1    **Scalable querying of human cell atlases via a foundational model reveals**

2    **commonalities across fibrosis-associated macrophages**

3

4

5    Graham Heimberg[1,2,*,✉], Tony Kuo[3,*], Daryle DePianto[4], Tobias Heigl[5], Nathaniel Diamant[2], Omar

6    Salem[3], Gabriele Scalia[2], Tommaso Biancalani[2], Shannon Turley[4,5], Jason Rock[5], Héctor Corrada Bravo[6],

7    Josh Kaminker[1,§,✉], Jason A. Vander Heiden[1,4,§,✉], Aviv Regev[7,§,✉]

8

9    **Affiliations:**

10   [1] Department of OMNI Bioinformatics, Genentech, South San Francisco, CA 94080, USA.

11   [2] Department of Machine Learning for Biology, Genentech, South San Francisco, CA 94080,

12   USA.

13   [3] Roche Informatics, F. Hoffmann-La Roche Ltd., Mississauga, Canada

14   [4] Department of Immunology Discovery, Genentech, South San Francisco, CA 94080, USA.

15   [5] Department of Regenerative Medicine, Genentech, South San Francisco, CA 94080, USA.

16   [6] Data Science and Statistical Computing, Genentech, South San Francisco, CA 94080, USA.

17   [7] Research and Early Development, Genentech, South San Francisco, CA 94080, USA.

18   [*] These authors contributed equally

19   [§] These authors contributed equally

20   ✉ email: heimberg@gene.com, regev.aviv@gene.com, vanderheiden.jason@gene.com,

21   kaminker@gene.com

## Abstract

Single-cell RNA-seq (scRNA-seq) studies have profiled over 100 million human cells across diseases, developmental stages, and perturbations to date. A singular view of this vast and growing expression landscape could help reveal novel associations between cell states and diseases, discover cell states in unexpected tissue contexts, and relate *in vivo* cells to *in vitro* models. However, these require a common, scalable representation of cell profiles from across the body, a general measure of their similarity, and an efficient way to query these data. Here, we present SCimilarity, a metric learning framework to learn and search a unified and interpretable representation that annotates cell types and instantaneously queries for a cell state across tens of millions of profiles. We demonstrate SCimilarity on a 22.7 million cell corpus assembled across 399 published scRNA-seq studies, showing accurate integration, annotation and querying. We experimentally validated SCimilarity by querying across tissues for a macrophage subset originally identified in interstitial lung disease, and showing that cells with similar profiles are found in other fibrotic diseases, tissues, and a 3D hydrogel system, which we then repurposed to yield this cell state *in vitro*. SCimilarity serves as a foundational model for single cell gene expression data and enables researchers to query for similar cellular states across the entire human body, providing a powerful tool for generating novel biological insights from the growing Human Cell Atlas.

## INTRODUCTION

Characterizing the contexts in which cells employ different expression programs is critical for deciphering their functional role in health and disease. To date, well over 100 million individual cells have been profiled using single-cell or single-nucleus RNA-seq (sc/snRNA-seq) across homeostatic, disease, and perturbed conditions[1]. Individually, even the largest multi-tissue scRNA-seq atlases[2–4] capture only a relatively small portion of cell states across human tissues; collectively, these atlases provide a vast, pan-human view of cell and disease biology that has the potential to address fundamental questions about human biology[1]. By aggregating across atlases, we may uncover biological insights and enable investigations into cell states that are common across multiple studies of the same organ and conditions (*e.g.*, similar neural progenitor populations across independent studies of brain development), different organs and conditions (*e.g.*, inflammatory fibroblasts in both ulcerative colitis and cancer); or between the human body and *in vitro* lab models (*e.g.*, regulatory T cells genetically perturbed to recapitulate *in vivo* cells from diseased tissue).

Despite this promise and the rapid growth in data, our ability to realize the potential of cross-datasets, pan-body, analyses remains limited and hampered by the need for laborious manual curation, harmonization, and dataset aggregation by expert analysts, as well as the painstaking process of selecting datasets, standardizing cell type annotations, and finding a common low-dimensional representation. As a result, most aggregation efforts have been limited in their biological scope and number of datasets, with some recent notable exceptions focused on genes rather than cell representation[5–8].

64  To leverage and query the massive scale and richness of available single-cell atlases, we need both

65  (1) a foundational model of cell states with an effective representation for single-cell profiles

66  across different cell types and conditions that can be used across many applications without

67  retraining; and (2) a measure of cell similarity that is robust to technical noise, scales to hundreds

68  of millions of cells, and accurately generalizes to datasets and cell states not observed in the

69  training. Established unsupervised methods to learn low dimensional representations of scRNA-

70  seq profiles, such as Principal Component Analysis (PCA) or autoencoders[3,9–11], faithfully preserve

71  information from the input[3,9–11] and may even eliminate technical variation for explicitly defined

72  batches. However, they do not learn general features that encode relationships between cells

73  needed to represent and query new data sets in the context of cross-study, pan-tissue biological

74  variation.

75

76  Machine learning methods, metric learning in particular, have successfully learned representations

77  for diverse entities and a measure of similarity between them, especially in image analysis. For

78  example, metric learning models for facial recognition are explicitly trained to embed images of

79  the same person closer together than images of different people, by exploiting visual features that

80  are critical to distinguish individuals[12]. Once trained, images are embedded into a low-dimensional

81  space, where distances between images represent a measure of similarity based on the learned

82  features. Users can then query with an image not in the training set to find additional similar images

83  that are nearby in the latent space and depict the same person. We reasoned that, analogously,

84  metric learning could provide a meaningful representation of and similarity metric for cell profiles.

85  By training a model using annotated scRNA-seq data, we can learn a low dimensional

86  representation that places similar cells near each other and dissimilar cells farther apart. If learned

87   from a sufficient diversity of cell profiles, such a representation would, in turn, provide a

88   foundational model of cells and would allow efficient searches for cells with similar expression

89   states (**Fig. 1a**).

90

91   Here, we introduce SCimilarity, a class of deep metric learning models that quantify similarity

92   between single-cell expression profiles (SCimilarity score) and provide a single-cell gene

93   expression foundational reference model to systematically query for comparable cell states across

94   tissues and diseases. SCimilarity uses a training set of diverse author-annotated cell profiles to

95   learn a universal representation and distance metric that facilitates efficient searches across a

96   massive reference meta-atlas for the most expression-similar cells. To train a foundational model

97   that can be broadly applied to many applications across tissues and studies, we built a

98   programmatic pipeline for massive data import and automated standardized curation and used it to

99   assemble a corpus of 22,699,774 cells from 399 datasets spanning a broad range of organs, systems

100   and conditions across the human body. After training and testing on a subset of 66 studies and

101   7.9M single-cell profiles, the learned models generalize well, representing and quantifying

102   similarities between 14.9M cells from another 347 studies excluded from training. By tuning a

103   single parameter during SCimilarity's training, we yield models optimized for either data

104   integration and visualization of millions of cells across hundreds of studies, or for fast and efficient

105   (millisecond) queries of a new cell state across tens of millions of cells. Finally, we illustrate the

106   power of SCimilarity by querying for a fibrosis-associated macrophage (FMΦ) subset previously

107   identified in interstitial lung disease (ILD), finding comparable cell populations (but with different

108   annotated names and signatures) in other ILD studies, as well as in new contexts, including

109   COVID-19, different tumors including pancreatic ductal adenocarcinoma (PDAC), and even

5

110  healthy lung (at low abundance). Surprisingly, SCimilarity recovered FMΦ-like cells among

111  PMBCs cultured and stimulated in a 3D hydrogel system *in vitro*, which we experimentally

112  validated, producing an FMΦ-like *in vitro* cell system for future functional studies of a tissue-

113  resident cell state. Overall, SCimilarity preserves expression diversity across cells in an integrated

114  foundational model of a human cell atlas, and allows a novel scaled cell search across organs,

115  systems, and conditions, as a powerful framework for generating biological insights and

116  experimentally testable hypotheses.

117

118  **RESULTS**

119  **SCimilarity: novel similarity metrics and representations for single-cell expression profiles**

120  SCimilarity is a family of models that blend unsupervised representation learning and supervised

121  metric learning, through simultaneously optimizing two objectives : (1) a supervised triplet loss

122  function, which is used to embed expression profiles from matching cell types close together,

123  effectively integrating cells of the same type across studies[13–15], and (2) an unsupervised mean

124  squared error (MSE) reconstruction loss function, which encourages the model to preserve

125  variation from the input expression profiles, capturing subtler differences in expression patterns

126  within cells of the same type, such as those related to tissue residency of immune cells (**Fig. 1b,**

127  **Methods**). The balance of these two objectives, set by a single hyperparameter, $\beta$, determines the

128  properties of the representation (**Methods**). Increasing the relative weight of the triplet loss

129  function improves dataset integration, while increasing the relative weight of the reconstruction

130  loss improves querying performance. Therefore, different loss function weightings within the same

131  model architecture can address different applications.

6

132

133     We train SCimilarity with tens of millions of cell triplets sampled from data with author-provided

134     standardized cell type annotations from the Cell Ontology[16] (**Fig. 1b**, **Methods**). Specifically, each

135     training triplet consists of similar anchor and positive cells (*i.e.*, same cell type) from different

136     studies, while anchor and negative cells are dissimilar (*i.e.*, distinct cell types; from the same or a

137     different study). However, even with standardized Cell Ontology terms, some cell type

138     comparisons are ambiguous due to arbitrary differences in annotation granularity across studies

139     (*e.g.*, it is ambiguous if cells annotated as "T cell" in one study and "CD4$^+$ T cell" in another are

140     similar or dissimilar). To address this, SCimilarity excludes cell pairings with such vertical

141     ancestor-descendant Cell Ontology relationships from training triplets, and learns only from cells

142     that are either explicitly similar or unambiguously dissimilar (**Fig. 1b, Methods**). By sampling

143     only unambiguous triplets we eliminate the need to manually flatten or harmonize every cell type

144     annotation and are able to seamlessly scale the training set across dozens of studies.

145

146     **A learned SCimilarity representation of 22.7M cells across dozens of tissues and disease**

147     **datasets collated by an automated curation and processing pipeline**

148     To test SCimilarity models, we assembled a compendium of sc/snRNA-seq datasets across human

149     biology. We focused on studies generated with one experimental platform (10x Genomics

150     Chromium droplet-based scRNA-seq) and data publicly available on the Gene Expression

151     Omnibus (GEO)[17] or CELLxGENE[18]. These data capture much of the published scRNA-seq data,

152     and were generated with similar library preparation protocols and computational pipelines[19]. There

153     were 753 human sc/snRNA-seq datasets matching our search criteria and keywords as of March

154     23$^{rd}$, 2021 (with Biopython Entrez[20], **Methods**). The number of samples and cells matching our

155    criteria has at least doubled every 6 months between December 2018 and March 2021; (**Extended**

156    **Data Fig. 1a,b**). We programmatically downloaded 13,401,599 cell profiles from 333 of the

157    identified studies with their respective GEO metadata and unnormalized gene count matrices

158    (**Methods, Extended Data Table 1**). We manually ingested another 66 well-annotated studies

159    from either the CELLxGENE portal[18] or from large studies and consortia not available through

160    GEO that passed the same dataset filtering criteria (**Methods**). Overall, we assembled a corpus of

161    399 studies comprising 22,699,774 cells from 33,815 tissue samples with 184 unique Tissue

162    Ontology terms[21], 132 Disease Ontology terms[22], and 204 Cell Ontology cell type terms[16], with

163    each Cell Ontology term appearing in at least two separate datasets (**Fig. 2a, Extended Data Fig.**

164    **1c, Extended Data Table 1**).

165

166    We trained SCimilarity models with a training set of 7,913,892 single-cell profiles from 52 studies

167    with Cell Ontology author annotations that reflected a diversity of conditions and tissues

168    (**Extended Data Fig. 1d**, **Extended Data Table 1**), sampling 50,000,000 of the most informative

169    triplets (**Methods**). We withheld 14 studies comprising 1,384,283 cells with Cell Ontology

170    annotations for testing the learned representation and metric (**Fig. 2a**). We excluded tumor, cell

171    lines, and iPSC-derived samples from the training and test sets, because cell identity of tumor cells

172    and cell lines can be ambiguous.

173

174    **Tuning of SCimilarity's reconstruction and triplet loss functions yields models optimized for**

175    **integration *vs*. cell search tasks**

176    We examined six different blends for SCimilarity's objective function, varying the relative

177    weighting of the reconstruction and triplet loss functions, and finding that the two loss function

8

178  components gave rise to different behaviors in a trained model. Briefly, we assessed the models

179  on two tasks – data integration and searching for cells similar to a query profile – using studies

180  entirely held out from training. To evaluate data integration, we quantified how coherently cells of

181  each type are clustered and how distinct each cell type cluster is from other clusters. To this end,

182  we created an ontology-aware variation of average silhouette width[23] to quantify integration

183  capabilities across datasets without harmonizing cell type annotations (**Methods**). To evaluate our

184  cell search distance metric, we compared searches with SCimilarity to gene signature scoring

185  (**Methods**). The higher the correlation between these two quantities, the more our similarity metric

186  corresponds to traditional signature-based similarity to represent a cell state of interest.

187

188  Models with higher triplet loss weighting scored higher on integration benchmarks, while models

189  with higher reconstruction loss weighting encoded distances between cells in a manner better

190  correlated with differences in representative expression signature scores (**Extended Data Fig.**

191  **2a,b**). Pure triplet loss, which is calculated at the level of cell type labels, does not reliably preserve

192  subtle cell state differences, such as tissue specificity or disease response within cells of the same

193  type. Mean squared error reconstruction loss complements this by preserving more subtle gene

194  expression patterns, while the triplet loss ensures that cells of the same type are embedded closely

195  together. Based on the biological question, a user can tune this balance to yield the highest utility.

196  We thus pursued two SCimilarity models: an integration model, optimized for the task of learning

197  a low dimensional representation that groups cells by type rather than by study; and a cell search

198  model that is optimized for the task of retrieving cells with an expression state similar to that of a

199  query cell across hundreds or thousands of scRNA-seq datasets (**Extended Data Fig. 2a**,

200  **Methods**).

**SCimilarity's latent space representation filters outlier cells and integrates test datasets**

**without batch correction**

We next benchmarked if SCimilarity's latent space representation from the integration model generalizes well to cells from entire datasets held out of training compared to other methods. In a low-dimensional embedding, unannotated cell profiles from nine lung studies (7 training set, 2 test set) visually intermix well when embedded into SCimilarity's learned 128-dimensional space (**Extended Data Fig. 2a**). SCimilarity's data integration model scored higher than Harmony, scVI, and scArches on integration tasks by the ontology-adjusted ASW measure of cluster coherence, but scored lower for normalized mutual information (NMI) and adjusted rand index (ARI), which measure the extent of study mixing within each cluster (**Extended Data Fig. 2b**). Thus, without directly training on the full data set or performing additional batch correction, the integration model clusters cells by type rather than study at a level that is competitive with existing methods trained directly on the data. This demonstrates that the triplet loss learns features that capture meaningful biology, while reducing technical sources of noise and avoiding overfitting to the training set.

