# Efficient and precise single-cell reference atlas mapping with Symphony

Joyce B. Kang[1-5], Aparna Nathan[1-5], Nghia Millard[1-5], Laurie Rumker[1-5], D. Branch Moody[3], Ilya Korsunsky[1-5]**, Soumya Raychaudhuri[1-6]**

[1] Center for Data Sciences, Brigham and Women's Hospital, Boston, MA, USA

[2] Division of Genetics, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA

[3] Division of Rheumatology, Inflammation, and Immunity, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA

[4] Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

[5] Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

[6] Versus Arthritis Centre for Genetics and Genomics, Centre for Musculoskeletal Research, Manchester Academic Health Science Centre, The University of Manchester, Manchester, UK

** These authors jointly supervised this work.

Correspondence to:

Ilya Korsunsky

Harvard New Research Building

77 Avenue Louis Pasteur

Boston, MA 02115

ikorsunskiy@bwh.harvard.edu


Soumya Raychaudhuri

Harvard New Research Building

77 Avenue Louis Pasteur, Suite 250

Boston, MA 02115

soumya@broadinstitute.org

Ph: 617-525-4484 Fax: 617-525-4488

# Abstract

Recent advances in single-cell technologies and integration algorithms make it possible to construct large, comprehensive reference atlases from multiple datasets encompassing many donors, studies, disease states, and sequencing platforms. Much like mapping sequencing reads to a reference genome, it is essential to be able to map new query cells onto complex, multimillion-cell reference atlases to rapidly identify relevant cell states and phenotypes. We present Symphony, a novel algorithm for building compressed, integrated reference atlases of $\geq 10^6$ cells and enabling efficient query mapping within seconds. Based on a linear mixture model framework, Symphony precisely localizes query cells within a low-dimensional reference embedding without the need to reintegrate the reference cells, facilitating the downstream transfer of many types of reference-defined annotations to the query cells. We demonstrate the power of Symphony by (1) mapping a query containing multiple levels of experimental design to predict pancreatic cell types in human and mouse, (2) localizing query cells along a smooth developmental trajectory of human fetal liver hematopoiesis, and (3) harnessing a multimodal CITE-seq reference atlas to infer query surface protein expression in memory T cells. Symphony will enable the sharing of comprehensive integrated reference atlases in a convenient, portable format that powers fast, reproducible querying and downstream analyses.

## Introduction

Advancements in single-cell RNA-sequencing (scRNA-seq) have launched an era in which individual studies can routinely profile $10^4$-$10^6$ cells[1–3], and multimillion-cell datasets are already emerging[4,5]. Single-cell resolution enables the discovery and refinement of cell states across diverse clinical and biological contexts[6–11]. To date, most studies redefine cell states from scratch, making it difficult to compare results across studies and thus hampering reproducibility. Coordinated large-scale efforts, exemplified by the Human Cell Atlas (HCA)[12], aim to establish comprehensive and well-annotated reference datasets comprising millions of cells that capture the broad spectrum of cell states. Building these reference datasets requires integrating multiple datasets which may have been collected under different technical and biological conditions. Hence, reference construction requires application of one of many recently developed single-cell integration algorithms[13–19]. Our group previously developed Harmony[15], a fast, accurate, and well-reviewed method[20] that is able to explicitly model complex study design, a property that makes it suitable for integrating complex datasets into reference atlases[21–24]. Once such atlases are constructed, powerful mapping algorithms will make it possible to rapidly and reproducibly map new single-cell datasets onto the reference and automatically annotate them by transferring information from nearby reference cells.

Fast mapping of query cells against a large, stable reference is a well-recognized open problem[25] and active area of research[18,26,27]. One inefficient but accurate approach to project reference and query cells into a joint embedding is to integrate both sets of cells together *de novo*, resulting in what might be considered a "gold standard" embedding. While this is a reasonable approach for relatively small reference datasets, the strategy is intractable for atlas-sized references with millions of cells. It requires users to "rebuild" the reference for each analysis, and requires potentially cumbersome and administratively challenging exchanges of large-scale datasets. Furthermore, it may corrupt the reference embedding once a reference is carefully constructed and annotated. It is instead preferable to freeze the reference when mapping new query cells onto it.

3

71  High-quality reference mapping requires both a fast and accurate mapping algorithm and a framework

72  to efficiently store a reference dataset. An ideal reference mapping algorithm must meet four key

73  requirements: handle complex study design in both the reference and query, scale to large datasets,

74  map with high accuracy, and enable inference of diverse query cell annotations based on reference

75  cells. Here, we present Symphony, a novel algorithm to compress a large, integrated reference and

76  map query cells to a precise location in the reference embedding within seconds. Through multiple real-

77  world dataset analyses, we show that Symphony can enable accurate downstream inference of cell

78  type, developmental trajectory position, and protein expression, even when the query itself contains

79  complex confounding technical and biological effects.

80  # Results

81  **Symphony compresses an integrated reference for efficient query mapping**

82  Symphony comprises two main algorithms: reference compression and mapping (**Methods, Fig. S1a**).

83  Symphony *reference compression* captures and structures information from multiple reference datasets

84  into an integrated and concise format that can subsequently be used to map query cells (**Fig. 1a-b**).

85  Symphony builds upon the same linear mixture model framework as Harmony[17]. Briefly, in a low-

86  dimensional embedding, such as principal component analysis (PCA), the model represents cell states

87  as soft clusters, in which a cell's identity is defined by probabilistic assignments across one or more

88  clusters. For *de novo* integration of the reference, cells are iteratively assigned soft cluster

89  memberships, which are used as weights in a linear mixture model to remove unwanted covariate-

90  dependent effects. To store the reference efficiently without saving information on individual reference

91  cells, Symphony computes summary statistics learned in the low-dimensional space (**Fig. 1b,**

92  **Methods**), returning computationally efficient data structures containing the "minimal reference

93  elements" needed to map new cells. These include the means and standard deviations used to scale

94  the genes, the gene loadings from PCA (or another low dimensional projection, e.g. canonical

95  correlation analysis [CCA]) on the reference cells, soft-cluster centroids from the integrated reference,
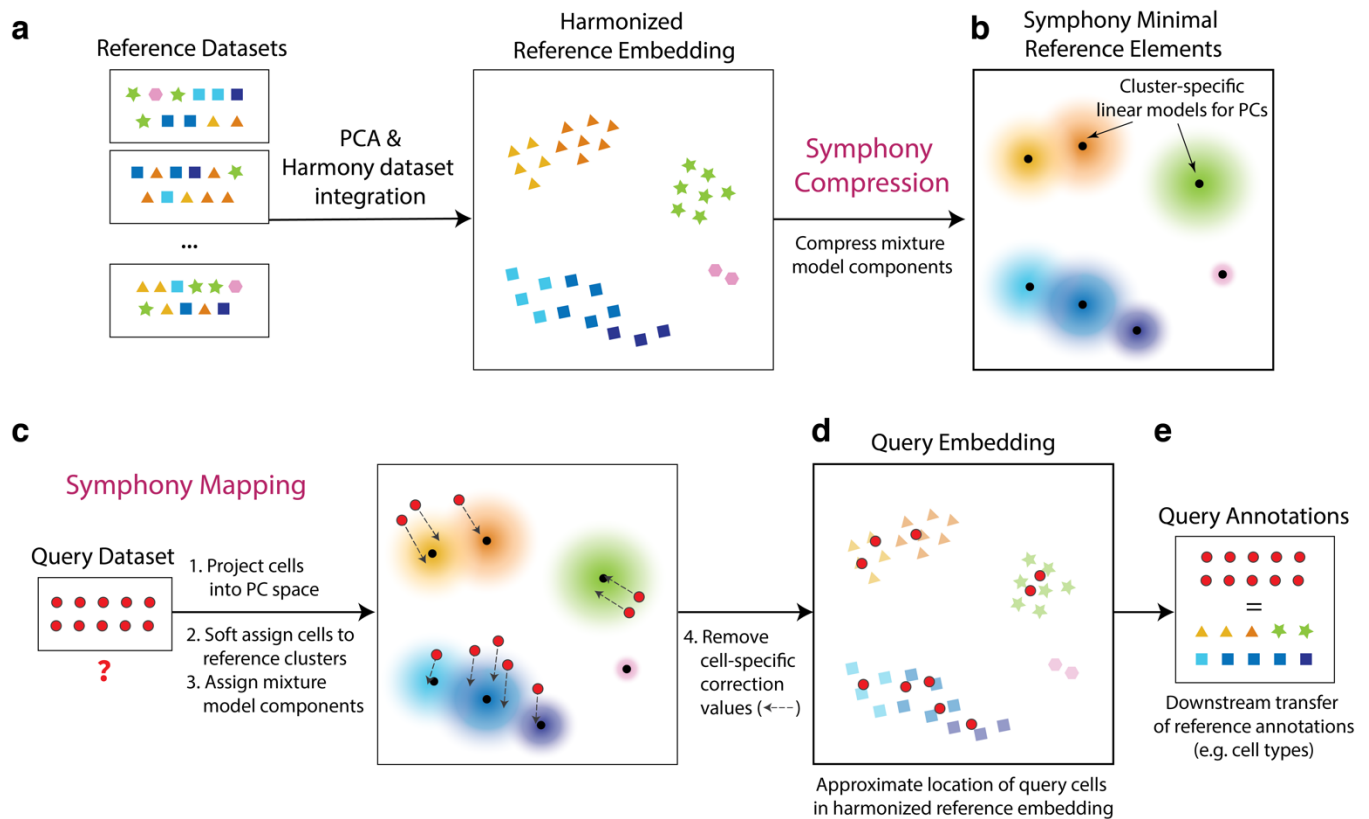
4

**Figure 1 | Overview of Symphony framework and algorithm.** Symphony comprises two algorithms: Symphony compression **(a-b)** and Symphony mapping **(c-d)**. **(a)** To construct a reference, cells from multiple datasets are embedded in a lower-dimensional space (e.g. PCA), in which dataset integration (Harmony) is performed to remove dataset-specific effects. Shape indicates distinct cell types, and color indicates finer-grained cell subtypes or states. **(b)** Symphony compression represents the information captured within the harmonized reference in a concise, portable format based on computing summary statistics for the reference-dependent components of the linear mixture model. Symphony returns the minimal reference elements needed to efficiently map new query cells to the reference. **(c)** Given an unseen query dataset and compressed reference, Symphony mapping localizes the query cells to their appropriate locations within the integrated reference embedding **(d)**. The reference cell locations do not change during mapping. **(e)** The shared Symphony feature embedding can be used for downstream transfer of reference-defined annotations to the query cells.

96    and two "compression terms" (a $k$ x 1 vector and $k$ x $d$ matrix, where $k$ is the number of clusters and $d$ is

97    the dimensionality of the embedding) (**Methods, Supplementary Equations, Fig. S1b**).

98    To map new query cells to the compressed reference, we apply Symphony *mapping*. The algorithm

99    approximates integration of reference and query cells *de novo* (**Methods**), but uses only the minimal

100   reference elements to compute the mapping (**Fig. S1c**). First, Symphony projects query gene

101   expression profiles into the same uncorrected low-dimensional space as the reference cells (e.g. PCs),

102   using the saved scaling parameters and reference gene loadings (**Fig. 1c**). Second, Symphony

103   computes soft cluster assignments for the query cells based on proximity to the reference cluster

104   centroids. Finally, to correct unwanted user-specified technical and biological effects in the query data,

105   Symphony assumes the soft cluster assignments from the previous step and uses stored mixture model

106   components to regress out the query batch effects (**Fig. 1d**). Importantly, the reference cell embedding

107   remains stable during mapping. Embedding the query within the reference coordinates enables

108   downstream transfer of annotations from reference cells to query cells, including discrete cell type

109   classifications, quantitative cell states (e.g. position along a trajectory), or expression of missing genes

110   or proteins (**Fig. 1e**).

**Symphony approximates *de novo* integration without reintegration of the reference cells**

112   As we demonstrate in the **Methods**, Symphony is equivalent to running *de novo* Harmony integration if

113   three conditions are met: (I) all cell states represented in the query data set are captured by the

114   reference dataset, (II) the number of query cells is much smaller than the number of reference cells,

115   and (III) the query dataset has a design matrix that is independent of reference datasets (i.e. non-

116   overlapping batches in reference and query). As the scope of available single-cell atlases continues to

117   grow, it is reasonable to assume that reference datasets are large and all-inclusive, making conditions

118   (I) and (II) well-supported. Condition (III) is also typically met if the query data was generated in

119   separate experiments from the reference.

120   To demonstrate that Symphony mapping closely approximates running *de novo* integration on all cells,

121   we applied Symphony to 20,792 peripheral blood mononuclear cells (PBMCs) assayed with three

122    different 10x technologies: 3'v1, 3'v2, and 5'. We performed three mapping experiments. For each, we

123    built an integrated Symphony reference from two technologies, then mapped the third technology as a

124    query. The resulting Symphony embeddings were compared to a gold standard embedding obtained by

125    running Harmony on all three datasets together. Visually, we found that the Symphony embedding for

126    each mapping experiment (**Fig. 2a**) closely reproduced the overall structure and cell type information of

127    the gold standard embedding (**Fig. 2b**). To quantitatively assess the degrees of dataset mixing we use

128    the Local Inverse Simpson's Index (LISI)[26–30] metric; higher LISI scores correspond to better mixing of

129    cells across batches. LISI scores in Symphony embeddings (mean LISI 2.16, 95% CI [2.16, 2.17]) and

130    *de novo* integration embeddings (mean LISI 2.14, 95% CI [2.13, 2.15]) were nearly identical (**Fig. 2c**,

131    **Methods**).

