

1 CAFE: An Integrated Web App for High-Dimensional Analysis and Visualization in Spectral Flow  
2 Cytometry

3 Md Hasanul Banna Siam<sup>1</sup>, Md Akkas Ali<sup>1</sup>, Donald Vardaman III<sup>1</sup>, Satwik Acharyya<sup>2,3</sup>, Mallikarjun  
4 Patil<sup>1</sup>, Daniel J. Tyrrell<sup>1</sup>

5 <sup>1</sup>Department of Pathology, Heersink School of Medicine, University of Alabama at Birmingham,  
6 Birmingham, AL, 35205 USA

7 <sup>2</sup>Department of Biostatistics, University of Alabama at Birmingham, Birmingham, AL, 35233 USA

8 <sup>3</sup>Department of Biomedical Informatics and Data Science, University of Alabama at Birmingham,  
9 Birmingham, AL, 35294 USA

10

11 *Corresponding author:*

12 Daniel J. Tyrrell, PhD

13 University of Alabama at Birmingham

14 PBMR2 #504

15 901 19<sup>th</sup> St. S.

16 Birmingham, AL 35205

17 Email: [danieltyrrell@uabmc.edu](mailto:danieltyrrell@uabmc.edu)

18 Running Title: Spectral Flow Cytometry Analysis Platform

19 **Abstract**

20 Spectral flow cytometry provides greater insights into cellular heterogeneity by simultaneous  
21 measurement of up to 50 markers. However, analyzing such high-dimensional (HD) data is  
22 complex through traditional manual gating strategy. To address this gap, we developed CAFE as  
23 an open-source Python-based web application with a graphical user interface. Built with Streamlit,  
24 CAFE incorporates libraries such as Scanpy for single-cell analysis, Pandas and PyArrow for  
25 efficient data handling, and Matplotlib, Seaborn, Plotly for creating customizable figures. Its robust  
26 toolset includes density-based down-sampling, dimensionality reduction, batch correction,  
27 Leiden-based clustering, cluster merging and annotation. Using CAFE, we demonstrated analysis  
28 of a human PBMC dataset of 350,000 cells identifying 16 distinct cell clusters. CAFE can generate  
29 publication-ready figures in real time via interactive slider controls and dropdown menus,  
30 eliminating the need for coding expertise and making HD data analysis accessible to all. CAFE is  
31 licensed under MIT and is freely available at <https://github.com/mhbsiam/cafe>.

32

33 *Keywords:* Cytometry, Bioinformatics, Protocol, Leiden,

34 *Word count:* 4633

35

## 36 **Introduction:**

37 Flow cytometry is a widely used technique in immunology to identify and quantify immune cells  
38 based on specific surface markers<sup>1</sup>. The development of spectral flow cytometry (SFCM) has  
39 further expanded immunophenotyping capabilities allowing the simultaneous analysis of a greater  
40 number of parameters through the complete emission spectra of fluorophores<sup>1</sup>. Compared to  
41 conventional flow cytometry, SFCM uses spectral unmixing algorithms to deconvolute the  
42 overlapping signals and achieves enhanced resolution and sensitivity to distinguish between  
43 different cell populations<sup>2</sup>. SFCM can incorporate broader range of antibodies with up to 50 colors  
44 in a single panel improving upon the conventional FCM where the number of parameters is limited  
45 by the instrument constraints<sup>3</sup>. Incorporating more parameters substantially increases the  
46 complexity in gating strategy which largely relies on established convention and prior  
47 knowledge<sup>4,5</sup>. Additional gating steps and combinations of markers used to subset cells  
48 complicate the interpretation of such high-dimensional data. Several clustering methods are  
49 available to identify cell populations such as FlowSOM<sup>6</sup>, xShift<sup>7</sup>, SPADE<sup>8</sup> and Phenograph<sup>9</sup>.  
50 SPADE and FlowSOM utilize hierarchical clustering with the latter employing self-organizing maps  
51 (SOMs) to cluster cells, whereas xShift detects clusters based on shifts in local cell density<sup>6-8</sup>.  
52 Phenograph, by contrast, constructs a K-nearest neighbor graph and applies the Louvain  
53 algorithm to identify cell clusters, but Louvain can produce poorly connected or disconnected  
54 communities<sup>9,10</sup>.

55 Recently, the Leiden clustering algorithm has emerged as a faster and more accurate alternative  
56 to improve community detection in networks<sup>10</sup>. Single-cell RNA sequencing (scRNA-seq) tools:  
57 Seurat<sup>11</sup> (R) and Scanpy<sup>12</sup> (Python) have integrated Leiden algorithms for community detection.  
58 However, running Leiden within Seurat resulted in drawbacks including higher memory usage,  
59 longer calculation time and random crashes in docker containers<sup>13</sup>. Scanpy resolves these issues,  
60 and unlike Seurat, Scanpy improves visualization quality by using consistent KNN and SNN  
61 graphs for both clustering and uniform manifold approximation and projection (UMAP)<sup>13,14</sup>. In  
62 February 2020, Phenograph version 1.5.3 was released, which incorporated an option to use  
63 Leiden for clustering; however, the default parameter is set to Louvain through the latest release  
64 (v.1.5.7). In our previous work, we showed that the use of Leiden algorithm in community detection  
65 for SFCM data provides superior result to Phenograph (Louvain), FlowSOM, and xShift<sup>15</sup>.  
66 Currently, there is a scarcity of open-source tools to utilize Leiden algorithm for SCFM data  
67 analysis<sup>16</sup>.