SCimilarity quantifies a confidence level for each cell's representation, providing both outlier detection and an assessment of the representation's relevance in the context of new data. When computing the representation of a new cell, the further outside the scope of model training it is, the harder it is for the model to accurately represent it. Using SCimilarity's score to quantify how distant a query cell is from the training data distribution provides a heuristic about the quality and scope of the representation – a cell scoring as highly similar to cells seen during training can be confidently represented by the model. Overall, 79.5% of *in vivo* holdout cells had high

224    representation confidence. Tissue samples with particularly low representation confidence, such

225    as stomach (n = 0 training studies), fetal gut (n=1), and bladder (n=0) were either absent or poorly

226    represented in training (**Methods, Extended Data Fig. 2c**), suggesting that more labeled training

227    datasets from those tissues could improve the model's representation. Similarly, 43.8% of *in vitro*

228    cell profiles were considered low confidence due to poor matching to the training set (which

229    excluded *in vitro* samples).

230

231    We combined SCimilarity's ability to generalize to new datasets and its confidence-based filtering

232    to systematically generate meta-atlases for 21 different human tissues without labor-intensive

233    dataset harmonization and no additional training (**Fig. 2b**). If datasets have already been embedded

234    using SCimilarity, this task only requires concatenation of cells of interest and standard

235    visualization workflows.

236

237    **SCimilarity assigns an unannotated query cell to a cell type by finding similar cells in a**

238    **labeled reference**

239    We next used SCimilarity to find the cells in the annotated reference that are most similar to an

240    unannotated query cell profile, and then annotate the query cell accordingly (**Fig. 3a, Methods**).

241    This approach is distinct from established annotation methods in that it (1) relies on a large, pan-

242    human annotated cell repository, (2) employs a measure of expression similarity, and (3) classifies

243    at the single cell rather than cluster level, providing greater transparency into the classification

244    itself. Thus, users can see which individual cells, studies, and tissues are driving the classification

245    decision. Moreover, since each cell is annotated independently, no clustering or associated

11

246  parameter selection, such as the number and resolution of clusters, are required. A user can choose

247  to annotate a cell's profile by comparing it either to a desired subset of cell types (*e.g.*, for a tissue-

248  specific query) or to the entire annotated cell reference. Because SCimilarity is built using metric

249  learning, finding the most similar cells is the same as retrieving the query cell's nearest neighbors.

250  This operation is extremely efficient with the hnswlib algorithm[24], where searching a precomputed

251  approximate nearest neighbor index of all the annotated reference cells in SCimilarity's latent

252  space takes just 20 milliseconds (**Methods**). Low SCimilarity scores to reference cells flag an

253  outlier query cell, which may be either a cell type that is not within the reference or a query cell of

254  low quality.

255

256  SCimilarity quickly and accurately assigned cell types for entire datasets held out from training,

257  as well as for the rest of the 22.7M cell corpus. When limiting potential cell types to author-selected

258  labels, 94.5% of SCimilarity's predicted labels from healthy kidney samples[25] match the author-

259  provided cell type annotations (**Fig. 3b-d**, **Methods**). In some cases, where SCimilarity's

260  predictions did not match author-provided annotations, SCimilarity's predictions were more

261  accurate or granular. For example, 94% of the cells that the authors[25] annotated as $CD4^+$ T cells

262  but SCimilarity annotated as $CD8^+$ T cells express *CD8A* or *CD8B* (and none express *CD4*),

263  supporting SCimilarity's annotation (**Fig. 3e-h**). Separately, when allowing cells to be annotated

264  as any cell type in the repository, 6.3% of the author-annotated $CD4^+$ T cells were reannotated by

265  SCimilarity as regulatory T cells ($T_{regs}$) (**Extended Data Fig. 3a**), most of which (85.2%)

266  expressed at least one $T_{regs}$ marker (FOXP3, IL2RA, or IKZF2, **Extended Data Fig. 3b-d**).

267  Similarly, 1.8% of author-annotated mesenchymal stem cells (**Fig. 3b**) were reassigned by

268  SCimilarity as myofibroblasts (**Extended Data Fig. 3a**) and 93% of those express the

12

269    myofibroblast-associated gene *ACTA2* (**Extended Data Fig. 3e**). Cell type prediction was rapid,

270    taking 3-5 seconds to embed and annotate 10,000 cells from a dataset. Overall, across all 14 test

271    datasets spanning 78 Cell Ontology terms, 71% of the cell populations had high agreement (>85%

272    of the cell population) between author and SCimilarity annotations (**Fig. 3i**). SCimilarity

273    performed poorly on one dataset (Cano-Gomez et al.[26]), due to fine granularity and redundancy of

274    author labels (*e.g.*, CD4$^+$ αβ T cells, helper T cell, memory T cell, naive T cell, and regulatory T

275    cell).

276

277    We used SCimilarity's cell type assignment to rapidly annotate all 22.7M cell profiles in one

278    common model, newly-annotating 13,401,599 profiles and reannotating 9,298,175 author-

279    annotated profiles (**Methods**) to a single set spanning 74 cell type labels (21 coarser lineages) from

280    25 simplified tissue categories (**Fig. 4a,b, Extended Data Fig. 3f**). A consistent annotation across

281    datasets facilitates cross-study and cross-tissue analyses of one cell type or lineage, as SCimilarity

282    can extract cells from hundreds of studies, aggregating vast biological diversity across one cell

283    type. For example, we readily aggregated 1,172,325 fibroblasts and myofibroblasts (**Extended**

284    **Data Fig. 3g**) and 2,507,879 monocytes and macrophages (**Extended Data Fig. 3h**) from

285    hundreds of studies profiling different primary tissue samples.

286

287    **SCimilarity's representations comprise of interpretable biological features**

288    To interpret SCimilarity's annotations, we quantified the importance of each gene for cell type

289    annotations assigned by the foundational query model using Integrated Gradients, a method that

290    identifies the impact on model predictions of small disturbances to the input expression profiles

291    (**Methods**). For example, the top gene attributions that distinguish lung alveolar type 2 (AT2) cells

13

292   are surfactant genes *SFTPA2*, *SFTPA1*, *SFTPB*, and *SFTPC*, consistent with known AT2 cell

293   function[27]. SCimilarity learned these without prior knowledge of cell type specific genes,

294   signatures, or highly variable genes. Overall, SCimilarity's top importance genes agreed well with

295   differentially expressed marker genes for 17 different matched types[3] with the exception of rare

296   neuroendocrine cells (average AUC=0.84, **Extended Data Table 2, Extended Data Fig. 3i**).

297   Thus, SCimilarity's representation captured known and validated biological markers within its

298   features.

299

300   **Cell search identifies fibrosis-associated macrophages across tissues and diseases**

301   With a single representation and common definition of cell types, we hypothesized that

302   SCimilarity could help elucidate the role of tissue-resident immune cells. As a case study, we

303   focused on macrophages, given their remarkable plasticity in cell states and their important

304   specialized roles in tissue repair, regeneration, and fibrosis[28,29]. Recent scRNA-seq studies in

305   fibrotic diseases, including lung fibrosis, cancer, obesity, and COVID-19 have reported seemingly-

306   similar $SPP1^+$ fibrosis-associated macrophage (FMΦ) populations[30–38]. However, because each

307   study identified them independently, using different nomenclatures and marker gene signatures to

308   define subsets, it is unclear how similar these cell states are. Moreover, it is unknown how broadly

309   associated such cell states are with other diseases, especially those with prominent fibrosis. We

310   reasoned that SCimilarity's cell search should allow us to query our corpus with an FMΦ cell

311   profile from one study to identify similar cells across other tissues and conditions, thereby

312   clarifying the cell identity of similarly-described cells and the conditions in which FMΦ arise (**Fig.**

313   **5a**).

314

14

315    We queried our model with the FMΦ cell profile, searching for similar cells across 2,578,221 cells

316    annotated by SCimilarity as monocytes or macrophages in the 22.7M cell corpus (**Fig. 5a**).

317    SCimilarity queries can use either an individual cell profile or a centroid of multiple cell profiles.

318    Here, we input the centroid profile of a macrophage cell subset from Adams *et al.*[30] that we defined

319    using a gene signature consisting of the extracellular matrix remodeling and fibrosis-associated

320    genes *SPP1*, *TREM2*, *GPNMB*, *MMP9*, *CHIT1*, and *CHI3L1* (**Methods**). In two seconds,

321    SCimilarity exhaustively computed the pairwise similarity of our query profile to each of the 2.6M

322    *in vivo* profiles of the cells it annotated as monocytes or macrophages in our corpus (**Fig. 5b** and

323    **Extended Data Fig. 4a**). Alternatively, simply identifying the 10,000 cells with the highest

324    SCimilarity score takes 0.05 seconds (**Methods**). By comparison, a more conventional approach

325    that scores each cell in the corpus with a literature-defined FMΦ gene signature took 2 hours and

326    46 minutes (**Extended Data Fig. 4b**). The gene signature and SCimilarity scores are broadly

327    correlated ($r = 0.50$, $p < 10^{-300}$, **Extended Data Fig. 4a-c**), showing that the granular cell state,

328    not just the cell type, is well-represented in SCimilarity query score and embedding.

329

330    The SCimilarity search showed that FMΦs are common in ILD lung samples in our compendium,

331    as well as present in some cancers, including uveal melanoma, pancreatic ductal adenocarcinoma

332    (PDAC), and colon cancer (**Fig. 5c-e, Extended Data Table 3**). Of the top 1% of monocytes and

333    macrophages most similar to our query, 99.1% were from lung tissue and 87.2% from ILD and

334    COVID-19 lung samples. The prevalence of FMΦ-like cells in the lung varied by disease: the

335    proportion of monocytes and macrophages that were FMΦ-like was 20% and 4% in two systemic

336    sclerosis (SSc) studies, 6.1% on average (SE = 1.4%) across 13 ILD studies (excluding SSc), 1.2%

337    on average across seven COVID-19 lung studies (SE = 0.5%, 0% in non-lung COVID-19 data)

338    and 0.4% in 19 studies annotated as "healthy", "normal" or with no disease annotation (SE =

339    0.2%). While abundant in SSc lung, FMΦ-like cells were much rarer (0.14% of myeloid cells) in

340    SSc skin[39]. There were some FMΦ-like cells in other fibrotic diseases and tissues, such as one

341    primary pancreatic ductal adenocarcinoma (PDAC) tumor[40] (0.85% of 1,171 myeloid cells) and

342    one liver metastasis[41] of PDAC (0.5% of 1,199 cells). Thus, while our query FMΦ profile was

343    derived from IPF samples, it uncovered FMΦ-like cells in many contexts, including SSc-ILD,

344    COVID-19 lung and PDAC. These results confirm previous observations of FMΦs in lung

345    injury[38,42] and suggest a role for FMΦ-like cells across other organs and diseases.

346    **Integrated gradients analysis reveals commonalities between SCimilarity score and**

347    **established gene signatures**

348    Because FMΦ-like cells are detected by SCimilarity across many ILD studies, we hypothesized

349    that the cells captured by different marker genes and nomenclature in different studies refer to the

350    same biological cell state. To test this, we applied integrated gradients to quantify each gene's

351    importance when SCimilarity distinguishes FMΦs from randomly sampled monocytes and

352    macrophages (**Methods**). The genes identified as important for distinguishing FMΦs are enriched

353    in key fibrotic processes, including extracellular matrix remodeling (*MMP7*, *MMP9*, *FN1*, *SDC2*,

354    *SPARC*, *SPP1*), lipid metabolism and lipoprotein clearance (*APOC1*, *APOE*, *LPL*, *LIPA*), and

355    damage-associated molecular pattern recognition (*MARCO*, *MSR1*) (**Fig. 5f**, **Extended Data Fig.**

356    **4d,e**, **Extended Data Table 4**). While SCimilarity found many FMΦ marker genes that were

357    already discussed in the literature, such as TREM2 (**Extended Data Fig. 4f**), it also identified

358    novel genes elevated in FMΦs such as HLA-DQA1 and RGS1 (**Extended Data Fig. 4g,h**).

359

360   The genes with the highest importance scores in the SCimilarity embedding of FMΦs significantly

361   overlap ($p<6.7 \times 10^{-13}$) with published gene signatures describing similar macrophage populations

362   or with genes whose differential expression defined each study's macrophage population of

363   interest (**Extended Data Table 5**). While cell signatures from IPF lung had a high signature match

364   ($AUC \geq 0.95$), the negative control signatures of M1 and M2 macrophages[43] had lower ones at the

365   bottom three ($AUC = 0.85$ ($2.65 \times 10^{-2}$) and $AUC=0.92$ ($p<4.92 \times 10^{-6}$), respectively; **Fig. 5f**,

366   **Extended Data Fig. 4d**).

367

368   **FMΦ-like cells identified among *ex vivo* stimulated peripheral blood mononuclear cells**

369   **(PBMCs) help establish a novel human cell model**

370   Research to understand the role of a novel cell state or subtype in disease, such as FMΦs, benefits

371   greatly from the ability to model, perturb, and study the cells *in vitro*. However, there is currently

372   no systematic way to identify *in vitro* culture conditions that generate cells that match cells

373   identified *in vivo*. To accelerate development of an *in vitro* FMΦ system, we used SCimilarity to

374   search for FMΦ-like cells across *in vitro* stimulated samples with the goal of identifying previously

375   employed experimental conditions that might resemble the tissue cell state. We filtered our full

376   reference cell collection for *in vitro* and *ex vivo* studies containing at least 50 monocytes or

377   macrophages, resulting in 41,926 monocytes and macrophages across 40 samples from 17 such

378   studies. These span diverse and complex conditions, such as lung organoids infected with SARS-

379   CoV-2[44], *ex vivo* treated acute myeloid leukemia samples[45], or PBMCs stimulated with morphine

380   and lipopolysaccharide[46].

381

17

382    The cells most similar to our query FMΦ expression profile were monocytes grown as part of a

383    heterogenous PBMC culture for 5 days in a 3D hydrogel culture system that was designed for

384    expansion of hematopoietic stem cells (HSCs) from PBMCs[47] (**Fig. 6a, Extended Data Table 6**).

385    This study is unrelated to lung biology and its authors did not report any results for myeloid cells.

386    Nevertheless, while no FMΦ-like cells were present among myeloid cells on day 0, 15% of cells

387    grown for five or more days in this system were highly similar to FMΦs (SCimilarity score >25)

388    and expressed *TREM2*, *GPNMB*, *CCL18* and *MMP9* (**Fig. 6b-e**). This was a surprising result,

389    because of the seeming irrelevance of the study to fibrosis or macrophage biology and the rarity

390    of FMΦ-like cells in PBMC samples *in vivo*.

