

1 **Scupa: Single-cell unified polarization assessment of immune cells using the single-cell foundation** 2 **model**

3 Wendao Liu^{1,2}, Zhongming Zhao^{1,2,*}

4 ¹The University of Texas MD Anderson Cancer Center UTHealth Houston Graduate School of Biomedical
5 Sciences, Houston, TX, USA.

6 ²Center for Precision Health, McWilliams School of Biomedical Informatics, The University of Texas Health
7 Science Center at Houston, Houston, TX, USA.

8 *Corresponding Author. E-mail: zhongming.zhao@uth.tmc.edu

9

10 **Abstract**

11 Immune cells undergo cytokine-driven polarization in respond to diverse stimuli. This process significantly
12 modulates their transcriptional profiles and functional states. Although single-cell RNA sequencing (scRNA-seq)
13 has advanced our understanding of immune responses across various diseases or conditions, currently there lacks
14 a method to systematically examine cytokine effects and immune cell polarization. To address this gap, we
15 developed **Single-cell unified polarization assessment (Scupa)**, the first computational method for comprehensive
16 immune cell polarization analysis. Scupa is trained on data from the Immune Dictionary, which characterizes 66
17 cytokine-driven polarization states across 14 immune cell types. By leveraging the cell embeddings from the
18 Universal Cell Embeddings model, Scupa effectively identifies polarized cells in new datasets generated from
19 different species and experimental conditions. Applications of Scupa in independent datasets demonstrated its
20 accuracy in classifying polarized cells and further revealed distinct polarization profiles in tumor-infiltrating
21 myeloid cells across cancers. Scupa complements conventional single-cell data analysis by providing new insights
22 into immune cell polarization, and it holds promise for assessing molecular effects or identifying therapeutic
23 targets in cytokine-based therapies.

24

25 **Introduction**

26 Immune cells detect and respond to a variety of stimuli, ensuring the body can effectively combat infections and
27 other environmental or biological threats. Cytokines as crucial signaling molecules that facilitate communication
28 between immune cells. These cytokines can induce significant changes in the transcriptional profiles and
29 functional states of immune cells, a phenomenon known as immune cell polarization^{1,2}. Through polarization,
30 immune cells adapt their responses to better address specific challenges, enhancing the overall effectiveness of the
31 immune system.

32 Recently, single-cell RNA sequencing (scRNA-seq) has been widely applied in immunological studies to
33 investigate the responses of immune cells under various conditions. The production and response to cytokines,
34 which play critical regulatory roles, have been extensively studied in numerous diseases, including COVID-19^{3,4},
35 rheumatoid arthritis⁵, and cancers^{6,7}. Importantly, a recent study systematically characterized the responses of 14
36 immune cell types to each of 86 cytokines and summarized the results as the Immune Dictionary⁸. A total of 66
37 cytokine-driven cell polarization states were identified in that study, serving as a valuable reference for assessing
38 immune cell polarization in other scRNA-seq studies.

39 Currently, the investigation of immune cell polarization in scRNA-seq data is limited and there lacks consistent
40 standards across studies. Most studies identified polarized cells based on the expression of certain signature genes

41 from previous findings. However, this empirical approach suffers from technical noises, such as dropout effect
42 and batch effect^{9,10}, as well as biological variations including tissue or disease variations. For example, many
43 studies analyzed macrophage polarization using M1 and M2 signature genes, but there was no consensus in the
44 signature gene lists and the expression of these genes varied by conditions¹¹⁻¹³. To address this challenge, recent
45 advances in single-cell foundation models offer unprecedented potential. These models were trained on vast
46 amounts of scRNA-seq data with millions of cells across all the representative human organs and learned cell
47 representations within a unified biological latent space¹⁴⁻¹⁶. As they have been successfully demonstrated in
48 multiple downstream tasks like cell type annotation and batch correction, we hypothesized that single-cell
49 foundation models could also effectively represent immune cell polarization.

50 To facilitate the analysis of immune cell polarization in scRNA-seq data, we developed the method Scupa for
51 **Single-cell unified polarization assessment**. After being trained on scRNA-seq data from the Immune Dictionary,
52 Scupa learns the representations of immune cell polarization within the latent space of Universal Cell Embeddings
53 (UCEs)¹⁶. As the first computational method for systematic immune cell polarization analysis, Scupa enables the
54 assessment of individual cell polarization across various predefined cytokine-driven cell polarization states in any
55 scRNA-seq dataset.

56

57 **Results**

58 **Overview of Scupa framework and immune cell polarization states**

59 Immune cells undergo transcriptional and phenotypical changes in cytokine-driven polarization. Scupa uses the
60 immune cell polarization states in the Immune Dictionary as the reference, and trains machine learning models to
61 distinguish between polarized cells from cytokine-treated samples and unpolarized cells from phosphate buffered
62 saline (PBS)-treated control samples. The specific polarized states of a cell type are usually driven by the specific
63 cytokines, with each state exhibiting a unique transcriptional profile. Instead of relying on gene expression
64 features, which can vary across species, tissues, and conditions, Scupa utilizes cell embeddings from the single-
65 cell foundation model UCE for assessing immune cell polarization¹⁶. We chose UCE over other single-cell
66 foundation models because of its advantages in multiple-species compatibility, and ‘zero-shot’ capability without
67 the need for additional fine-tuning. We further reduced the dimension of cell embeddings using principal
68 component analysis (PCA), and trained support vector machine (SVM) models to classify polarized cells and
69 unpolarized cells based on principal components (PCs). The models learned the representation of each polarized
70 state in the latent space of UCEs, allowing for transferability to other datasets (Fig. 1a). We found that SVM
71 outperformed several other machine learning models in this task (Supplementary Table 1).

72 For any new scRNA-seq dataset containing immune cells presented in the Immune Dictionary, Scupa examines
73 whether the cells have the similar transcriptional changes as reference polarized cells, thereby inferring their
74 polarization states and received cytokines. The polarization of each cell is assessed based on its cell embeddings
75 from UCE. The trained SVM models predict the polarization scores of each cell using the learned representation
76 in the unified latent space of UCEs. According to the cell type, Scupa assigns a score to each individual cell for
77 every polarization state, ranging from 0 (unpolarized) to 1 (fully polarized). Additionally, Scupa outputs a p-value
78 by comparing the polarization score of a cell with the null distribution of polarization scores from the reference
79 unpolarized cells, facilitating the identification of significantly polarized cells. Scupa is designed for integration
80 into the widely used Seurat pipeline for comprehensive single-cell data analyses¹⁷, enabling the output scores and
81 p-values to be readily visualized in multiple formats (Fig. 1a).

82 Scupa supports the analysis of 66 polarization states in 14 immune cell types from the Immune Dictionary (Fig.
83 1b, Supplementary Fig. 1). Among these states, the ‘a states’ of all cell types represent those driven by some type-

84 I interferons (IFN- α 1, IFN- β , IFN- ϵ , IFN- κ). The ‘b states’ represent those driven by IFN- γ and interleukins
85 inducing IFN- γ expression (IL-2, IL-12, IL-15, IL-18). The ‘c states’ are for those driven by two proinflammatory
86 cytokines IL-1 α and IL-1 β , and the ‘d states’ are mainly driven by TNF- α . In contrast, the driving cytokines of the
87 ‘e states’ and ‘f states’ in various cell types tend to vary. Identification of polarized cells using Scupa suggests the
88 potential presence of one or more driving cytokines in the tissue, thereby facilitating the understanding of immune
89 cell environment, communication and response in scRNA-seq data.

