean expression value in the rest of the cells  
 1106 in the dataset ( $y$ -axis). Genes were grouped into 20 bins of equal size ( $x$ -axis, on the basis of their ranking relative  
 1107 to the target cell on the MCA space (**Figure 1**). High ranks (on the left side of the  $x$ -axis) correspond to small gene-  
 1108 to-cell distances, indicating that gene expression is highly specific to the target cell, whereas lower ranks (on the  
 1109 right) reflect a lack of specificity to the target cell. The gene expression values in the target cell were not  
 1110 considered in the assessment of such log-fold changes. For each bin, two boxplots are shown, summarizing (i)  
 1111 the distribution of log-fold changes in expression across genes for which the target cell presented a non-zero  
 1112 count (red); and (ii) the distribution of log-fold changes in expression across genes for which the target cell  
 1113 presented a zero count (blue). **(B)** Analogous figure to **(A)**, showing the generalization of patterns to all individual  
 1114 cells in a randomly selected dataset. Each individual cell was first independently assessed as in **(A)**, and its per-  
 1115 bin median values were extracted to plot a distribution across cells (boxplots in **B**). Figures analogous to **(A)** and  
 1116 **(B)** are shown when the per-cell gene rankings displayed on the  $x$ -axis through bins of equal size were obtained  
 1117 from (i) a naïve approach based on the log-fold changes in gene expression observed in a cell relative to all the  
 1118 other cells in the dataset (**C, D**) or (ii) from highest-to-lowest expression values within a cell, with random ranks  
 1119 for ties, as previously described (AUCell, **E, F**).  
 1120



1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135

**Supplementary Figure 3. Automatic Cell-ID cell-type identification of human peripheral blood mononuclear cells with pre-established signatures.** UMAP representation of 7488 peripheral blood mononuclear cells (PBMCs) profiled with a REAP-seq protocol, with color-coding of the cells according to blood cell type classification based on (A) single-cell protein marker levels as provided, and (B) Cell-ID predictions based on a reference collection of well-established blood cell signatures. (C) Heatmap representing, for each individual cell (displayed in columns), the  $-\log_{10}$  transformed  $p$  value obtained by Cell-ID testing of the association of the gene signature with each of the evaluated pre-established signatures (displayed in rows). The heatmap color scale extends from dark blue ( $p$  value=1) to yellow ( $p$  value =  $10^{-2}$ ) to dark red ( $p$  value =  $10^{-10}$ ), with  $p$  values  $<10^{-10}$  fixed at this value). Non-significant associations ( $p$  value  $>10^{-2}$  after Benjamini Hochberg correction for the number of gene signatures tested) are shown in blue. The columns in the heatmaps were grouped by the reference cell-type label, as indicated by the colored bands at the top and in the associated legend.