132    To directly assess similarity of the local neighborhood structures, we computed the correlation between

133    the local neighborhood adjacency graphs generated by Symphony and *de novo* integration. We define

134    a new metric called k-nearest-neighbor correlation (k-NN-corr), which quantifies how well the local

135    neighborhood structure in a given embedding is preserved in an alternative embedding by looking at

136    the correlation of neighbor cells sorted by distance (**Fig. S2a-e**). Anchoring on each query cell, we

137    calculate (1) the pairwise similarities to its *k* nearest reference neighbors in the gold standard

138    embedding and (2) the similarities between the same query-reference neighbor pairs in the alternate

139    embedding (**Methods**), then calculate the Spearman correlation between (1) and (2). k-NN-corr ranges

140    from -1 to +1, where +1 indicates a perfectly preserved sorted ordering of neighbors. We find that for

141    *k*=500, the Symphony embeddings produce a k-NN-corr >0.4 for 77.3% of cells (and positive k-NN-corr

142    for 99.9% of cells), demonstrating that Symphony not only maps query cells to the correct broad cluster

143    but also preserves the distance relationships between nearby cells in the same local region (**Fig. 2d**).

144    As a comparison, we calculated k-NN-corr for a simple PC projection of the query cells (with no

145    correction step) using the original reference gene loadings prior to integration and observed

146    significantly lower correlations (Wilcoxon signed-rank $p<2.2e-16$), with k-NN-corr >0.4 for 39.9% of cells
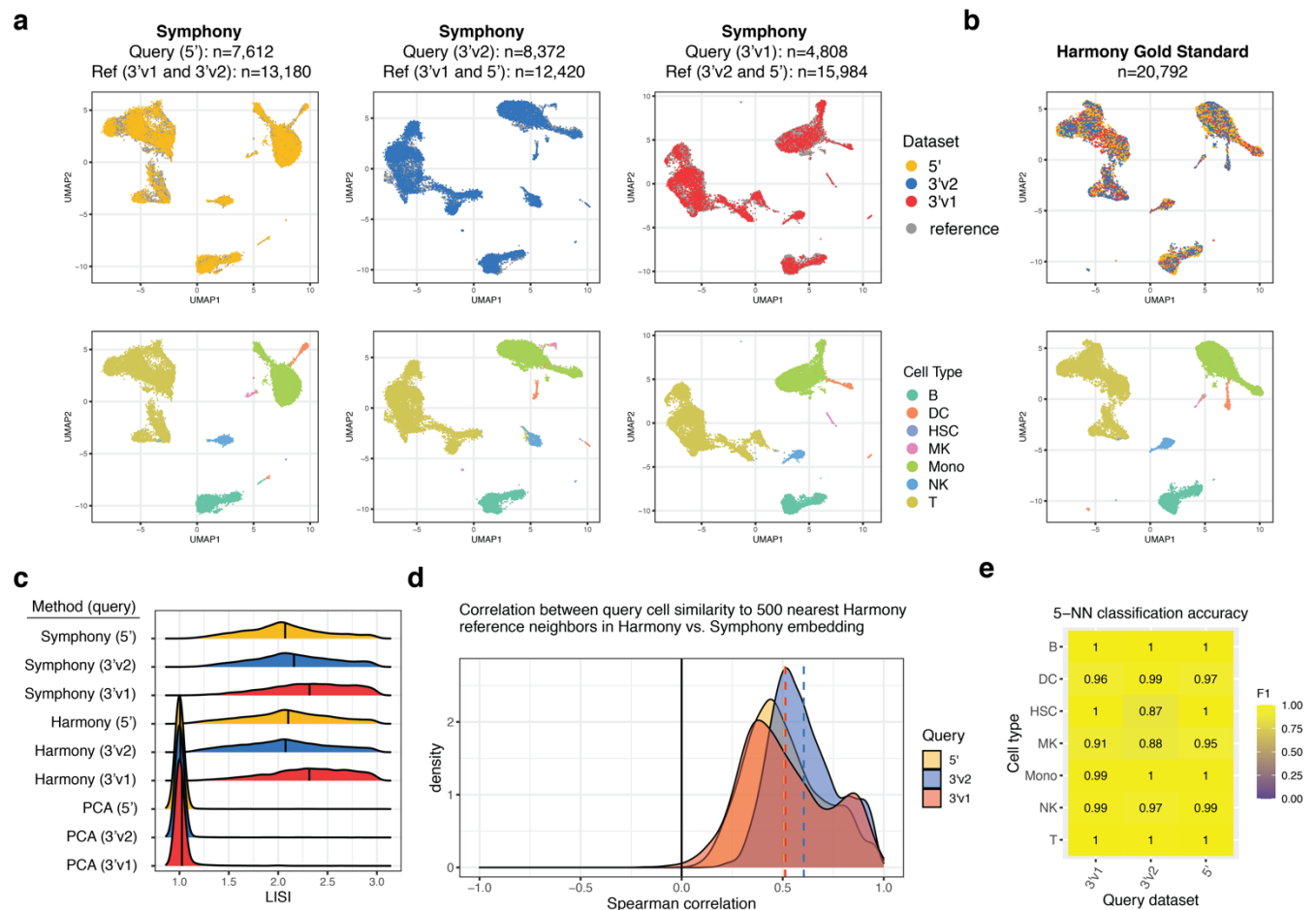
147    (**Fig. S2f**).

**Figure 2 | Symphony approximates *de novo* integration without reintegration of the reference.** Three PBMC datasets were sequenced with different 10x protocols: 5' (yellow, n=7,697 cells), 3'v2 (blue, n=8,380 cells), and 3'v1 (red, n=4,809 cells). We ran Symphony three times, each time mapping one dataset onto the other two. **(a)** Symphony embeddings generated across the three mapping experiments (columns). Top row: cells colored by query (yellow, blue, or red) or reference (gray), with query cells plotted in front. Bottom row: cells colored by cell type: B cell (B), dendritic cell (DC), hematopoietic stem cell (HSC), megakaryocyte (MK), monocyte (Mono), natural killer cell (NK), or T cell (T), with query cells plotted in front. **(b)** For comparison, gold standard *de novo* Harmony embedding colored by dataset (top) and cell type (bottom). **(c)** Distribution of LISI scores for query cells in the Symphony embeddings, gold standard, and a standard PCA pipeline on all cells. **(d)** Distribution of k-NN-corr (Spearman correlation between the similarities between the neighbor-pairs in the Harmony embedding and the similarities between the same neighbor-pairs in the Symphony embedding) across query cells for k=500, colored by query dataset. **(e)** Classification accuracy as measured by cell type F1 scores for downstream query cell type annotation using 5-NN on the Symphony embedding.

148 **Symphony enables accurate cell type classification**

149 If Symphony is effective, then cells should be mapped close to cells of the same cell type, enabling

150 accurate cell type classification. To test this, we performed post-mapping query cell type classification

151 in the 10x PBMCs example from above. We used a 5-NN classifier to annotate query cells across 7 cell

152 types based on the nearest reference cells in the harmonized embedding and compared the predictions

153 to the ground truth labels assigned *a priori* with lineage-specific marker genes (**Methods, Table S2**).

154 Across all three experiments, predictions using the Symphony embeddings achieved 99.5% accuracy

155 overall, with a median cell type F1-score (harmonic mean of precision and recall, ranging from 0 to 1) of

156 0.99 (**Fig. 2e, Table S3**). This indicates that Symphony appropriately localizes query cells in

157 harmonized space to enable the accurate transfer of cell type labels.

158 Automatic cell type classification represents an open area of research[28–32]. Existing supervised

159 classifiers assign a limited set of labels to new cells based on training data and/or marker genes. To

160 benchmark Symphony-powered downstream inference against existing classifiers, we followed the

161 same procedure as a benchmarking analysis in Abdelaal et al. (2019)[28]. The benchmark compared 22

162 cell type classifiers on the PbmcBench dataset consisting of two PBMC samples sequenced using 7

163 different protocols[33]. For each protocol train-test pair (42 experiments) and donor train-test pair

164 (additional 6 experiments) (**Methods**), we built a Symphony reference from the training dataset then

165 mapped the test dataset. We used the resulting harmonized feature embedding to predict query cell

166 types using three downstream models: 5-NN, SVM with radial kernel, and multinomial logistic

167 regression. The Symphony-based classifiers achieve consistently high cell type F1-scores (average

168 median F1 of 0.79-0.83) comparable to the top three supervised classifiers for this benchmark

169 (scmapcell, singleCellNet, and SCINA, average median F1 of 0.77-0.83) (**Fig. S3a**). Notably, as

170 discussed in Abdelaal et al., the median F1-score alone can be misleading given that some classifiers

171 (including SCINA) leave low-confidence cells as "unclassified", whereas we used Symphony to assign a

172 label to every cell. This benchmark is also arguably suboptimal in that the reference in each experiment

173 is comprised of a single dataset (no reference integration involved).

7

174 **Symphony maps against a large reference within seconds**

175 To demonstrate scalability to large reference atlases, we evaluated Symphony's computational speed.

176 We downsampled a large memory T cell dataset[34] to create benchmark reference datasets with 20,000,

177 50,000, 100,000, 250,000, and 500,000 cells (from 12, 30, 58, 156, and 259 donors, respectively).

178 Against each reference, we mapped three different-sized queries: 1,000, 10,000, and 100,000 cells

179 (from 1, 6, and 64 donors) and measured total elapsed runtime (**Fig. S4, Table S4**). The speed of the

180 reference building process is comparable to that of running *de novo* integration since they both start

181 with expression data and require a full pipeline of scaling, PCA, and Harmony integration. However, a

182 reference need only be built and saved once in order to map all subsequent query datasets onto it. For

183 instance, initially building a 500,000-cell reference with Symphony took 5,163 seconds (86.1 min) and

184 mapping a subsequent 10,000-cell query onto it took only 0.99 secs, compared to 4,806 secs (80.1

185 mins) for *de novo* integration on all cells. Symphony offers a 5000x speedup in this application. These

186 results show that Symphony scales efficiently to map against multimillion-cell references, enabling it to

187 power potential web-based queries within seconds.

188 Importantly, Symphony mapping time does not depend on the number of cells or batches in the

189 reference since the reference cells are modeled post-batch correction (**Methods**); however, it does

190 depend on the reference complexity (number of centroids $k$ and dimensions $d$) and number of query

191 cells and batches (**Table S4**) since the query mapping algorithm solves for the query batch coefficients

192 for each of the reference-defined clusters.

193 **Mapping a query dataset with multi-level experimental design in human and mouse**

194 **pancreas**

195 Symphony is designed to handle query datasets with multiple batches, technologies, individuals, or

196 other structures. To demonstrate, we used Symphony in a scenario in which both the reference and the

197 query have complex experimental designs (**Fig 3a**). The reference contained 6,177 pancreatic islet

198 cells from 32 human donors across four independent studies[35–38], each employing a different plate-

199 based scRNA-seq technology (CEL-seq, CEL-seq2, Smart-seq2, and Fluidigm C1; **Fig 3b**). We
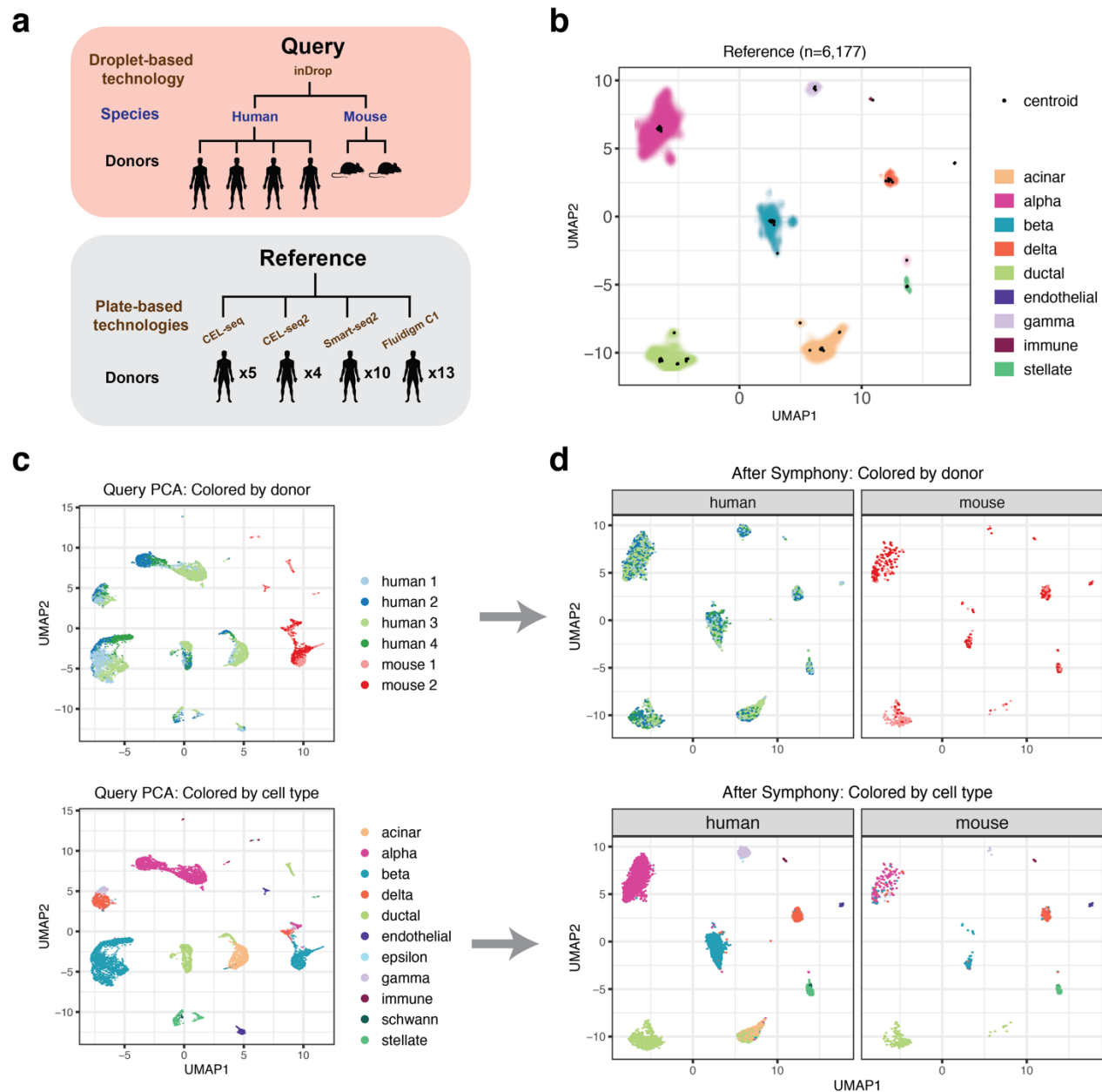
8