391

392    To validate SCimilarity's prediction of an FMΦ-like cell culture condition, we used a similar

393    protocol to replicate the 3D hydrogel system[47], followed by scRNA-seq to assess the yield of

394    FMΦ-like cells (**Fig. 6b,c,f**). While relative cellular abundances differed between the original day

395    5 data (Xu et al, 2022) and our day 8 replication of the same conditions (**Methods**), 10.1% of all

396    cells in the Day 8 experiment were predicted as HSCs by SCimilarity (**Fig. 6g**). Moreover, 41.5%

397    of the myeloid cells in day 8 validation experiments from three donors were predicted as FMΦ-

398    like macrophages (**Fig. 6b,f**, 37.1%, 42.5%, and 44.9%; SCimilarity score > 25). Furthermore,

399    FMΦ hallmark genes, such as *CCL18*, *GPNMB*, *SPP1*, and *TREM2*, were enriched in the myeloid

400    compartment of our replicate experiment compared to day 0 conditions (**Fig. 6c**). This experiment

401    validates that an FMΦ-like population can be generated from PBMCs in culture conditions. Taken

402    together, these results demonstrate SCimilarity's ability to interrogate publicly available data at

403    scale, query a reference of *in vivo* and *in vitro* data for biologically similar conditions, and help

404    identify experimental conditions to reproduce those results in laboratory settings.

405

**DISCUSSION**

406

407     To date, more than a hundred million human cells have been profiled across tissues in health and

408     disease, and such data continue to grow exponentially. This growing human cell atlas should be

409     the starting point for researchers aiming to readily search, query and compare cell states of interest

410     across different protocols, treatments, tissues, and diseases.

411

412     SCimilarity systematically annotates and repurposes tens of millions of expression profiles from

413     hundreds of studies, to create an integrated, searchable and queryable foundational model of pan-

414     human cellular diversity. SCimilarity is comprised of three key features: (1) a 22.7M cell human

415     scRNA-seq data repository (at present), (2) a foundational model for single cell gene expression

416     with a generalizable embedding and similarity metric (which could readily be retrained for larger

417     datasets), and (3) methods to efficiently query across this entire pan-body human cell atlas.

418     Together, these provide new context, capabilities, and workflows for extracting insights from new

419     and existing scRNA-seq datasets in the human cell atlas and other atlases. SCimilarity's

420     framework architecture can easily accommodate quick updates as data continue to grow.

421

422     Because SCimilarity can generalize to cells and datasets not seen in the training, cell profiles can

423     be added as entirely new studies or removed by applying new cell filters without recomputing the

424     low dimensional representations. This flexibility allows us to change the analysis' scope at any

425     point without redoing work, enabling modularized workflows for scRNA-seq analysis.

426     Downstream tasks, such as cell type annotation, cell queries, and gene signature derivation all are

427     simplified using SCimilarity's generalized low dimensional representation and can be applied to

428    cells not seen during training without informing the model about the importance or variability of

429    specific genes during training. Outlier detection helps both filter out technical errors and highlight

430    potentially novel cell subsets. Although generalized models that do not require recomputing low

431    dimensional representations would alleviate time and expertise barriers that currently impede

432    researchers, to the best of our knowledge, generalization has rarely been optimized in single-cell

433    expression analysis.

434

435    There is no single objective measure of similarity, or dissimilarity, between cellular profiles.

436    Curated gene signatures are useful when a small number of explanatory genes are sufficient to

437    define a cell state. SCimilarity uses the full expression profile of a cell as its query, defined by

438    either a single representative cell or the centroid of a set of cell profiles. Thus, SCimilarity's cell-

439    based search bypasses the manual curation requirements and biases inherent in defining a gene

440    signature. In cases where such a gene signature is desired, SCimilarity can compute a robust

441    signature for a cell state across studies.

442

443    Exploration of transcriptionally-similar populations across a vast atlas of human scRNA-seq data

444    provides critical context to a cell population of interest. First, observing a query population across

445    many similar studies shows that the original observation was reproducible, a key for subsequent

446    scientific research[48]. Second, SCimilarity queries can connect results from independent studies.

447    While one study may find a cell population in a disease, another may show similar cells with

448    functional characterization, allowing us to formulate a new hypothesis on the functional properties

449    of disease-associated cells.

450

451     This is illustrated by how SCimilarity allowed us to search for and identify FMΦ-like cells across

452     tissues and disease states, construct a cross-study set of explanatory marker genes, and uncover a

453     cell culture system that elicits a similar FMΦ-like state *in vitro*. Modeling FMΦs from readily

454     available PBMCs is exceptionally valuable, because isolation of cells from human lung explants

455     is prohibitive for many functional assays. Surprisingly, in addition to fibrotic lung, FMΦs were

456     present in multiple tumor types, particularly PDAC, a heavily fibrotic cancer, where macrophages

457     play an important role in mediating the associated fibrosis and have been linked to tumor

458     progression[49]. The identification of a common FMΦ state across fibrosis, cancer, and infection

459     suggests a broader role for these cells in the damage response and tissue remodeling processes

460     across diseases. Moreover, SCimilarity's search identified FMΦ-like cells in an *in vitro* study– an

461     observation that could not have been gleaned by reviewing the paper or based on the description

462     of the culture system – but that we validated in the lab. The variations we observed between the

463     original and replicate *in vitro* experiment may be attributed to differences in culture duration, cell

464     extraction from the hydrogel, lymphocyte proportions, or other batch effects. Furthermore, these

465     results invite new hypotheses, such as whether the 3D hydrogel provides key ECM-like

466     environmental cues that promote an FMΦ-like state and induction of remodeling genes, such as

467     *MMP9* and *SPP1*, and which factors can be added to drive an even stronger FMΦ phenotype. Thus,

468     SCimilarity provides a powerful framework to iteratively generate and validate such experimental

469     hypotheses.

470

471     SCimilarity is not appropriate for all applications and will need further improvements to continue

472     to scale with exponential data growth and to more comprehensively span human biology as the

473     Human Cell Atlas continues to grow. Training SCimilarity requires Cell Ontology labels.

474   Fortunately, scRNA-seq data sharing practices are increasingly relying on using the Cell Ontology

475   for standardization. However, the Cell Ontology itself is a large, yet incomplete, effort. Cell states

476   are only considered in training if they are recognized in the Cell Ontology, and the number of these

477   states is growing rapidly. Furthermore, while we trained SCimilarity on vast amounts of data,

478   cancer cells and cell lines were deliberately withheld from training due to lack of clear cell type

479   identity and therefore may not be well represented. In addition, in our experience, we see poor

480   performance on fetal samples, likely due to most of the training data being sourced from adult

481   tissues.

482

483   The current data integration and cell search models provide generalizable representations of 22.7M

484   single-cell profiles across the human body, and include a Python API for querying cell profiles of

485   interest. Future improvements to SCimilarity could include pre-training on the massive amounts

486   of unlabeled data, effectively exposing the model to more cell states and more technical variability

487   during training. With effective representations we can more easily combine embeddings to include

488   other species or data modalities. We believe that SCimilarity brings a new framework to single-

489   cell genomics, enabling re-use of rich public data resources through instantaneous queries and

490   demonstrates how this can be used to provide novel biological insights.

491

492

493   **Contributions:**

494   GH conceived of the method with input from AR, JVH, HCB, and JK. GH and TK performed

495   data ingest and model implementation with input from JVH, ND, GS, TB, JK and AR. Python

496   API was developed by TK with help from JVH, OS, and GH. Interpretability was developed by

22

497  ND and GS with input from HCB and TB.  JVH conceived of biological application of method

498  with input from GH, ST, JR and DD. DD and TH performed experimental validation with

499  guidance from JR and ST. GH wrote the manuscript with input from JVH, JK, AR and HCB.

500

501  **Acknowledgements:**

502  We thank Anupriya Tripathi for coming up with the name "SCimilarity" and Jenna Collier,

503  Gokcen Eraslan, John Marioni, and Jake Freimer for their suggestions that strengthened the

504  manuscript.

505

506  **Competing interests:**

507  All authors are employees of Genentech or Roche. A.R. is a co-founder and equity holder of

508  Celsius Therapeutics, an equity holder in Immunitas, and until July 31, 2020 was an S.A.B.

509  member of Thermo Fisher Scientific, Syros Pharmaceuticals, Neogene Therapeutics and

510  Asimov.

511

512

513  **References**

514

515  1.  Rood, J.E., Maartens, A., Hupalowska, A., Teichmann, S.A., and Regev, A. (2022). Impact
516      of the Human Cell Atlas on medicine. Nat. Med. *28*, 2486–2496.

517  2.  Domínguez Conde, C., Xu, C., Jarvis, L.B., Rainbow, D.B., Wells, S.B., Gomes, T., Howlett,
518      S.K., Suchanek, O., Polanski, K., King, H.W., et al. (2022). Cross-tissue immune cell analysis
519      reveals tissue-specific features in humans. Science *376*, eabl5197.

520  3.  Eraslan, G., Drokhlyansky, E., Anand, S., Fiskin, E., Subramanian, A., Slyper, M., Wang, J.,
521      Van Wittenberghe, N., Rouhana, J.M., Waldman, J., et al. (2022). Single-nucleus cross-tissue
522      molecular reference maps toward understanding disease gene function. Science *376*,
523      eabl4290.

4.  Tabula Sapiens Consortium*, Jones, R.C., Karkanias, J., Krasnow, M.A., Pisco, A.O., Quake, S.R., Salzman, J., Yosef, N., Bulthaup, B., Brown, P., et al. (2022). The Tabula Sapiens: A multiple-organ, single-cell transcriptomic atlas of humans. Science *376*, eabl4896.

5.  Rosen, Y., Brbić, M., Roohani, Y., Swanson, K., Li, Z., and Leskovec, J. (2023). Towards Universal Cell Embeddings: Integrating Single-cell RNA-seq Datasets across Species with SATURN. bioRxiv. 10.1101/2023.02.03.526939.

6.  Cui, H., Wang, C., Maan, H., Pang, K., Luo, F., and Wang, B. (2023). scGPT: Towards Building a Foundation Model for Single-Cell Multi-omics Using Generative AI. bioRxiv, 2023.04.30.538439. 10.1101/2023.04.30.538439.

7.  Theodoris, C.V., Xiao, L., Chopra, A., Chaffin, M.D., Al Sayed, Z.R., Hill, M.C., Mantineo, H., Brydon, E.M., Zeng, Z., Liu, X.S., et al. (2023). Transfer learning enables predictions in network biology. Nature. 10.1038/s41586-023-06139-9.

8.  Shen, H., Shen, X., Hu, J., Liu, J., Zhang, C., Wu, D., Feng, M., Yang, M., Li, Y., Yang, Y., et al. (2022). Generative pretraining from large-scale transcriptomes: Implications for single-cell deciphering and clinical translation. bioRxiv, 2022.01.31.478596. 10.1101/2022.01.31.478596.

9.  Heimberg, G., Bhatnagar, R., El-Samad, H., and Thomson, M. (2016). Low Dimensionality in Gene Expression Data Enables the Accurate Extraction of Transcriptional Programs from Shallow Sequencing. Cell Systems *2*, 239–250.

10. Lopez, R., Regier, J., Cole, M.B., Jordan, M.I., and Yosef, N. (2018). Deep generative modeling for single-cell transcriptomics. Nat. Methods *15*, 1053–1058.

11. Lotfollahi, M., Naghipourfar, M., Luecken, M.D., Khajavi, M., Büttner, M., Wagenstetter, M., Avsec, Ž., Gayoso, A., Yosef, N., Interlandi, M., et al. (2022). Mapping single-cell data to reference atlases by transfer learning. Nat. Biotechnol. *40*, 121–130.

12. Schroff, F., Kalenichenko, D., and Philbin, J. (6/2015). FaceNet: A Unified Embedding for Face Recognition and Clustering. 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 815–823.

13. Simon, L., Wang, Y.-Y., and Zhao, Z. (2021). Integration of millions of transcriptomes using batch-aware triplet neural networks. Nature Machine Intelligence *3*, 1–11.

14. Yang, M., Yang, Y., Xie, C., Ni, M., Liu, J., Yang, H., Mu, F., and Wang, J. (2022). Contrastive learning enables rapid mapping to multimodal single-cell atlas of multimillion scale. Nature Machine Intelligence *4*, 696–709.

15. Yu, X., Xu, X., Zhang, J., and Li, X. (2023). Batch alignment of single-cell transcriptomics data using deep metric learning. Nat. Commun. *14*, 960.

16. Diehl, A.D., Meehan, T.F., Bradford, Y.M., Brush, M.H., Dahdul, W.M., Dougall, D.S., He, Y., Osumi-Sutherland, D., Ruttenberg, A., Sarntivijai, S., et al. (2016). The Cell Ontology

560  2016: enhanced content, modularization, and ontology interoperability. J. Biomed. Semantics
561  *7*, 44.

562  17.  Edgar, R., Domrachev, M., and Lash, A.E. (2002). Gene Expression Omnibus: NCBI gene
563      expression and hybridization array data repository. Nucleic Acids Res. *30*, 207–210.

564  18.  Chan Zuckerberg CELLxGENE Discover (2022). Cellxgene Data Portal.

565  19.  Zheng, G.X.Y., Terry, J.M., Belgrader, P., Ryvkin, P., Bent, Z.W., Wilson, R., Ziraldo, S.B.,
566      Wheeler, T.D., McDermott, G.P., Zhu, J., et al. (2017). Massively parallel digital
567      transcriptional profiling of single cells. Nat. Commun. *8*, 14049.

568  20.  Cock, P.J.A., Antao, T., Chang, J.T., Chapman, B.A., Cox, C.J., Dalke, A., Friedberg, I.,
569      Hamelryck, T., and Kauff, F. (2009). Biopython: freely available Python tools for
570      computational molecular biology and bioinformatics. Bioinformatics.

571  21.  Gremse, M., Chang, A., Schomburg, I., Grote, A., Scheer, M., Ebeling, C., and Schomburg,
572      D. (2011). The BRENDA Tissue Ontology (BTO): the first all-integrating ontology of all
573      organisms for enzyme sources. Nucleic Acids Res. *39*, D507-13.

574  22.  Schriml, L.M., Mitraka, E., Munro, J., Tauber, B., Schor, M., Nickle, L., Felix, V., Jeng, L.,
575      Bearer, C., Lichenstein, R., et al. (2019). Human Disease Ontology 2018 update:
576      classification, content and workflow expansion. Nucleic Acids Res. *47*, D955–D962.

577  23.  Luecken, M.D., Büttner, M., Chaichoompu, K., Danese, A., Interlandi, M., Mueller, M.F.,
578      Strobl, D.C., Zappia, L., Dugas, M., Colomé-Tatché, M., et al. (2022). Benchmarking atlas-
579      level data integration in single-cell genomics. Nat. Methods *19*, 41–50.

580  24.  Malkov, Y.A., and Yashunin, D.A. (2020). Efficient and Robust Approximate Nearest
581      Neighbor Search Using Hierarchical Navigable Small World Graphs. IEEE Trans. Pattern
582      Anal. Mach. Intell. *42*, 824–836.