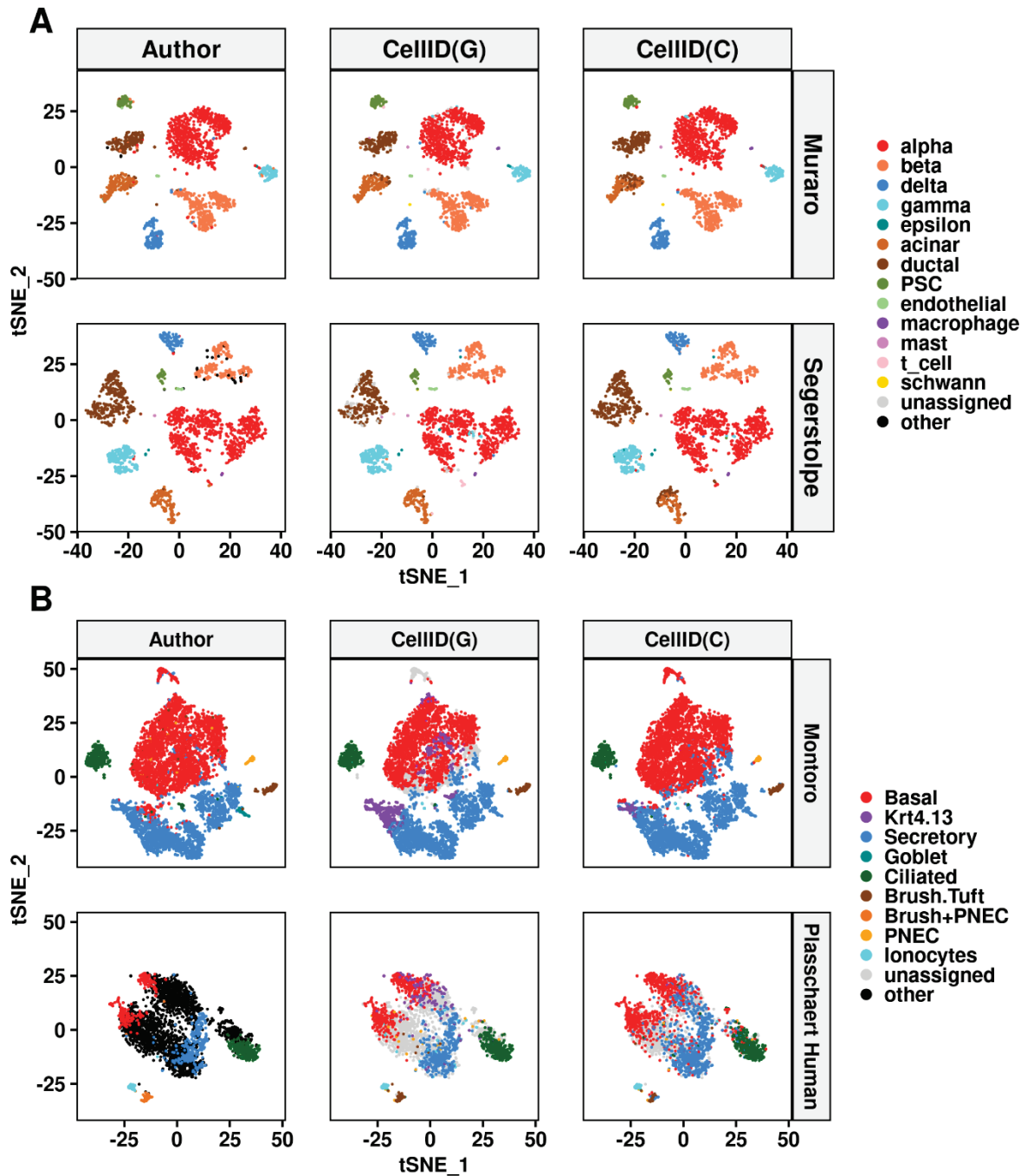
1136



1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145

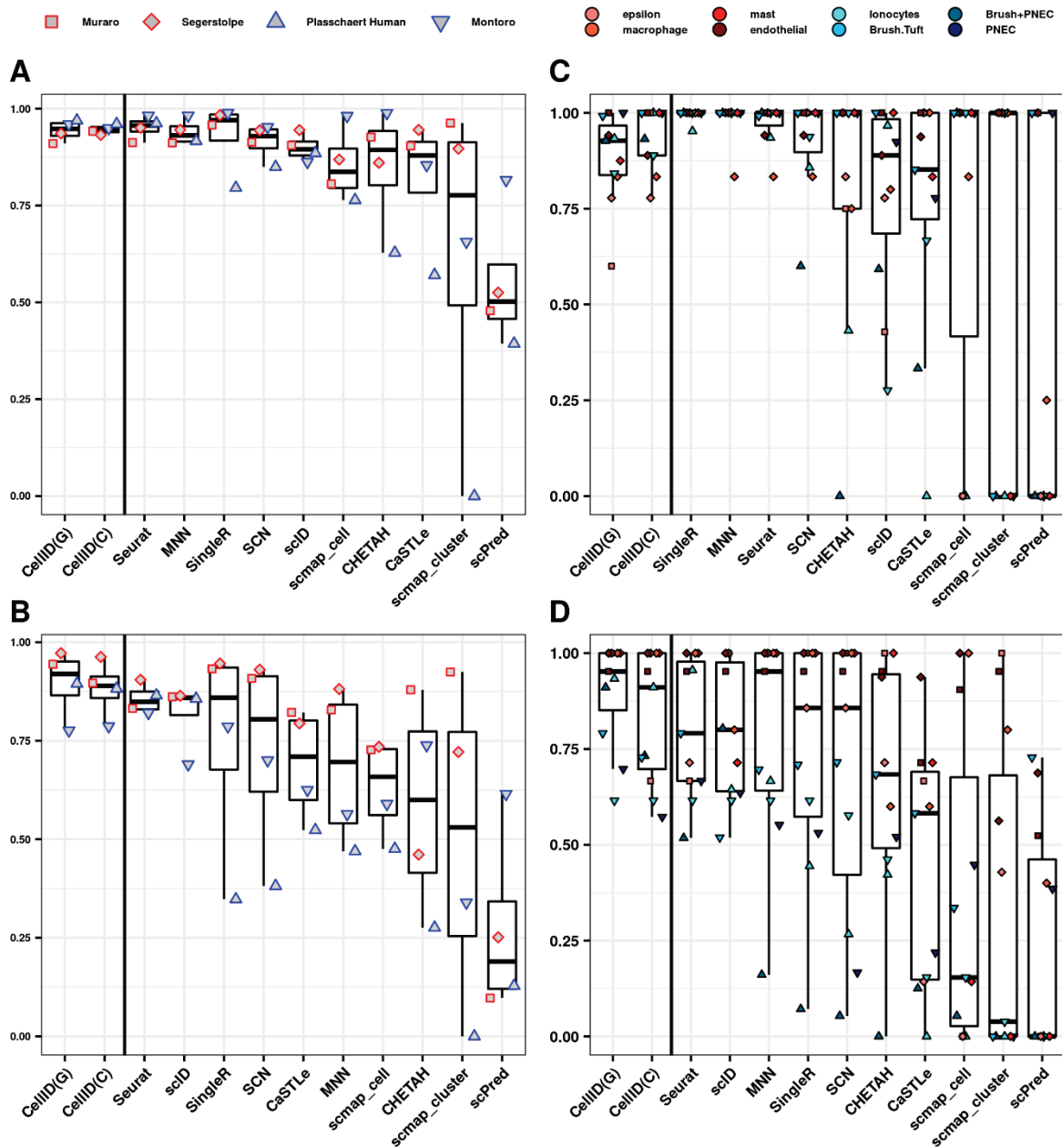
**Supplementary Figure 4.** Performance measured through Precision, recall, and F1 scores achieved by Cell-ID, AUCell and SCINA cell type predictions on (A) CBMCs Cite-Seq data, and (B) PBMCs Reap-Seq data, for each of the blood cell types reported in the original publications. Boxplots summarize the corresponding scores for each method. The numbers above boxplots denote the global performance (macro F1 score, upper digits) and its standard deviation (lower digits), where the maximum and minimum values across methods are colored in black and grey, respectively.