**Figure 3 | Symphony maps pancreas cells with multi-level query. (a)** Symphony mapping can model and remove multiple potentially nested sources of batch variation in the query, such as technology, species, and donor. In this example, the query dataset (n=10,455 cells, from 4 human donors and 2 mouse donors) was sequenced on a new technology (inDrop) previously unseen in the reference. **(b)** A Symphony reference (n=6,177 cells, 32 donors) was built from four human pancreas datasets, each assayed using different technologies, and broad cell types were annotated using canonical marker genes. Contour fill represents density of reference cells. Black points represent soft-cluster centroids in the Symphony mixture model. **(c)** Without mapping, a standard PCA pipeline shows that query cells exhibit strong species and donor effects. **(d)** Query cells are mapped against the reference by simultaneously removing the effect of technology, species, and donor in the query such that the cells group by cell type with mixing between species and among donors. Top row colored by donor; bottom row colored by cell type as previously defined by Baron et al. (2016).

200 integrated across donors and technologies, defined clusters, and manually annotated cell types using

201 cluster-specific marker genes (**Methods, Table S5**). The query contained 8,569 pancreatic islet cells

202 from 4 human donors and 1,866 cells from 2 mice, all profiled with inDrop, a droplet-based scRNA-seq

203 technology absent in the reference[39].

204 PCA of the query dataset alone revealed large sources of variation from both species and donor

205 identity (**Fig. 3c**). Symphony mapped the multi-level droplet-based query onto the plate-based

206 reference by simultaneously modeling and removing the effects of technology, species, and donor

207 within the query (**Fig. 3d**). By removing all three nested sources of variation, we accurately predicted

208 query cell types with a 5-NN classifier in the harmonized embedding: median cell type F1-scores of

209 0.97 (overall accuracy 96.4%) for human and median cell type F1 of 0.94 (overall accuracy 86.4%) for

210 mouse cells, with ground truth labels defined by the original publication[39] (**Table S6**). By mapping

211 against a reference, Symphony is able to overcome strong species effects and simultaneously map

212 analogous cell types between mouse and human.

213 **Localizing query cells along a reference-defined trajectory of human fetal liver**

214 **hematopoiesis**

215 A successful mapping method should position cells not only within cell type clusters but also along

216 smooth transcriptional gradients, commonly used to model differentiation and activation processes over

217 time (**Fig. 4a**). To test Symphony in a gradient mapping context, we built and mapped to a reference

218 atlas profiling human fetal liver hematopoiesis, containing 113,063 liver cells from 14 donors spanning

219 7-17 post-conceptional weeks of age and 27 author-defined cell types, sequenced with 10x 3' chemistry

220 (**Fig. 4b, Fig. S5a**)[40]. Trajectory analysis of immune populations with the force directed graph (FDG)

221 algorithm[40] highlights relationships among progenitor and differentiated cell types (**Fig. 4c**). Notably, the

222 hematopoietic stem cell and multipotent progenitor population branches into three major trajectories,

223 representing the lymphoid, myeloid, and megakaryocyte-erythroid-mast (MEM) lineages. This reference

224 contains two forms of annotation for downstream query inference: discrete cell types and positions

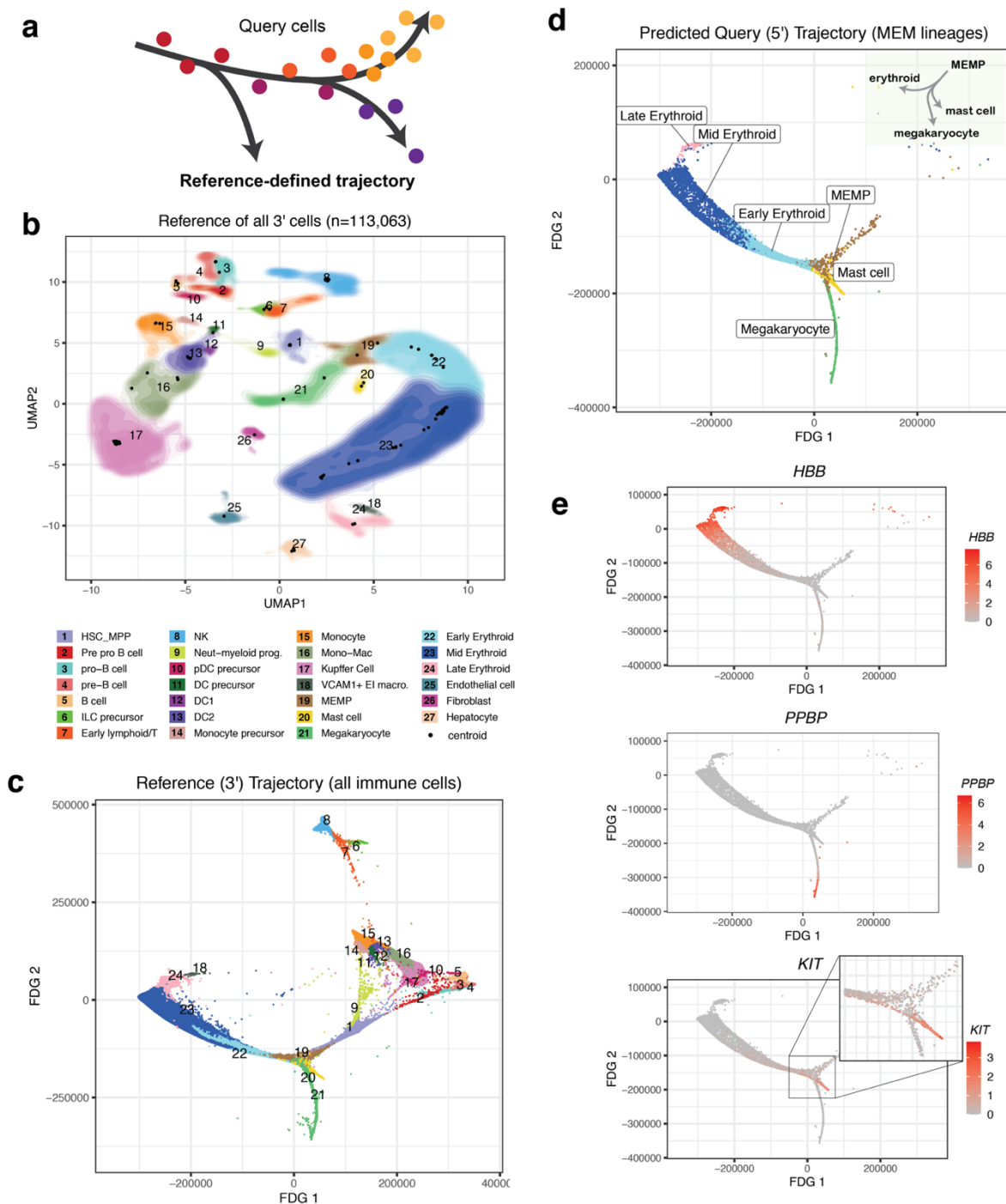225 along differentiation gradients.

9

**Figure 4 | Localizing query cells along a trajectory of fetal liver hematopoiesis. (a)** Symphony can precisely place query cells along a reference-defined trajectory. The reference (n=113,063 cells, 14 donors) was sequenced using 10x 3' chemistry, and the query (n=25,367 cells, 5 donors) was sequenced with 10x 5' chemistry. **(b)** Symphony reference colored by cell types as defined by Popescu et al. (2019). Contour fill represents density of cells. Black points represent soft-cluster centroids in the Symphony mixture model. **(c)** Reference developmental trajectory of 3'-sequenced immune cells (FDG coordinates obtained from original authors). Query cells in the MEM lineages (n=5,141 cells) were mapped against the reference and query coordinates along the trajectory were predicted with 10-NN **(d)**. The inferred query trajectory preserves branching within the MEM lineages, placing terminally differentiated states on the ends. **(e)** Expression of lineage marker genes (*PPBP* for megakaryocytes, *HBB* for erythroid cells, and *KIT* for mast cells). Cells colored by log-normalized expression of gene.

226    We mapped a query consisting of 21,414 new cells from 5 of the original 14 donors, sequenced with

227    10x 5' chemistry. We first inferred query cell types with k-NN classification (**Methods**) and confirmed

228    accurate cell type assignment based on the authors' independent query annotations[40] (median cell type

229    F1=0.92 across 14 held-out donor experiments within 3' dataset only, median cell type F1=0.83 for the

230    5'-to-3' experiment; **Fig. S6**). To evaluate query trajectory inference, we used the Symphony joint

231    embedding to position query cells from the MEM lineage (n=5,141) in the reference-defined trajectory

232    by averaging the 10 nearest reference cell FDG coordinates. The inferred query trajectory (**Fig. 4d**)

233    recapitulated known branching from MEM progenitors (MEMPs, brown) into distinct megakaryocyte

234    (green), erythroid (blue, pink), and mast cell (yellow) lineages. Moreover, transitions from MEMPs to

235    differentiated types were marked by gradual changes in canonical marker genes (**Fig. 4e**): *PPBP* for

236    megakaryocytes, *HBB* for erythrocytes, and *KIT* for mast cells. These gradual expression patterns are

237    consistent with correct placement of query cells along differentiation gradients.

238    **Inferring query surface protein marker expression by mapping to a reference assayed**

239    **with CITE-seq**

240    Recent technological advances in multimodal single-cell technologies (e.g., CITE-seq) make it possible

241    to simultaneously measure mRNA and surface protein expression from the same cells using

242    oligonucleotide-tagged antibodies[41,42]. With Symphony, we can construct a reference from these data,

243    map query cells from experiments that measure only mRNA expression, and infer surface protein

244    expression for the query cells to expand possible analyses and interpretations (**Fig. 5a**).

245    To demonstrate this, we used a CITE-seq dataset that measures the expression of whole-transcriptome

246    mRNA and 30 surface proteins on 500,089 peripheral blood memory T cells from 271 samples[43]. We

247    leveraged both mRNA and protein features to build a multimodal reference from 80% of samples

248    (n=217) and map the remaining 20% of samples (n=54). Instead of using PCA, which is best for one

249    modality[44], we used canonical correlation analysis (CCA) to embed reference cells into a space that

250    leverages both. Specifically, CCA constructs a pair of correlated low-dimensional embeddings, one for