68 Here we present CAFE, Cell Analyzer for Flow Experiment, a user-friendly web application  
69 developed in Python that works across Windows, MacOS, and Linux. The app is lightweight and  
70 can perform high-dimensional SFCM data analysis using a standard computing machine (i.e.,  
71 Apple M1 chip with 16gb RAM), and it provides the flexibility to be deployed on HPC clusters for  
72 enhanced scalability. Once installed, the tool runs entirely offline and does not require an active  
73 internet connection to load files. This also enables users to maintain compliance with data  
74 security, especially with protected health information (PHI) when analyzing patient samples.  
75 CAFE can be used to process data, reduce dimension, batch correction, run Leiden clustering,  
76 perform statistics and generate a wide range of figures. Figures can be adjusted and viewed within  
77 the tool in real time. Additionally, the tool offers Kernel Density Estimation (KDE)-based data  
78 downsampling, advanced clustering with predefined markers, cluster quality evaluation, merging  
79 subclusters into metaclusters, and cell type annotation. Designed as an open-source interactive  
80 data analysis platform, CAFE enables biologists with no-coding experience to analyze SFCM data  
81 and create publication quality visualization with customizable parameters. CAFE is freely  
82 available to download at: <https://github.com/mhbsiam/caffe>

83

## 84 **Methods:**

### 85 ***Implementation***

86 The CAFE webtool was developed using Python programming language due to its compatibility  
87 with Scanpy library<sup>12</sup>. Figure 1 illustrates the components and workflows of CAFE. Streamlit  
88 (streamlit.io) library was used to develop the web interface that provides dynamic updates based  
89 on user inputs in the graphical user interface (GUI) without writing or editing code directly.  
90 Streamlit was chosen due to its simplicity in development and compatibility with other Python  
91 libraries across operating systems. Streamlit v1.39.0 is compatible with any modern HTML5 web  
92 browser. For data loading and processing, we relied on Pandas v2.2.3 with PyArrow v18.0.0 which  
93 achieves faster data loading and processing compared to Pandas alone. We used NumPy v1.26.4  
94 for data type and range selection, RGB array creation for color handling, and grid setup for  
95 subplots. Seaborn v0.13.2, Matplotlib v3.9.2, and Plotly v4.24.1 libraries were used for data  
96 visualization and users are provided with options to adjust parameters: plot size, color profile, and  
97 output formats in PNG, JPG, SVG, or PDF. CAFE integrates AnnData<sup>12</sup>, a widely used framework  
98 in single-cell RNA sequencing analysis that allows for efficient storage and manipulation of both

99 sparse and dense matrices along with metadata. CAFE outputs an AnnData object which can be  
100 used outside of CAFE if users wish to deposit their data with analysis or perform custom analyses  
101 using other tools.

102 The Scanpy library was used to perform key analyses including dimension reduction, batch  
103 correction, and Leiden clustering. The user has the option to reduce dimension (`sc.tl.pca`) of the  
104 data through Principal Component Analysis (PCA) or skip it. PCA is a linear dimensionality  
105 reduction technique that retains the global structure of the data by capturing the variances across  
106 all dimensions. Because the app performs PCA through Scanpy library, by default, the number of  
107 components retained is limited to the lesser of the two values: the number of cells or the number  
108 of markers. Also, the Singular Value Decomposition (SVD) solver was set to “auto” that chooses  
109 the most appropriate solver based on the size of the dataset; however, users have options to set  
110 a percentage of variances they want to retain, and the type of solver used. The reduced dataset  
111 is stored and can be further processed for batch correction (`sc.pp.combat`) using ComBat  
112 (Combined Batch)<sup>17</sup>. This is particularly useful if a user has collected samples in different batches  
113 as the algorithm standardizes the data by making it comparable and removing unwanted  
114 variability.

115 To group cells into distinct clusters based on marker expression profiles, Leiden clustering is run  
116 (`sc.tl.leiden`) and users can select either ‘iGraph’ or ‘leidenalg’ algorithm flavor<sup>10,18</sup>. To define  
117 clustering resolution, a user can choose from 0.01 to 2.0 where the lowest value provides the  
118 lowest number of clusters. The user can fine tune Leiden calculation by altering the number of  
119 neighbors and minimum distance values in Uniform Manifold Approximation and Projection  
120 (UMAP) calculation within the app. CAFE generates AnnData object (H5AD) file, CSV outputs,  
121 and visual outputs including UMAP plots, dot-heatmaps, expression pots, and barplots as high-  
122 resolution images and provides download buttons to save them to a desired folder. The app allows  
123 various visualization settings, with changes made and displayed immediately within the app. The  
124 generated AnnData object can be further used to perform a range of statistical analyses.