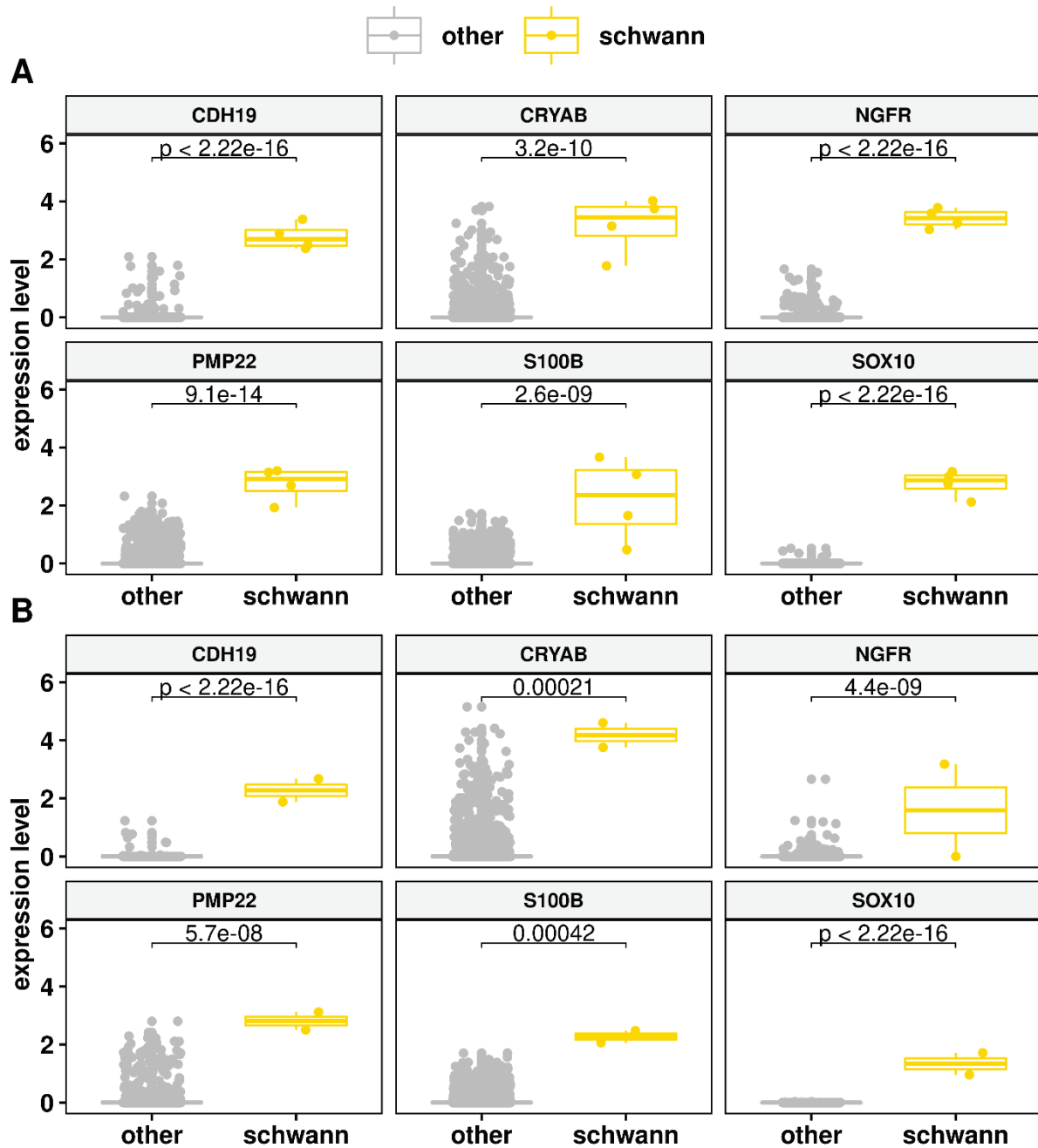
1146  
1147

1148 **Supplementary Figure 5.** tSNE representation of (A) human pancreatic islet cells from the Muraro (top panels)  
1149 and Segerstolpe (bottom panels) datasets, and (B) airway epithelium cells from the Montoro's mouse dataset  
1150 (top panels) and Plasschaert's human dataset (bottom panels). Cells are color-coded according to the cell type  
1151 labels annotated in the original publications (left panels), as well as by the cell type predictions from Cell-ID(g)  
1152 (middle panels), and CellID(c) (right panels).  