251    mRNA and one for protein features, each with a linear projection function akin to gene loadings in PCA.
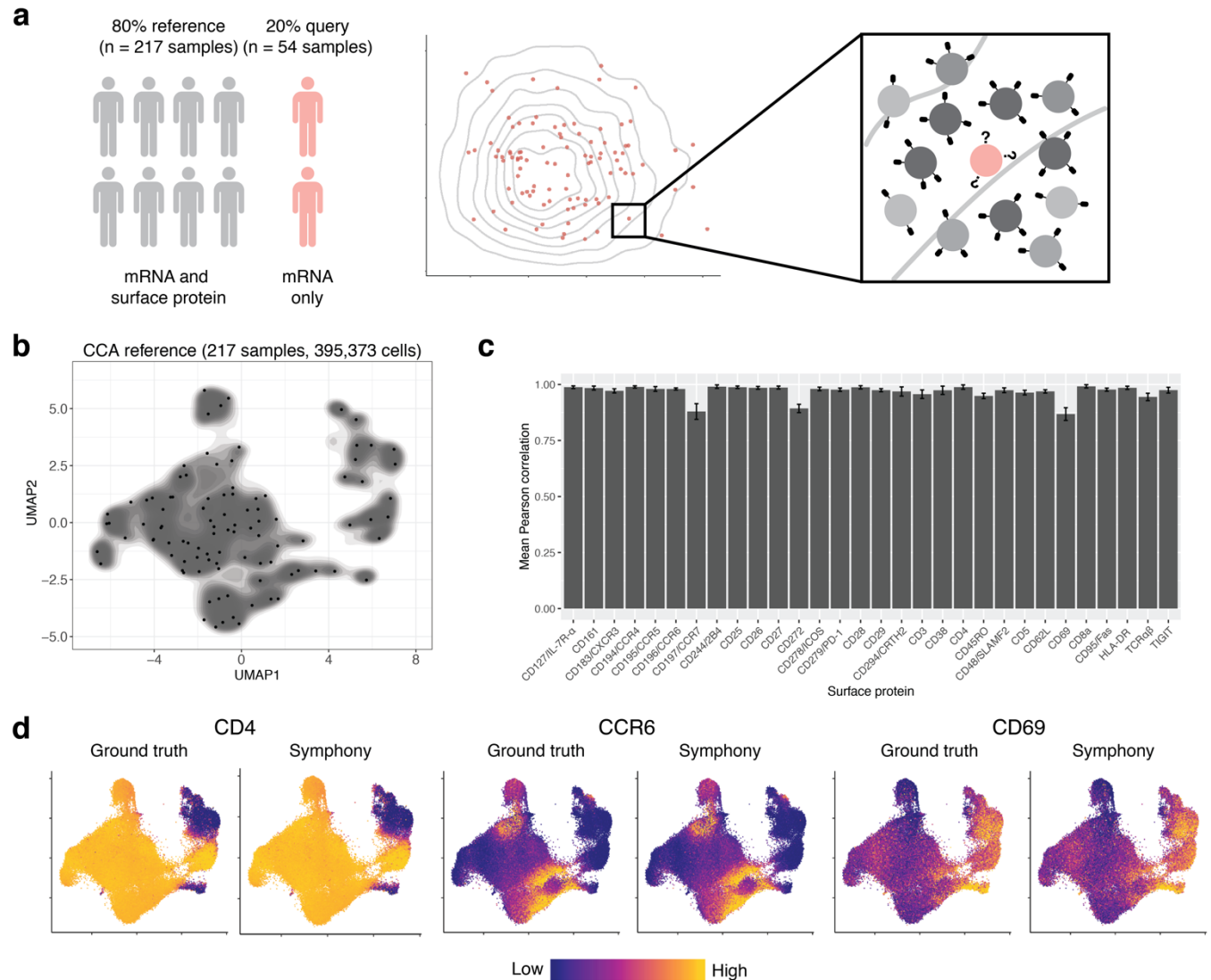
**Figure 5 | Mapping onto a multimodal reference to infer query surface protein expression in memory T cells. (a)** Schematic of multimodal mapping experiment. The dataset was divided into training and test sets (80% and 20% of samples, respectively). The training set was used to build a Symphony reference, and the test set was mapped onto the reference to predict surface protein expression in query cells (pink) based on 50-NN reference cells (gray). **(b)** Symphony reference built from mRNA/protein CCA embedding. Contour fill represents density of reference cells. Black points represent soft-cluster centroids in the Symphony mixture model. **(c)** We measured the accuracy of protein expression prediction with the Pearson correlation between predicted and ground truth expression for each surface protein across query cells in each donor. Bar height represents the average per-donor correlation for each protein, and error bars represent standard deviation. **(d)** Ground truth and predicted expression of CD4, CCR6, and CD69 based on CCA reference. Ground truth is the 50-NN-smoothed expression measured in the CITE-seq experiment. Colors are scaled independently for each marker from minimum (blue) to maximum (yellow) expression.

252    We corrected reference batch effects in CCA space with Harmony and built a Symphony reference

253    (**Fig. 5b**), saving the gene loadings for the CCA embedding from mRNA features. Then, we mapped

254    the held-out query using only mRNA expression to mimic a unimodal scRNA-seq experiment, reserving

255    the measured query protein expression as a ground truth for validation. We accurately predicted the

256    surface protein expression of each query cell using the 50-NN average from the reference cells in the

257    harmonized embedding. For all proteins, we found strong concordance between predicted and (50-NN

258    smoothed) measured expression (Pearson r: 0.88-0.99, **Fig. 5c,d**). For all but three proteins, we

259    achieved comparable results with as few as 5 or 10 nearest neighbors (**Fig. S7a**).

260    We note that it is also possible to conduct the same analysis with a unimodal PCA-based reference

261    built from the cells' mRNA expression only. This approach has slightly worse performance for some

262    proteins (Pearson r: 0.65-0.97, **Fig. S7b-d**), demonstrating that a reference built jointly on both mRNA

263    and protein permits better inference of protein expression than an mRNA-only reference, which is

264    consistent with previous observations that mRNA expression is not fully representative of protein

265    expression[41,42]. This analysis highlights how users can start with a low-dimensional embedding other

266    than PCA, such as CCA, to better capture rich multimodal information in the reference.

267    # Discussion

268    We frame reference mapping as a specialized case of integration, between one dataset and a second

269    larger, more comprehensive, and previously integrated dataset. Because the reference is already

270    integrated, it is natural to use the same mathematical framework from the integration to perform

271    mapping. For instance, the scArches[26] algorithm uses an autoencoder-based framework to map to

272    references built with autoencoder-based integration algorithms trVAE[46] and CVAE. Similarly, Symphony

273    uses the mixture modeling framework to map to references built with Harmony mixture modeling

274    integration. Symphony compresses the reference by extracting relevant reference-derived parameters

275    from the mixture model to map query cells in seconds. With this compression, references can be

276    distributed without the need to share raw expression data or donor-level metadata, which enables data

277    privacy[28–32]. Symphony compression greatly reduces the size of a reference dataset: for the memory T

11

278    cell dataset of 500,089 cells, the raw expression matrix is 8.9 GB, whereas the Symphony minimal

279    reference elements are 1.3 MB.

280    Useful reference atlases contain annotations not present in the query, such as cell type labels (**Fig. 3**),

281    trajectory coordinates (**Fig. 4**), or multimodal measurements (**Fig. 5**). Transfer of these annotations

282    from reference to query is an open area of research that includes algorithms for automated cell type

283    classification[47]. We approach annotation transfer in two steps. We first learn a predictive model in the

284    reference embedding, then map query cells and use their reference coordinates to predict query

285    annotations. In this two-step approach, Symphony mapping provides a feature space but is otherwise

286    independent from the choice of downstream inference model. In PBMC type prediction (**Fig. S3**), we

287    used Symphony embeddings to train multiple competitive classifiers: k-NN, SVM, and logistic

288    regression. In our analyses, we were encouraged to find that a simple k-NN classifier can achieve high

289    performance with only 5-10 neighbors. In practice, users can choose more complex inference models if

290    it is warranted for certain annotation types. Moreover, we expect prediction results to improve with more

291    accurate and reproducible annotation methods, such as consistent cell type taxonomies provided by

292    the Cell Ontology[48] project and better modeling of multimodal expression data[17].

293    We defined three conditions under which Symphony and *de novo* integration with Harmony yield

294    equivalent results. In subsequent examples, we showed that Symphony still performs well when the last

295    two conditions are relaxed. The pancreas query contains more cells than its reference (**condition II**),

296    while the liver hematopoiesis reference and query overlap in donors (**condition III**). Condition I, which

297    requires comprehensive cell type coverage in the reference, is less flexible. When the query contains a

298    brand new cell type, it will be aligned to its most transcriptionally similar reference cluster. Note that

299    condition I only pertains to cell types and not clinical and biological contexts. For instance, we

300    successfully mapped mouse pancreas query to an entirely human pancreas reference (**Fig. 3**),

301    because the same pancreatic cell types are shared in both species. Mapping novel cell types is a

302    current limitation and important direction for future work. For now, we advise users interested in novel

303    cell type discovery to supplement a Symphony analysis with *de novo* analyses of the query alone.

12

304    Instead of one monolithic reference, we expect the proliferation of multiple, well-annotated specialized

305    references. For instance, the memory T cell reference (**Fig. 5**) will be useful to annotate fine-grained T

306    cell states, while an unsorted PBMC reference (**Fig. 2**) would better suit annotation of more diverse

307    immune populations. Similarly, a reference with only healthy individuals is useful for annotation of cell

308    types, while a reference with both healthy and diseased individuals is useful for annotation of cell types

309    and pathological cell states.

310    As large-scale tissue and whole-organism single-cell reference atlases become available in the near

311    future, Symphony will enable investigators to leverage the rich information in these references to

312    perform integrative analyses and rapidly transfer reference coordinates and diverse annotations to new

313    datasets.

13

314    <u>Methods</u>

315    **1. Symphony**

316    1.1 Symphony overview

317    The goal of single-cell reference mapping is to embed newly assayed query cells into an existing

318    comprehensive reference atlas, facilitating the automated transfer of annotations from the reference to

319    the query. The optimal mapping method needs to be able to operate at various levels of resolution,

320    capture continuous intermediate cell states, and scale to multimillion cells[17]. Consider a scenario in

321    which we wish to map a query of $m$ cells against reference datasets with $n$ cells, where $m<<n$.

322    Unsupervised integration of measurements across donors, studies, and technological platforms is the

323    standard way to compare single cell datasets and identify cell types. Hence, a "gold standard"

324    reference mapping strategy might be to run Harmony integration on all $m+n$ cells *de novo*. However,

325    this approach is impractical because it is cumbersome and time-intensive to process all the cell-level

326    data for the reference datasets every time a user wishes to reharmonize it with a query. Instead, we

327    envision a pipeline where a reference atlas need only be carefully constructed and integrated once, and

328    all subsequent queries can be rapidly mapped into the same stable reference embedding.

329    Symphony is a reference mapping method that efficiently places query cells in their precise location

330    within an integrated low-dimensional embedding of reference cells, approximating *de novo*

331    harmonization without the need to reintegrate the reference cells. Symphony is comprised of two

332    algorithms: reference compression and mapping. Expanding upon the linear mixture model framework

333    introduced in Harmony[18], Symphony compression takes in an integrated reference and faithfully

334    compresses it by capturing the components of the model into efficient data structures. The output of

335    reference compression is the minimal set of elements needed for mapping **(Fig. S1b)**. The Symphony

336    mapping algorithm takes as input a new query dataset as well as minimal reference elements and

337    returns the appropriate locations of the query cells within the integrated embedding **(Fig. S1c)**.

14

338 Once a harmonized reference is constructed and compressed using Symphony, subsequent mapping

339 of query cells executes within seconds (**Fig. S4**). Efficient implementations of Symphony are available

340 as part of an R package at https://github.com/immunogenomics/symphony, along with several

341 precomputed references constructed from public scRNA-seq datasets. The following sections introduce

342 the Symphony model, then describes Symphony compression and mapping in terms of the underlying

343 data structures and algorithms. We also provide **Supplementary Equations** containing more detailed

344 derivations for reference compression terms.

345 *Glossary*

346 We define all symbols for data structures used in the discussion of Symphony below, including their

347 dimensions and possible values. Dimensions are in terms of the following parameters:

348    • *n:* the number of reference cells

349    • *m:* the number of query cells

350    • *N:* the total number of cells (*n* + *m*)

351    • *g:* the number of genes in the reference after any gene selection

352    • *d:* the dimensionality of the embedding (e.g. PCs). *d* applies to both reference and query.

353    • *b:* the number of batches in the reference

354    • *c:* the number of batches in the query

355    • *k:* the number of clusters in the mixture model for reference integration (representing latent cell

356      states)

357 **Reference-related symbols:**

| | |
|---|---|
| $G_r \in \mathbb{R}^{g \times n}$ | Input reference gene expression matrix, prior to scaling. |
| $G_{rs} \in \mathbb{R}^{g \times n}$ | Scaled reference gene expression matrix. |
| $X_r \in \{0,1\}^{b \times n}$ | One-hot design matrix assigning reference cells (columns) to batches (rows). |
| $X_r' \in \{0\}^{c \times n}$ | Zero matrix assigning reference cells (columns) to *query* batches (rows). All values are 0 because reference cells do not belong to query batches. This term is used in the derivation for the reference compression terms. |

| | |
|---|---|
| $\mu \in \mathbb{R}^{g \times 1}$ | Reference gene means used to center each gene for PCA. |
| $\sigma \in \mathbb{R}^{g \times 1}$ | Reference gene standard deviations used to scale each gene for PCA. |
| $U \in \mathbb{R}^{g \times d}$ | Gene loadings from the original PCA (before Harmony integration). |
| $Z_r \in \mathbb{R}^{d \times n}$ | Original (non-harmonized) PC embedding for reference cells. |
| $\hat{Z}_r \in \mathbb{R}^{d \times n}$ | Integrated embedding for reference cells in harmonized PC (hPC) space, as output by Harmony. |
| $R_r \in [0,1]^{k \times n}$ | Soft cluster assignment of reference cells (columns) to clusters (rows), as output by Harmony. Each column is a probability distribution that sums to 1. |
| $Y_{cos} \in \mathbb{R}^{d \times k}$ | Cluster centroid locations in the harmonized embedding, L2 normalized. |
| $B_r \in \mathbb{R}^{k \times (1+b) \times d}$ | 3D tensor of the estimated parameters (betas and intercepts) of the linear mixture model for each of $k$ clusters for the reference cells. |
| $N_r \in \mathbb{R}^{k \times 1}$ | First reference compression term. Vector containing the size of each of the $k$ clusters, effectively the number of reference cells contained within them. |
| $C \in \mathbb{R}^{k \times d}$ | Second reference compression term. |
| $Ref = \{\mu, \sigma, U, Y_{cos}, N_r, C\}$ | Symphony minimal reference elements comprising $\mu, \sigma, U, Y_{cos}, N_r, C$. |

358 **Query-related symbols:**

| | |
|---|---|
| $G_q \in \mathbb{R}^{g \times m}$ | Input query gene expression matrix, prior to scaling. |
| $G_{qs} \in \mathbb{R}^{g \times m}$ | Query gene expression matrix, scaled by *reference* gene means $\mu$ and standard deviations $\sigma$. |
| $X_q \in \{0,1\}^{c \times m}$ | Design matrix assigning query cells (columns) to query batches (rows). |
| $Z_q \in \mathbb{R}^{d \times m}$ | Query cell locations in original (non-harmonized) PC embedding. |
| $\hat{Z}_q \in \mathbb{R}^{d \times m}$ | Approximate query cell locations in integrated embedding (hPC space). Output of Symphony reference mapping. |
| $R_q \in [0,1]^{k \times m}$ | Soft cluster assignment of query cells (columns) to clusters (rows). Each column is a probability distribution that sums to 1. |
| $B_q \in \mathbb{R}^{k \times (1+c) \times d}$ | 3D tensor of the estimated parameters (betas and intercepts) of the linear mixture model for each of $k$ clusters. |

359 1.2 Symphony model and conditions for equivalence to Harmony integration

16

360  Symphony and Harmony both use a linear mixture model framework, but the two methods perform

361  different tasks: Harmony integrates a reference, whereas Symphony compresses the reference and

362  enables efficient query mapping. To motivate the Symphony model, it is helpful to first briefly review the

363  mixture model, which serves as the basis. Harmony integrates scRNA-seq datasets across batches

364  (e.g. multiple donors, technologies, studies) and projects the cells into a harmonized embedding where

365  cells cluster by cell type rather than batch-specific effects. Harmony takes as input a low-dimensional

366  embedding of cells ($Z$) and design matrix with assignments to batches ($X$) and outputs a harmonized

367  embedding ($\hat{Z}$) with batch effects removed. Briefly, Harmony works by iterating between two

368  subroutines—maximum diversity clustering and linear mixture model correction—until convergence. In

369  the clustering step, cells are probabilistically assigned to soft clusters with a variant of soft $k$-means with

370  a diversity penalty favoring clusters represented by multiple datasets rather than single datasets. In the

371  correction step, each cluster learns a cluster-specific linear model that explains cell locations in PC

372  space as a function of a cluster-specific intercept and batch membership. Then, cells are corrected by

373  cell-specific linear factors weighted by cluster membership to remove batch-dependent effects. The full

374  algorithm and implementation are detailed in Korsunsky et al. (2019)[47].