125 The app includes advanced functionalities for clustering and cluster evaluation. Because setting  
126 up appropriate values for Leiden resolution and UMAP parameters is central to obtaining quality  
127 clustering results, a user can leverage CAFE’s Cluster Evaluation tab to generate multiple  
128 AnnData files with various combinations of these parameters and compare UMAP plots as well  
129 as Silhouette score, Calinski-Harabasz score, Davies-Bouldin score, and Elbow method to assess  
130 clustering results. Besides, CAFE provides clustering with pre-selected markers, merging

131 subclusters into metaclusters, and annotation of clusters directly within the app. The advanced  
132 Downsampling tab offers to downsample (e.g. 20,000 events per sample) data using a PCA-KDE  
133 based method. This approach combines PCA and KDE (Kernel Density Estimation) based  
134 algorithm from `scipy.stats` using Gaussian kernel function and silverman bandwidth<sup>19</sup>. KDE is  
135 applied to the PCA-transformed data to estimate the density of data points. Based on these  
136 density estimates, the code probabilistically downsamples data, thus reducing sampling bias and  
137 preserving original data distribution. This method offers an informed approach compared to simple  
138 random downsampling, and it can be used to filter out noise while retaining meaningful biological  
139 information in a smaller dataset.

140

### 141 ***Statistical analysis***

142 In the visualization tab under statistical analyses section, users can perform different statistical  
143 tests and generate plots. The Shapiro-Wilk test from “`scipy.stats`” is used to determine if marker  
144 expression within each group or cluster follows a normal distribution. Based on these results, the  
145 app recommends either parametric (T-test) or non-parametric (Mann-Whitney U test) tests using  
146 “`scipy.stats`”. For comparing multiple clusters, the app allows users to perform ANOVA  
147 (`scipy.stats.f_oneway`) or Kruskal-Wallis (`scipy.stats.kruskal`) tests. To reduce statistical artifacts,  
148 multiple testing correction is applied using the Benjamini-Hochberg False Discovery Rate (FDR)  
149 through “`statsmodels.stats.multitest`”. Additionally, effect size measures are computed to  
150 complement statistical p-values, with users able to choose between parametric tests (Cohen’s d)  
151 or non-parametric tests (Cliff’s Delta). Cohen’s d is calculated using basic functions from Numpy,  
152 while Cliff’s Delta is computed with the “`cliffs_delta`” package. To assess associations between  
153 clusters and groups, we used Chi-square testing from “`scipy.stats.chi2_contingency`” and  
154 contingency tables with “`pandas.crosstab`”. Residual calculations were displayed in Streamlit as  
155 tables to help users understand which clusters are more prevalent within certain groups.