1153



1154  
1155  
1156  
1157  
1158  
1159  
1160

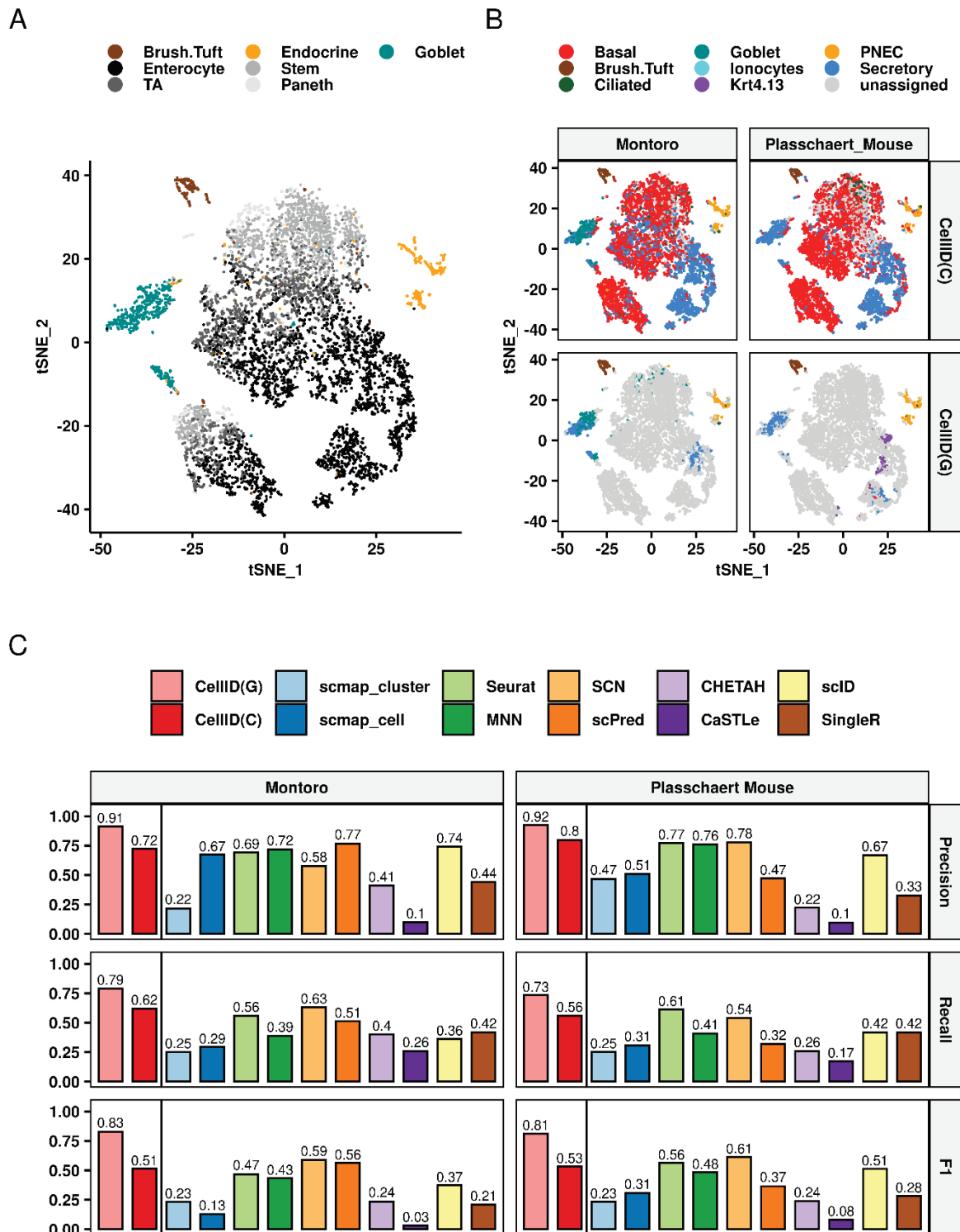
**Supplementary Figure 6. Precision and recall of Cell-ID cell-to-cell matching across independent scRNA-seq datasets from the same or different tissue of origin, within and across species.** Performance achieved by Cell-ID(g), Cell-ID(c) and 10 alternative state-of-the-art methods (x-axis), measured through Precision (**A and C**), and Recall (**B and D**). Panels are analogous to the ones described in **Figure 3 (A-B)** legend.