583  25.  Young, M.D., Mitchell, T.J., Custers, L., Margaritis, T., Morales-Rodriguez, F., Kwakwa, K.,
584      Khabirova, E., Kildisiute, G., Oliver, T.R.W., de Krijger, R.R., et al. (2021). Single cell
585      derived mRNA signals across human kidney tumors. Nat. Commun. *12*, 3896.

586  26.  Cano-Gamez, E., Soskic, B., Roumeliotis, T.I., So, E., Smyth, D.J., Baldrighi, M., Willé, D.,
587      Nakic, N., Esparza-Gordillo, J., Larminie, C.G.C., et al. (2020). Single-cell transcriptomics
588      identifies an effectorness gradient shaping the response of CD4+ T cells to cytokines. Nat.
589      Commun. *11*, 1801.

590  27.  Beers, M.F., and Moodley, Y. (2017). When Is an Alveolar Type 2 Cell an Alveolar Type 2
591      Cell? A Conundrum for Lung Stem Cell Biology and Regenerative Medicine. Am. J. Respir.
592      Cell Mol. Biol. *57*, 18–27.

593  28.  Wynn, T.A., and Vannella, K.M. (2016). Macrophages in Tissue Repair, Regeneration, and
594      Fibrosis. Immunity *44*, 450–462.

595    29.  Lis-López, L., Bauset, C., Seco-Cervera, M., and Cosín-Roger, J. (2021). Is the Macrophage
596          Phenotype Determinant for Fibrosis Development? Biomedicines *9*.
597          10.3390/biomedicines9121747.

598    30.  Adams, T.S., Schupp, J.C., Poli, S., Ayaub, E.A., Neumark, N., Ahangari, F., Chu, S.G.,
599          Raby, B.A., DeIuliis, G., Januszyk, M., et al. (2020). Single-cell RNA-seq reveals ectopic
600          and aberrant lung-resident cell populations in idiopathic pulmonary fibrosis. Science
601          Advances *6*, eaba1983.

602    31.  Ayaub, E.A., Poli, S., Ng, J., Adams, T., Schupp, J., Quesada-Arias, L., Poli, F., Cosme, C.,
603          Robertson, M., Martinez-Manzano, J., et al. (2021). Single Cell RNA-seq and Mass
604          Cytometry Reveals a Novel and a Targetable Population of Macrophages in Idiopathic
605          Pulmonary Fibrosis. 10.1101/2021.01.04.425268.

606    32.  Jaitin, D.A., Adlung, L., Thaiss, C.A., Weiner, A., Li, B., Descamps, H., Lundgren, P.,
607          Bleriot, C., Liu, Z., Deczkowska, A., et al. (2019). Lipid-Associated Macrophages Control
608          Metabolic Homeostasis in a Trem2-Dependent Manner. Cell *178*, 686-698.e14.

609    33.  Morse, C., Tabib, T., Sembrat, J., Buschur, K.L., Bittar, H.T., Valenzi, E., Jiang, Y., Kass,
610          D.J., Gibson, K., Chen, W., et al. (2019). Proliferating SPP1/MERTK-expressing
611          macrophages in idiopathic pulmonary fibrosis. Eur. Respir. J. *54*. 10.1183/13993003.02441-
612          2018.

613    34.  Mulder, K., Patel, A.A., Kong, W.T., Piot, C., Halitzki, E., Dunsmore, G., Khalilnezhad, S.,
614          Irac, S.E., Dubuisson, A., Chevrier, M., et al. (2021). Cross-tissue single-cell landscape of
615          human monocytes and macrophages in health and disease. Immunity *54*, 1883-1900.e5.

616    35.  Ramachandran, P., Dobie, R., Wilson-Kanamori, J.R., Dora, E.F., Henderson, B.E.P., Luu,
617          N.T., Portman, J.R., Matchett, K.P., Brice, M., Marwick, J.A., et al. (2019). Resolving the
618          fibrotic niche of human liver cirrhosis at single-cell level. Nature *575*, 512–518.

619    36.  Reyfman, P.A., Walter, J.M., Joshi, N., Anekalla, K.R., McQuattie-Pimentel, A.C., Chiu, S.,
620          Fernandez, R., Akbarpour, M., Chen, C.-I., Ren, Z., et al. (2019). Single-Cell Transcriptomic
621          Analysis of Human Lung Provides Insights into the Pathobiology of Pulmonary Fibrosis. Am.
622          J. Respir. Crit. Care Med. *199*, 1517–1536.

623    37.  Wendisch, D., Dietrich, O., Mari, T., von Stillfried, S., Ibarra, I.L., Mittermaier, M., Mache,
624          C., Chua, R.L., Knoll, R., Timm, S., et al. (2021). SARS-CoV-2 infection triggers profibrotic
625          macrophage responses and lung fibrosis. Cell *184*, 6243-6261.e27.

626    38.  Gao, X., Jia, G., Guttman, A., DePianto, D.J., Morshead, K.B., Sun, K.-H., Ramamoorthi, N.,
627          Vander Heiden, J.A., Modrusan, Z., Wolters, P.J., et al. (2020). Osteopontin Links Myeloid
628          Activation and Disease Progression in Systemic Sclerosis. Cell Reports Medicine *1*, 100140.

629    39.  Mirizio, E., Tabib, T., Wang, X., Chen, W., Liu, C., Lafyatis, R., Jacobe, H., and Torok, K.S.
630          (2020). Single-cell transcriptome conservation in a comparative analysis of fresh and
631          cryopreserved human skin tissue: pilot in localized scleroderma. Arthritis Res. Ther. *22*, 263.

26

40. Lin, W., Noel, P., Borazanci, E.H., Lee, J., Amini, A., Han, I.W., Heo, J.S., Jameson, G.S., Fraser, C., Steinbach, M., et al. (2020). Single-cell transcriptome analysis of tumor and stromal compartments of pancreatic ductal adenocarcinoma primary tumors and metastatic lesions. Genome Med. *12*, 80.

41. Kemp, S.B., Steele, N.G., Carpenter, E.S., Donahue, K.L., Bushnell, G.G., Morris, A.H., The, S., Orbach, S.M., Sirihorachai, V.R., Nwosu, Z.C., et al. (2021). Pancreatic cancer is marked by complement-high blood monocytes and tumor-associated macrophages. Life Science Alliance *4*, e202000935.

42. Bhattacharya, M. (2022). Insights from Transcriptomics: CD163+ Profibrotic Lung Macrophages in COVID-19. Am. J. Respir. Cell Mol. Biol. *67*, 520–527.

43. Martinez, F.O., Gordon, S., Locati, M., and Mantovani, A. (2006). Transcriptional Profiling of the Human Monocyte-to-Macrophage Differentiation and Polarization: New Molecules and Patterns of Gene Expression1. The Journal of Immunology *177*, 7303–7311.

44. Salahudeen, A.A., Choi, S.S., Rustagi, A., Zhu, J., van Unen, V., de la O, S.M., Flynn, R.A., Margalef-Català, M., Santos, A.J.M., Ju, J., et al. (2020). Progenitor identification and SARS-CoV-2 infection in human distal lung organoids. Nature *588*, 670–675.

45. Duy, C., Li, M., Teater, M., Meydan, C., Garrett-Bakelman, F.E., Lee, T.C., Chin, C.R., Durmaz, C., Kawabata, K.C., Dhimolea, E., et al. (2021). Chemotherapy Induces Senescence-Like Resilient Cells Capable of Initiating AML Recurrence. Cancer Discov. *11*, 1542–1561.

46. Karagiannis, T.T., Cleary, J.P., Gok, B., Henderson, A.J., Martin, N.G., Yajima, M., Nelson, E.C., and Cheng, C.S. (2020). Single cell transcriptomics reveals opioid usage evokes widespread suppression of antiviral gene program. Nat. Commun. *11*, 2611.

47. Xu, Y., Zeng, X., Zhang, M., Wang, B., Guo, X., Shan, W., Cai, S., Luo, Q., Li, H., Li, X., et al. (2022). Efficient expansion of rare human circulating hematopoietic stem/progenitor cells in steady-state blood using a polypeptide-forming 3D culture. Protein Cell. 10.1007/s13238-021-00900-4.

48. Regev, A., Teichmann, S.A., Lander, E.S., Amit, I., Benoist, C., Birney, E., Bodenmiller, B., Campbell, P., Carninci, P., Clatworthy, M., et al. (2017). The Human Cell Atlas. Elife *6*, e27041.

49. Liou, G.-Y. (2017). Inflammatory Cytokine Signaling during Development of Pancreatic and Prostate Cancers. J Immunol Res *2017*, 7979637.

# Methods

## SCimilarity model architecture and loss function

## Model architecture

The SCimilarity model consists of one fully connected encoder and one decoder stage and reuses the same encoding network three times per training triplet, such that updates to the model after each batch are shared equally for each subsequent batch of training triplets. The decoder stage is not part of the conventional triplet loss architecture, but is included to compute a mean squared error  (MSE) reconstruction loss.

Expression profiles are reduced through an encoder network, starting from 28,231 genes through three hidden layers with dimensions 1,024, 1,024, and 128. The 128-dimensional outputs are unit length normalized, forcing all low dimensional cell representations to lie on the surface of a hypersphere. During training, the input layer is subjected to 40% dropout, zeroing out many gene expression values at random, and each hidden layer is subjected to 50% dropout rates for maximum regularization [1].

While hyperspheric spaces have been infrequently used for representation of single-cell profiles [2], the triplet loss model often uses hypersphere embeddings to ensure consistency between the model hyperparameters [3]. During triplet loss training, the objective is to place cells of different types sufficiently far apart. The minimum desired distance between cells of different types is called the margin. By fixing the volume of the embedding space to the surface of a unit length

64-dimensional hypersphere, the margin is interpreted consistently between model runs. Without normalization, cells can be placed up to an infinite distance apart, rendering the margin meaningless.

**Triplet loss training**

To learn features that place data points considered similar near each other, the loss function depends on distances between data points embedded in a learned low dimensional latent space, described with:

$$d(x, y) = ||f(x) - f(y)||_2^2$$

where $x$ and $y$ are two high dimensional vectors (here, cell profiles), passed through a neural network encoder $f()$.

The triplet loss model learns from three vectors at a time: the anchor $(x_i^a)$, positive $(x_i^p)$ and negative $(x_i^n)$. The anchor and positive vectors are considered similar, whereas the anchor and negative are dissimilar.

The model parameters are iteratively updated to decrease the number of triplets where the distance between the anchor and negative data vectors is insufficiently large relative to the distance between the anchor and the positive points, thus minimizing the triplet loss function:

$$L_{triplet} = \frac{\sum_i^N max\left(d\left(x_i^a, x_i^p\right) - d\left(x_i^a, x_i^n\right) + \alpha, 0\right)}{N}$$

where $\alpha$ is the margin, which denotes how much further the negatives should be from the anchor than the positives, and $i$ is the index of the triplet.

## Reconstruction loss training

The reconstruction loss is computed on the anchor cell only, because each anchor cell is used only once as an anchor within a batch. The reconstruction loss is defined as:

$$L_{MSE} = \frac{\sum_i^N \left\| x_i^a - g(f(x_i^a)) \right\|_2^2}{N}$$

where $N$ is the number of anchor cells in a batch, set to $N=1000$ in SCimilarity, and $g()$ is the function learned by the neural network decoder stage.

## Combined loss function

Adding a reconstruction loss to classification models has been shown to improve generalization [4] through a regularization effect. The SCimilarity loss function combines the triplet loss and reconstruction loss functions as follows:

$$L = (1 - \beta) * L_{MSE} + \beta * L_{triplet}$$

where $\beta$ is a weighting term in $[0, 1]$. $\beta = 0$ corresponds to a conventional autoencoder, and $\beta = 1$ corresponds to a pure triplet loss model. Empirically, $\beta = 0.001$ performed best on the cell search task (query model) and $\beta = 1$ performed best on batch integration (integration model) (**Extended Data Fig. 2a**).

## Use of Cell Ontology terms and relationships

Authors may annotate cell types at different granularities, which confounds triplet sampling by introducing cell type annotations with hierarchical relationships that cannot be unambiguously defined as either similar or dissimilar. As such, cell type annotations used for training are defined using standardized Cell Ontology terms and valid triplets are restricted to cells without vertical

3

Cell Ontology relationships between members of the triplet. A vertical relationship is defined as any directed path of one or more ancestor-descendant relationships in the Cell Ontology network. Thus, there are three binary relations defined for annotation: (1) similar pairs with identical annotations (*e.g.*, "T cell" and "T cell"), (2) dissimilar pairs with non-vertical ontology relationships (*e.g.*, "CD4-positive, alpha-beta T cell" and "CD8-positive, alpha-beta T cell"), and (3) ambiguous pairs with vertical relationships (*e.g.*, "T cell" and "CD4-positive, alpha-beta T cell"). Positives are drawn from cells similar to the anchor, negatives are drawn from cells dissimilar to the anchor, and cells that are ambiguous to the anchor are excluded from sampling.

**GEO data aggregation**

334 human scRNA-seq datasets were obtained from the Gene Expression Omnibus (GEO)[5]. Multiple filtering steps were used to restrict the datasets analyzed to samples from human tissue, that were generated using the 10x Chromium platform, and which reported unnormalized gene count data that could be automatically processed. To select appropriate datasets, search criteria were designed for the Biopython Entrez search tool (Cock et al., 2019) to find GEO studies that had specific properties, such as metadata keywords, file formats, and species. Then, using GEOparse[6], the GEO text metadata was downloaded for each sample and searched for blacklisted words in the metadata or download URLs (*e.g.*, "smartseq", "trizol", and "fasta") to further filter out samples that were not generated using 10x Chromium. Data for samples and corresponding download links that passed the metadata filter stage were automatically downloaded. No datasets were realigned. 753 studies were identified for download. A set of import functions was designed for the most common file type formats (.mtx, .h5ad, and gene

expression matrices in .tsv or .csv). Any dataset that could not be successfully downloaded or read in was discarded. Once read in, each sample was automatically tested for count data and gene names that match a reference gene list or gene name mapper before saving each file in a uniform h5ad format for later processing. This resulted in a total of 334 published studies that were not duplicates of studies found in CELLxGENE [7] for use in our analysis.

## Data preprocessing

All UMI count data were natural log normalized per-cell with a scaling factor of 10,000 using the scanpy.pp.normalize_to_target(adata, 10000) and scanpy.pp.log1p(adata) functions from scanpy[8].

## Manual data aggregation, normalization and filtering

Datasets with author-provided cell type annotations used for training were obtained from Tabula Sapiens[9], 10x Genomics[10], the single nucleus cross-tissue atlas[11], and the human lung cell atlas[12] and subjected to the same preprocessing procedures as programmatically-downloaded datasets. Cell type annotations were manually converted into terms contained within the Cell Ontology. Cells that with annotations that did not clearly map to the Cell Ontology were not included in training.

Cell profiles previously annotated as doublets, scored as doublets by infer_doublets from Pegasus[13], had >20% total UMI counts aligned to mitochondrial genes, or had <500 total genes detected were removed.