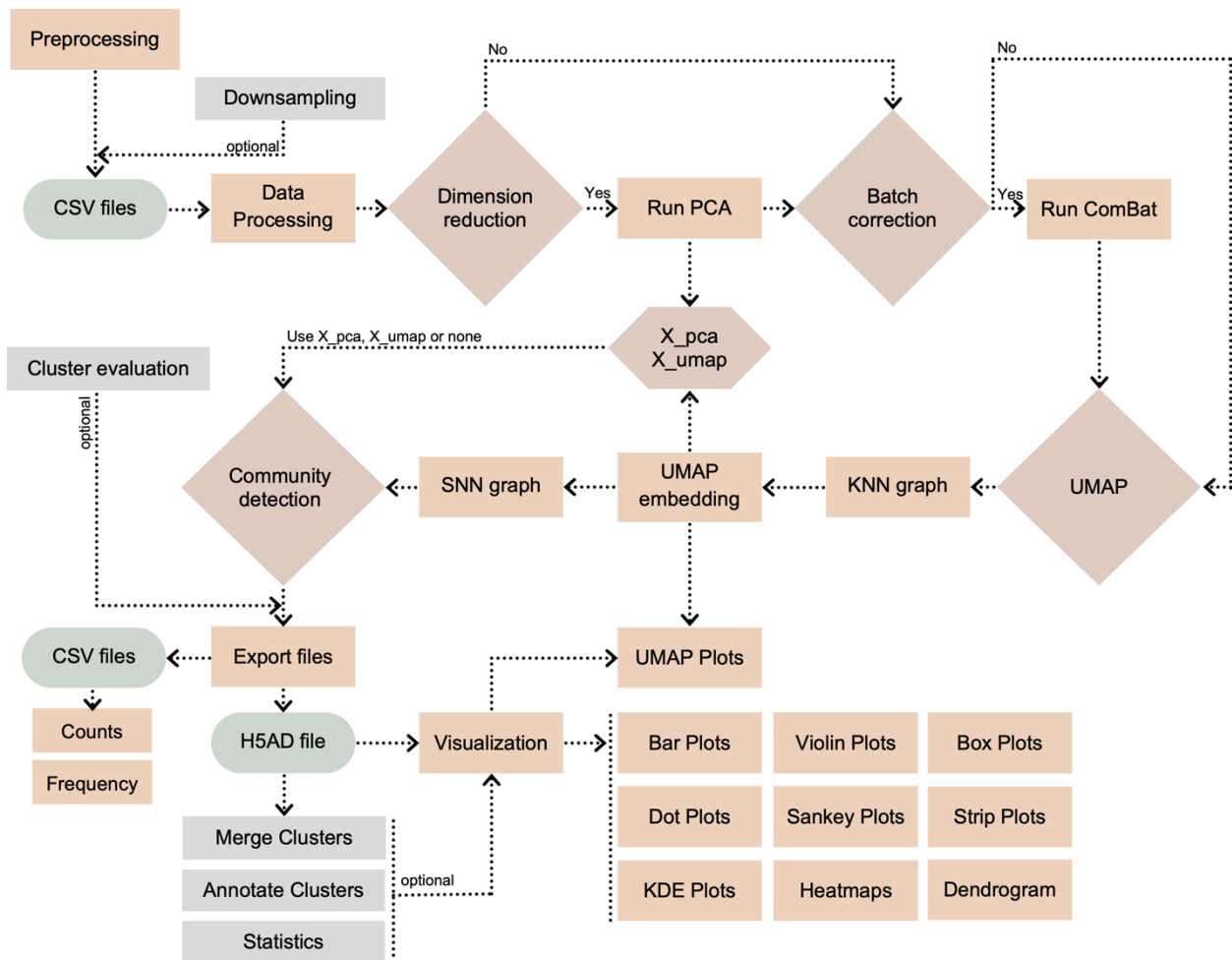
156

### 157 ***Performance and reproducibility***

158 We have set a global setting for Scanpy (`sc.settings.n_jobs = -1`) to use all available CPU cores.  
159 For advanced clustering, multi-threading was achieved using Python’s `joblib` library. Two other  
160 libraries were used, `watchdog v5.0.2` and `iGraph v0.10.8`<sup>18</sup>. `Watchdog` helps in monitoring file

161 change events and improves performance of Streamlit by providing real-time feedback. iGraph is  
 162 designed to handle complex networks and graph operations and is used by Scanpy as part of the  
 163 Leiden clustering to perform graph operations. We recommend 'iGraph' over 'leidenalg' as iGraph  
 164 is implemented in C and achieves advantages in performance compared to high-level interpreted  
 165 languages such as Python. To export Plotly figures, we have used Kaleido engine v0.2.1. We  
 166 have tested the app with various datasets using an Apple M3 Pro System with 18GB of random-  
 167 access memory (RAM). CAFE is primarily intended to be used using local computer; however, it  
 168 can be scaled up using any High-Performance Computing (HPC) system that supports an HTML5  
 169 web browser. We also provided scripts in our GitHub page to generate AnnData with dimension  
 170 reduction and Leiden clustering through command-line interface (CLI) based HPC systems.

171



172

173 **Figure 1: The flowchart outlines steps and components of CAFE's workflow.**

174 Preprocessing includes compensation, data scaling/transformation using a standard FCM  
 175 software and scaled CSV files are then exported and renamed as Sample\_Group.csv. Data  
 176 processing performs error checks and concatenation of CSV files into an AnnData object/H5AD

177 file. Major steps requiring user input include dimension reduction, batch correction, UMAP  
178 (UMAP Uniform Manifold Approximation and Projection) and community detection. Outputs are  
179 downloadable as CSV, H5AD, PNG, JPG, SVG and PDF files.  
180

## 181 **Results:**

182 To demonstrate the functionality of the app, we have analyzed 35-color spectral flow cytometry  
183 data (Publicly available at FlowRepository: FR-FCM-Z3WR) of human peripheral blood  
184 mononuclear cells (PBMC) obtained from COVID-19 hospitalized patients and healthy controls<sup>20</sup>.  
185 A total of 10 samples were analyzed with 5 from each group. For best practices, we installed and  
186 ran CAFE through Pixi package manager. Users can also install and run the app using Anaconda  
187 package manager as described in our Github documentation. Once initiated through a terminal  
188 (*pixi run cafe, or python cafe.py*), a web browser opens with the CAFE app at localhost on port  
189 8501. The default data loading limit is set to 3GB, but a user can change the value from the  
190 *cafe.py* script if necessary.

## 191 ***Data processing***

192 The uploaded public data were available as doublets-debris removed and CD45+ gated; so we  
193 obtained the CSV files just by exporting scaled values from FlowJo v10.10.0. Data scaling is  
194 generally recommended for high resolution clustering but there may be instances where users  
195 may use raw values. Data can be similarly exported from other flow cytometry software such as  
196 FCS Express. It is required that flow cytometry data have proper compensation. We recommend  
197 manual inspection of flow cytometry data and removal of debris, dead cells, and doublets prior to  
198 exporting the scaled files. A user can also gate on appropriate cell type and export the data to  
199 obtain more focused clustering results. To streamline downstream analysis, we have implemented  
200 a naming convention for the CSV files. Each CSV file name must begin with a unique  
201 "SampleName" followed by "GroupName", separated by an underscore; for instance,  
202 "Sample01\_Control.csv" and "Sample02\_Treatment.csv". After loading the data, the app will  
203 import the required libraries and perform initial checks for data structure and incorporate  
204 SampleID and GroupName into the dataframe based on the CSV file names. Within the  
205 dataframe, rows containing any missing values are skipped and anomalies in data structure are  
206 reported. In this study, we used the advanced KDE-based downsampling option in CAFE to  
207 downsample data to 35,000 cells per sample for a total of 350,000 cells and 12.25M data points  
208 (number of cells multiplied by number of markers). This is an optional step prior to data

209 processing. After loading the files, the app processed (7.8 sec) and combined the expression data  
210 and metadata without errors to create an AnnData object.