90 **Scupa learns the representation of cell polarization in various immune cell types**

91 Cytokines are key regulators of intracellular signaling and gene expression, acting as messengers that mediate and
92 modulate immune responses. They activate signaling cascades that lead to the phosphorylation and activation of
93 various transcription factors, which then translocate to the nucleus to modulate the transcription of genes.
94 Consequently, the cytokine-driven immune cell polarization states are characterized by unique transcriptional
95 profiles⁸. UCE can effectively capture these transcriptional changes and represent them as variations in cell
96 embeddings. For example, we found that the CD8⁺ T cells treated with driving cytokines of each polarization state
97 tended to have cell embeddings shifting away from the unpolarized cells, as visualized by uniform manifold
98 approximation and projection (UMAP, Fig. 2a). By filtering cells based on cosine similarity and the expression of
99 top marker genes, we obtained fully polarized cells in each state, which displayed distinctly different cell
100 embedding distributions from unpolarized cells (Fig. 2b).

101 For each polarization state, we trained an SVM model to classify polarized cells and unpolarized cells, and
102 quantified the polarization using a score derived from the trained models. With this approach, the gradients of
103 polarization scores in the unified cell embedding space represent the directions of cell polarization (Fig. 2c).
104 Among all CD8⁺ T cells, we found that the polarization scores of some states are highly correlated (Fig. 2d). For
105 example, the Spearman correlation coefficient between T8-e and T8-f state scores is 0.96. This indicates a high
106 similarity in transcriptional changes between these two states, which aligns with the findings in the original
107 study⁸. We then calculated the p-values for each polarization state by comparing the polarization scores with the
108 distribution of scores from unpolarized cells. This metrics is useful because the commonly applied p-value
109 threshold of 0.05 can serve as a good criterion for identifying fully polarized cells across most polarization states
110 (Fig. 2e). Importantly, in addition to CD8⁺ T cells, we observed similar patterns of polarization scores and p-
111 values across all other 13 cell types and 61 polarization states. This demonstrated that Scupa could effectively
112 learn the representation of all polarization states in the unified cell embedding space (Supplementary Figs. 2-14).

113 Next, we evaluated the performance of Scupa across all polarization states. Using a random 70% of the data for
114 training and the remaining 30% for testing, we repeated this process 20 times and calculated the area under the
115 receiver operating characteristic curve (AUROC) for testing set. The median AUROC values were above 0.95 in
116 56 out of the 66 polarization states (Fig. 2f). For polarization states with lower performance, such as Neu-e
117 (AUROC = 0.836), the primary factor was the insufficient number of cells leading to poor model fitting
118 (Supplementary Fig. 15). Overall, Scupa achieves superior performance in learning the representations of
119 polarization states in scRNA-seq data.

120 **Scupa characterizes cytokine-driven immune cell polarization in both *in vitro* and *in vivo* datasets**

121 To evaluate Scupa’s performance in independent datasets, we collected two scRNA-seq datasets generated from
122 cytokine-treated samples. The first dataset comprises human peripheral blood mononuclear cell (PBMC) samples
123 treated with IFN- β or left untreated *in vitro*¹⁸. IFN- β treatment dramatically induces the transcription of interferon-
124 stimulated genes, causing the cell embeddings of treated cells to diverge significantly from those of untreated
125 cells across all immune cell types (Fig. 3a). We applied Scupa to analyze the immune cell polarization in this
126 dataset. Considering that IFN- β is the driving cytokine of -a states in all immune cell types (Fig. 1b), we
127 examined the polarization scores of these states. We found that all cells from the IFN- β -treated sample had

128 polarization scores close to 1, while those cells from the untreated sample had scores close to 0. This clearly
129 indicates a tremendous difference (Fig. 3b). Additionally, nearly all cells from the IFN- β -treated sample had
130 polarization p-values below 0.05, whereas those cells from the untreated sample had p-values above 0.05 (Fig. 3c,
131 Supplementary Fig. 16). The ROC curves further confirmed the near-perfect performance of using polarization
132 scores to classify cells from treated and untreated samples, with AUROCs above 0.99 across all cell types (Fig.
133 3d). Importantly, although Scupa was trained on mouse scRNA-seq data, it demonstrated excellent performance
134 on this human scRNA-seq data without any adaptation, suggesting that Scupa learned unified representations of
135 immune cell polarization across species.

136 In the second study, the mice chronically infected with lymphocytic choriomeningitis virus (LCMV) were treated
137 with IL-2, anti-PD-L1, or a combination therapy¹⁹. The virus-specific CD8⁺ T cells from the mouse spleen were
138 sorted for scRNA-seq. We clustered cells from three treatment groups and control group using UCE cell
139 embeddings (Fig. 3e), and applied Scupa to analyze the CD8⁺ T cell polarization in this dataset. As IL-2 is the
140 driving cytokine for two polarization states, T8-b and T8-e, we examined the polarization scores and p-values of
141 these two states (Fig. 3f, g). A small number of cells were polarized to T8-b state and enriched in cluster 3, while
142 more cells were polarized to T8-e state and enriched in cluster 2. To verify the two polarization states with shared
143 driving cytokines, we examined the expression of polarization state marker genes identified in the Immune
144 Dictionary⁸. We found that the expression of some marker genes correlated with polarization scores, such as *Stat1*
145 with T8-b scores and *Ptma* with T8-e scores. This result suggested that Scupa could distinguish these two
146 polarization states. On the other hand, we observed the ubiquitous high or low expression of some marker genes
147 (e.g., *Igtp*, *Gbp4*, *Ncl*, and *Npm1*) in all CD8⁺ T cells, highlighted the variability of gene expression in different
148 datasets (Fig. 3h, i). Therefore, the feature of independence on marker genes in Scupa makes it more applicable to
149 various scRNA-seq data from diverse sources.

150 Since IL-2 was administrated *in vivo*, it was likely that only a subset of spleen T cells were polarized by IL-2,
151 while other T cells were polarized by different cytokines or remained unpolarized. To analyze the effect of IL-2
152 treatment, we compared the percentage of cells with significant T8-e polarization (polarization p-value<0.05)
153 among groups. This percentage significantly increased in the IL-2 treatment group when compared to the control
154 group (p<0.0001), and increased in the combination therapy treatment group when compared to the anti-PD-L1
155 treatment group (p<0.001, Fig. 3j). We also compared the percentage of cells in cluster 2, which was highly
156 enriched with the cells polarized to T8-e state. Similar increase was observed with IL-2 treatment (Fig. 3j),
157 confirming the CD8⁺ T cell polarization to IL-2 driving polarization states.

158 As above, Scupa was validated using scRNA-seq datasets generated from both *in vitro* and *in vivo* samples. It
159 superbly classified stimulated and unstimulated cells from *in vitro* samples, and also revealed the increases of cell
160 polarized to a cytokine-driving polarization state in *in vivo* samples.

161 **Scupa reveals polarization states and proinflammatory responses of myeloid cells across cancer types**

162 Myeloid cells play a crucial role in the tumor microenvironment, influencing cancer progression and response to
163 therapy. These cells, which include macrophages, monocytes, dendritic cells (DCs), and neutrophils, can exhibit
164 either pro-tumor or anti-tumor functions depending on their polarization state and the cytokine milieu^{20,21}.
165 Understanding the dual roles of myeloid cells in cancers is essential for developing therapeutic strategies that
166 modulate their function to inhibit tumor progression and improve patient outcomes. In order to systematically
167 analyze the myeloid cell polarization in multiple cancer types, we applied Scupa to a pan-cancer single-cell atlas
168 of tumor infiltrating myeloid cell dataset¹¹.