1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170

**Supplementary Figure 7.** Distribution of log-transformed normalized gene expression levels (y-axis) of Schwann cell markers (SOX10, S100B, CRYAB, NGFR, CDH19, and PMP22) across pancreatic islet cells from (A) Segerstolpe and (B) Muraro datasets. Values their associated boxplots are represented in yellow for the  $n=4$  and  $n=2$  cells in the Muraro and Segerstolpe datasets respectively that were identified by both Cell-ID(c) and Cell-ID(g) as putative Schwann cells, and in grey for the rest of cells. Brackets and  $p$  values indicate the results of two-sided Wilcoxon rank sum tests comparing the distribution across putative Schwann cells versus the other cells in the dataset.

1171



1172

1173

1174

1175

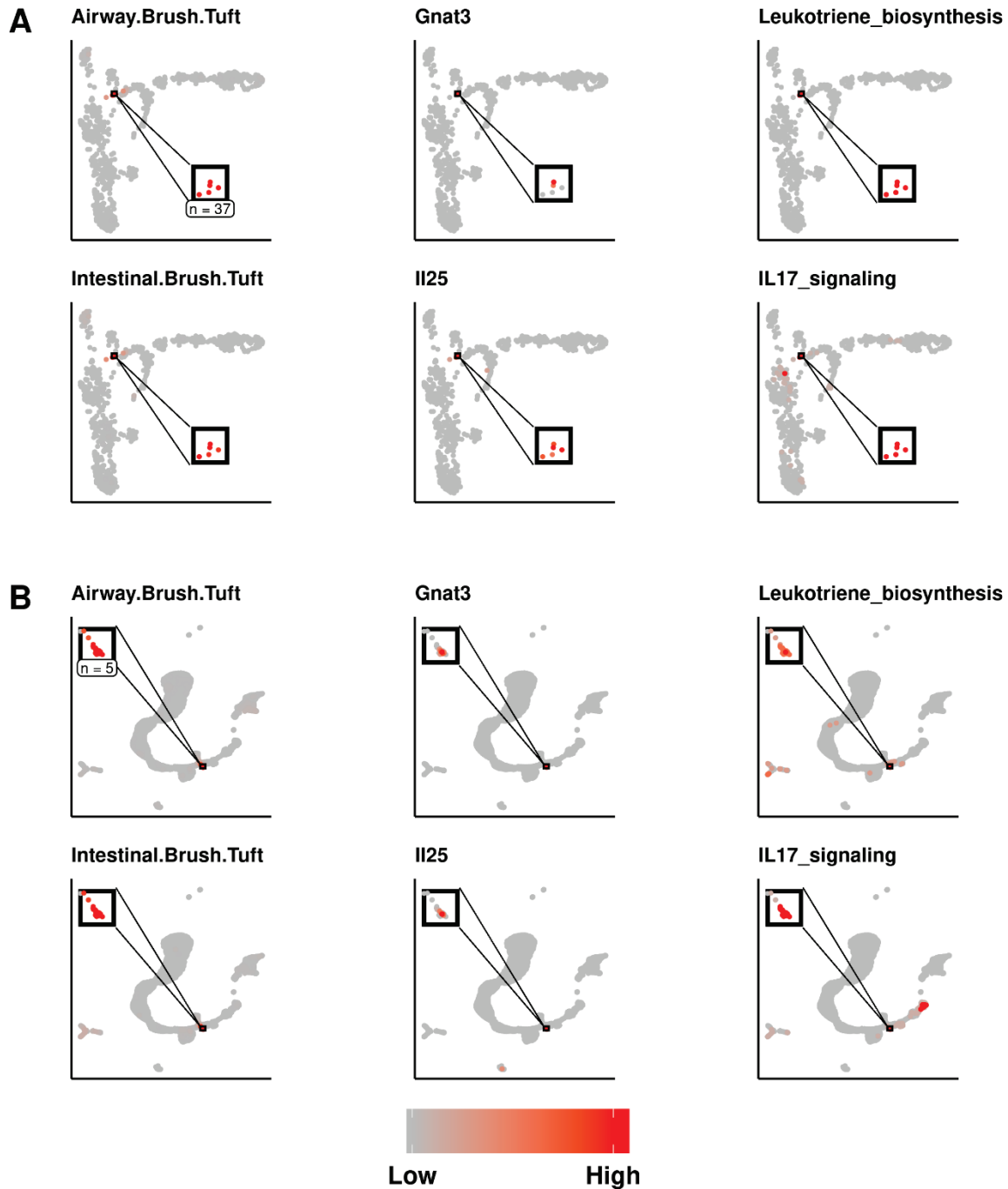
1176

1177

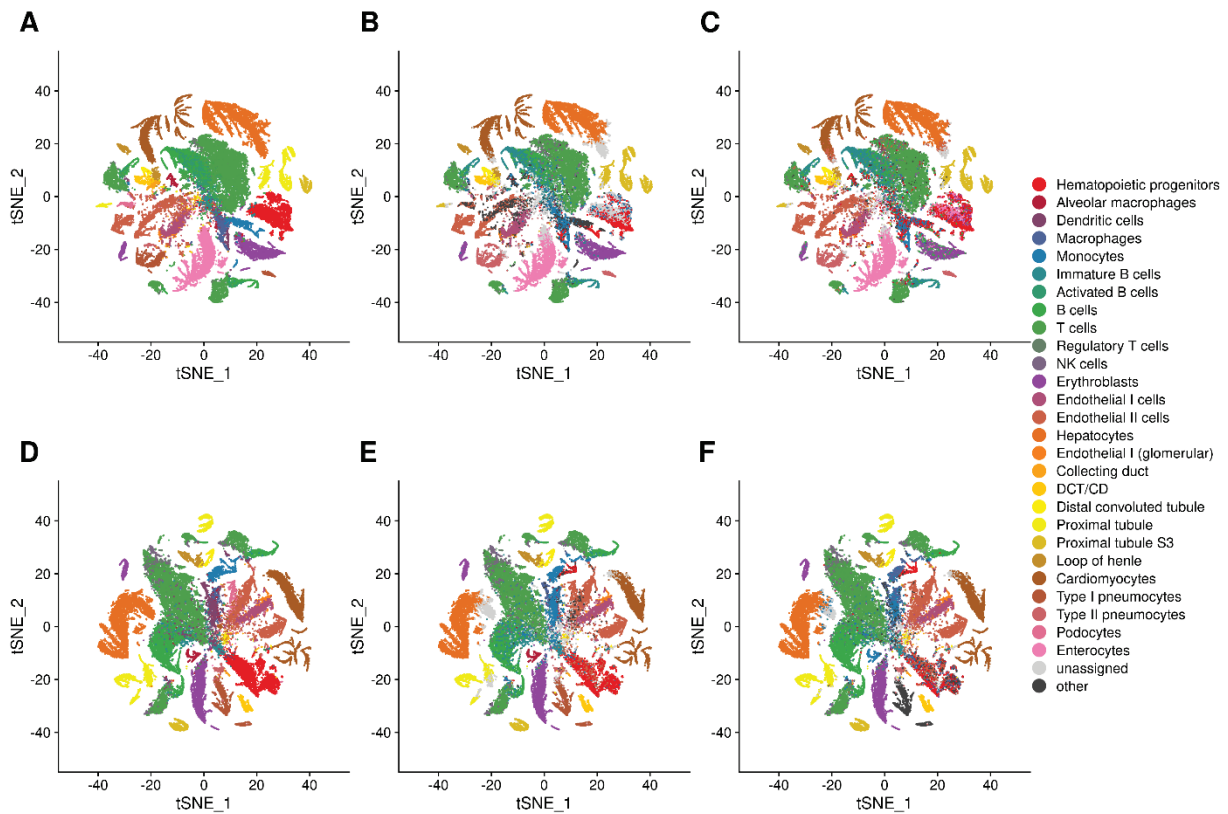
1178

1179

**Supplementary Figure 8.** t-SNE representation of 7216 cells from mouse small intestinal epithelium, where dots representing cells are color coded according to (A) manual cell type annotations provided in Haber et al. and (B) Cell-ID(c) and Cell-ID(g) cell type predictions (top and bottom panels, respectively), using as a reference the mouse airway epithelial gene signatures extracted from Montoro and Plasschaert datasets (left and right panels, respectively). Cells with no significant enrichments were left unassigned and are displayed in grey. (C) Precision, Recall and F1 score (y-axis), achieved by Cell-ID(g), Cell-ID(c) and 10 alternative state-of-the-art methods (x-axis), for the label transferring results depicted in (B).