## Preparation of training and test data

Training and test sets were chosen such that entire studies were held out of training (rather than holding out a subset of cells from each dataset) (**Extended Data Table 1**); there were 52 and 14 datasets in the training and test sets, respectively. This presents a harder generalization challenge and reflects how users are likely to use SCimilarity. Test datasets were selected to reflect the tissue diversity within the training sets.

**Selection of Cell Ontology terms for training**

Cell Ontology terms were selected for training if they were observed in at least two separate studies in the training set. Terms that appeared in only one study were not used because SCimilarity is trained by comparing cells across studies. To rescue single-study terms, the immediate parent terms were inspected across studies. If a single-study term's parent was observed in at least two other datasets then the original cell type annotation was replaced with the coarser parent term (**Extended Data Table 1**). Otherwise, all cells with this annotation were removed from training. As the size or annotation quality of training data grows, the number of Cell Ontology terms meeting the inclusion criteria are expected to increase.

**Triplet sampling and semi-hard triplet mining**

During training, batches of 1,024 cells are sampled from the training datasets. This sampling is weighted by study and cell type to have a similar number of observations per cell type from each study per batch.

Because of the *maximum* operation within the loss function, not all viable triplets contribute to the gradient, and are categorized as easy, semi-hard or hard, based on their contribution to the gradient.

Easy negatives are defined as:

$$\left\lVert f(x_i^a) - f(x_i^p) \right\rVert_2^2 < \left\lVert f(x_i^a) - f(x_i^n)) \right\rVert_2^2 + \alpha$$

Easy negatives provide no information to the gradient because the distances between the cells in the low dimensional embedding already satisfy the objective, such that the *maximum* operation returns 0 to the triplet loss sum. Because there are many easy triplets after training a small number of batches, randomly sampling triplets does not train models effectively. To accelerate training, triplets are mined to search for training triplets that are especially informative for model training[3].

Hard negatives are defined as:

$$\left\lVert f(x_i^a) - f(x_i^p) \right\rVert_2^2 > \left\lVert f(x_i^a) - f(x_i^n)) \right\rVert_2^2 + \alpha$$

Hard negatives contribute the largest quantity to the loss function, because they do not fit and are far from fitting the desired latent relationships. In practice, hard triplets are rarely useful for training, because they contribute to model collapse during training[3,14]. Hard negatives may be enriched for incorrectly annotated cells.

Semi-hard negatives are defined as:

$$\left\lVert f(x_i^a) - f(x_i^n)) \right\rVert_2^2 - \left\lVert f(x_i^a) - f(x_i^p) \right\rVert_2^2 < \alpha$$

Semi-hard negatives contribute small amounts to the loss function because they nearly satisfy the desired distances between cells in low dimensional space. Meaning, the negative cell profile is further from the anchor cell than the positive cell, but by a less than desired distance α. Semihard negatives are often used in triplet loss models[3].

Overall, we chose to train SCimilarity using only semi-hard negative triplets.

**Explainability framework and marker gene identification**

An explainability framework was used to identify genes whose variation leads to the most significant variations of the learned features and, in turn, affects the relative distance between different cells.

An explanation for a pair of cells is defined as those genes which have the greatest impact on the relative distance between those cells in latent space. Given $d(x, y) = ||f(x) - f(y)||_2^2$, the distance between two cell profiles $x$ and $y$ in latent space $f$, the integrated gradient approach (Sundararajan et al. 2017) was extended to compute the importance of each gene $i$ in the comparison between cell profiles $x$ and $y$ as:

$$Importance_i(x) = \left| \max((x_i - y_i), 0) \times \int_{\alpha=0}^{1} \frac{\partial d(y + \alpha \times (x - y), y)}{\partial x_i} \right|$$

High values of $Importance_i(x)$ correspond to genes that are highly expressed in $x$, and their modification (i.e., gradient) affects $d(x, y)$ more. Intuitively, the expression of each gene in $y$ is gradually increased to match $x$ along the trajectory from $x$ to $y$. Through this trajectory, the rate

of change of $d(x, y)$ is computed for each gene, aggregating the results. The score is scaled by $(x_i - y_i)$. In order to identify genes that are up regulated in a subset of interest, genes $i$ with expression $x_i < y_i$ are ignored.

This approach differs in several key ways from the standard integrated gradient approach, because: (1) gradients are computed with respect to a learned distance instead of output features, (2) attributions where $x_i < y_i$ are ignored and (3) the sign of the integral is ignored due to the complex interactions between features.

To identify important genes for a cell type $t$, a set of cells $T \in \{t_1, ..., t_N\}$ with cell type $t$ and a set of cells $B \in \{b_1, ..., b_N\}$ with cell types different from $t$ are randomly sampled. Pairwise importances are computed for each pair of cells $t_i$ in $T$ and $b_j$ in $B$ and aggregated to obtain a signature that characterizes cell type $t$ as:

$$Signature_i(t) = \frac{1}{N} \sum_{c=1}^{N} Importance_i(t_c, b_c)$$

Since the pairwise comparisons are averaging relative comparisons, the sampling of $\{b_1, ..., b_N\}$ impacts the signature scoring. To obtain general cell type markers, a background of all cell types is sampled. To obtain a cell state specific signature, a background of cells in other states of the same type are sampled. Confidence intervals for each gene $i$ are computed as the standard error of the mean. This results in an attribution score for each gene.

9

**Training and evaluation metrics**

**SCimilarity score**

The SCimilarity score is defined as the inverse of the cosine distance of two embedded cell profiles:

$$SCimilarity\ score\ = \frac{1}{1 - c_i \bullet c_j}$$

where $c_i$ and $c_j$ are the embeddings of the $i$th and $j$th cell profiles with unit length, respectively and

$i \neq j$. The threshold for similarity varies in practice by question and cell types.

**Ontology-aware modified average silhouette width**

Average silhouette width (ASW) has been used to assess the performance of data integration tasks on multiple scRNA-seq studies[15] by quantifying how coherently grouped each cell type is across studies. The silhouette width of cell profile $i$ of cell type $t$ typically compares the average intra-cell type distances $a(i)$ and the average inter-cell type distances $b(i)$ between cells of type $t$ and cells of the nearest cell type, defined as:

$$a(i) = \frac{1}{|C_I| - 1} \sum_{j \in C_I, i \neq j} d(i, j)$$

$$b(i) = \min_{J \neq I} \frac{1}{|C_J| - 1} \sum_{j \in C_J} d(i, j)$$

where, typically, $C_I$ is the set of cells of author-annotated type $t$ and $C_J$ are the cells of all other cell types.

However, the ASW as typically formulated does not account for differences in granularity of cell type annotations across studies. To address those, a modified formulation is used where $C_I$ contains cell type label $t$ and all of its ontological descendants and $C_J$ is the set of all other cell types, except cells of type $t$ and any of its ontological descendants or ancestors. For example, if computing $a(i)$ for a T cell, distances between all types of T cell terms ("CD4-positive", "alpha-beta T cell", "CD8-positive", "alpha-beta T cell and CD4-positive", "CD25-positive", "alpha-beta regulatory T cell", etc) are members of the "T cell" term. Ancestor terms of T cells, such as the term "Lymphocytes", are not members of the T cell class (nor a T cell subset) but are excluded from the summation indices in the calculations of $a(i)$ and $b(i)$.

**Correlation to predefined gene signatures**

To test how the SCimilarity distance represents distance between predefined cell states, a signature-based definition of cell state was correlated with the SCimilarity score (above).

For each cell in the test set, both the signature score[16] and a SCimilarity score *vs*. the cell query are calculated, yielding two vectors, and Pearson's correlation coefficient is calculated between the vectors.

**Selection of models for downstream analysis**

Models were run in triplicate along 6 different β parameters ranging from [0,1] and one query model and one integration model were selected based on two criteria. First, query performance was tested by how well cell similarities to a query FMΦ profile correlated with a signature defining that same state (*TREM2, GPNMB, SPP1, CCL18, MMP9, CTSK, APOE, CHIT1, LIPA,*

11

*CHI3L1*, *CD14*, *APOC1*). Second, ontology aware ASW was used to quantify how well the cells of the same type from different studies intermixed in SCimilarity's representation. The query model was selected as the model with the highest query test performance. The integration model was selected from the β=1 models. Since the three replicates had nearly identical integration scores, we picked the model with the highest query test score as it performed much better on the query task than the other high integration models. (**Extended Data Fig. 2a**). The selected integration model had more study mixing than the query model according to the study (NMI) and study adjusted rand index (ARI)[15].

## Benchmarking vs. integration methods

SCimilarity's integration and cell search models were each compared to three batch integration methods: Harmony[17], scVI[18], and scArches[19]. A test dataset of 34,713 cells was created by sampling cells from lung tissue studies with uniform probability across studies. The modified ASW (above), adjusted Rand index (ARI) and normalized mutual information (NMI) were calculated as integration benchmark metrics. Harmony was run using the wrapper in Pegasus[13] following the workflow described in https://pegasus-tutorials.readthedocs.io/en/latest/_static/tutorials/batch_correction.html. scVI and scArches were run using the scvi-tools workflow described in https://docs.scvi-tools.org/en/stable/tutorials/notebooks/harmonization.html and https://docs.scvi-tools.org/en/stable/tutorials/notebooks/scarches_scvi_tools.html, respectively. As the scArches workflow requires a reference dataset, 101,133 cell profiles were sampled across all training datasets with uniform probability across studies for use as the reference.

**Cell type annotation**

Cell type assignments were performed by *k*-nearest neighbors (*k*-NN) classification combined with an annotated reference set. SCimilarity's reduced dimensionality latent space was used to determine *k*=50 nearest neighbors in the reference data set to a query cell *t*, and the query cell was annotated either by tallying votes based each cell's annotation with either equal weights,

$$\text{Celltype}(t) = \arg \max_t \left( \sum_{i \in t} \frac{1}{n} \right)$$

or with weights by distance in SCimilarity's reduced dimensionality latent space:

$$\text{Celltype}(t) = \arg \max_t \left( \sum_{y \in t} \frac{1}{d(x,y)} \right)$$

To allow users to annotate new datasets from a restricted list of cell types of interest, excluding (blocklisting) or limiting to (safelisting) specified cell type annotations is used, and is recommended when feasible to improve interpretability and reduce spurious annotations. However, extensive blocklisting or safelisting can slow the annotation process significantly, because the pre-built *k*-NN indices are not optimized for a modified target cell type list.

**_k_NN parameters for annotation and querying**

Two separate *k*NN indices were used for efficient and accurate queries. For cell type annotation, a 7.9M cell *k*-NN index was built using hnswlib[20] with ef_construction = 1000 and M = 80. Searching this *k*-NN found the 50 nearest neighbors (default behavior) for cell type annotation (k=50) and ef=100.

Cell query relied on a separate 22.7M cell $k$-NN index also built using hnswlib. This index was constructed with the following parameters: ef_construction=400 and M=50. The search parameters are set by the user's request for how many similar cells to return. Default behavior is set to $k$=1000 and ef=$k$, but in practice $k$ can vary widely depending on the use case.

**Outlier filtering**

To filter outlier cells prior to visualization and downstream analysis, SCimilarity's score is used to flag cells that are out of distribution. Cells with a SCimilarity score < 50 from the nearest cell in the training set were removed prior to further analysis. Many of these cells were from immortalized cell lines, and reflect their difference from primary cells (and absence in the training). Note that if out of distribution cells are not removed, these cells won't be accurately annotated and can confound visualization.

**Macrophage query preprocessing**

To prepare a cell query for FMΦ cells, a public dataset[21] (GSE136831 and https://www.ipfcellatlas.com) was preprocessed with the same steps for all ingested data and scored use Scanpy's scanpy.tl.score_genes function with a gene signature of *SPP1*, *TREM2*, *GPNMB*, *MMP9*, *CHIT1*, and *CHI3L1* Scanpy[8]. The average profile of the top 50 scoring cell was embedded using SCimilarity and used as the input query to SCimilarity's cell search model and used throughout analyses in **Fig. 5** and **6**.

**Important genes and pathway enrichment**

Important genes were identified using SCimilarity's attribution score method. This method requires two cell groups to compare, identifying which genes differ between them. Here we used 1,000 cells that were considered similar to the average FMΦ profile calculated from Adams et al. as the FMΦ-like group. This query excluded any cells from the Adams et al. dataset. To compare to the FMΦ-like group comparison, 1,000 dissimilar monocytes and macrophages were randomly sampled (any monocyte or macrophage that was not within the top 10,000 most FMΦ similar results).

Reactome pathways enriched for the 100 genes with the top importance scores for FMΦ were determined using the method provided in the ReactomePA[22] R package, with multiple hypothesis correction using the Benjamini-Hochberg method and the background gene universe restricted to the ~28,000 genes included in SCimilarity. Pathways were considered significant if they met the criteria of adjusted p-value (Q) ≤ 0.05 and gene count ≥ 5.

**3DCS culture of PBMC**

Peripheral blood was sourced from healthy volunteers at Genentech that were consented as per IRB. Samples were collected in heparin collection tubes and subsequently diluted 1:1 with a solution of PBS containing 2% FBS and 1mM EDTA. 30 ml of diluted blood was overlayed onto 15 ml of Lymphoprep (STEMCELL Technologies) in a 50ml tube and centrifuged at 3,000 rpm for 20 minutes at 4°C. PBMCs were isolated from the interphase after centrifugation and diluted

with PBS containing 2% FBS and 1 mM EDTA and centrifuged at 300 x g for 10 minutes at 4°C.

Cell pellet was washed again with PBS containing 2% FBS and 1mM EDTA. Red blood cell

lysis was performed on the cell pellet by resuspending in RBC Lysis Buffer (Cell Signaling

Technology) for 5 minutes at room temperature, followed by inactivation with addition of RPMI

media containing 10% FBS. Cells were pelleted by centrifugation at 300 x g for 10 minutes at

4°C and subsequently washed with PBS containing 2% FBS and 1 mM EDTA. Cells were then

resuspended in a 10% sucrose solution at a concentration of 2 x 106 cells/ml right before plating

into 3D hydrogel culture. Puramatrix hydrogel (Corning) was vortexed for 30 seconds and

diluted 1:1 with a 20% sucrose solution. 250 µl of diluted Puramatrix hydrogel was mixed with

250 µl of resuspended PBMCs and plated in a 24-well tissue culture plate. To induce gelation,

RPMI media was overlaid onto the hydrogel/PBMC mixture and incubated for 5 minutes in a

37°C incubator with 5% $CO_2$. Overlayed media was aspirated off of the 3D hydrogel and washed

twice with RPMI media, after which 600 µl of 3DCS media, formulated as previously described

(Xu, Y. et al., Protein & Cell 2022, 13:808-824) was overlaid onto the hydrogel. Cells were

cultured in a 37°C incubator with 5% CO2 for 8 days, with media exchanges every other day. On

day 8, culture cells were recovered from the 3D hydrogel for scRNA-seq.