169 We first revisited the macrophage polarization using Scupa as it has been extensively investigated in cancer
170 research²². We calculated the polarization scores of five polarization states for all macrophage clusters and
171 compared these scores across seven cancer types. Among the five states, Mac-b, Mac-c, and Mac-d are all M1-

172 like states driven by proinflammatory cytokines, while Mac-e is a M2-like state driven by cytokines that induce
173 M2 polarization (Fig. 1b)². Several macrophage clusters displayed consistent polarization profiles across cancer
174 types. *CIQC*⁺ and *LYVE1*⁺ macrophages generally had low polarization scores of all polarization states, indicating
175 that they were mainly unpolarized. *ISG15*⁺ macrophages generally had high Mac-a, Mac-b, Mac-c, and Mac-d
176 polarization scores, suggesting that they were primarily polarized by type-I interferons and proinflammatory
177 cytokines. *NLRP3*⁺ and *INHBA*⁺ macrophages also had high Mac-b, Mac-c, and Mac-d scores. This might indicate
178 polarization by proinflammatory cytokines. In particular, *SPPI*⁺ macrophages showed significant variation across
179 these cancer types. They had highest Mac-e polarization scores in thyroid carcinoma (THCA), but had much
180 lower scores in pancreatic adenocarcinoma (PAAD) or uterine corpus endometrial carcinoma (UCEC). When
181 comparing macrophage polarization in different cancer types, we found that PAAD and kidney cancer had the
182 lowest overall polarization scores in all polarization states. This is potentially due to lower cytokine production in
183 these tumors (Fig. 4a).

184 Next, we analyzed the relationship between polarization scores of macrophage polarization states and the
185 expression of M1 and M2 signature genes¹¹. The mean expression of M1 signature genes was positively
186 correlated with Mac-b, Mac-c, and Mac-d polarization scores across all cancer types. In contrast, there was no
187 consistent correlation between the mean expression of M2 signature genes and Mac-e polarization scores (Fig.
188 4b). This lack of correlation was likely attributed to the observation of that the M2 signature genes did not overlap
189 with the Mac-e marker genes from the Immune Dictionary, thus resulting in strong variation in different studies.
190 Additionally, we included angiogenesis and phagocytosis signature genes in the correlation analysis¹¹. The mean
191 expression of angiogenesis signature genes was moderately positively correlated with Mac-b (Spearman's
192 correlation coefficient $\rho=0.22$), Mac-c ($\rho=0.42$), Mac-d ($\rho=0.35$), and Mac-e ($\rho=0.35$) polarization scores,
193 highlighting macrophages' simultaneous contributions to inflammation and angiogenesis²³. Conversely, the mean
194 expression of phagocytosis signature genes was negatively correlated with polarization scores of all polarization
195 states, suggesting that phagocytotic macrophages were generally less polarized by cytokines (Fig. 4c).
196 Considering the high correlation between Mac-b, Mac-c, and Mac-d polarization scores, and proinflammatory
197 cytokines as driving cytokines for these states, we defined a macrophage proinflammatory state score as the
198 maximum of these three scores in each single cell. This score summarizes macrophage proinflammatory
199 polarization and its distribution indicates the proinflammatory activity in different cancers. We found that
200 lymphoma (LYM) was characterized with the highest overall proinflammatory scores, followed by THCA and
201 esophageal carcinoma (ESCA). PAAD and kidney cancer exhibited the lowest overall proinflammatory scores
202 (Fig. 4d).

203 In addition to macrophages, we examined the polarization of monocytes and different DC populations. These cell
204 populations displayed stronger variation in polarization across cancer types than macrophages (Fig. 4e, f). In
205 LYM, *CD14*⁺ monocytes, *CD16*⁺ monocytes, cDC2, and pDCs had high polarization scores of all polarization
206 states, supporting the strong effect and important role of multiple cytokines in the cancer²⁴. In ESCA, UCEC, and
207 THCA, monocytes, cDC1, and pDCs were also polarized to multiple states. Notably, the overall polarization of
208 macrophages, monocytes, and DCs showed a consistent trend in all seven cancer types. All these myeloid cell
209 populations were more polarized in LYM, THCA, UCEC, and THCA, but less polarized in PAAD, LYE, and
210 kidney cancer. This trend likely reflects the variation in cytokine milieu in the tumor microenvironment of these
211 cancer types (Fig. 4d, e, f). Taken together, these results indicated the distinct cytokine environments, myeloid
212 cell polarization patterns, and proinflammatory responses in different cancers, demonstrating Scupa's capability
213 to complement the conventional scRNA-seq analysis with unique cell polarization analysis.

214

215 Discussion

216 In this study, we introduce Scupa, the first computational method designed to assess immune cell polarization
217 from scRNA-seq data. Scupa is design to complement to the conventional scRNA-seq analysis pipelines by
218 providing additional perspectives in the cytokine environment and immune cell polarization. The method
219 leverages the recently released Immune Dictionary, which systematically characterized the responses of 14
220 immune cell types to 86 cytokines and then identified 66 cytokine-driven polarization states. Unlike traditional
221 approaches that rely on predefined signature genes, Scupa utilizes cell embeddings from the single-cell foundation
222 model, Universal Cell Embeddings, to capture the nuanced transcriptional changes associated with different
223 polarization states. Our results clearly indicated that Scupa could effectively classify polarized and unpolarized
224 cells by training machine learning models on cell embeddings. The method was validated using independent
225 datasets, including human PBMCs treated with IFN- β and CD8⁺ T cells from mice spleen treated with IL-2 and
226 anti-PD-L1. Scupa accurately identified polarized cells and revealed the cytokine-driven polarization states within
227 these datasets. With UCE's multiple-species compatibility, even though our method was trained on the mouse
228 data of Immune Dictionary, it could effectively be applied to human dataset without additional fine tuning. While
229 this needs further evaluation, it suggested the robustness of our model across species.

230 Additionally, we apply Scupa to a pan-cancer single-cell atlas to investigate the polarization of myeloid cells
231 across seven cancer types. The analysis reveals distinct polarization profiles and proinflammatory responses in
232 macrophages, monocytes, and DC populations. Notably, LYM, THCA, UCEC exhibit higher polarization scores
233 when compared to PAAD, LYE, and kidney cancer, reflecting distinct cytokine environments in different cancer
234 types.

235 Macrophage polarization has been one of the top research interests in immunology due to its significant role in
236 health and disease. Traditionally, macrophages have been categorized into two main polarization states. The first
237 state is M1, or classically activated macrophages that are induced by proinflammatory cytokines like IFN- γ and
238 TNF- α . These macrophages are associated with antimicrobial and tumoricidal activities. The second state is M2,
239 or alternatively activated macrophages that are stimulated by cytokines such as IL-4 and IL-13. These
240 macrophages are involved in tissue repair and immune regulation^{2,22}. However, this dichotomous classification
241 has been increasingly recognized as an oversimplification, with emerging evidence suggesting a spectrum of
242 intermediate states influenced by a variety of cytokines and environmental cues²⁵. In many scRNA-seq studies,
243 the lack of consistent marker genes and standards further makes the polarization inference arbitrary. Scupa
244 provides a signature gene-free approach for analyzing macrophage polarization to five cytokine-driving states.
245 Using Scupa, we identified the *CIQC*⁺ or *LYVE1*⁺ unpolarized macrophage subpopulations and multiple
246 macrophage subpopulations polarized to M1-like polarization states. Of note, Scupa analysis revealed *SPPI*⁺
247 macrophages polarized to Mac-e state, a M2-like polarization state, suggesting its pro-tumor roles. This result was
248 further supported by previous findings that the worse clinical outcomes were associated with higher *SPPI*
249 expression¹¹.