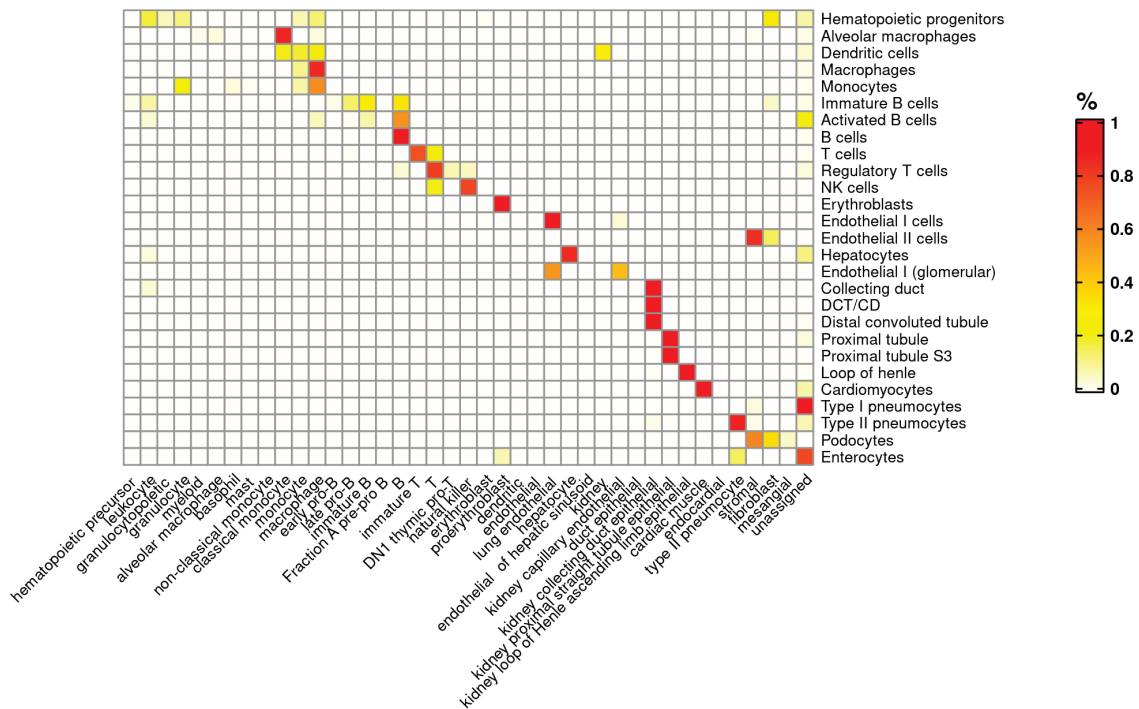


1180  
1181 **Supplementary Figure 9.** UMAP representation of olfactory epithelium cells from (A) Fletcher et al. and (B) Wu  
1182 et al. datasets. The inset panels focus on cells identified by Cell-ID as putative SCCs. Cells represented by dots are  
1183 color-coded according to (i) their Cell-ID(g)  $-\log_{10}$  enrichment  $p$ -value using as a reference the mouse brush/tuft  
1184 gene signatures extracted from mouse airway epithelium, and from small intestinal epithelium, as indicated in  
1185 the panel titles; (ii) the log-normalized expression of two solitary chemosensory cells markers, i.e. interleukin 25  
1186 (IL25) and G protein subunit alpha Transducin 3 (GNAT3); and (iii) The  $-\log_{10}$  enrichment  $p$ -values for 2 functional  
1187 terms, i.e.: “leukotriene biosynthetic process” (Gene Ontology term GO:0019370) and “interleukin-17-mediated  
1188 signalling pathway” (GO:0097400). The color scale throughout panels extends from gray (indicating a non-  
1189 significant  $p$  value or a lack of detection of gene expression) to red, corresponding to a significant  $p$  value or high  
1190 level of gene expression.  
1191



1192  
1193 **Supplementary Figure 10.** (A-C) t-SNE representation of 57370 cells profiled with scATAC-seq from the Mouse  
1194 ATAC atlas corresponding to the eight tissues for which Tabula Muris SmartSeq scRNA-seq datasets are available,  
1195 i.e.: heart, kidney, liver, lung, bone marrow, spleen, thymus and large intestine. (D-E) t-SNE representation of  
1196 50284 cells profiled with scATAC-seq from the Mouse ATAC atlas corresponding to the seven tissues for which  
1197 Tabula Muris 10X Genomics scRNA-seq datasets are available, i.e.: heart, kidney, liver, lung, bone marrow,  
1198 spleen, thymus. Cells are colored according to the manually annotated cell types provided by the atlas,  
1199 regrouped by the categories represented in the legend and described in **Supplementary Table 5**, (A and D), as  
1200 well as by the cell type predictions from Cell-ID(g) (B and E), and CellID(c) (C and F) using the group gene  
1201 signatures extracted from (A and D, respectively).  
1202