### 211 ***Dimension reduction and batch effects***

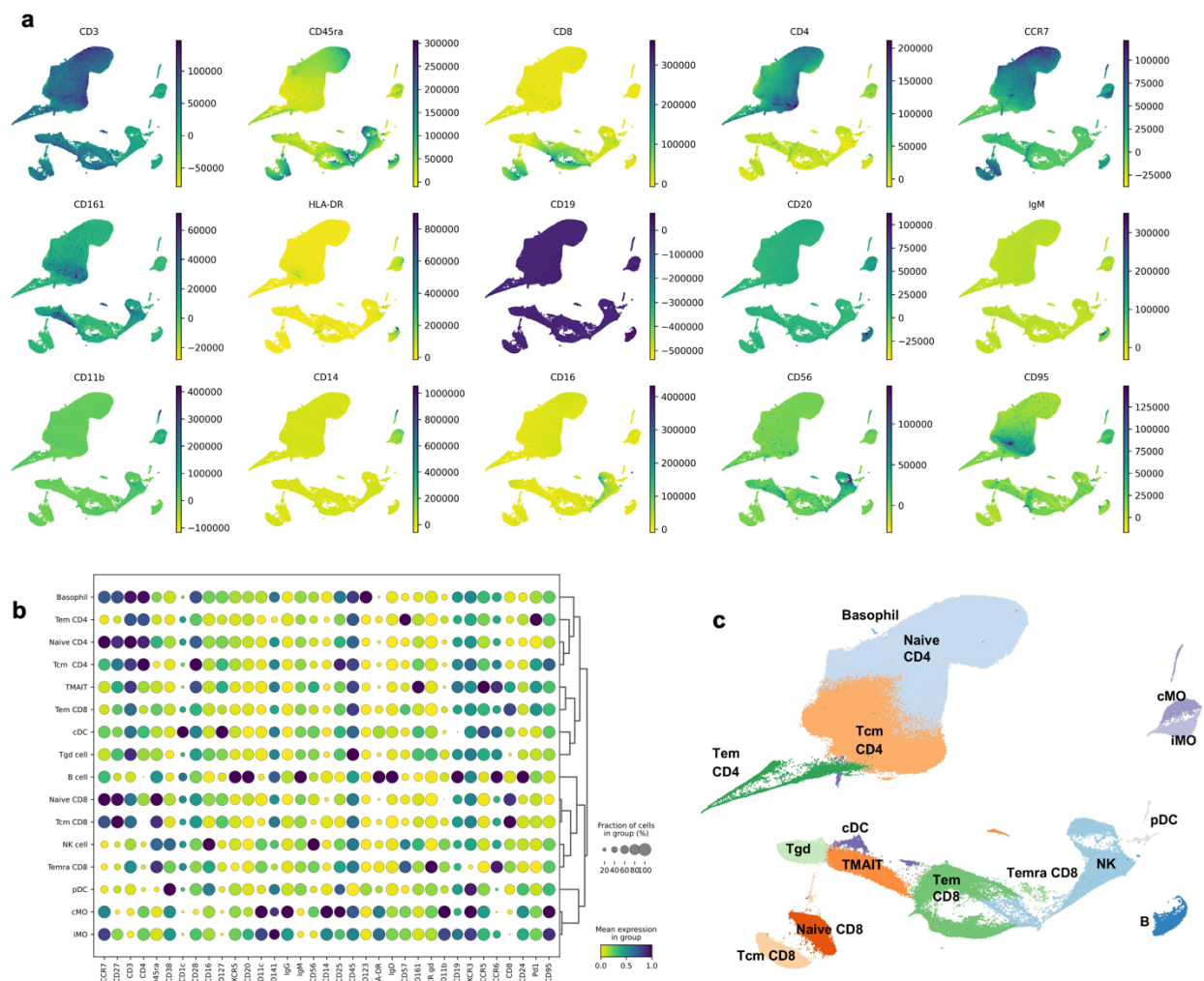
212 After generating the AnnData object (H5AD file), we selected dimension reduction using PCA with  
213 SVD solver set to auto and retained components with 95% variance. The app ran PCA (2.16 sec),  
214 kept 12 components and generated (PC1 x PC2) graphs by Groups. Depending on the data size,  
215 a user can choose from auto, full, arpack, and randomized SVD solver. Randomized, for example,  
216 is better suited for larger datasets as it provides a balance between speed and accuracy. For  
217 batch correction, we applied ComBat (1.06 sec) and proceeded to Leiden clustering.

### 218 ***Leiden clustering and metaclustering***

219 For the dataset, we applied Leiden resolution of 1.0 with flavor set to iGraph, UMAP n\_neighbors  
220 to 15, min\_dist to 0.1 and distance calculation method as Euclidean. A user has the option to use  
221 a slider control to choose from resolution values 0.01 to 2.0. To find the optimal resolution, we  
222 initially made use of Advanced Cluster Evaluation option in CAFE to generate a series of AnnData  
223 files with varied Leiden resolution and n\_neighbor values and observed the UMAPs to find distinct  
224 clusters that are biologically meaningful for the dataset. With Leiden resolution of 1.0, we initially  
225 obtained a total of 30 clusters for the PBMC dataset which took 11.5 minutes for calculation. Once  
226 clustering was completed, CAFE generated a frequency table of each sample by Leiden cluster  
227 for the number of cells, frequency of cells, and median fluorescence intensity of each marker for  
228 each cluster. Using these 3 tables, users can perform statistical analyses to compare cluster count  
229 and frequency by groups and expression of marker proteins within clusters by group. Using the  
230 Advanced Cluster Merging option, we merged the subclusters with similar profile into  
231 corresponding metaclusters resulting in a new total of 16 clusters.