250 Cytokine therapies have emerged as a powerful strategy for treating various diseases, leveraging the potent
251 regulatory effects of cytokines to modulate immune responses and target disease processes. These therapies
252 harness the ability of cytokines to influence cellular behavior, enhancing or suppressing immune responses as
253 needed. For instance, cytokines such as interferons have been used to treat viral infections and certain cancer
254 types^{26,27}, while interleukins have shown promise in enhancing immune responses in cancer immunotherapy^{19,28}.
255 Scupa's ability to systematically assess immune cell polarization based on scRNA-seq data offers a valuable tool
256 for assessing the molecular effects of cytokine therapies (Fig. 3e-k). By evaluating how cytokine treatments
257 impact the polarization and functional states of immune cells, Scupa can provide insights into the therapeutic
258 mechanisms at a granular level. It also holds the potential to identify therapeutic targets for cytokine therapies,
259 facilitating the development and optimization of cytokine-based treatments.

260 There are some limitations in this work. First, there lacks high-quality experimentally generated immune cell
261 polarization states. Although we included all predefined polarization states from the Immune Dictionary, there
262 may be additional, unidentified polarization states due to the constraints of experimental design. For example,
263 some chemokines have been found to induce macrophage polarization^{29,30}, but the chemokine family was not
264 included in the Immune Dictionary. Second, there are noncytokine pathways of immune cell polarization, such as
265 hypoxia and lactate for macrophage polarization^{31,32}. Identification of new immune cell polarization states will
266 require further investigations comparable to the Immune Dictionary, and necessary immunological expertise.
267 Such data has not been available yet, but Scupa can be updated to include newly identified polarization states in
268 future studies. Third, while Scupa has been demonstrated robust in identifying polarization and states in both mice
269 and humans, more evaluation will be needed to enhance the models across species.

270 In conclusion, we introduced Scupa, the first method for comprehensive immune cell polarization assessment
271 using scRNA-seq data. It is broadly applicable to the studies of various diseases involving immune cell populations
272 and is particularly useful in contexts where cytokines play important roles in disease pathogenesis, progression,
273 and treatment.

274

275 **Methods**

276 **Generating cell embeddings using UCE and dimension reduction**

277 We used single-cell foundation model, Universal Cell Embeddings (<https://github.com/snap-stanford/UCE>), to
278 generate cell embeddings for all scRNA-seq datasets used in this study. The pretrained 4-layer model was
279 employed with a batch size of 50. The UCE cell embeddings are 1,280-dimensional, representing cells in the
280 unified latent space. However, this high dimensionality poses challenges for training machine learning models on
281 most polarization states with a limited number of cells.

282 To address this issue, we performed PCA to reduce the dimensionality of UCE cell embeddings for each cell type.
283 Principal components (PCs) are linear combinations of vector bases in the cell embedding space, with the top PCs
284 representing directions with the largest variation for that cell type. By default, Scupa uses the first 20 PCs as
285 features for machine learning training and prediction. Additionally, we generated two-dimensional UMAPs using
286 the first 20 PCs for data visualization.

287 **Identifying fully polarized cells in the Immune Dictionary**

288 According to both our analysis and the original study⁸, only a subset of immune cells were polarized after
289 cytokine treatment in the *in vivo* experiments. This was likely due to variable cytokine concentration, receptor
290 expression, and cellular status of different cells. The cell embeddings of some cells from cytokine-treated samples
291 were closer to those of unpolarized cells from the PBS-treated samples than other cells from cytokine-treated
292 samples, likely suggesting an unpolarized state or mildly polarization state (Fig. 2e, Supplementary Fig. 2-14).
293 Therefore, we first identified fully polarized cells of each polarization state for training machine learning models.

294 We identified the fully polarized cells of each polarization state based on following three criteria. 1) The cell is
295 from a sample treated with one of the driving cytokines. 2) The mean expression of top marker genes of the
296 polarization in the cell is higher than that of most other cells. 3) The UCE cell embeddings of the cell are similar
297 to those of the other cells from the samples treated with driving cytokines.

298 For criterion 2, we used a consistent threshold of the 90th quantile, i.e., the mean expression of top marker genes
299 in the cell was required to be higher than that in 90% cells of the same cell type. For criterion 3, we found that the
300 UCE cell embeddings of all cells of the same cell type were highly correlated and, thus, not informative for

301 identifying fully polarized cells. To overcome this issue, we calculated the ‘embedding shift’ as the vector
302 difference between the cell embedding of each cell and the cell embedding of unpolarized cell center, which
303 represented the cell embedding change from the unpolarized state. We then calculated the cosine similarity
304 between the embedding shift of each cell with the rest cells from the samples treated with driving cytokines. Fully
305 polarized cells had to secure a minimum mean cosine similarity, ranging from 0.08 to 0.2 among different cell
306 types. Those cells that satisfied with all criteria were considered fully polarized cells and then were used for
307 training machine learning models alongside unpolarized cells from PBS-treated samples.

308 **Training and testing machine learning models**

309 When training machine learning models to classify unpolarized cells and polarized cells, we tested several models
310 including: 1) logistic regression using the ‘glm’ function from R package ‘stats’, 2) SVM using the ‘svm’ function
311 from R package ‘e1071’, 3) random forest using the ‘randomForest’ function from R package ‘randomForest’,
312 and 4) semi-supervised learning approach. For each cell type, 70% cells were randomly selected for training and
313 the remaining 30% for testing. During training, unpolarized cells were labeled with a polarization score of 0,
314 while the fully polarized cells identified in the previous step were labeled with a polarization score of 1. For
315 binary classification models, the predicted probability to the fully polarized state was used as the polarization
316 score. For regression models, the output prediction was clamped to a range of 0 to 1 and used as the polarization
317 score. We repeated the training and testing for 20 times and calculated the mean AUROC values for each machine
318 learning model. SVM showed the best performance with the highest mean AUROC values across all polarization
319 states. The final SVM models in Scupa were trained using the ‘svm’ function from R package ‘e1071’, with a
320 linear kernel and eps-regression type³³.

321 For the semi-supervised learning approach, we included the cells from the samples treated with cytokines other
322 than the polarization state-driving cytokines as unlabeled data. We first trained supervised machine learning
323 models (logistic regression, SVM, random forest) on labeled data: unpolarized cells and fully polarized cells from
324 the previous step. The trained models were then used to classify unlabeled cells as either unpolarized or polarized.
325 In the end, the final machine learning models were trained on the combined data from initial identification and
326 following prediction. In our comparison of the testing results from supervised models with semi-supervised
327 models, we found that the semi-supervised models generally had slight worse performance compared to the
328 corresponding supervised models, despite that they improved the performance on some polarization states with
329 small cell numbers (Supplementary Table 1). Therefore, we did not use the trained semi-supervised models for
330 prediction in the Scupa package.

331 **Estimating statistical significance of polarization**

332 When calculating the p-value of a cell for being polarized to a polarization state, we set the null hypothesis as H₀:
333 the observed cell is an unpolarized cell, and the alternative hypothesis as H₁: the observed cell is a polarized cell.
334 The polarization score serves as the test statistic. We used the polarization score distribution from unpolarized
335 cells of a specific cell type in the Immune Dictionary as the null distribution. The p-value is calculated as the
336 probability of obtaining a polarization score equal or greater than the observed polarization score from the null
337 distribution. This hypothesis testing is a non-parametric, with no assumptions about the cell embedding
338 distributions of the input dataset.