375  In the scenario of mapping $m$ query cells against $n$ reference cells, the *de novo* integration strategy

376  would model all cells as in (1), where the $H$ subscript denotes the Harmony solution, in contrast to the

377  Symphony model which is presented in (2). Let $X_H \in \{0,1\}^{(c+b)\times(m+n)}$ represent the one-hot encoded

378  design matrix assigning all cells across batches. $X_H^*$ denotes $X_H$ augmented with a row of 1s for the

379  batch-independent intercept term: $X_H^* = 1||X_H$. The intercept terms represent cluster centroids (location

380  of "experts" in the mixture of experts model). $Z_H$ represents the low-dimensional PCA embedding of all

381  cells. $R_H$ represents the probabilistic assignment of cells across $k$ clusters, and $diag(R_{Hk}) \in \mathbb{R}^{N\times N}$

382  denotes the diagonalized $k$th row of $R_H$. For each cluster $k$, the parameters of the linear mixture model

383  $B_k \in \mathbb{R}^{(1+c+b)\times d}$ can therefore be solved for as in (1), using ridge regression with ridge penalty

384  hyperparameter $\lambda$. Note that we do not penalize the batch-independent intercept term: $\lambda_0 = 0$,

385  $\forall_{a\in[1:(c+b)]}\lambda_a = 1$.

386 ### *De novo* **Harmony model:**

$$B_k = (X_H^* \, diag(R_{Hk}) X_H^{*T} + \lambda I)^{-1} X_H^* diag(R_{Hk}) \, Z_H^T \tag{1}$$

387 The goal of Symphony mapping is to add new query cells to the model in order to estimate and remove

388 the query batch effects. Symphony mapping approximates *de novo* Harmony integration on all cells,

389 except the reference cell positions in the harmonized embedding do not change. In order for Symphony

390 mapping to be equivalent to *de novo* Harmony, several conditions must be met:

391     I.   All cell states represented in the query dataset are captured by the reference datasets—i.e.

392        there are no completely novel cell types in the query.

393     II.   The number of reference cells is much larger than the query ($m$<<$n$).

394     III.   The query dataset is obtained independent of the reference datasets—i.e. the reference

395        batch design matrix ($X_r$) has no interaction with the query batch design matrix ($X_q$).

396 We consider these to be fair assumptions for large-scale reference atlases, allowing Symphony to

397 make three key approximations:

398   (1) With a large reference, the reference-only PCs approximate the PCs for the combined reference

399     and query datasets. This allows us to project the query cells into the pre-harmonized reference

400     PCA space using the reference gene loadings ($U$).

401   (2) The cluster centroids ($Y$) for the integrated reference cells approximate the cluster centroids

402     from harmonizing all cells.

403   (3) The reference cell cluster assignments ($R_r$) remains approximately stable with the addition of

404     query cells.

405 Given these approximations, we can thereby harmonize the reference cells *a priori* and save the

406 reference-dependent portions of the Harmony mixture model (**Supplementary Equations**). In

407 Symphony, we model the reference cells as already harmonized with batch effects removed, so we can

408 thereafter ignore the reference design matrix structure. The Symphony design matrix $X \in [0,1]^{c \times N}$

409 assigns all cells (reference and query) to *query* batches only. $X^*$ denotes $X$ augmented with a row of 1s

18

410  $(X^*_{[0,\cdot]})$ corresponding to the batch-independent intercepts (we model the intercepts for all cells). The

411  remaining $c$ rows $(X^*_{[1:c,\cdot]})$ represent the one-hot batch assignment of the cells among the $c$ query

412  batches. Note that for the reference cell columns, these values are all 0 since the reference cells do not

413  belong to any *query* batches. The parameters $(B_{qk} \in \mathbb{R}^{(1+c)\times d})$ of the model for each cluster $k$ can

414  then be solved for as in (2). Similar to Harmony, we use ridge regression penalizing the non-intercept

415  terms, where $\lambda_0 = 0, \forall_{a\in[1:c]}\lambda_a = 1$.

**Symphony model:**

$$B_{qk} \approx (X^* \, diag(R_k) \, X^{*T} + \lambda I)^{-1} X^* \, diag(R_k) \, Z^T \qquad (2)$$

417  The matrix $R \in \mathbb{R}^{k\times N}$ denotes the assignment of query and reference cells (columns) across the

418  reference clusters (rows). $Z \in \mathbb{R}^{d\times N}$ denotes the horizontal matrix concatenation of the uncorrected

419  query cells in original PC space $(Z_q)$ and corrected reference cells in harmonized space $(\hat{Z}_r)$. For each

420  cluster $k$, let matrix $B_{qk} \in \mathbb{R}^{(1+c)\times d}$ represent the query parameters to be estimated. The first row of

421  $B_{qk}$ represents the batch-independent intercept terms, and the remaining $c$ rows of $B_{qk}$ represent the

422  query batch-dependent coefficients, which can be regressed out to harmonize the query cells with the

423  reference. Note that the intercept terms from Symphony mapping should equal the cluster centroid

424  locations from the integrated reference since the harmonized reference cells are modeled only by a

425  weighted average of the centroid locations for the clusters over which it belongs (and a cell-specific

426  residual). Hence, the reference cell positions should not change when removing query batch effects.

427  The matrices $X^*$, $R_k$, and $Z$ in (2) can be partitioned into query and reference-dependent portions. In the

428  **Supplementary Equations**, we show in detail how the reference-dependent portions can be further

429  simplified into a $k$ x 1 vector and $k$ x $d$ matrix ($N_r$ and C), which we call "reference compression terms."

430  Intuitively, the vector $N_r$ contains the size (in cells) of each reference cluster. The matrix $C = R_r \hat{Z}_r^T$ does

431  not have as intuitive an explanation but follows from the derivation (**Supplementary Equations**). These

432  terms can be computed at the time of reference building and saved as part of the minimal reference

433  elements to reduce the necessary computations during mapping.

19

## 1.3 Reference building and compression

434

435    Reference compression is the key idea that allows for the efficient mapping of new query cells onto the

436    harmonized reference embedding without the need to reintegrate all cells. To construct a Symphony

437    reference with minimal elements needed for mapping, reference cells are first harmonized in a low-

438    dimensional space (e.g. PCs) to remove batch-dependent effects. Symphony then compresses the

439    Harmony mixture model components to be saved for subsequent query mapping.

440    **Data structures**

441    Symphony takes as input a gene expression matrix for reference cells ($G_r$) and corresponding one-hot-

442    encoded design matrix ($X_r$) containing metadata about assignment of cells to batches. It outputs a set

443    of data structures, referred to as the Symphony minimal reference elements, that captures key

444    information about the reference embedding that can be subsequently used to efficiently map previously

445    unseen query cells (**Algorithm 1**). These components include the gene mean ($\mu$) and standard

446    deviation ($\sigma$) used to scale the genes, the PCA gene loadings ($U$), the final L2-normalized cluster

447    centroid locations ($Y_{cos}$), and precomputed values which we call the "reference compression terms" ($N_r$

448    and $C$) that expedite the correction step of query mapping (**Supplementary Equations**). These

449    elements are a subset of the components available once Harmony integration is applied to the

450    reference cells. Note that other input embeddings, such as canonical correlation analysis (CCA), may

451    be used in place of PCA as long as the gene loadings to perform query projection into those

452    coordinates are saved.

453    **Table 1** lists the Symphony minimal reference elements required to perform mapping. **Table 2** shows

454    additional components of a "full" Harmony reference that are not included in the Symphony reference

455    elements. Importantly, the dimensions of the Symphony data structures do not require information on

456    the *n* individual reference cells and hence do not scale with the raw number of reference cells. Rather

457    the components scale with the biological complexity captured (i.e. number of clusters *k* and

458    dimensionality of embedding *d*). Conversely, the Harmony data structures store information on a per-

459    cell basis (*n*). Note that in practice the integrated embedding of reference cells ($\hat{Z}_r$) listed in **Table 2** is

20

460    needed to perform downstream transfer of annotations from reference to query cells (e.g. k-NN), but it

461    is not required during any computations of the mapping step.

462    **Table 1: Symphony minimal reference elements**

| | |
|---|---|
| $\mu \in \mathbb{R}^{g \times 1}$ | Reference gene means used to center each gene for PCA. |
| $\sigma \in \mathbb{R}^{g \times 1}$ | Reference gene standard deviations used to scale each gene for PCA. |
| $U \in \mathbb{R}^{g \times d}$ | Gene loadings to project from expression to PCA (or CCA) space |
| $Y_{cos} \in \mathbb{R}^{d \times k}$ | Cluster centroid locations in harmonized PC space, L2 normalized. |
| $N_r \in \mathbb{R}^{k \times 1}$ | First reference compression term. Vector containing the size of each of the $k$ clusters, effectively the number of reference cells contained within them. |
| $C \in \mathbb{R}^{k \times d}$ | Second reference compression term. |

463

464    **Table 2: Additional components of Harmony reference**

| | |
|---|---|
| $G_r \in \mathbb{R}^{g \times n}$ | Input reference gene expression matrix, prior to scaling. |
| $X_r \in \{0, 1\}^{b \times n}$ | Design matrix assigning reference cells (columns) to reference batches (rows). |
| $B_r \in \mathbb{R}^{k \times (1+b) \times d}$ | 3D tensor of the estimated parameters (betas and intercepts) of the linear mixture model for each of $k$ clusters for the reference cells. |
| $\hat{Z}_r \in \mathbb{R}^{d \times n}$ | Integrated embedding for reference cells in harmonized PC ("hPC") space, as output by Harmony. |
| $R_r \in [0, 1]^{k \times n}$ | Soft cluster assignment of reference cells (columns) to clusters (rows), as output by Harmony. Each column is a probability distribution that sums to 1. |

465

466    **Algorithm**

467    Starting from reference cell gene expression, we first perform within-cell library size normalization (if not

468    already done) and variable gene selection to obtain $G_r$, scaling of the genes to have mean 0 and

469    variance 1 (saving $\mu$ and $\sigma$ for each gene), and PCA to embed the reference cells in a low-dimensional

470    space, saving the gene loadings ($U$) (**Implementation Details**). Then, the PCA embedding ($Z_r$) and

471    batch design matrix ($X_r$) are used as input to Harmony integration to harmonize over batch-dependent

472    sources of variation. Given the resulting harmonized embedding ($\hat{Z}_r$) and final soft assignment of

21

473 reference cells to clusters ($R_r$), the locations of the final reference cluster centroids $Y \in \mathbb{R}^{d \times k}$ can be

474 calculated as in (3) and saved.

$$Y = \hat{Z}_r R_r^T \tag{3}$$

475 Symphony then computes the reference compression terms $N_r$ (intuitively, the number of cells per

476 cluster) and $C$, which does not have an intuitive explanation but can be directly computed as $C = R_r \hat{Z}_r^T$.

477 Refer to the **Supplementary Equations** for a complete mathematical derivation of the compression

478 terms. Symphony reference building ultimately returns the minimal reference elements: $\mu, \sigma, U, Y_{cos}, N_r$,

479 and $C$ (**Fig. S1a**).