232





233

234 **Figure 2: Profiling of Human PBMCs Reveals Distinct Immune Subpopulations and Marker**  
235 **Expression Patterns.** (a) UMAP plots showing selected marker expression intensities across all  
236 cells in the UMAP space to highlight lineage-specific marker distribution. (b) Dot plot of all marker  
237 expression across all identified PBMC cell types. Dendrogram highlighted distinct marker-based  
238 groupings. (c) UMAP visualization showing 16 distinct clusters with annotated cell types including  
239 Naive CD4 and CD8 T cells, central memory CD4 and CD8 T cells (Tcm), effector memory CD8  
240 T cells (Tem), terminally differentiated effector memory CD8 T cells (Temra), mucosal-associated  
241 invariant (MAIT) T cells, classical monocytes (cMO), intermediate monocytes (iMO), B cells, NK  
242 cells, gamma delta ( $\gamma\delta$ ) T cells (Tgd), conventional dendritic cell (cDC) and plasmacytoid  
243 dendritic cell (pDC).

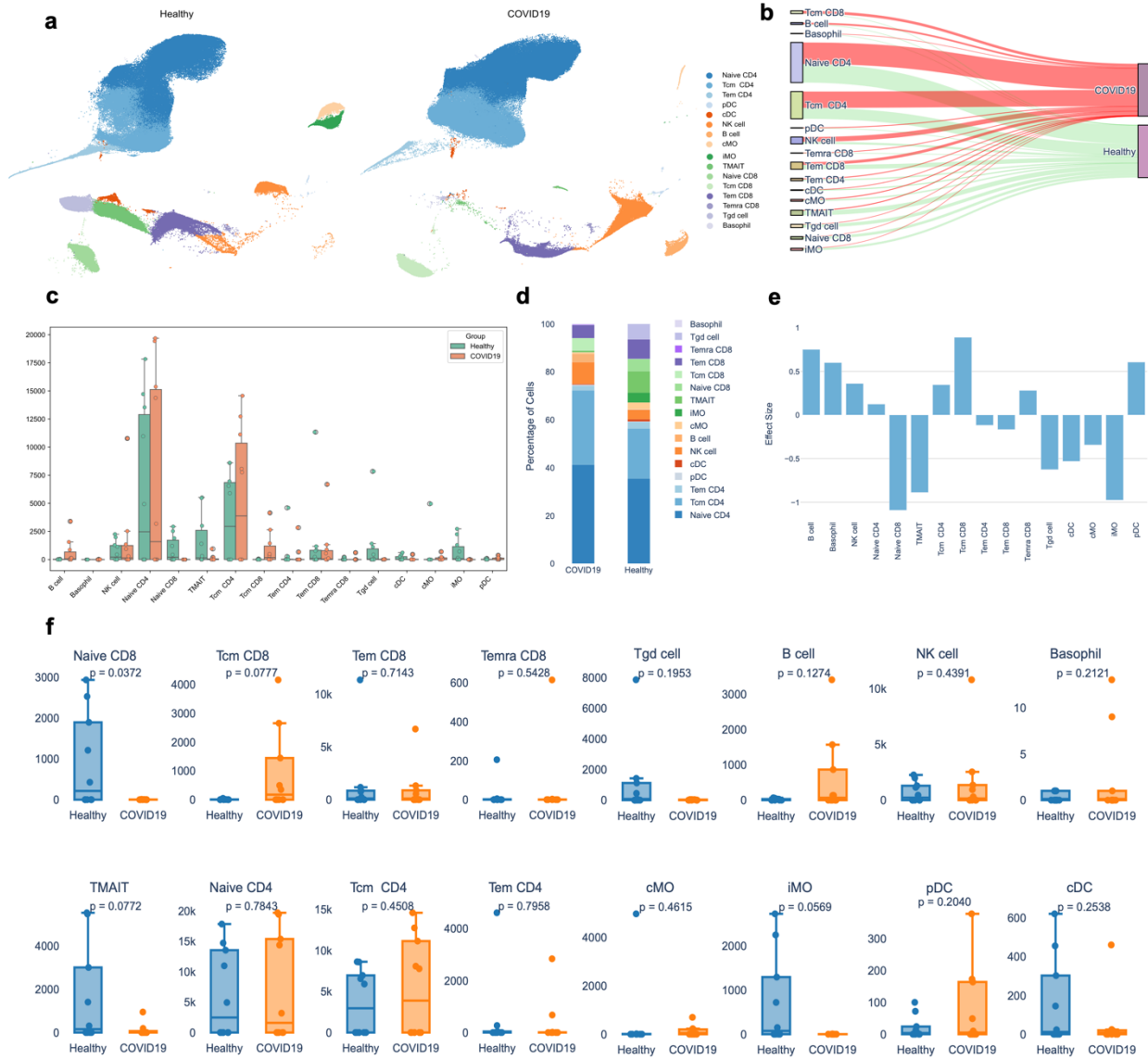
#### 244 ***Characterization of PBMC Subpopulations***