339 Scupa adjusts the p-values using the Benjamini-Hochberg procedure for the input dataset. We found that both
340 unadjusted p-values and adjusted p-values accurately classified polarized cells in the samples from *in vitro*
341 experiments (Supplementary Fig. 16). However, there were significantly fewer cells with adjusted p-values <0.05
342 in the samples from *in vivo* experiments. This discrepancy is likely due to varying extents of immune cell
343 polarization across different studies, and possibly more dynamic natures in the living cells. Multiple experimental
344 factors can influence the cytokine concentration in the tissue, thereby affecting the extent of immune cell

345 polarization. Consequently, using a stringent adjusted p-value threshold for identifying fully polarized cells would
346 result in few positive cells if cytokine concentration were low. Based on our analysis of four datasets included in
347 this study, it is recommended to use unadjusted p-values < 0.05 for identifying polarized cells in the samples from
348 *in vivo* experiments.

349 **Cross-dataset batch effect correction**

350 As a single-cell foundation model, UCE is robust to dataset and batch-specific artifacts, though cross-dataset
351 batch effects may still persist between the Immune Dictionary and other datasets. To enhance Scupa's
352 transferability to diverse datasets, we provide a straightforward and effective approach for cross-dataset batch
353 effect correction. UCE's capability allows us to represent cross-dataset batch effects as the difference in
354 unpolarized cell embeddings between two datasets. When there are untreated control and treated samples, cells
355 from the control samples could be specified as unpolarized cells. Scupa first calculates the center of reference
356 unpolarized cells' UCE cell embeddings \mathbf{c}_{ref} , and the center of input unpolarized cells' UCE cell embeddings
357 \mathbf{c}_{in} . For a cell k with UCE cell embeddings ($\mathbf{emb}_{k,original}$) from the input dataset, its cell embeddings are
358 adjusted to:

$$359 \quad \mathbf{emb}_{k,adjusted} = \mathbf{emb}_{k,original} - \mathbf{c}_{in} + \mathbf{c}_{ref}$$

360 This adjustment allows the learned representations of immune cell polarization to be applied to the adjusted cell
361 embeddings, bypassing complicated data integration processes and thus, preserving polarization information that
362 might be lost with scRNA-seq data integration methods^{34,35}. In the Scupa package, we implement this batch effect
363 correction approach and provide an optional parameter for users to specify unpolarized cells in the input dataset,
364 suitable for experimental designs with untreated healthy controls.

365 Regarding the two cytokine treatment datasets in Scupa evaluation, we used this batch effect correction approach
366 when analyzing immune cell polarization. In the IFN- β treated human PBMC scRNA-seq dataset, the cells from
367 the untreated sample were specified as unpolarized cells. Similarly, in the IL-2 treated mouse spleen scRNA-seq
368 dataset, the cells from the untreated mouse were specified as unpolarized cells. We found only slight differences
369 in the polarization analysis results when specifying unpolarized cells or not in both datasets. This evaluation
370 indicated the robustness of UCE and Scupa to cross-dataset batch effect correction. For the pan-cancer myeloid
371 cell dataset, we did not indicate unpolarized cells due to the absence of an untreated healthy control sample in the
372 dataset.

373 **Statistical analysis**

374 The calculation of p-values and adjusted p-values in Scupa is described in subsection "Deriving statistical
375 significance of polarization" above. When comparing the cell proportions between two conditions in the IL-2
376 treated mouse spleen scRNA-seq dataset, we used two-sided Fisher exact test. All statistical analyses were
377 performed in R (v 4.1.3).

378 **Data availability**

379 All data used in the study are from published studies. The Immune Dictionary scRNA-seq dataset was
380 downloaded from Single Cell Portal (https://singlecell.broadinstitute.org/single_cell/study/SCP2554/). The IFN- β
381 treated human PBMC scRNA-seq dataset with cell type annotation was downloaded using SeuratData
382 (<https://github.com/satijalab/seurat-data>). The rest datasets were downloaded from Gene Expression Omnibus
383 (GEO) database with following accession numbers: IL-2 treated mouse spleen scRNA-seq dataset (GEO206732)
384 and pan-cancer infiltrating myeloid cell scRNA-seq data (GSE154763). The processed datasets with generated
385 UCE cell embeddings are available at <https://zenodo.org/doi/10.5281/zenodo.13312248>.

386 **Code availability**

387 Scupa is available as a R package. The source code and vignettes of Scupa are freely available at
388 <https://github.com/bsml320/Scupa>. The code for generating the results is available at
389 <https://zenodo.org/doi/10.5281/zenodo.13312248>.

390 **Acknowledgements**

391 W.L. is a CPRIT Predoctoral Fellow in the Biomedical Informatics, Genomics and Translational Cancer Research
392 Training Program (BIG-TCR) funded by Cancer Prevention & Research Institute of Texas (CPRIT RP210045).
393 Z.Z. was partially supported by NIH grants (R01LM012806, U01AG079847, and R01CA276513) and Cancer
394 Prevention and Research Institute of Texas (CPRIT RP180734). We thank the authors of Immune Dictionary for
395 providing the high-quality reference dataset that was used for training Scupa, and thank the authors of UCE for
396 providing the foundation model that Scupa depends on.

397 **Author contributions**

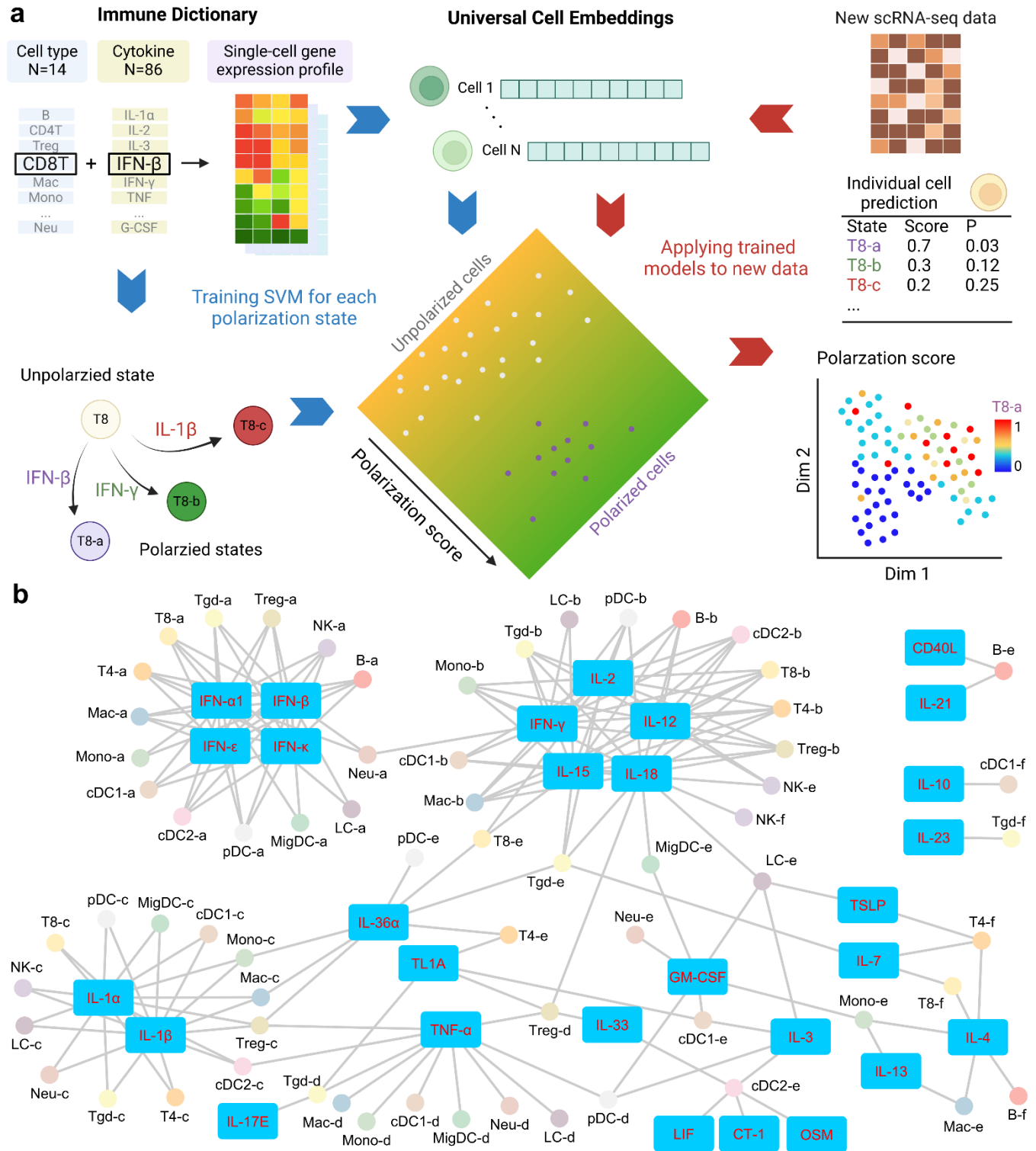
398 W.L. and Z.Z. conceived the project and wrote the manuscript. W.L. collected data, developed the method, and
399 performed the analysis.