**Single cell RNA-Seq from 3DCS cultures**

Wells of 3D hydrogel culture were washed with PBS, followed by recovery of the hydrogel and

cells by gentle pipetting in PBS buffer. This solution was centrifuged for 5 minutes at 750 x g

and the hydrogel/PBMC pellet was resuspended in TrypLE solution (ThermoFisher Scientific)

and incubated at 37°C for 10 minutes. RPMI media with 10% FBS was added and the solution

was centrifuged for 5 minutes at 750 x g. The resultant pellet was washed twice with PBS to

16

remove hydrogel matrix debris. PBMCs were resuspended in PBS and passed through a 40 μM filter, pelleted by centrifugation at 300 x g for 5 minutes, and resuspended in RPMI media with 10% FBS. The cell solution was subjected to FACS to isolate cells from any remaining hydrogel debris and recovered cells were concentrated to 1,000 cells/μl in RPMI media with 10% FBS for downstream profiling by scRNA-seq.

ScRNA-seq was performed using the Chromium Single Cell 3' Library and Gel bead kit v3 (10x Genomics), following manufacturer's user guide. Briefly, cell density and viability of single-cell suspension were determined by Vi-CELL XR cell counter (Beckman Coulter). Cell density was used to impute the volume of single cell suspension needed in the reverse transcription (RT) master mix, aiming to achieve ~10,000 cells per sample. cDNAs and libraries were prepared following the manufacturer's user guide (10x Genomics). Libraries were profiled by Bioanalyzer High Sensitivity DNA kit (Agilent Technologies) and quantified using Kapa Library Quantification Kit (Kapa Biosystems). Libraries were sequenced on a NovaSeq 6000 (Illumina) following the manufacturer's specifications with 28+90 bp paired-end reads at a depth of 101M mate-pair reads. Sequencing reads were aligned to the GENCODE 27 Basic gene model on the human genome assembly GRCh38 using Cell Ranger v6.0 (10x Genomics, Pleasanton, CA, USA).

Individual samples were genetically demultiplexing using the singularity container provided with Souporcell 2.0 [23]. No genotype information was provided to the pipeline. Since PBMCs were provided from 3 donors, a k of 3 was used to cluster the samples into 3 genotypes. These samples

were pre-processed consistently with the previously ingested samples and then embedded using SCimilarity to enable direct comparisons to Xu et al as well as the rest of the public datasets.

SCimilarity cell type classification was applied to both public and validation cells using SCimilarity with the following safelist: B cell, CD4-positive, alpha-beta T cell, CD8-positive, alpha-beta T cell, conventional dendritic cell, hematopoietic stem cell, macrophage, monocyte, natural killer cell, plasma cell, plasmacytoid dendritic cell.

**Code performance benchmarking**

Benchmarks were run on servers with 8 Intel Xeon E5-2650 v4 CPUs with 2.20GHz cores and a total of 128 GB of RAM.

Query runtimes, using the pre-built approximate $k$-NN index[20] to find the top $n$ most similar cells, had an average runtime of 50 milliseconds. Some API functions use the query and summarize the metadata within one function call. That function timing is dominated by summarizing metadata and computing statistics from the query results, which requires an additional 3.3 seconds. This performance differs from an exhaustive comparison (**Fig. 5b**), where the query was directly compared against 2.58M monocytes and macrophages with a runtime of 2 seconds.

Cell signatures were calculated using scanpy.tl.score_genes. The scanpy score_genes function was applied to the already normalized data. This runtime totalled 2 hours, 46 minutes and 20

seconds when it was applied across each h5ad file (one file per tissue sample). Even though h5ad files were not stored with any compression, file reading was a dominant factor in runtime.

## Code availability

Code and tutorials are available at https://github.com/Genentech/scimilarity.

## Licensing

- Code license: Apache 2.0

- Pretrained model weights, kNN and pre-built indices license: CC-BY-SA 4.0

## References

1. Baldi, P., and Sadowski, P. (2013). Understanding dropout. In Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2 NIPS'13. (Curran Associates Inc.), pp. 2814–2822.

2. Ding, J., and Regev, A. (2021). Deep generative model embedding of single-cell RNA-Seq profiles on hyperspheres and hyperbolic spaces. Nat. Commun. 12, 2554.

3. Schroff, F., Kalenichenko, D., and Philbin, J. (2015). FaceNet: A Unified Embedding for Face Recognition and Clustering. 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 815–823.

4. Le, L., Patterson, A., and White, M. (2018). Supervised autoencoders: improving generalization performance with unsupervised regularizers. In Proceedings of the 32nd International Conference on Neural Information Processing Systems NIPS'18. (Curran Associates Inc.), pp. 107–117.

5. Edgar, R., Domrachev, M., and Lash, A.E. (2002). Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. Nucleic Acids Res. 30, 207–210.

6.  Gumienny, R. GEOparse: Python library to access Gene Expression Omnibus Database (GEO).

7.  Chan Zuckerberg CELLxGENE Discover (2022). Cellxgene Data Portal.

8.  Wolf, F.A., Angerer, P., and Theis, F.J. (2018). SCANPY: large-scale single-cell gene expression data analysis. Genome Biol. 19, 15.

9.  Tabula Sapiens Consortium*, Jones, R.C., Karkanias, J., Krasnow, M.A., Pisco, A.O., Quake, S.R., Salzman, J., Yosef, N., Bulthaup, B., Brown, P., et al. (2022). The Tabula Sapiens: A multiple-organ, single-cell transcriptomic atlas of humans. Science 376, eabl4896.

10. Zheng, G.X.Y., Terry, J.M., Belgrader, P., Ryvkin, P., Bent, Z.W., Wilson, R., Ziraldo, S.B., Wheeler, T.D., McDermott, G.P., Zhu, J., et al. (2017). Massively parallel digital transcriptional profiling of single cells. Nat. Commun. 8, 14049.

11. Eraslan, G., Drokhlyansky, E., Anand, S., Fiskin, E., Subramanian, A., Slyper, M., Wang, J., Van Wittenberghe, N., Rouhana, J.M., Waldman, J., et al. (2022). Single-nucleus cross-tissue molecular reference maps toward understanding disease gene function. Science 376, eabl4290.

12. Travaglini, K.J., Nabhan, A.N., Penland, L., Sinha, R., Gillich, A., Sit, R.V., Chang, S., Conley, S.D., Mori, Y., Seita, J., et al. (2020). A molecular cell atlas of the human lung from single-cell RNA sequencing. Nature 587, 619–625.

13. Li, B., Gould, J., Yang, Y., Sarkizova, S., Tabaka, M., Ashenberg, O., Rosen, Y., Slyper, M., Kowalczyk, M.S., Villani, A.-C., et al. (2020). Cumulus provides cloud-based data analysis for large-scale single-cell and single-nucleus RNA-seq. Nat. Methods 17, 793–798.

14. Wu, C.-Y., Manmatha, R., Smola, A.J., and Krähenbühl, P. (2017). Sampling Matters in Deep Embedding Learning. arXiv [cs.CV].

15. Luecken, M.D., Büttner, M., Chaichoompu, K., Danese, A., Interlandi, M., Mueller, M.F., Strobl, D.C., Zappia, L., Dugas, M., Colomé-Tatché, M., et al. (2022). Benchmarking atlas-level data integration in single-cell genomics. Nat. Methods 19, 41–50.

16. Satija, R., Farrell, J.A., Gennert, D., Schier, A.F., and Regev, A. (2015). Spatial reconstruction of single-cell gene expression data. Nat. Biotechnol. 33, 495–502.

17. Korsunsky, I., Millard, N., Fan, J., Slowikowski, K., Zhang, F., Wei, K., Baglaenko, Y., Brenner, M., Loh, P.-R., and Raychaudhuri, S. (2019). Fast, sensitive and accurate integration of single-cell data with Harmony. Nat. Methods 16, 1289–1296.

18. Lopez, R., Regier, J., Cole, M.B., Jordan, M.I., and Yosef, N. (2018). Deep generative modeling for single-cell transcriptomics. Nat. Methods 15, 1053–1058.

19. Lotfollahi, M., Naghipourfar, M., Luecken, M.D., Khajavi, M., Büttner, M., Wagenstetter, M., Avsec, Ž., Gayoso, A., Yosef, N., Interlandi, M., et al. (2022). Mapping single-cell data to reference atlases by transfer learning. Nat. Biotechnol. 40, 121–130.

20. Malkov, Y.A., and Yashunin, D.A. (2020). Efficient and Robust Approximate Nearest Neighbor Search Using Hierarchical Navigable Small World Graphs. IEEE Trans. Pattern Anal. Mach. Intell. 42, 824–836.

21. Adams, T.S., Schupp, J.C., Poli, S., Ayaub, E.A., Neumark, N., Ahangari, F., Chu, S.G., Raby, B.A., DeIuliis, G., Januszyk, M., et al. (2020). Single-cell RNA-seq reveals ectopic and aberrant lung-resident cell populations in idiopathic pulmonary fibrosis. Science Advances 6, eaba1983.

22. Yu, G., and He, Q.-Y. (2016). ReactomePA: an R/Bioconductor package for reactome pathway analysis and visualization. Mol. Biosyst. 12, 477–479.