1203



1204

1205

1206

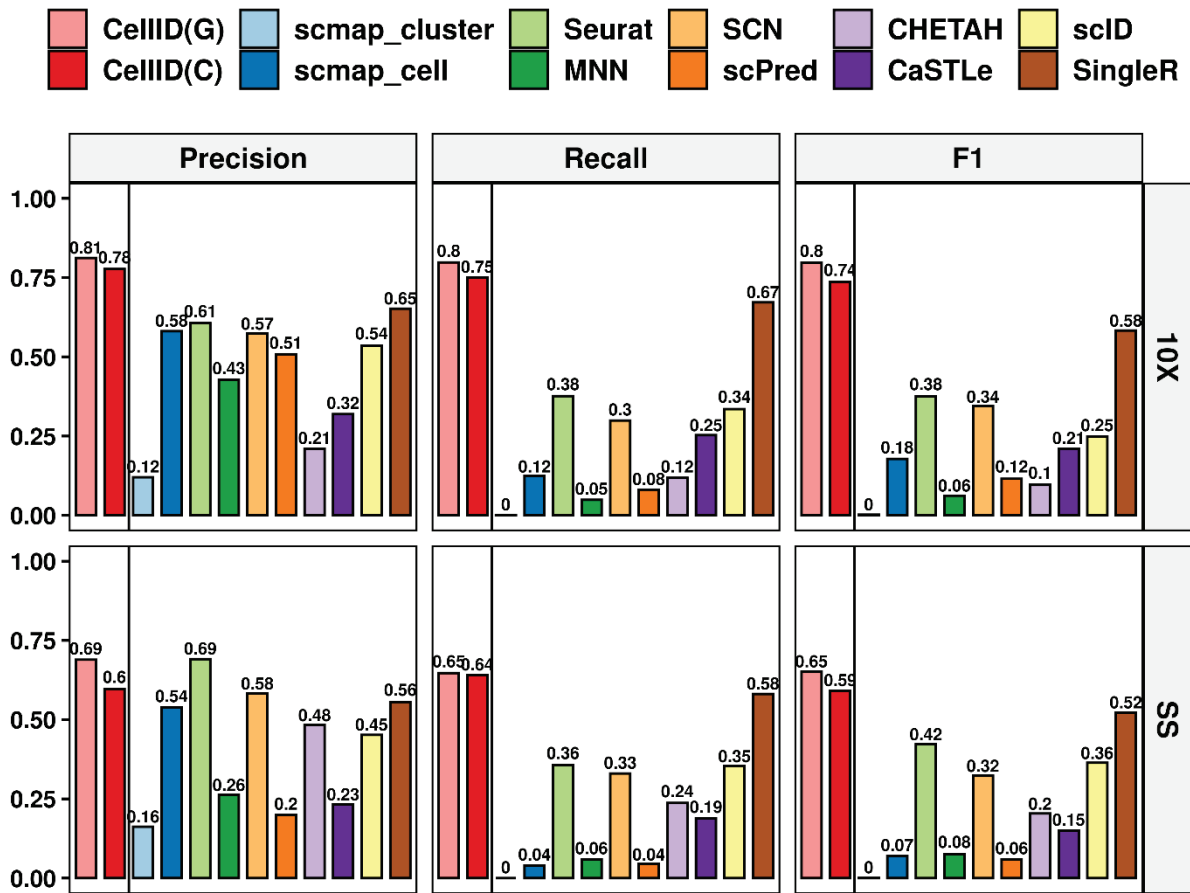
1207

1208

1209

1210

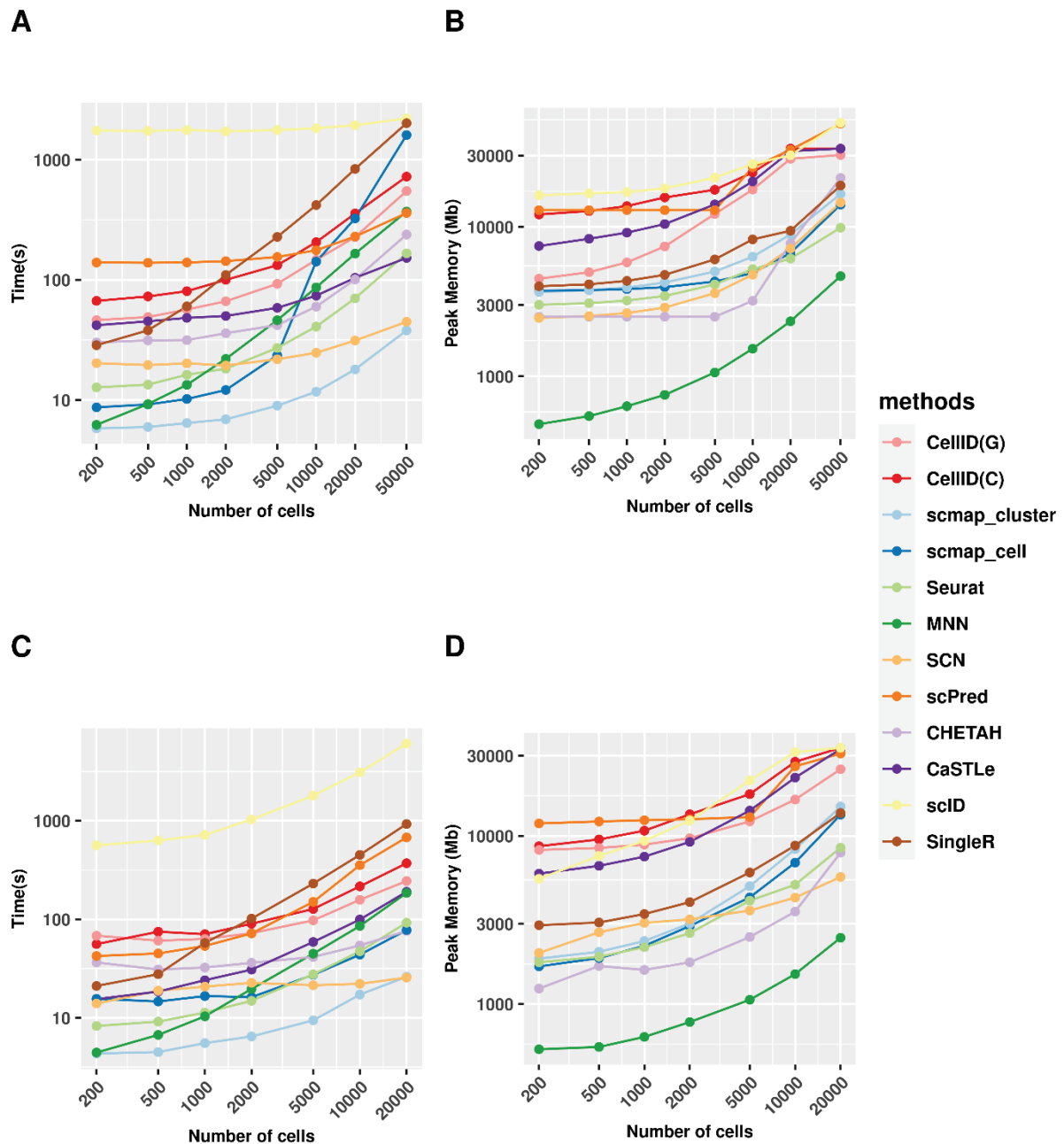
**Supplementary Figure 11.** Heatmap representing the confusion matrix between the manually curated cell type annotations from the Mouse ATAC atlas (displayed per rows) and the Cell-ID(c) cell type predictions (displayed per columns) using the gene-signatures extracted from the manually annotated cell types provided by the Tabula Muris mouse cell atlas. The color code in the heatmap represents the ratio  $r$  of the cell types displayed per rows that are allocated in the cell types represented per columns, ranging from white ( $r = 0$ ) to red ( $r = 1$ ).