400 **Competing interests**

401 The authors declare no competing interests.

402

403 **Figures**

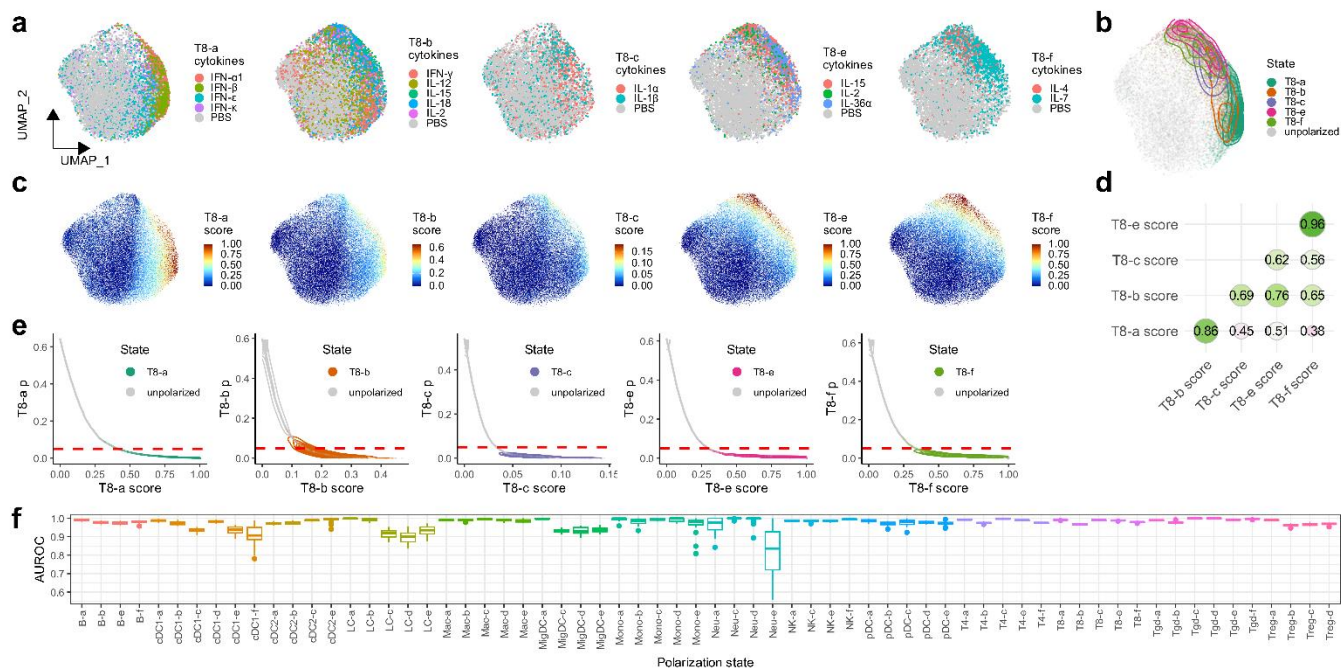


404

405 **Fig. 1. The framework of Scupa and cytokine-driven cell polarization states.** a. Scupa uses the Immune
 406 Dictionary as the reference and training set for measuring immune cell polarization. Major immune cell
 407 polarization states were defined in the Immune Dictionary according to driving cytokines and top marker genes.
 408 Universal Cell Embeddings (UCEs) were generated for the Immune Dictionary and any new scRNA-seq dataset
 409 for Scupa. Scupa trained support vector machine (SVM) models for various cell polarization states using the

410 Immune Dictionary, and then it applies the trained models to predict polarization in a new dataset. It outputs the
 411 polarization score and p-value for each individual cell. Created with BioRender.com. b. A network showing the
 412 driving cytokines and immune cell polarization states.