---

480 **Algorithm 1** Build Symphony reference

---

481     **function** BUILDREFERENCE($G_r, X_r$)

482         $\mu, \sigma, G_{rs} \leftarrow$ **SCALE**($G_r$)

483         $U, Z_r \leftarrow$ **PCA**($G_{rs}$)

484         $\hat{Z}_r, R_r \leftarrow$ **HARMONIZE**($Z_r, X_r$)

485         $Y \leftarrow \hat{Z}_r R_r^T$

486         $Y_{cos} \leftarrow {Y_{[\cdot,i]}} \Big/ {\left\|Y_{[\cdot,i]}\right\|_2}$            $\triangleright$ $L_2$ *normalize cluster centroids*

487         $N_r \leftarrow rowSums(R_r)$            $\triangleright$ *First compression term*

488         $C \leftarrow R_r \hat{Z}_r^T$            $\triangleright$ *Second compression term*

489         $Ref \leftarrow (\mu, \sigma, U, Y_{cos}, N_r, C)$

490         **return** $Ref$            $\triangleright$ *Return minimal reference elements*

---

491

492 ## 1.4 Symphony mapping

493 The Symphony mapping algorithm localizes new query cells to their appropriate locations in the

494 harmonized embedding without the need to run integration on the reference and query cells altogether.

495 The joint embedding of reference and query cells can be used for downstream analyses, such as

496 transferring cell type annotations from the reference cells to the query cells.

22

**Data structures**

498 Symphony mapping takes as input the gene expression matrix for query cells ($G_q$), query design matrix

499 assigning query cells to batches ($X_q$), and the precomputed minimal elements for a reference ($Ref$). It

500 outputs a query object containing the locations of query cells in the integrated reference embedding

501 ($\hat{Z}_q$; **Algorithm 2**). **Table 3** lists the components of the query object that is returned by Symphony.

502 **Table 3: Components of Symphony query**

| $G_q \in \mathbb{R}^{g \times m}$ | Input query gene expression matrix, prior to scaling. |
|---|---|
| $X_q \in \{0,1\}^{c \times m}$ | Design matrix assigning query cells (columns) to query batches (rows). |
| $Z_q \in \mathbb{R}^{d \times m}$ | Query cell locations in original (non-harmonized) PC embedding. |
| $\hat{Z}_q \in \mathbb{R}^{d \times m}$ | Approximate query cell locations in integrated embedding (hPC space). |
| $R_q \in [0,1]^{k \times m}$ | Soft cluster assignment of query cells (columns) to clusters (rows). Each column is a probability distribution that sums to 1. |
| $B_q \in \mathbb{R}^{k \times (1+c) \times d}$ | 3D tensor of the estimated parameters (betas and intercepts) of the linear mixture model for each of $k$ clusters. |

503

504 **Algorithm**

505 The input to the query mapping procedure is a gene expression matrix ($G_q$) and design matrix ($X_q$) for

506 query cells, and the output is the locations of the cells in the harmonized embedding ($\hat{Z}_q$). At a high

507 level, the mapping algorithm first projects the query cells into the original, non-harmonized PC space as

508 the reference cells using the reference gene loadings ($U$) and assigns probabilistic cluster membership

509 across the reference cluster centroid locations. Then, the query cells are modeled using the Symphony

510 mixture model and corrected to their approximate locations in the integrated embedding by regressing

511 out the query batch-dependent effects (**Algorithm 2**).

512 ***Projection of query cells into pre-harmonized PC Space***

513 Symphony projects the query cells into the same original PCs ($Z_r$) as the reference. Symphony

514 assumes that, given a much smaller query compared to the reference ($m<<n$), the PCs will remain

23

515    approximately stable with the addition of query cells. To project the query cells, we first subset the

516    query expression data by the same variable genes used in reference building and scale the normalized

517    expression of each gene by the same mean and standard deviations used to scale the reference cells

518    $(\mu, \sigma)$. Let $G_{qs}$ denote the query gene expression matrix scaled by the reference gene means and

519    standard deviations. We can then use the reference gene loadings $(U)$ to project $G_{qs}$ into reference PC

520    space. In (4), $Z_q \in \mathbb{R}^{d \times m}$ denotes the PC embedding for the query cells. Note that if an alternate

521    starting embedding (e.g. CCA) is used instead of PCA, the gene loadings must be saved to enable this

522    query projection step.

$$Z_q = U^T G_{qs} = \Sigma_q V_q^T \tag{4}$$

523    ***Soft assignment across reference clusters***

524    Once the query cells are projected into PC space, we soft assign the cells to the reference clusters

525    using the saved reference centroid locations $(Y_{cos})$. Symphony assumes that the reference cluster

526    centroid locations remain approximately stable with the addition of a much smaller query dataset since

527    the query contains no novel cell types. Under these conditions, we use a previously published objective

528    function for soft $k$-means clustering (5), which includes a distance term and an entropy regularization

529    term over $R$ weighted by hyperparameter $\sigma$. This is the same objective function as the clustering step of

530    Harmony, except it does not include the diversity penalty term. In Harmony, the purpose of the diversity

531    term is to penalize clusters that are only represented by one or a few datasets (suggesting they do not

532    represent true cell types). In contrast, Symphony does not require the use of a diversity penalty

533    because the reference centroids have already been established. Furthermore, the query cell types can

534    comprise a subset of a larger set of reference cell types, and therefore not all clusters are necessarily

535    expected to be represented in the query. We can solve for $R_q$, the optimal probabilistic assignment for

536    query cells across each of the $k$ reference clusters (**Implementation Details**).

$$\min_{R,Y} \sum_{i,k} R_{ki} \|Z_i - Y_k\|^2 + \sigma R_{ki} \log R_{ki} \tag{5}$$

24

$$\text{s.t. } \forall_i \forall_k R_{ki} > 0, \forall_i \sum_{k=1}^{K} R_{ki} = 1$$

### Mixture of experts correction

537

538 The final step in Symphony mapping is to model then remove the query batch effects to obtain $\hat{Z}_q$, the

539 approximate location of query cells in the harmonized reference embedding. In equation (2), we

540 modeled the reference and query cells together and wish to solve for the query parameters $B_{qk} \in$

541 $\mathbb{R}^{(1+c) \times d}$ for each cluster $k$. The reference-dependent terms in (2) were previously computed and

542 saved in compressed form ($N_r$ and $C$). With $R_q$ and $Z_q$ calculated from query cell projection and

543 clustering, we can finally solve for $B_{qk}$. Similar to the correction step of Harmony, we obtain cell-specific

544 correction values for the query cells by removing the batch-dependent terms captured in $B_{qk[1:c,\cdot]}$. Note

545 that the reference batch terms are neither modeled nor corrected during reference mapping, so the

546 harmonized reference cells do not move.

547 The final locations of the query cells in the harmonized embedding are estimated by iterating over all $k$

548 clusters and subtracting out the non-intercept batch terms for each cell weighted by cluster membership

549 (6). Intuitively, the query centroids are moved so that they overlap perfectly with the reference centroids

550 in the harmonized embedding. $\hat{Z}_{q[i]}$ denotes the approximate location in harmonized PC space for

551 query cell $i$.

$$Z_{q[i]} = \sum_k R_{q[k,i]} \left[ B_{qk[0,\cdot]}^T + B_{qk[1:c,i]}^T X_q \right] + \varepsilon$$

$$\hat{Z}_{q[i]} = Z_{q[i]} - \sum_k R_{q[k,i]} B_{qk[1:c,\cdot]}^T X_q \tag{6}$$

$$\hat{Z}_{q[i]} = \sum_k R_{q[k,i]} B_{qk[0,\cdot]}^T + \varepsilon$$

---

552 **Algorithm 2** Map query cells onto reference

---

553 **function** QUERYMAPPING($G_q, X_q, Ref$)

554 $\quad G_{qs} \leftarrow$ **SCALE**($G_q, Ref\$\mu, Ref\$\sigma$) $\qquad \triangleright \$ \text{ denotes accessing a component of } Ref$

555      $Z_q \leftarrow$ **PCAPROJECTION**$(G_{qs}, Ref\$U)$

556      $R_q \leftarrow$ **CLUSTER**$(Z_q, Ref\$Y_{cos})$

557      $\hat{Z}_q \leftarrow Z_q$

558      **for** $k \leftarrow 1 \dots k$ **do**

559        $E \leftarrow X_q^* R_q^{(k)} X_q^{*T}$           ▷ $X_q^*$: query design matrix augmented with row of 1s

560        $E_{[0,0]} \leftarrow E_{[0,0]} + Ref\$N_{r(k)}$

561        $F \leftarrow X_q^* R_q^{(k)} Z_q^T$

562        $F_{[0,\cdot]} \leftarrow F_{[0,\cdot]} + Ref\$C_{[k,\cdot]}$

563        $B_{qk} \leftarrow (E + \lambda I)^{-1}(F)$

564        $B_{qk[0,\cdot]} \leftarrow 0$           ▷ Do not correct the intercept terms

565        $\hat{Z}_q \leftarrow \hat{Z}_q - B_{qk}^T X_q^* R_q^{(k)}$

566      **return** $\hat{Z}_q$           ▷ Return query locations

567

## 568   1.5 Implementation details

### 569   **Reference building and compression**

#### 570   *Variable gene selection and scaling*

571   Starting with the gene expression matrix for reference cells, we perform log(CP10K) library size

572   normalization of the cells (if not already done), subset by the top $g$ variable genes by the vst method

573   (as provided in Seurat[49]), which fits a line to the log(variance) and log(mean) relationship using local

574   polynomial regression, then standardizes the features by observed mean and expected variance,

575   calculating gene variance on the standardized values, which is re-implemented as a standalone

576   function at https://github.com/immunogenomics/singlecellmethods. The data is scaled such that the

577   expression of each gene has a mean expression of 0 and variance of 1 across all cells.

#### 578   *PCA*

579   We perform dimensionality reduction on the scaled gene expression $G_{rs}$ using principal component

580   analysis (PCA). PCA projects the data a low-dimensional, orthonormal embedding that retains most of

581     the variation of gene expression in the dataset. Singular value decomposition (SVD) is a matrix

582     factorization method that can calculate the PCs for a dataset. Here, we use SVD (irlba package in R[48])

583     to perform PCA. SVD states that matrix $G_{rs}$ with dimensions $g \times n$ can be factorized as:

$$G_{rs} = U \Sigma V^T \qquad\qquad (7)$$

584     In (7), $\Sigma V^T = Z_r$ (dimensions $d \times n$) represents the embedding of reference cells in PC space, after

585     truncating the matrix on the first $d$ (by default, $d = 20$) PCs. The gene loadings ($U \in \mathbb{R}^{g \times d}$) are saved.

586     Note that an alternative embedding, such as canonical correlation analysis (CCA) may be used in place

587     of PCA, as long as the gene loadings are saved.

588     ***Harmony integration***

589     The PCA embedding ($Z_r$) is then input to Harmony for dataset integration. By default, Symphony uses

590     the default parameters for the cluster diversity enforcement ($\theta = 2$), the entropy regularization

591     hyperparameter for soft *k*-means ($\sigma = 0.1$), and the number of clusters $k = \min\left(100, \frac{n}{30}\right)$. We save the

592     L2-normalized cluster centroid locations $Y_{cos}$ to the reference object since query mapping employs a

593     cosine distance metric. If the reference has a single-level batch structure, no integration is performed,

594     and the clusters are defined using soft k-means.

595     **Query mapping**

596     ***Normalization and scaling***

597     The gene expression for query cells are assumed to be library size normalized in the same manner that

598     was used to normalize the reference cells (e.g. log(CP10K)). During scaling, the query data is subset

599     by the same variable genes from the reference datasets, and query gene expression is scaled by the

600     *reference* gene means and standard deviations. Any genes present in the query but not the reference

601     are ignored, and any genes present in the reference but not the query have scaled expression set to 0.

602     ***Clustering step uses cosine distance***

603    As in Harmony, in practice we use cosine distance rather than Euclidean distance in the clustering step.

604    For the computation of the distance term, we L2-normalize the columns (cells) of $Z$ and columns

605    (centroids) of $Y_k$ such that the squared values sum to 1 across each column. Let the terms $Z_{q\_cos\,[\cdot,i]}$ and

606    $Y_{\cos\,[\cdot,k]}$ represent the L2-normalized locations of query cell $i$ and the reference centroid for cluster $k$ in

607    PC space, respectively. We compute the cosine distance between the cells and centroids. Since all

608    $Z_{q\_cos\,[\cdot,i]}$ and $Y_{\cos\,[\cdot,k]}$ each have unity norm, the squared Euclidean distance $\left\| Z_{q\_cos\,[\cdot,i]} - Y_{\cos\,[\cdot,k]} \right\|^2$ is

609    equivalent to the cosine distance $2\big(1 - \cos(Y_{\cos\,[\cdot,k]}, Z_{q\_cos\,[\cdot,i]})\big) = 2(1 - Y^T_{\cos\,[k,\cdot]} Z_{q\_cos\,[\cdot,i]})$. Therefore, the

610    objective function for query assignment to centroids becomes:

$$\min_{R,Y} \sum_{i,k} 2R_{q[k,i]}(1 - Y^T_{\cos\,[k,\cdot]} Z_{q\_cos\,[\cdot,i]}) + \sigma\, R_{q[k,i]} \log R_{q[k,i]} \tag{8}$$

$$\text{s.t. } \forall_i \forall_k R_{q[k,i]} > 0, \forall_i \sum_{k=1}^{K} R_{q[k,i]} = 1$$

611    We can solve the optimization problem using an expectation-maximization framework. Following the

612    same strategy as Korsunsky et al. (2019), we calculate $R_i$, the optimal probabilistic assignment for each

613    query cell $i$ across each of the $k$ reference clusters. In (9), we can interpret $R_{q[k,i]}$ as the probability that

614    query cell $i$ belongs to cluster $k$. The denominator term simply ensures that for any given cell $i$, the

615    probabilities across all $k$ clusters sum to one.

$$R_{q(k,i)} = \frac{\exp\left(-\dfrac{2}{\sigma}(1 - Y^T_{\cos\,[k,\cdot]} Z_{q\_cos\,[\cdot,i]})\right)}{\sum_{k=1}^{K} \exp\left(-\dfrac{2}{\sigma}(1 - Y^T_{\cos\,[k,\cdot]} Z_{q\_cos\,[\cdot,i]})\right)} \tag{9}$$

## 2. Analysis details

617    2.1 10x PBMCs and pancreas examples

618    ***Preprocessing scRNA-seq data***

619    As the three 10x PBMCs and reference pancreas datasets were previously preprocessed by our group

620    as part of the Harmony publication, we used the same log(CP10K) normalized expression data, filtered

621    as described in Korsunsky et al. (2019)[50]. The PBMCs consist of cells from three technologies: 3'v1

622    (n=4,808 cells), 3'v2 (8,372 cells), and 5' (7,612 cells). The pancreas reference datasets were each

623    sequenced with a different technology: Fluidigm C1 (n=638 cells), CEL-seq (946 cells), CEL-seq2

624    (2,238 cells), Smart-seq2 (2,355 cells). The pancreas query dataset (inDrop, n=8,569 human and 1,886

625    mouse cells) along with author-defined cell type labels were downloaded from https://hemberg-

626    lab.github.io/scRNA.seq.datasets/human/pancreas/.