23. Heaton, H., Talman, A.M., Knights, A., Imaz, M., Gaffney, D.J., Durbin, R., Hemberg, M., and Lawniczak, M.K.N. (2020). Souporcell: robust clustering of single-cell RNA-seq data by genotype without reference genotypes. Nat. Methods 17, 615–620.
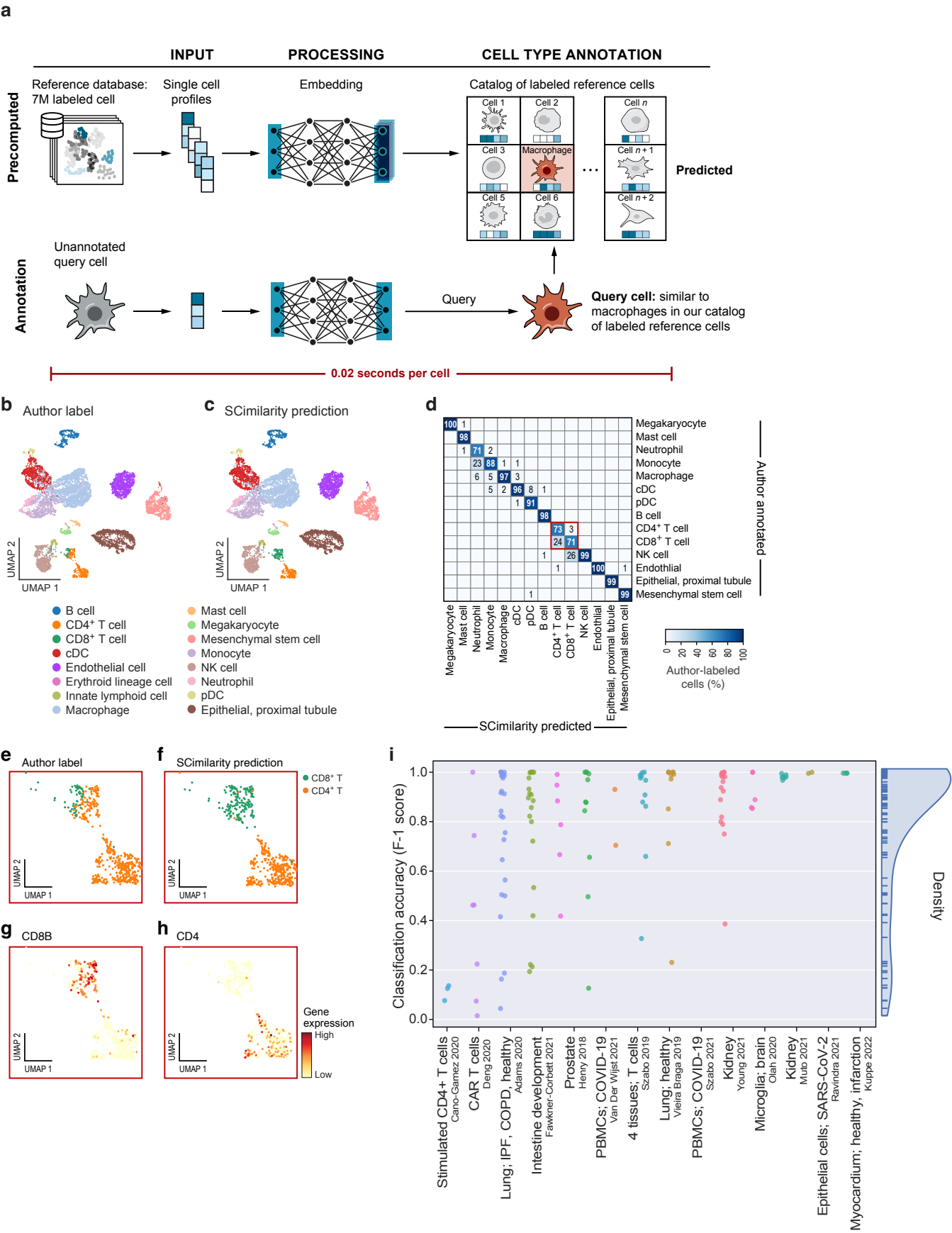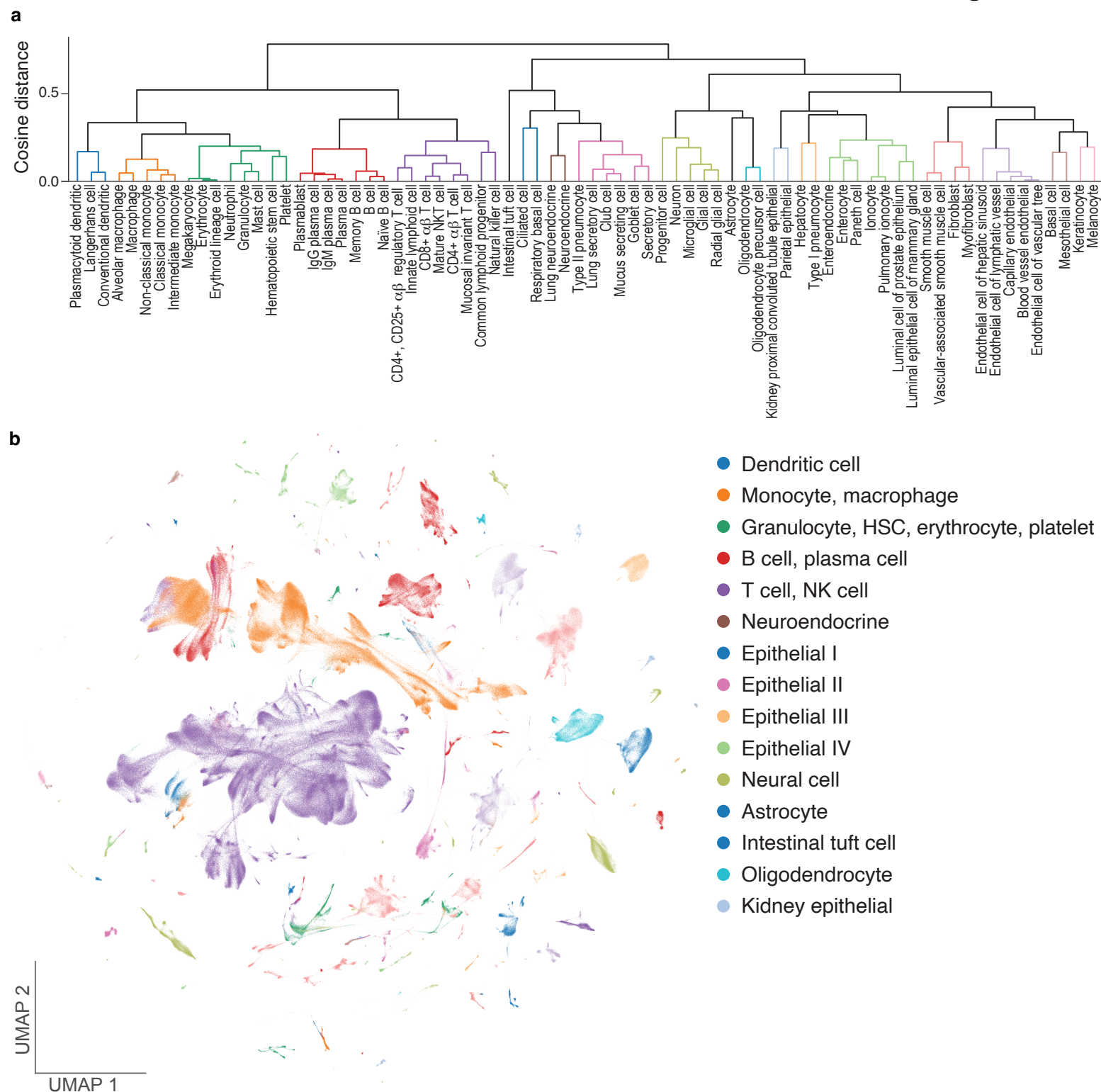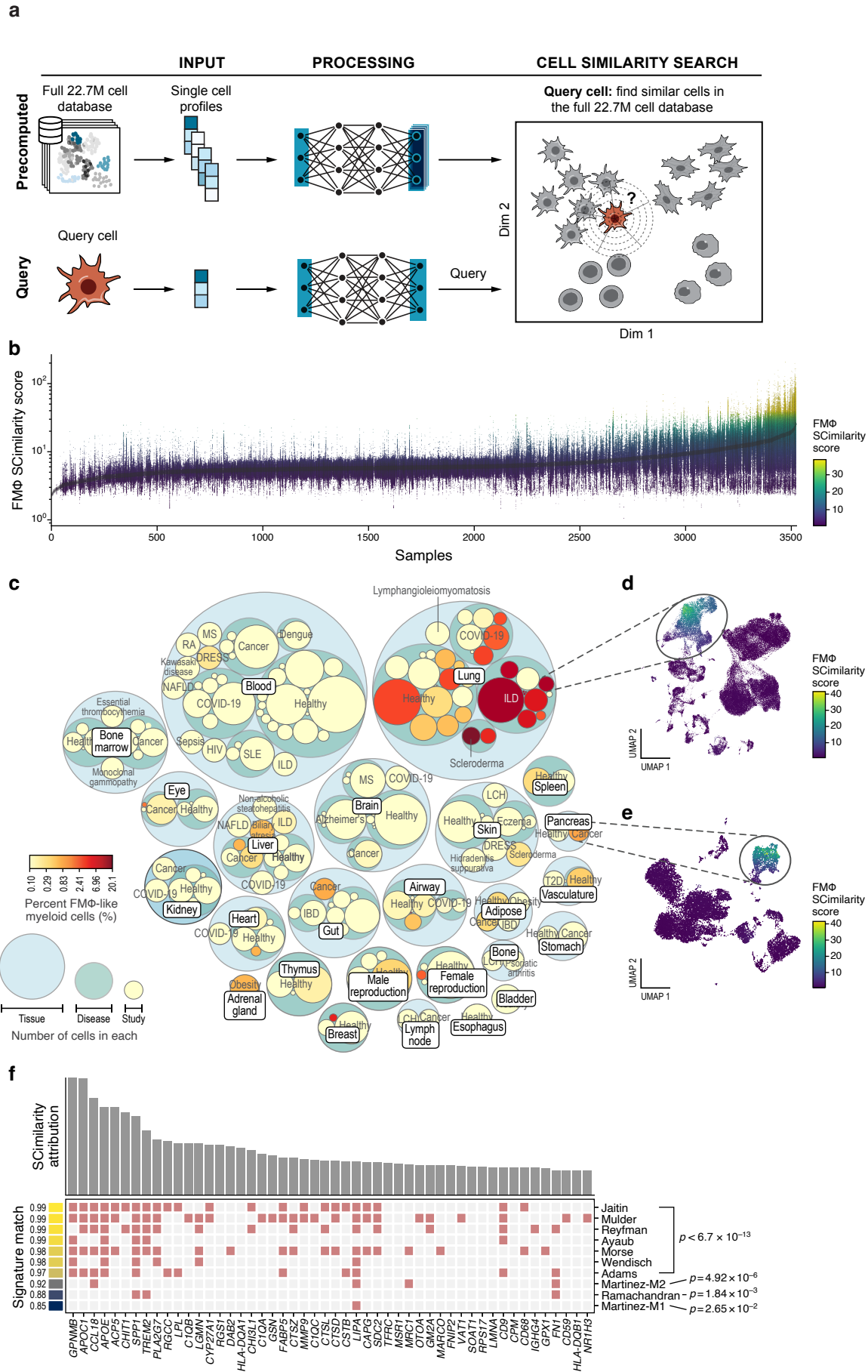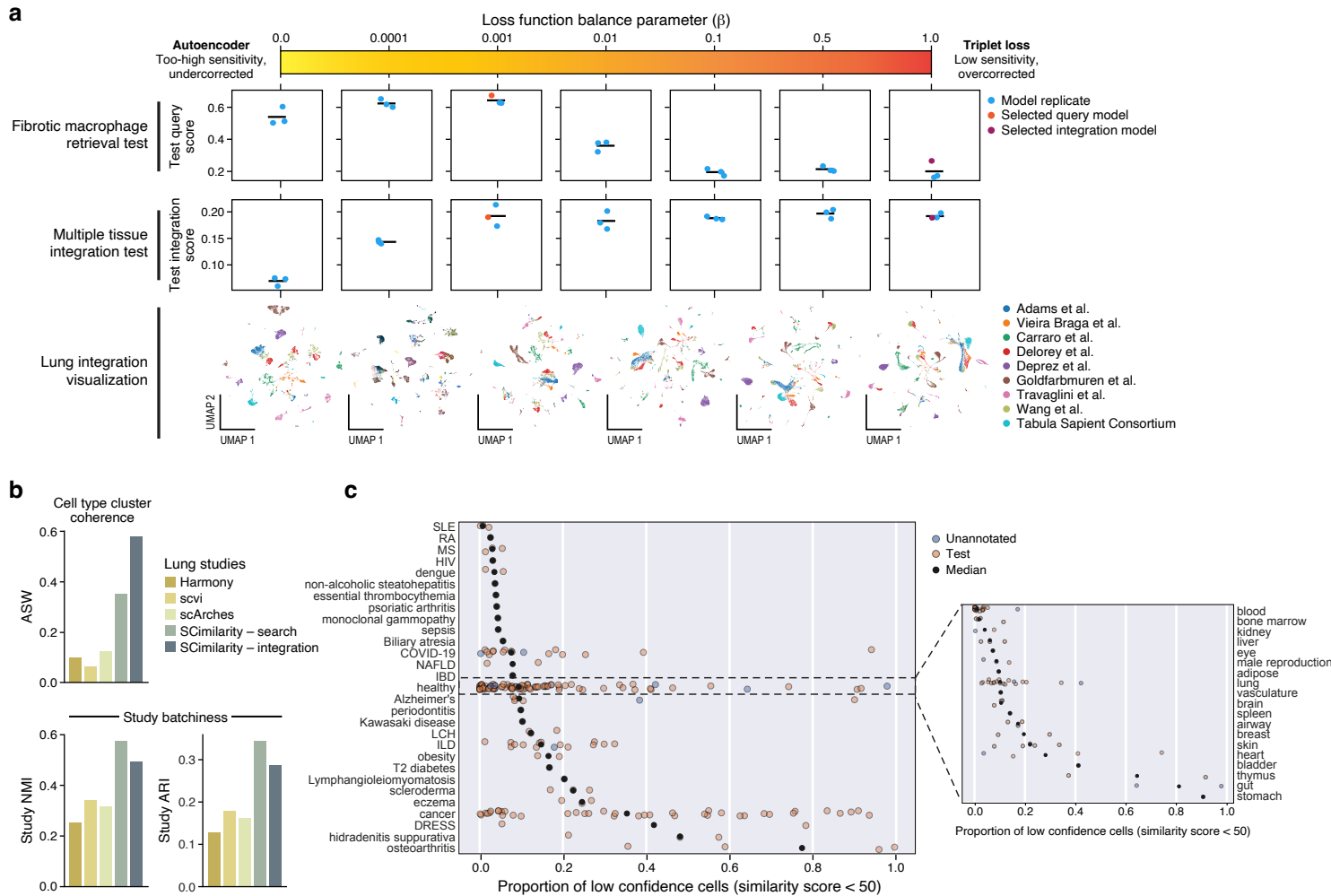
Figure 1

Figure 2

Figure 3

Figure 4

Figure 5
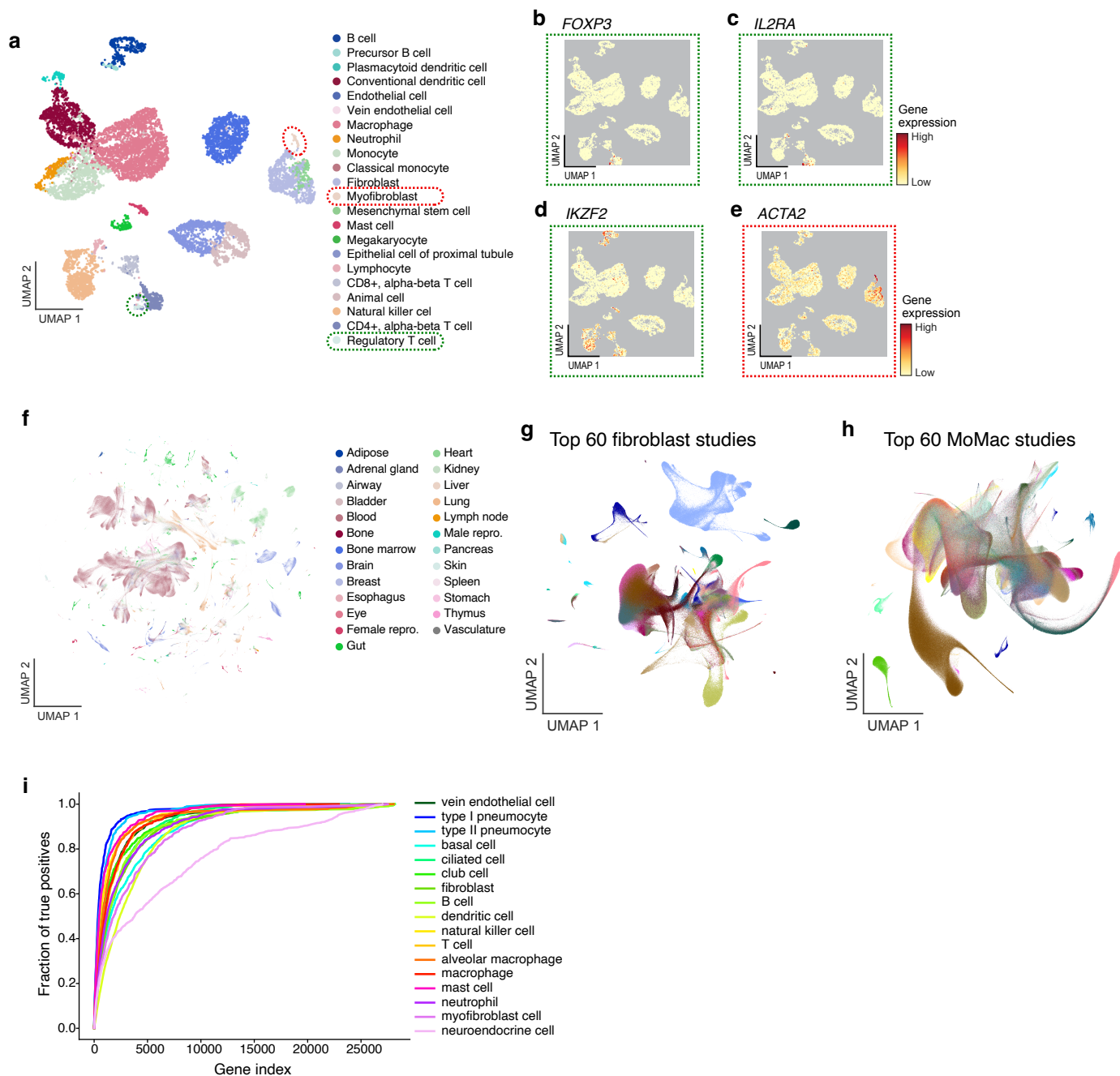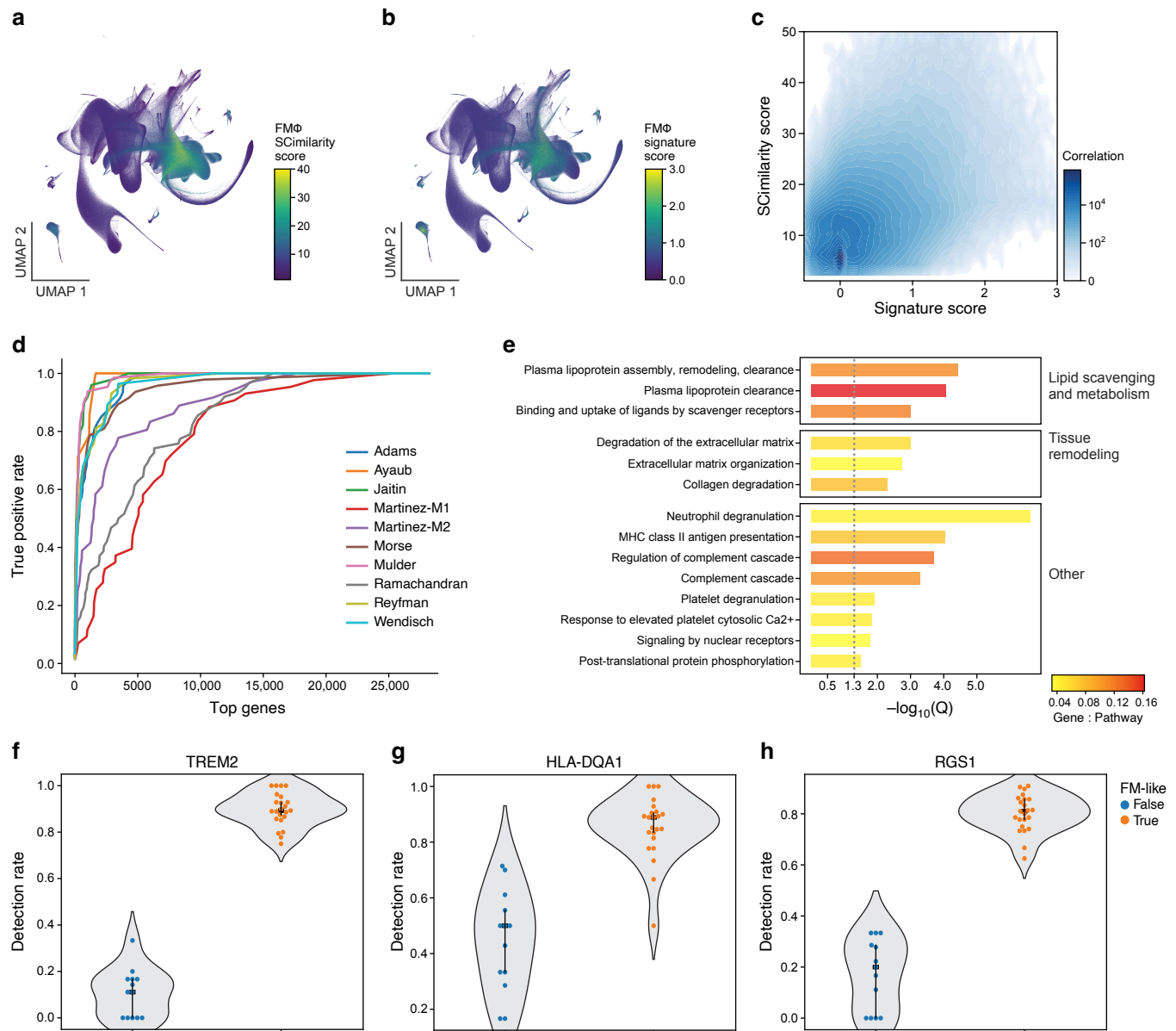
Figure 6