245 To characterize the phenotypic properties, we examined the expression of surface markers for  
246 each identified cluster by protein expression UMAP plots (Figure 2a). We found high CD3  
247 expression in T cells clusters with CD4 and CD8 expression showing corresponding T cell

248 subtypes. CD8<sup>+</sup> effector memory (Tem) and central memory (Tcm) subsets were differentiated by  
249 high CCR7 and CD27 expression in Tcm. We used CD45RA expression to identify terminally  
250 differentiated Tem cells (Temra). Monocyte clusters were identified by CD14 and CD16  
251 expression, distinguishing classical monocytes (cMO) from other monocyte subsets, while natural  
252 killer (NK) cells showed high levels of CD56, corresponding with CD56<sup>Bright</sup> NK cells. Based on  
253 shared marker profiles and hierarchical ranking, T cell subsets (Tem, Tcm, and Temra) formed a  
254 distinct grouping separate from B cells and myeloid-derived cells, reflecting the differential  
255 expression of lineage-specific markers.

256 Based on the expression profiles of marker proteins, we annotated the clusters using CAFE's  
257 Advanced Annotation tab and classified them into 16 distinct cell types. We also used the dotplot  
258 to confirm annotations of the cell types (Figure 2b). For instance, the B cell-specific marker CD19  
259 and CD20 were used to identify the B cell cluster, the CD14 marker to identify monocytes, and  
260 the CD16 marker to identify NK cells. High CD20 expression in B cell cluster indicated their mature  
261 stage in immune response. Our annotated UMAP (Figure 2c) shows well-defined clusters that  
262 correspond to PBMC lineages, including Naive CD4<sup>+</sup> and CD8<sup>+</sup> T cell and Tcm for both CD4<sup>+</sup>  
263 and CD8<sup>+</sup> subsets. We also identified Tem and Temra cells, as well as mucosal associated  
264 invariant CD8<sup>+</sup> T cell (MAIT). We also identified cMO, intermediate monocytes (iMO), B cells, NK  
265 cells, Gamma delta ( $\gamma\delta$ ) T cells (Tgd) and dendritic cell (DC) types (cDC and pDC).

266



267

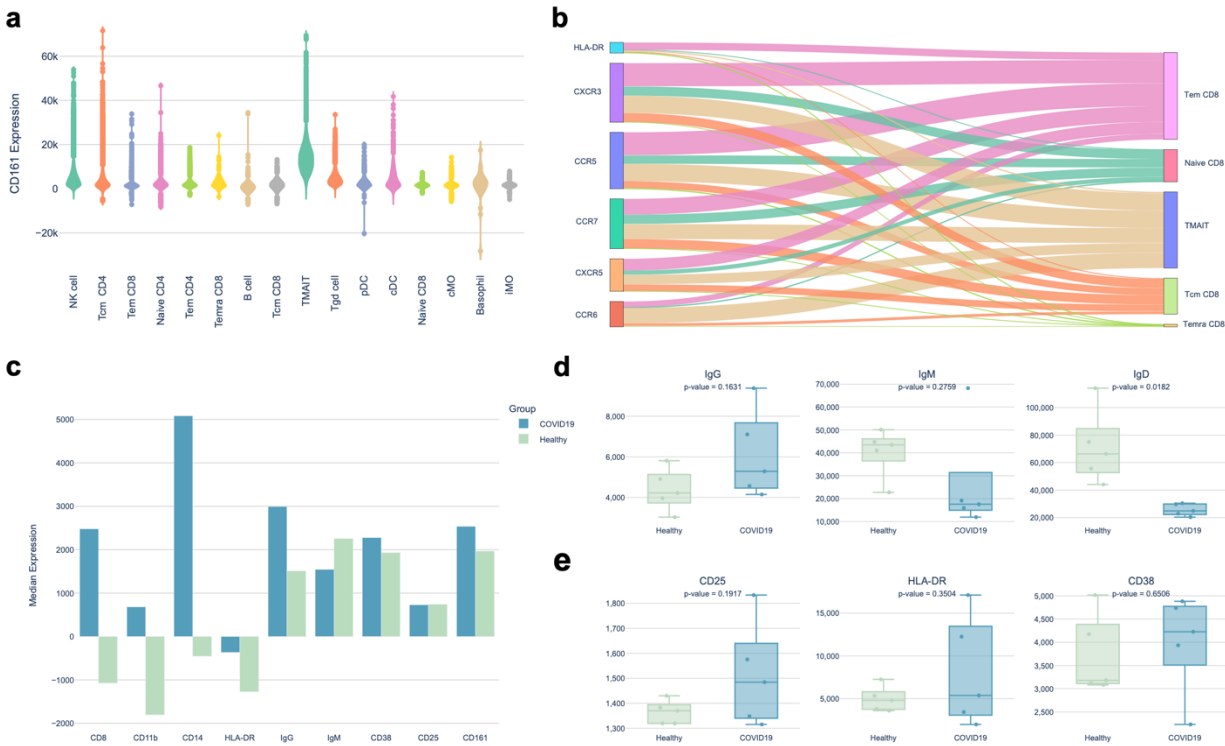
268 **Figure 3: Comparative Analysis of Immune Cell Subpopulations in Healthy and COVID-19**  
 269 **Individuals.** (a) UMAP plots displaying distinct clustering patterns and differential distribution of  
 270 cell types in PBMC across healthy and COVID-19 group. (b) Sankey diagram illustrates the  
 271 distribution of cells across groups, with thicker flow indicating more cells. (c) Composite bar-strip  
 272 plot summarizing cell count distribution across cell subpopulations. Dots represent each individual  
 273 samples colored by group. (d) Stacked bar chart showing distribution of cells in percentage across  
 274 two groups. (e) Effect size calculated using Cohen's d indicating changes in the number of cells  
 275 in COVID-19 compared to reference healthy control. (f) Comparison of individual cell type  
 276 frequencies between healthy and COVID-19 groups with p-values for statistical significance.  
 277 N=9/group