627    ***Constructing the pancreas query with mouse and human***

628    For the pancreas query (Baron et al., 2016), we downloaded both the human and mouse expression

629    matrices. In order to combine the two matrices into a single aggregated query, we "humanized" the

630    mouse expression matrix by mapping mouse genes to their orthologous human genes. This mapping

631    was computed using the biomaRt R package[47], mapping `mgi_symbol` from the

632    `mmusculus_gene_ensembl` database to `hgnc_symbol` from the `hsapien_gene_ensembl`

633    database. We represented this map as a matrix, with mouse genes as rows, human genes as columns,

634    and values in {0,1} assigned to denote whether a mouse gene maps to a human gene. We then

635    normalized the matrix to have each column sum to one, effectively creating a count-preserving

636    probabilistic map from d mouse to D human genes $M \in R^{D \times d}$. Mapping from mouse to human genes is

637    then performed with matrix multiplication: $U_{human} = MU_{mouse}$. Note that while the mouse gene expression

638    matrix $U_{mouse}$ contains only integers ($U_{mouse} \in Z^{d \times N}$), the many-to-many mapping means that the mapped

639    human gene expression matrix $U_{human}$ may contain non-integers ($U_{human} \in R^{D \times N}$). For any human

640    orthologs that were missing in the mouse expression data, we filled in the expression with zeroes. The

641    mouse mapping was based on 2,140 overlapping genes (of the 3,000 variable genes used for

642    reference building) for which the mouse genes had human orthologs in the reference.

643    ***Symphony mapping experiments***

644    To construct each reference for query mapping, we aggregated all reference datasets into a single

645    normalized expression matrix and identified the top $g$ variable genes across all cells ($g$=2000 for

646    PBMCs, $g$=3000 for pancreas) using the variance stabilizing transformation (vst) procedure[49]. We

647    scaled the genes (mean = 0, variance = 1), performed PCA[18], and ran Harmony on the top 20 PCs and

648    default 100 clusters. For PBMCs, we harmonized over 'technology' with default parameters. For

649    pancreas, we harmonized over 'donor' ($\theta$ = 2) and 'technology' ($\theta$ = 4), with $\tau$ = 5. During Symphony

650    mapping, we specified query 'technology' covariate for PBMCs and query 'donor', 'species', and

651    'technology' covariates for pancreas.

652    ***Constructing gold standard embedding***

653    To construct the gold standard *de novo* Harmony embedding, we concatenated the reference and

654    query datasets together into a single expression matrix, subsetted by the top $g$ variable genes over all

655    (both reference and query) cells ($g$=2000 for PBMCs, $g$=3000 for pancreas) and ran Harmony

656    integration on the top 20 PCs[49]. For PBMCs, we harmonized over 'technology' with default parameters.

657    For pancreas, we harmonized over 'donor' ($\theta$ = 2) and 'technology' ($\theta$ = 4), with $\tau$ = 5.

658    ***Assigning ground truth cell types***

659    We clustered the cells in the gold standard embedding using the Louvain algorithm as implemented in

660    the Seurat functions *BuildSNN* and *RunModularityClustering*[51]. For PBMCs, we used nn_k = 5 (to

661    capture rare HSCs), nn_eps = 0.5, and resolution = 0.8. For pancreas, we used the same parameters

662    except nn_k = 30. We labeled clusters with ground truth cell types according to expression of canonical

663    lineage marker genes (**Table S2,5**). PBMCs were assigned across 7 types: T (*CD3D*), NK (*GNLY*), B

664    (*MS4A1*), Monocytes (*CD14, FCGR3A*), DCs (*FCER1A*), Megakaryocytes (*PPBP*), and HSCs (*CD34*).

665    Pancreas cells were assigned across 9 types: alpha (*GCG*), beta (*MAFA*), gamma (*PPY*), delta (*SST*),

666    acinar (*PRSS1*), ductal (*KRT19*), endothelial (*CDH5*), stellate (*COL1A2*), and immune (*PTPRC*).

667    Clusters were labeled if the AUC (calculated using presto[28]) for the corresponding lineage marker was

668    >0.62. For clusters that did not express a specific lineage marker, we manually assigned a cell type

669    based on the top differentially expressed genes (**Table S2,5**). In the PBMCs, cluster 20 was identified

670   as low-quality cells (high in mitochondrial genes; **Table S2**). We removed all cells in this cluster (n=94)

671   from further analyses. The final ground truth labels were used in downstream analyses and cell type

672   classification accuracy evaluation.

673   ***Evaluation of cell type classification accuracy***

674   We predicted query cell types by transferring reference cell type annotations using the *knn* function in

675   the 'class' R package (k=5). Note for the pancreas human cell type classification, we excluded query

676   epsilon and Schwann cells from the accuracy metrics because those cell types are not present in the

677   reference. We calculated overall accuracy across all query cells and cell type F1 scores (the harmonic

678   mean of precision and recall, ranging from 0 to 1). Precision = TP/(TP+FP), recall = TP/(TP+FN), F1 =

679   (2 * precision * recall) / (precision + recall). Cell type F1 was the metric Abdelaal et al. recently used to

680   benchmark automated cell type classifiers.[40] We used their *evaluate.R* script to calculate confusion

681   matrices and F1 by cell type.

682   ***Quantifying local similarity between two embeddings***

683   k-NN-correlation (k-NN-corr) is a new metric that quantifies how well a given alternative embedding

684   preserves the local neighborhood structure with respect to a gold standard embedding. Anchoring on

685   each query cell, we calculate (1) the pairwise similarities to its *k* nearest reference neighbors in the gold

686   standard embedding and (2) the similarities between the same query-reference neighbor pairs in an

687   alternate embedding (**Methods**), then calculate the Spearman (rank-based) correlation between (1)

688   and (2). For similarity, we use the radial basis function kernel: *similarity*$(x,y) = \exp(-\|x-y\|^2/(2\sigma^2))$. For

689   each query cell, we obtain a single k-NN-corr value capturing how well the relative similarities to its *k*

690   nearest reference neighbors are preserved. Note that k-NN-corr is asymmetric with respect to which

691   embedding is selected as the gold standard and which is selected as the alternative because the

692   nearest neighbor pairs are fixed based on how they were defined in the gold standard. The distribution

693   of k-NN-corr scores for all query cells can measure the embedding quality, where higher k-NN-corr

694   indicates greater recapitulation of the gold standard. Lower values for *k* assess more local

695   neighborhoods, whereas higher *k* assesses more global structure.

31

696 We calculated k-NN-corr between the gold standard Harmony embedding and two alternative

697 embeddings: (1) the full Symphony mapping algorithm (projection, clustering, and correction) and (2)

698 PCA-projection only as a comparison to a batch-naïve mapping. PCA-projection refers to the first step

699 of Symphony mapping, where query cells are projected from gene expression to pre-harmonized PC

700 space: $Z_q = U^T G_q$.

## 2.2 Fetal liver hematopoiesis trajectory inference example

702 We obtained post-filtered, post-doublet removal data directly from the authors[52] along with author-

703 defined cell type annotations for 113,063 cells sequenced with 10x 3' end bias and a separate 25,367

704 cells sequenced with 10x 5' end bias. For building the harmonized reference from all 3' cells, we

705 followed the same variable gene selection procedures as the original authors, using the Seurat

706 variance/mean ratio (VMR) method with parameters min_expr = .0125, max_expr = 3, and

707 min_dispersion = 0.625 (resulting in 1,917 variable genes). For each of 14 held-out donor experiments

708 within the 3' dataset, we integrated the reference with Harmony on 13 donors ($\theta = 3$). During Symphony

709 mapping, we specified query 'donor' covariate. For mapping 5' cells against a 3' reference, we removed

710 two donors (F2 and F5, n=3,953) from the 5' query based on low library complexity (**Fig. S5b**), leaving

711 n=21,414 cells from 5 donors. We integrated the reference (all 14 donors sequenced with 3' end bias)

712 with Harmony over 'donor' ($\theta = 3$). During Symphony mapping, we specified both 'donor' and

713 'technology' as covariates. We predicted query cell types by transferring reference cell type annotations

714 using the *knn* function in the 'class' R package (k=30). We visualized the aggregated confusion matrix

715 across all 14 held-out donor experiments as well as the confusion matrix for the single 5'-to-3'

716 experiment using ComplexHeatmap R package[40].

717 For the trajectory inference analysis, we obtained trajectory coordinates from the force directed graph

718 (FDG) embedding of all 3'-sequenced cells from the original authors[53], forming a reference trajectory.

719 We restricted the trajectory to immune cell types only (excluding hepatocytes, fibroblasts, and

720 endothelial). We then mapped a subset of the query cells belonging to the MEM lineage (MEMPs,

721 megakaryocytes, mast cells, early-late erythroid; n=5,141) to the reference-defined trajectory by

722  averaging the FDG coordinates of the 10 reference immune cell neighbors in the Symphony

723  embedding. Note: in addition to the author-provided FDG trajectory, we explored building a trajectory

724  on the reference immune cells with DDRTree[43], but we found that the inferred trajectory was not as

725  clean as FDG.

## 2.3 Memory T cell surface protein inference example

727  We used a memory T cell CITE-seq dataset collected from a tuberculosis disease progression cohort of

728  259 individuals of admixed Peruvian ancestry[54]. The dataset includes expression of the whole

729  transcriptome (33,538 genes) and 30 surface protein markers from 500,089 memory T cells isolated

730  from PBMCs. Including technical replicates, 271 samples were processed across 46 batches.

731  To assess protein prediction accuracy using Symphony embeddings, we randomly selected 217

732  samples (411,004 cells), normalized the expression of each gene (log2(CP10K)) and built a Symphony

733  reference based on mRNA expression, correcting for donor and batch. The held-out 54 samples

734  comprised the query that we mapped onto the reference. We predicted the expression of each of the 30

735  surface proteins in each of the query cells by averaging the protein's expression across the cell's 50

736  nearest reference neighbors. Nearest neighbors were defined based on Euclidean distance in the

737  batch-corrected low-dimensional embedding. As a ground truth for each protein in each query cell, we

738  computed a smoothed estimate of the cells' measured protein expression by averaging the protein's

739  expression across the cell's 50 nearest neighbors in the batch-corrected complete PCA embedding of

740  all 259 donors. We did not use the cells' raw measured protein expression due to dropout. We

741  computed the Pearson correlation coefficient between our predicted expression and the ground truth

742  expression across all cells per donor for each marker.

743  To assess protein prediction accuracy based on mapping to a joint mRNA and protein-based

744  Symphony reference, we first built an integrated reference by using canonical correlation analysis

745  (CCA) to project cells into a low-dimensional embedding maximizing correlation between mRNA and

746  protein features. We randomly selected 217 samples (395,373 cells) to comprise this reference, and

747  normalized the expression of each gene (log2(CP10K)), selected the top 2,865 most variable genes,

33

748    and scaled (mean = 0, variance = 1) all mRNA and protein features. We computed 20 canonical

749    variates (CVs) with the *cc* function in the CCA R package[55] and corrected the mRNA CVs for donor and

750    batch effects with Harmony. Then, we used Symphony to construct a reference based on the batch-

751    corrected CVs, gene loadings on each CV, and mean and standard deviation used to scale each gene

752    prior to CCA. The held-out 54 samples comprised the query that we mapped onto the reference. As

753    described above, we predicted the expression of each of the 30 surface proteins in each of the query

754    cells based on the cell's 5, 10, or 50 nearest neighbors in the reference, estimated the smoothed

755    ground truth expression of each protein in each query cell (now based on the batch-corrected CCA

756    embedding of all 259 donors) and computed the Pearson correlation coefficient for each marker.

## 2.4 Visualization

758    For visualizing the embeddings using UMAP[28,33], we used the 'uwot' R package with the following

759    parameters: n_neighbors=30, learning_rate=0.5, init = 'laplacian', metric = 'cosine', min_dist=0.1

760    (except min_dist=0.3 for fetal liver example). For each Symphony reference, we saved the uwot model

761    at the time of UMAP using the *uwot::save_uwot* function and saved the path to the model file as part of

762    the Symphony reference object. Saving the reference UMAP model allows for the fast projection of new

763    query cells into reference UMAP space from the query embedding from Symphony mapping using the

764    function *uwot::transform*.