# Extended data Fig. 3

Extended data Fig. 4

1    **Fig. 1. SCimilarity metric learning enables cell search in large human scale atlases.**

2    **a,** Cell Querying with SCimilarity. Left: A query cell profile is compared to a searchable reference

3    collection of 22.7M profiles from 399 studies. Center: Sample with similar cells are identified and

4    returned with information about the original sample conditions, including unexpected tissue, *in*

5    *vitro* or diseases contexts. Right: A SCimilarity score is computed between the query cell and each

6    cell within a tissue sample. **b,** Triplet loss training. From left: 52 training and 14 test annotated (by

7    the Cell Ontology) datasets from across the body are sampled for cell triplets (an anchor, a

8    "positive" (anchor-similar), and a "negative" (anchor-dissimilar) cell; based on Cell Ontology

9    annotations) to train a neural network that embeds similar cells closer than dissimilar ones

10   (**Methods**).

1 **Fig. 2. SCimilarity learns a universal representation that generalizes to new datasets.**

2 **a,** A large-scale reference database of public gene expression datasets across tissues and diseases.

3 Number of cells (circle size) across tissues (outermost light blue circles) and disease states (middle

4 green circles) across individual studies (innermost circles) in the training (gold), test (pink) or

5 unannotated (purple) datasets. **b,** Integration between studies without feature selection or batch

6 correction. Uniform manifold approximation and projection (UMAP) embedding of cell profiles

7 (dots) generated on a 128-dimensional latent space from SCimilarity's integration model

8 (**Methods**) for cells from 21 tissues (panels) and 239 unique studies (color code). For tissues with

9 more than a million cell profiles, the UMAP embedding was computed on a random uniform

10 subsampling of 1 million cells from the studies for that tissue.

1 **Fig. 3. SCimilarity accurately annotates cell types across the human body**

2 **a**, SCimilarity cell annotation. A new unannotated cell (red, bottom left) is embedded in

3 SCimilarity's common low-dimensional space and compared against the precomputed reference

4 for cell type annotation (0.02 seconds per cell). **b-d**, SCimilarity annotation of a kidney scRNA-

5 seq dataset. (**b,c**) UMAP embedding of cell profiles (dots) from SCimilarity's latent representation

6 of a held out kidney dataset[1] colored by author provided (**b**) or SCimilarity predicted (**c**) cell type

7 annotations. (**d**) Percentage (color bar and number) of author-annotated cells (columns) with each

8 SCimilarity annotation (rows). **e-h,** SCimilarity-corrected author annotations. UMAP embedding

9 of cell profiles (dots) from SCimilarity's latent representation for cells either author-annotated or

10 predicted as CD4$^+$ or CD8$^+$ T cells, colored by author (**e**) or SCimilarity (**f**) annotations, and by

11 *CD8A* (**g**) or *CD4* (**h**) expression. **i,** Classification performance. F1 score (y axis) for SCimilarity

12 *vs*. author annotation for each cell population (dot) from each of 14 held out datasets (x axis).

13 Right: Distribution of F1 scores.

1    **Fig. 4. SCimilarity annotations scale to tens of millions of cells from hundreds of datasets**

2    **a,b** Predicted cell types group by biologically-accurate lineages. **a,** Hierarchical clustering

3    dendrogram of centroids of predicted cell types (leaves) in SCimilarity latent space, colored by

4    lineage. Clustering was performed using cosine distance and average linkage. **b,** UMAP of

5    2,000,000 embedded cells uniformly sampled from the 22.7M reference, colored by clusters (as

6    labeled in **a**).

1    **Fig. 5. SCimilarity cell search reveals FMΦs across ILD and other diseases**

2    **a,** SCimilarity cell search. A query cell profile (bottom left) is embedded into the learned

3    SCimilarity representation along with the reference of 22.7M cells, and its nearest neighbors,

4    determined by distance from the embedded query in the low dimensional space, are tabulated by

5    study, tissue and disease. **b-e,** Identification of FMΦs across tissues by SCimilarity cell search. (**b**)

6    SCimilarity scores (y axis, $\log_{10}$ scale, and color bar) against a FMΦ query profile for each

7    annotated monocyte or macrophage (dot) from 1,041 *in vivo* tissue samples from 143 studies (x

8    axis), ordered by mean SCimilarity score. (**c**) Number of cells (circle size) across tissues

9    (outermost light blue circles), disease states (middle green circles), and individual studies

10    (innermost circles, colored by fraction of monocytes and macrophages with SCimilarity scores

11    >95[th] percentile of all FMΦ SCimilarity scores (log scaled color bar)). Circle size for disease and

12    individual study are scaled relative to other diseases in the same tissue, or studies in the same

13    disease. (**d,e**) UMAP embeddings of cell profiles (dots) from the SCimilarity representation (query

14    model) from an ILD[2] (**d**) and PDAC[3] (**e**) studies, colored by FMΦ query SCimilarity scores (color

15    bar). **f,** Identification of FMΦs associated genes by importance. Integrated gradients attribution

16    scores (y axis, top) for the genes (x axis top, and columns, bottom) with top 50 scores for FMΦs

17    *vs.* lung macrophages (**Methods**), and their membership (red: presence; grey: absence) in

18    published macrophage signatures (bottom, rows). Left color bar: AUC of the ranking of each

19    published signature in SCimilarity attribution scores (AUC=1: all *n* signature genes are listed as

20    the top *n* genes by SCimilarity attribution scores for distinguishing FMΦ). Martinez-M1 and -M2:

21    macrophage states expected to be different from FMΦs. P-value: hypergeometric test

1  **Fig. 6. SCimilarity cell search identifies *in vitro* cells matching an *in vivo* FMΦ state and a**

2  **novel *in vitro* disease model.**

3  **a,** Identification of FMΦs-like cells across *in vitro* samples by SCimilarity cell search. SCimilarity

4  scores (y axis, $\log_{10}$ scale, and color bar) against a FMΦ query profile for each annotated myeloid

5  cell (dot) from 40 *in vitro* samples (x axis) (from 17 studies), ordered by mean SCimilarity score.

6  Gray boxes: Day 0 and Day 5 samples from a 3D-hydrogel culture system[4]. **b-f,** 3D conditions

7  yield FMΦ-like cells *in vitro* in validation experiments. **b,** SCimilarity scores (y axis, $\log_{10}$ scale,

8  and color bar) against a FMΦ query profile for each annotated myeloid cell (dot) from the original

9  3D-hydrogel culture system dataset[4] and from 3 donors in the validation experiment (x axis). **c,**

10  Mean expression (dot color) and proportion of expressing cells (dot size) of genes (rows) with high

11  SCimilarity attribution score for distinguishing FMΦs *in vivo* (as in **Fig. 5f**) for myeloid cells in

12  the original 3D-hydrogel culture system[4] and the validation experiment (columns). **d-f,** UMAP

13  embedding from SCimilarity's query model latent space of cell profiles (dots) from day 0 (**d**) or

14  day 5 (**e**) of the original 3D-hydrogel culture system[4] or from day 8 of the replication experiment

15  (**f**), colored by FMΦ SCimilarity score (color bar). **g,** replication of Xu et al.'s original finding of

16  HSC expansion. Proportion of HSCs between Xu et al.'s day 0, day 5 and our validation day 8

17  time points.

1 **Extended Data figure legends**

2 **Extended Data Fig. 1. Data compendium to assemble a pan-human reference.**

3 **a,b**, Cumulative number of cells (**a**, y axis) and samples (**b**, y axis) profiled by sc/snRNA-seq (and

4 matching our filters; **Methods**) over time (x axis). Doubling time is calculated based on the

5 publication date from the most recent 150 data points (dashed red line). **c,** Author-annotated cell

6 types used in training. Number of author-annotated cells (color bar) from each Cell Ontology type

7 (rows) and study (columns) used for SCimilarity model training. **d,** Tissues and diseases used in

8 training. Number of studies (heatmap tiles, text and color bar) and cells (margins, y or x axis) used

9 for model training from each tissue (rows, right y-axis) and disease (columns, top x-axis).

1    **Extended Data Fig. 2. SCimilarity training details and hyperparameter search.**

2    **a,** Impact of triplet and autoencoder loss mixing on model performance, where the leftmost column

3    is a traditional autoencoder and the rightmost column is exclusively triplet loss. The FMΦ retrieval

4    test quantifies how much correlation there is between signature scoring of FMΦs and SCimilarity

5    score to FMΦ. The multiple tissue integration test quantifies an ontology-aware average silhouette

6    width where a higher score denotes more coherent clusters for each cell type. The bottom row

7    shows UMAPs for each loss function mix for nine lung scRNA-seq datasets, colored by study. **b,**

8    Benchmarking SCimilarity to established data integration models. Ontology-aware average

9    silhouette width (ASW, y axis, top), normalized mutual information (NMI, y axis, bottom left) and

10   adjusted Rand index (ARI, y axis, bottom right) for SCimilarity's integration and search models

11   and for scVI, scArches, and Harmony (x axis), each applied to nine lung datasets. **c,** Outlier cells

12   from different types and tissues. Fraction (x axis) of cells from different disease (left) or healthy

13   (right) tissue samples with low similarity (SCimilarity score <50) to training data.

14

1  **Extended Data Fig. 3. Validation of large-scale integration and annotation.**

2  **a,** Unconstrained cell annotation. UMAP embedding of single cell profiles (dots) from

3  SCimilarity's latent representation of the held out scRNA-Seq kidney data[1] (as in **Fig. 3b,c,**

4  colored by cell annotation without constraining target labels to the scope of author-provided labels

5  in this studyor by expression of select marker genes of regulatory T cells (**b-d**) or myofibroblasts

6  (**e**). **f,** UMAP of 2,000,000 embedded cells uniformly sampled from the 22.7M reference, colored

7  by tissue (as in **Fig. 4b**). **g,h,** UMAP embedding of cell profiles predicted by SCimilarity as

8  fibroblast/myofibroblast (**g**) or monocytes/macrophages (**h**), colored by study (for the 60 studies

9  contributing most cells). **i,** SCimilarity cell-type important genes match cell-type specific

10  signatures. Fraction of cell type-specific differentially expressed genes (from Eraslan *et al.*[5]) (y

11  axis) captured by top-n important genes (x axis) for that cell type by SCimilarity's integrated

12  gradients attribution analysis.

1     **Extended Data Fig. 4. FMΦs among monocytes and macrophages.**

2     **a-c,** Agreement between SCimilarity and traditional FMΦ cell scores. **a,b** UMAP embedding of

3     2,578,221 monocyte and macrophage cell profiles (dots) from SCimilarity's latent space

4     representation colored by SCimilarity score using a prototypical FMΦ cellular profile defined from

5     Adams *at el.*[6] (**a**) or Scanpy's signature score for FMΦ associated genes (**b**). (**c**) Scanpy FMΦ gene

6     signature score (x axis) and FMΦ SCimilarity score (y axis) for each cell (shown as density). **d,**

7     Agreement between SCimilarity FMΦ important genes and published FMΦ signatures. ROC curve

8     of the fraction of each study's gene sets (y axis) captured within the top genes by SCimilarity

9     attribution ranking (x axis). **e,** FMΦ important genes are enriched for relevant pathways. False

10     discovery rate ($-\log_{10}$(q value), hypergeometric test, x axis) for enrichment of Reactome pathways

11     (y axis, $Q \leq 0.05$ and gene count $\geq 5$) with the 100 genes with the top integrated gradients

12     attribution scores for the FMΦ query (ranked by score). Color: ratio of important genes within a

13     Reactome pathway to the total size of the pathway. **f-h,** Expression of known and novel genes

14     associated with FMΦs. Pseudobulked gene expression values for ILD tissue samples for known

15     marker TREM2 (**f**) and enriched genes not previously described (**g,h**).

1.  Young, M. D. *et al.* Single cell derived mRNA signals across human kidney tumors. *Nat. Commun.* **12**, 3896 (2021).

2.  Morse, C. *et al.* Proliferating SPP1/MERTK-expressing macrophages in idiopathic pulmonary fibrosis. *Eur. Respir. J.* **54**, (2019).

3.  Lin, W. *et al.* Single-cell transcriptome analysis of tumor and stromal compartments of pancreatic ductal adenocarcinoma primary tumors and metastatic lesions. *Genome Med.* **12**, 80 (2020).

4.  Xu, Y. *et al.* Efficient expansion of rare human circulating hematopoietic stem/progenitor cells in steady-state blood using a polypeptide-forming 3D culture. *Protein Cell* (2022) doi:10.1007/s13238-021-00900-4.

5.  Eraslan, G. *et al.* Single-nucleus cross-tissue molecular reference maps toward understanding disease gene function. *Science* **376**, eabl4290 (2022).

6.  Adams, T. S. *et al.* Single-cell RNA-seq reveals ectopic and aberrant lung-resident cell populations in idiopathic pulmonary fibrosis. *Science Advances* **6**, eaba1983 (2020).