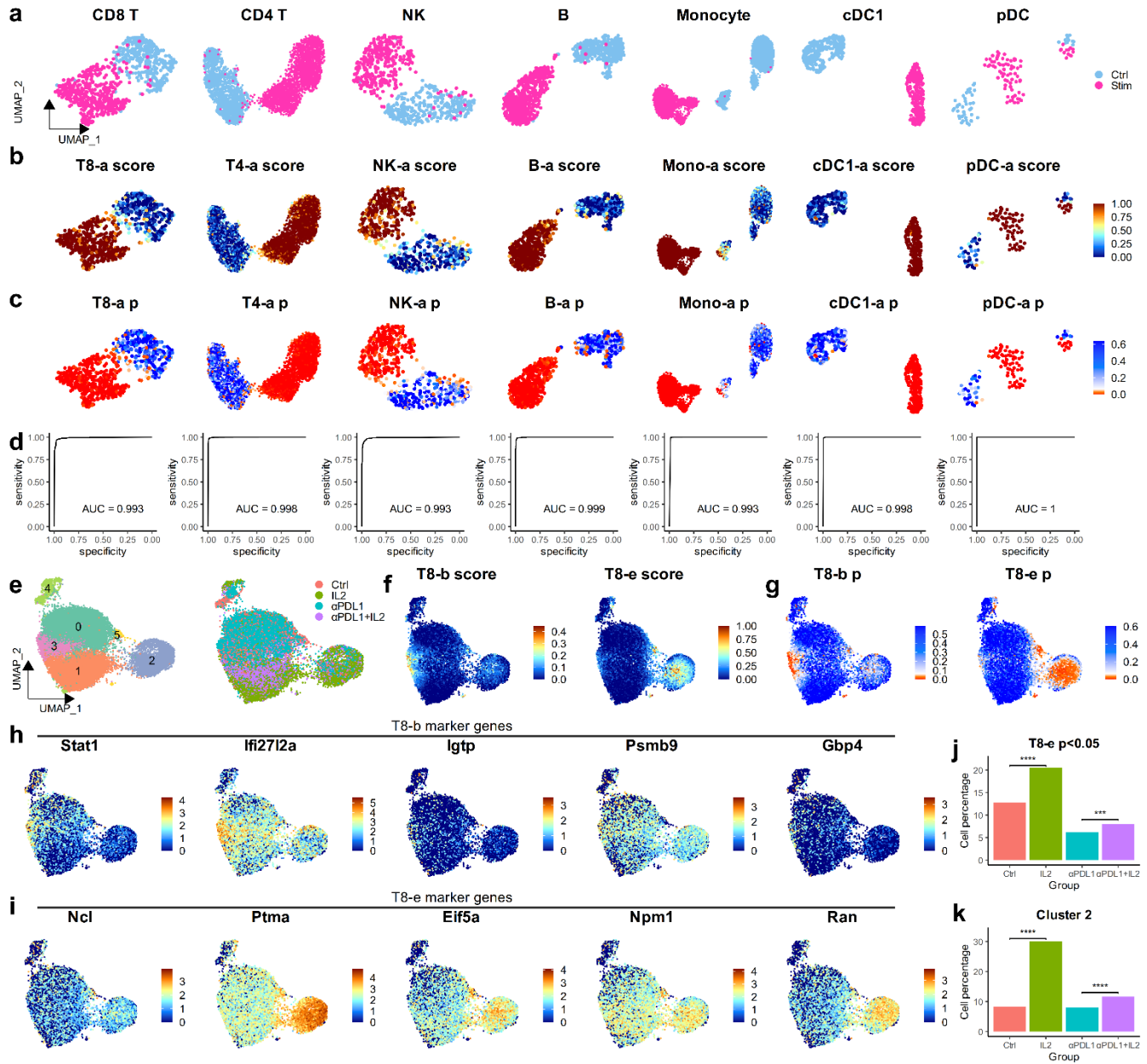
413



414

415 **Fig. 2. Scupa learns the representation of polarization states in CD8⁺ T cells and other cells.** a. Uniform
 416 manifold approximation and projection (UMAP) plots showing CD8⁺ T cells from the samples treated with
 417 driving cytokines of each polarization state or from the control samples treated with PBS. UMAPs are derived
 418 from UCE cell embeddings rather than gene expression. b. UMAP plot showing the distribution of fully polarized
 419 cells of each polarization state after filtering. c. UMAP plots showing polarization scores of each polarization
 420 state from Scupa prediction. d. The Spearman correlation coefficients between polarization scores from each two
 421 polarization states. e. The polarization scores and p-values of fully polarized cells of each polarization state and
 422 unpolarized cells. The red dashed line indicates p-value=0.05. f. Box plots showing the testing AUROC values in
 423 20 repeats across all polarization states.

424

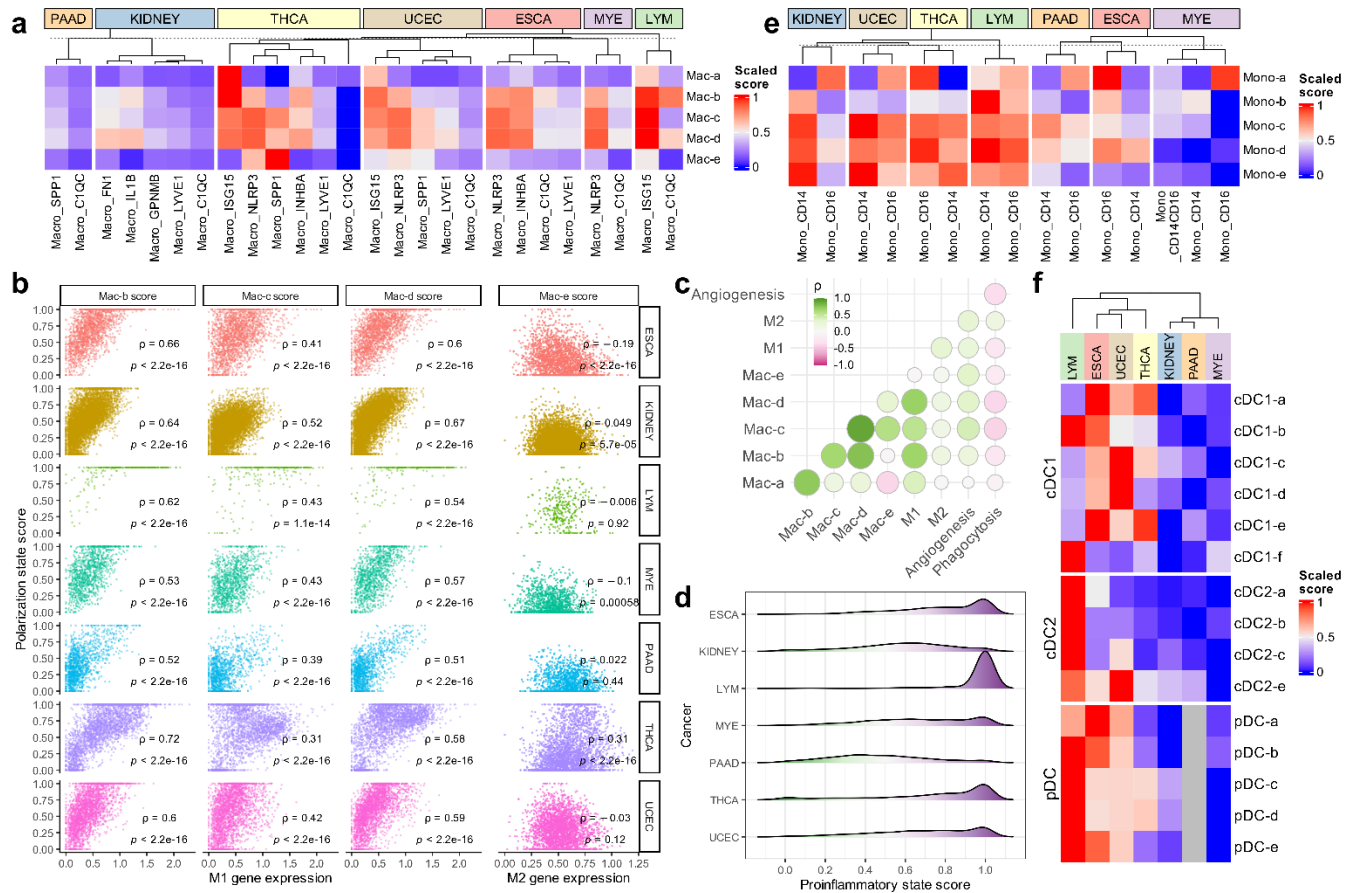


425

426 **Fig. 3. Validation of Scupa in *in vitro* and *in vivo* datasets.** a. UMAP plots showing the immune cells from the
 427 control and IFN- β -stimulated samples. UMAPs are derived from UCE cell embeddings rather than gene
 428 expression. Ctrl: untreated control sample. Stim: interferon-stimulated sample. b. UMAP plots showing the
 429 polarization scores in different immune cell types. The cells from IFN- β -stimulated samples generally had much
 430 higher scores. c. UMAP plots showing the polarization p-values in different immune cell types. The cells from
 431 IFN- β -stimulated samples generally had much lower p-values. d. ROC curves showing the superior performance
 432 of polarization scores for classifying immune cells from the control and IFN- β -stimulated samples. e. UMAP
 433 plots showing CD8+ T cells in different clusters and treatment groups. f. UMAP plots showing polarization scores
 434 of two IL-2-driving states, T8-b and T8-e. g. UMAP plots showing polarization p-values of two IL-2-driving
 435 states. h. The expression level of five T8-b marker genes. i. The expression level of five T8-e marker genes. j. The

436 percentage of cells with T8-e polarization p-value<0.05 among different treatment groups. k. The percentage of
 437 cells in cluster 2 among different treatment groups. ***: p<0.001. ****: p<0.0001.