765    To distinguish the reference plots from query plots, we visually present the reference embedding as a

766    contour density instead of individual cells. The density plots were generated using ggplot2 function

767    *stat_density_2d* with geom = 'polygon' and contour_var = 'ndensity'. We provide custom functions to

768    generate these plots as part of the Symphony package.

## 2.5 Benchmarking against automatic cell type classifiers

770    We downloaded the PbmcBench benchmarking dataset used by a recent comparison of automatic cell

771    type identification methods[28]. For each of 48 train-test experiments previously described[51], we used the

772    same evaluation metrics (median cell type F1 score) to evaluate Symphony in comparison to the 22

34

773     other classifiers. We obtained the numerical F1-score results for all other classifiers for each of the 48

774     experiments directly from the authors in order to determine Symphony's place within the rank ordering

775     of classifier performance.

776     During reference building, we tried two different gene selection methods: (1) unsupervised (top 2000

777     variable genes) and (2) supervised based on identifying the top 20 differentially expressed (DE) genes

778     per cell type. Option (2) was included to give Symphony the same information as prior-knowledge

779     classifiers (e.g. SCINA with 20 marker genes per cell type). We used the 'presto' package[56] for DE

780     analysis. No integration was performed because the reference had a single-level batch structure

781     (clusters were simply assigned using soft k-means). Onto each of 7 references (each representing 1

782     protocol for donor pbmc1), we mapped either a second protocol for donor pbmc1 (6 experiments) or the

783     same protocol for donor pbmc2 (1 experiment). Given the resulting Symphony joint feature

784     embeddings, we used three downstream classifiers to predict query cell types: 5-NN, SVM with a radial

785     kernel, and glm_net with ridge[34]. A total of 6 Symphony-based classifiers were tested (2 gene selection

786     methods * 3 downstream classifiers).

787     ## 2.6 Runtime analysis

788     We downsampled a large memory T cell dataset[34] to create benchmark reference datasets with 20,000,

789     50,000, 100,000, 250,000, and 500,000 cells. For each, we built a reference (20 PCs, 100 centroids)

790     integrating over 'donor' and mapped three different-sized queries: 1,000, 10,000, and 100,000 cells. To

791     isolate the separate effects of number of query cells and number of query batches on mapping time, we

792     mapped against the 50,000-cell reference: (1) varying the number of query cells (from 1,000 to 10,000

793     cells) while keeping the number of donors constant and (2) varying the number of query donors (6 to

794     120 donors) while keeping the number of cells constant (randomly sampling 10,000 cells). We also

795     performed separate experiments varying the number of reference centroids (25 to 400) and number of

796     dimensions (10 to 320 PCs) while keeping all other parameters constant. We ran all jobs on Linux

797     servers allotted 4 cores and 64 GB of memory (Intel Xeon E5-2690 v.3 processors) and used the

798     *proc.time* R function to measure elapsed time.

# Data availability

Datasets for all analyses were obtained from the links in **Table S1.** All datasets are publicly available except the memory T cell CITE-seq data, which will be available at GEO accession GSE158769.

# Code availability

We provide an implementation of Symphony along with prebuilt references from all examples at https://github.com/immunogenomics/symphony. Scripts reproducing results of this paper will be made available at https://gihub.com/immunogenomics/referencemapping.

# Acknowledgements

# Author contributions

I.K., J.B.K., and S.R. conceived the project. J.B.K. and I.K. developed the method and performed the analyses under the guidance of S.R. S.R., A.N., and D.B.M. contributed to generating the memory T

821 cell dataset. A.N. performed analysis of the memory T cell dataset. All authors participated in

822 interpretation and writing the manuscript.

## 823 Competing interests

## 825 References

826 1. Klein, A. M. & Treutlein, B. Single cell analyses of development in the modern era. *Development*

827 **146**, (2019).

828 2. Svensson, V., da Veiga Beltrame, E. & Pachter, L. A curated database reveals trends in single-cell

829 transcriptomics. *biorxv* (2019) doi:10.1101/742304.

830 3. Han, X. *et al.* Construction of a human cell landscape at single-cell level. *Nature* (2020)

831 doi:10.1038/s41586-020-2157-4.

832 4. Cao, J. *et al.* The single-cell transcriptional landscape of mammalian organogenesis. *Nature* **566**,

833 496–502 (2019).

834 5. Jerber, J. *et al.* Population-scale single-cell RNA-seq profiling across dopaminergic neuron

835 differentiation. *bioRxiv* 2020.05.21.103820 (2020) doi:10.1101/2020.05.21.103820.

836 6. Zhang, F. *et al.* Defining inflammatory cell states in rheumatoid arthritis joint synovial tissues by

837 integrating single-cell transcriptomics and mass cytometry. *Nat. Immunol.* **20**, 928–942 (2019).

838 7. Reyes, M. *et al.* An immune-cell signature of bacterial sepsis. *Nat. Med.* **26**, 333–340 (2020).

839 8. Kotliarov, Y. *et al.* Broad immune activation underlies shared set point signatures for vaccine

840 responsiveness in healthy individuals and disease activity in patients with lupus. *Nat. Med.* **26**,

841 618–629 (2020).

842 9. Schafflick, D. *et al.* Integrated single cell analysis of blood and cerebrospinal fluid leukocytes in

843 multiple sclerosis. *Nat. Commun.* **11**, 247 (2020).

844    10.  Smillie, C. S. *et al.* Intra- and Inter-cellular Rewiring of the Human Colon during Ulcerative Colitis.

845          *Cell* **178**, 714-730.e22 (2019).

846    11.  Litviňuková, M. *et al.* Cells of the adult human heart. *Nature* (2020) doi:10.1038/s41586-020-2797-

847          4.

848    12.  Rozenblatt-Rosen, O., Stubbington, M. J. T., Regev, A. & Teichmann, S. A. The Human Cell Atlas:

849          from vision to reality. *Nature* **550**, 451–453 (2017).

850    13.  Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-

851          sequencing data are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–

852          427 (2018).

853    14.  Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell transcriptomes

854          using Scanorama. *Nat. Biotechnol.* **37**, 685–691 (2019).

855    15.  Welch, J. D. *et al.* Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain

856          Cell Identity. *Cell* **177**, 1873-1887.e17 (2019).

857    16.  Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-

858          cell transcriptomics. *Nat. Methods* **15**, 1053–1058 (2018).

859    17.  Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat.*

860          *Methods* **16**, 1289–1296 (2019).

861    18.  Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21 (2019).

862    19.  He, Z., Brazovskaja, A., Ebert, S., Camp, J. G. & Treutlein, B. CSS: cluster similarity spectrum

863          integration of single-cell genomics data. *Genome Biol.* **21**, 224 (2020).

864    20.  Tran, H. T. N. *et al.* A benchmark of batch-effect correction methods for single-cell RNA

865          sequencing data. *Genome Biol.* **21**, 12 (2020).

866    21.  Zhang, Q. *et al.* Landscape and Dynamics of Single Immune Cells in Hepatocellular Carcinoma.

867          *Cell* **179**, 829-845.e20 (2019).

868    22.  Wei, K. *et al.* Notch signalling drives synovial fibroblast identity and arthritis pathology. *Nature* **582**,

869          259–264 (2020).

870   23.  Kirita, Y., Wu, H., Uchimura, K., Wilson, P. C. & Humphreys, B. D. Cell profiling of mouse acute

871         kidney injury reveals conserved cellular responses to injury. *Proc. Natl. Acad. Sci. U. S. A.* **117**,

872         15874–15883 (2020).

873   24.  Sandu, I. *et al.* Landscape of Exhausted Virus-Specific CD8 T Cells in Chronic LCMV Infection.

874         *Cell Rep.* **32**, 108078 (2020).

875   25.  Lähnemann, D. *et al.* Eleven grand challenges in single-cell data science. *Genome Biol.* **21**, 31

876         (2020).

877   26.  Lotfollahi, M. *et al.* Query to reference single-cell integration with transfer learning. *bioRxiv* (2020).

878   27.  Cao, Z.-J., Wei, L., Lu, S., Yang, D.-C. & Gao, G. Searching large-scale scRNA-seq databases via

879         unbiased cell embedding with Cell BLAST. *Nat. Commun.* **11**, 3458 (2020).

880   28.  Abdelaal, T. *et al.* A comparison of automatic cell identification methods for single-cell RNA

881         sequencing data. *Genome Biol.* **20**, 194 (2019).

882   29.  Zhang, Z. *et al.* SCINA: A Semi-Supervised Subtyping Algorithm of Single Cells and Bulk Samples.

883         *Genes*  **10**, (2019).

884   30.  Kiselev, V. Y., Yiu, A. & Hemberg, M. scmap: projection of single-cell RNA-seq data across data

885         sets. *Nat. Methods* **15**, 359–362 (2018).

886   31.  Alquicira-Hernandez, J., Sathe, A., Ji, H. P., Nguyen, Q. & Powell, J. E. scPred: accurate

887         supervised method for cell-type classification from single-cell RNA-seq data. *Genome Biol.* **20**, 264

888         (2019).

889   32.  Tan, Y. & Cahan, P. SingleCellNet: A Computational Tool to Classify Single Cell RNA-Seq Data

890         Across Platforms and Across Species. *Cell Syst* **9**, 207-213.e2 (2019).

891   33.  Ding, J. *et al.* Systematic comparative analysis of single cell RNA-sequencing methods. *bioRxiv*

892         632216 (2019) doi:10.1101/632216.

893   34.  Nathan, A. *et al.* Multimodal Profiling of 500,000 Memory T Cells from a Tuberculosis Cohort

894         Identifies Cell State Associations with Demographics, Environment, and Disease. (2020)

895         doi:10.2139/ssrn.3652337.

896   35. Segerstolpe, Å. *et al.* Single-Cell Transcriptome Profiling of Human Pancreatic Islets in Health and

897        Type 2 Diabetes. *Cell Metab.* **24**, 593–607 (2016).

898   36. Lawlor, N. *et al.* Single-cell transcriptomes identify human islet cell signatures and reveal cell-type-

899        specific expression changes in type 2 diabetes. *Genome Res.* **27**, 208–222 (2017).

900   37. Grün, D. *et al.* De Novo Prediction of Stem Cell Identity using Single-Cell Transcriptome Data. *Cell*

901        *Stem Cell* **19**, 266–277 (2016).

902   38. Muraro, M. J. *et al.* A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Syst* **3**, 385-

903        394.e3 (2016).

904   39. Baron, M. *et al.* A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals

905        Inter- and Intra-cell Population Structure. *Cell Syst* **3**, 346-360.e4 (2016).

906   40. Popescu, D.-M. *et al.* Decoding human fetal liver haematopoiesis. *Nature* **574**, 365–371 (2019).

907   41. Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells. *Nat.*

908        *Methods* **14**, 865–868 (2017).

909   42. Peterson, V. M. *et al.* Multiplexed quantification of proteins and transcripts in single cells. *Nat.*

910        *Biotechnol.* **35**, 936–939 (2017).

911   43. Nathan, A. *et al.* Multimodal memory T cell profiling identifies a reduction in a polyfunctional Th17

912        state associated with tuberculosis progression. *bioRxiv* 2020.04.23.057828 (2020)

913        doi:10.1101/2020.04.23.057828.

914   44. Luecken, M. D. & Theis, F. J. Current best practices in single-cell RNA-seq analysis: a tutorial. *Mol.*

915        *Syst. Biol.* **15**, e8746 (2019).

916   45. Lotfollahi, M., Naghipourfar, M., Theis, F. J. & Alexander Wolf, F. Conditional out-of-sample

917        generation for unpaired data using trVAE. *arXiv [cs.LG]* (2019).

918   46. Berger, B. & Cho, H. Emerging technologies towards enhancing privacy in genomic data sharing.

919        *Genome Biol.* **20**, 128 (2019).

920   47. Wang, S., Pisco, A. O., Karkanias, J. & Altman, R. B. Unifying single-cell annotations based on the

921        Cell Ontology. *bioRxiv* 810234 (2019) doi:10.1101/810234.

40

922    48.  Gayoso, A. *et al.* A Joint Model of RNA Expression and Surface Protein Abundance in Single Cells.

923          *biorxiv* 791947 (2019) doi:10.1101/791947.

924    49.  Baglama, J. & Reichel, L. Augmented Implicitly Restarted Lanczos Bidiagonalization Methods.

925          *SIAM Journal on Scientific Computing* vol. 27 19–42 (2005).

926    50.  Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the integration of

927          genomic datasets with the R/Bioconductor package biomaRt. *Nat. Protoc.* **4**, 1184–1191 (2009).

928    51.  Korsunsky, I., Nathan, A., Millard, N. & Raychaudhuri, S. Presto scales Wilcoxon and auROC

929          analyses to millions of observations. *bioRxiv* 653253 (2019) doi:10.1101/653253.

930    52.  Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in

931          multidimensional genomic data. *Bioinformatics* **32**, 2847–2849 (2016).

932    53.  Qi Mao, Li Wang, Tsang, I. W. & Yijun Sun. Principal Graph and Structure Learning Based on

933          Reversed Graph Embedding. *IEEE Trans. Pattern Anal. Mach. Intell.* **39**, 2227–2241 (2017).

934    54.  Leurgans, S. E., Moyeed, R. A. & Silverman, B. W. Canonical Correlation Analysis When the Data

935          are Curves. *J. R. Stat. Soc. Series B Stat. Methodol.* **55**, 725–740 (1993).

936    55.  McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for

937          Dimension Reduction. *arXiv [stat.ML]* (2018).

938    56.  Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear Models via

939          Coordinate Descent. *J. Stat. Softw.* **33**, 1–22 (2010).