1211  
1212  
1213  
1214  
1215

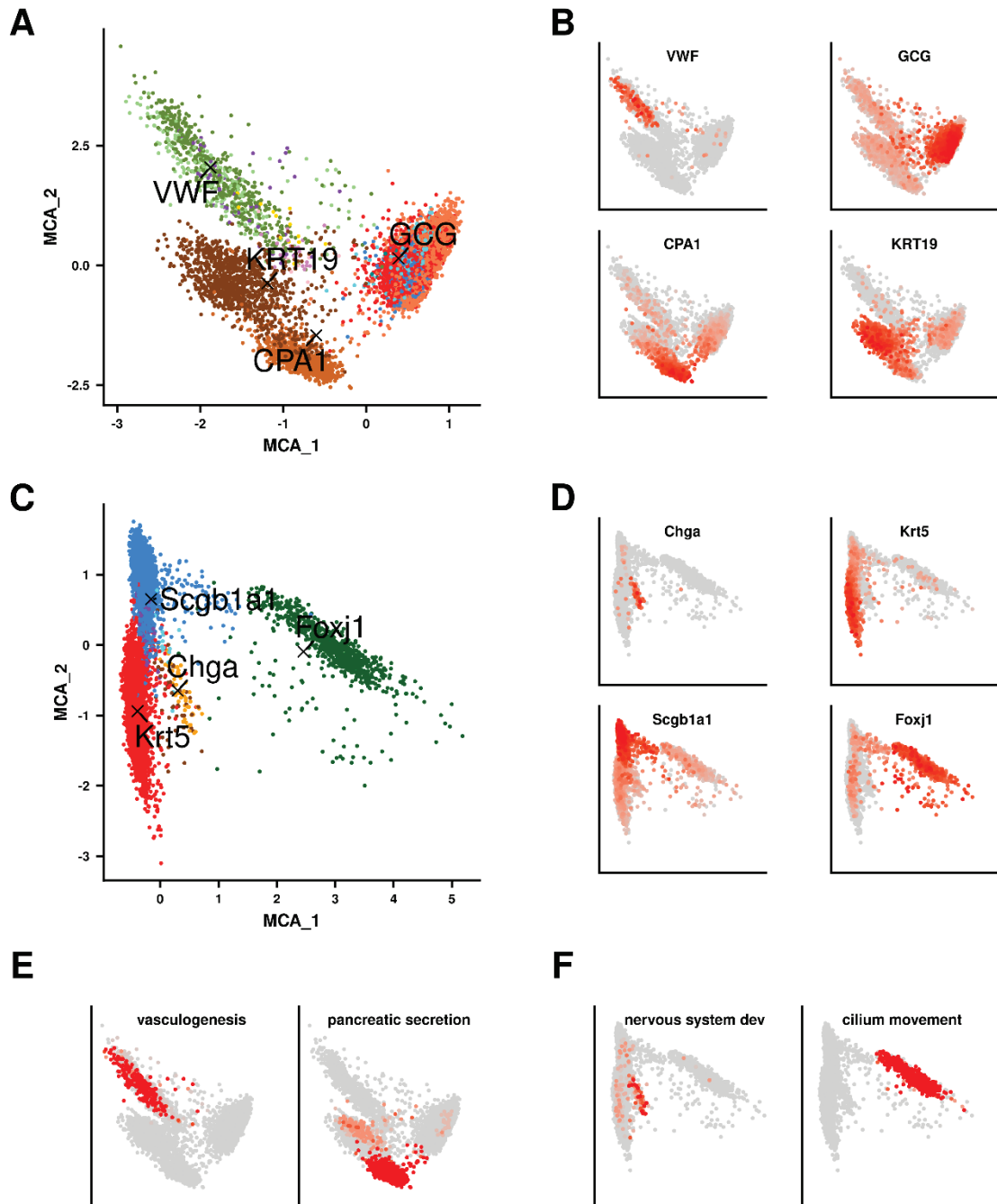
**Supplementary Figure 12.** Performance measured as F1 score (y-axis) achieved by Cell-ID(g), Cell-ID(c) and 10 alternative state-of-the-art methods (x-axis), in the cell type label transferring from Tabula Muris scRNAseq to scATAC Atlas for 10X (top panel) and smarts-seq (bottom panel), as depicted in Supplementary Figure 10.





1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227

**Supplementary Figure 13. Computational time and memory consumption of label transfer methods on simulated datasets.** (A and B) Line plots showing on the y axis (A) the computation time and (B) the total memory allocation as a function of the number of cells (x axis) randomly sampled from the Tabula Muris 10X data that were used as *query* dataset for cell type label transferring from a *reference* dataset of 5000 randomly sampled cells from the Tabula Muris SmartSeq dataset. Results are plotted for Cell-ID(g), Cell-ID(c) and the 10 state of the art methods evaluated. (C and D) Line plots showing on the y axis (C) the computation time and (D) the total memory allocation as a function of the number of cells (x axis) randomly sampled from the Tabula Muris SmartSeq data used as *reference* dataset for cell type label transferring on a *query* dataset of 5000 randomly sampled cells from the Tabula Muris 10X data. Results are plotted for Cell-ID(g), Cell-ID(c) and the 10 state of the art methods evaluated.



1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

**Supplementary Figure 14. Novel visualization options provided by Cell-ID for the explorative analysis of single-cell RNA-seq data.** Simultaneous MCA representation of prototypical marker genes (indicated by back cross symbols and gene name labels) on (A) Baron human pancreatic cells, and (C) Plasschaert mouse airway cells, color-coded according to the cell type labels annotated in the original publication. (B) and (D) Equivalent representation of cells in MCA space as depicted in (A) and (C) where cells are color coded according to their levels of expression for the corresponding selected markers. (E) and (F) Equivalent representation of cells in MCA space as depicted in (A) and (C) where cells are color coded according to their functional enrichment  $-\log_{10}$  p-value for a given functional term or biological pathway, as indicated in the figure titles.

## 1239 Supplementary Tables Legends

1240 **Supplementary Table 1.** Functional enrichment analysis results for putative Schwann cells identified in pancreas  
1241 datasets, and for putative solitary chemosensory cells identified in olfactory epithelium datasets.

1242 **Supplementary Table 2.** Summary of the datasets used in this article.

1243 **Supplementary Table 3.** Summary of external tools and packages used in this article.

1244 **Supplementary Table 4.** Blood cell-type gene markers from the XCell repository used in this study.

1245 **Supplementary Table 5.** Mapping scheme for the original cell population labels from query data to cell  
1246 population labels in the reference data.

1247 **Supplementary Table 6** Summary of the functional pathways used in this article.

1248 **Supplementary Table 7.** Performance metrics for cell type prediction using pre-established marker lists on blood  
1249 cells from CITE-seq and REAP-seq datasets. The table reports precision, recall and F1 score for CellID, AUCell and  
1250 SCINA for each of the evaluated cell types as well as for the overall assessment.

1251 **Supplementary Table 8.** Performance metrics for the prediction of rare cell types from pre-established immune  
1252 signatures on 100 simulated subsets of CITEseq.

1253 **Supplementary Table 9.** Performance metrics for cell type predictions represented in Figure 3A and  
1254 Supplementary Figure 6 A-B.

1255 **Supplementary Table 10.** Performance metrics for cell type predictions represented in Figure 3D and  
1256 Supplementary Figure 8C.

1257 **Supplementary Table 11.** Performance metrics for cell type predictions represented in Figure 4D and  
1258 Supplementary Figure 12.

1259

1260