438



439

440 **Fig. 4. Comparison of the polarization states of infiltrating myeloid cells across seven cancer types by**
 441 **Scupa.** a. The mean polarization scores in different macrophage subpopulations from different cancer types.
 442 Polarization scores are scaled across cancer types. b. The polarization scores of macrophage states Mac-b,
 443 Mac-c, and Mac-d were positively correlated with M1 signature gene expression. In contrast, the polarization scores of
 444 macrophage state Mac-e were not positively correlated with M2 signature gene expression. c. The correlation
 445 between macrophage state polarization scores, M1 signature gene expression, M2 signature gene expression,
 446 angiogenesis signature gene expression, and phagocytosis signature gene expression in all macrophages. d. LYM
 447 macrophages displays overall highest proinflammatory state scores. e. The mean polarization scores in different
 448 monocyte subpopulations from seven cancer types. f. The mean polarization scores in different DC
 449 subpopulations from seven cancer types. ESCA: esophageal carcinoma, KIDNEY: kidney cancer, LYM:
 450 lymphoma, MYE: myeloma, PAAD: pancreatic adenocarcinoma, THCA: thyroid carcinoma, UCEC: uterine
 451 corpus endometrial carcinoma.

452

453 **References**

454 1 Kourilsky, P. & Truffa-Bachi, P. Cytokine fields and the polarization of the immune response. *Trends*
 455 *Immunol* **22**, 502-509 (2001).

- 456 2 Murray, P. J. Macrophage Polarization. *Annu Rev Physiol* **79**, 541-566 (2017).
- 457 3 Ren, X. *et al.* COVID-19 immune features revealed by a large-scale single-cell transcriptome atlas. *Cell*
458 **184**, 1895-1913 e1819 (2021).
- 459 4 Liu, W. *et al.* Delineating COVID-19 immunological features using single-cell RNA sequencing.
460 *Innovation (Camb)* **3**, 100289 (2022).
- 461 5 Zhang, F. *et al.* Defining inflammatory cell states in rheumatoid arthritis joint synovial tissues by
462 integrating single-cell transcriptomics and mass cytometry. *Nat Immunol* **20**, 928-942 (2019).
- 463 6 Chung, W. *et al.* Single-cell RNA-seq enables comprehensive tumour and immune cell profiling in
464 primary breast cancer. *Nat Commun* **8**, 15081 (2017).
- 465 7 Garris, C. S. *et al.* Successful Anti-PD-1 Cancer Immunotherapy Requires T Cell-Dendritic Cell
466 Crosstalk Involving the Cytokines IFN-gamma and IL-12. *Immunity* **49**, 1148-1161 e1147 (2018).
- 467 8 Cui, A. *et al.* Dictionary of immune responses to cytokines at single-cell resolution. *Nature* **625**, 377-384
468 (2024).
- 469 9 Kharchenko, P. V., Silberstein, L. & Scadden, D. T. Bayesian approach to single-cell differential
470 expression analysis. *Nat Methods* **11**, 740-742 (2014).
- 471 10 Hicks, S. C., Townes, F. W., Teng, M. & Irizarry, R. A. Missing data and technical variability in single-
472 cell RNA-sequencing experiments. *Biostatistics* **19**, 562-578 (2018).
- 473 11 Cheng, S. *et al.* A pan-cancer single-cell transcriptional atlas of tumor infiltrating myeloid cells. *Cell* **184**,
474 792-809 e723 (2021).
- 475 12 Cochain, C. *et al.* Single-Cell RNA-Seq Reveals the Transcriptional Landscape and Heterogeneity of
476 Aortic Macrophages in Murine Atherosclerosis. *Circ Res* **122**, 1661-1674 (2018).
- 477 13 Theocharidis, G. *et al.* Single cell transcriptomic landscape of diabetic foot ulcers. *Nat Commun* **13**, 181
478 (2022).
- 479 14 Cui, H. *et al.* scGPT: toward building a foundation model for single-cell multi-omics using generative AI.
480 *Nat Methods* (2024).
- 481 15 Hao, M. *et al.* Large-scale foundation model on single-cell transcriptomics. *Nat Methods* (2024).
- 482 16 Rosen, Y. *et al.* Universal Cell Embeddings: A Foundation Model for Cell Biology. *bioRxiv*,
483 2023.2011.2028.568918 (2023).
- 484 17 Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587 e3529 (2021).
- 485 18 Kang, H. M. *et al.* Multiplexed droplet single-cell RNA-sequencing using natural genetic variation. *Nat*
486 *Biotechnol* **36**, 89-94 (2018).
- 487 19 Hashimoto, M. *et al.* PD-1 combination therapy with IL-2 modifies CD8(+) T cell exhaustion program.
488 *Nature* **610**, 173-181 (2022).
- 489 20 Engblom, C., Pfirschke, C. & Pittet, M. J. The role of myeloid cells in cancer therapies. *Nat Rev Cancer*
490 **16**, 447-462 (2016).
- 491 21 van Vlerken-Ysla, L., Tyurina, Y. Y., Kagan, V. E. & Gabrilovich, D. I. Functional states of myeloid cells
492 in cancer. *Cancer Cell* **41**, 490-504 (2023).
- 493 22 Najafi, M. *et al.* Macrophage polarity in cancer: A review. *J Cell Biochem* **120**, 2756-2765 (2019).
- 494 23 Ono, M. Molecular links between tumor angiogenesis and inflammation: inflammatory stimuli of
495 macrophages and cancer cells as targets for therapeutic strategy. *Cancer Sci* **99**, 1501-1506 (2008).
- 496 24 Skinnider, B. F. & Mak, T. W. The role of cytokines in classical Hodgkin lymphoma. *Blood* **99**, 4283-
497 4297 (2002).
- 498 25 Ma, R. Y., Black, A. & Qian, B. Z. Macrophage diversity in cancer revisited in the era of single-cell
499 omics. *Trends Immunol* **43**, 546-563 (2022).
- 500 26 Parker, B. S., Rautela, J. & Hertzog, P. J. Antitumour actions of interferons: implications for cancer
501 therapy. *Nat Rev Cancer* **16**, 131-144 (2016).
- 502 27 Kaufmann, S. H. E., Dorhoi, A., Hotchkiss, R. S. & Bartenschlager, R. Host-directed therapies for
503 bacterial and viral infections. *Nat Rev Drug Discov* **17**, 35-56 (2018).
- 504 28 Berraondo, P. *et al.* Cytokines in clinical cancer immunotherapy. *Br J Cancer* **120**, 6-15 (2019).
- 505 29 Ruytinx, P., Proost, P., Van Damme, J. & Struyf, S. Chemokine-Induced Macrophage Polarization in
506 Inflammatory Conditions. *Front Immunol* **9**, 1930 (2018).

- 507 30 Xuan, W., Qu, Q., Zheng, B., Xiong, S. & Fan, G. H. The chemotaxis of M1 and M2 macrophages is
508 regulated by different chemokines. *J Leukoc Biol* **97**, 61-69 (2015).
- 509 31 Colegio, O. R. *et al.* Functional polarization of tumour-associated macrophages by tumour-derived lactic
510 acid. *Nature* **513**, 559-563 (2014).
- 511 32 El Kasmi, K. C. *et al.* Adventitial fibroblasts induce a distinct proinflammatory/profibrotic macrophage
512 phenotype in pulmonary hypertension. *J Immunol* **193**, 597-609 (2014).
- 513 33 Meyer, D. & Wien, F. Support vector machines. *R News* **1**, 23-26 (2001).
- 514 34 Polanski, K. *et al.* BBKNN: fast batch alignment of single cell transcriptomes. *Bioinformatics* **36**, 964-
515 965 (2020).
- 516 35 Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902 e1821 (2019).
- 517