278

279 ***Distinct Cellular and Molecular Signatures Observed in COVID-19 Compared to Healthy***  
280 ***Controls***

281 UMAP analysis of COVID-19 hospitalized patients compared to healthy controls revealed distinct  
282 clustering patterns between groups, particularly among monocytes, NK cells, and CD8 T cells  
283 (Figure 3a). To understand changes between the two groups, CAFE offers varied visualization  
284 options, for instance, we used a Sankey diagram to demonstrate that MAIT cells and Tgd are  
285 much less abundant in COVID-19 patients compared to healthy controls (Figure 3b). We also  
286 found that CD8 Tcm and B cells were significantly expanded in COVID-19 patients. A composite  
287 bar-strip plot also demonstrates the distribution of cells in frequency where each dot represented  
288 each sample colored by specific group (Figure 3c). The total number of cells in iMO were largely  
289 reduced in COVID-19 patients compared to healthy controls (Figure 3d). These data may indicate  
290 a possible shift from an innate response towards an adaptive response. To quantify the effect size  
291 of changes observed, we compared cell types within the COVID-19 group to healthy controls as  
292 a reference and found changes in naive CD8, TMAIT, and iMO cells have a larger effect size,  
293 demonstrating a bigger difference between the two groups (Figure 3e). We further compared  
294 these cell types by plotting box plots for individual cell types (Figure 3f) which demonstrated a  
295 non-statistically significant increase in Tcm CD8 ( $p=0.0777$ ) and statistically significant decrease  
296 in naive CD8 cells ( $p=0.0372$ ) in COVID-19 patients compared to healthy controls.

297



298

299 **Figure 4: Marker expression and distribution differences between COVID-19 and healthy**  
 300 **individuals.** (a) Violin plot showed the median expression levels of CD161 across immune cell  
 301 subtypes. (b) Sankey diagram illustrated marker expression in CD8+ T cells, with thicker flows  
 302 indicating more cells expressing that marker. (c) Bar chart showed median expression of markers  
 303 across all cell types between COVID-19 and healthy individuals where positive values indicated  
 304 upregulation. Box plots displayed (d) the number of cells expressing IgG, IgM and IgD in B cells,  
 305 and (e) the number of CD8+ T cells expressing CD25, HLA-DR, and CD38.

306

### 307 ***Altered Marker Expression Profiles in COVID-19 Patients***

308 CAFE outputs a file with expression data for all markers on each cluster by sample for use in  
 309 external plotting and statistical software. In addition, further exploration of specific markers within  
 310 clusters can be performed within CAFE. We used this approach to identify that MAIT cells have  
 311 the greatest expression of CD161, as shown by violin plot (Figure 4a). We also examined marker  
 312 distribution in CD8 T cell populations and found that more cells within the Tem CD8 population  
 313 appeared activated based on greater HLA-DR, CXCR3, and CCR5 expression compared to other  
 314 CD8 T cell subsets (Figure 4b). We found that median expression of CD8, CD14, CD11b, IgG  
 315 were increased in COVID-19 patients compared to healthy controls across all clusters (Figure  
 316 4c). These reflect the overall differences in some of the cell populations we observed between

317 groups. We examined marker expression within the B cell cluster and found that more B cells  
318 expressed IgG in COVID-19 samples compared to healthy controls although the difference was  
319 not statistically significant ( $p=0.1631$ ), while IgD expression was statistically significantly reduced  
320 ( $p=0.0162$ ) in COVID-19 samples compared to healthy controls (Figure 4d). We also examined  
321 activation markers in CD8+ T cells and found that COVID-19 patients had more CD8+ T cells that  
322 expressed CD25 and HLA-DR than healthy controls (Figure 4e).

323

324

## 325 **Discussion:**

326

327 As research in immunology increasingly relies on high-dimensional cytometric data, there is a  
328 growing need for a user-friendly analysis tools for everyday use. Here, we present CAFE as a  
329 free and open-source tool designed to address the analytical and accessibility issues posed by  
330 SFCM data. CAFE uses a GUI and interactive controls to enable immunologists to analyze  
331 complex data without needing specialized coding knowledge.

332

333 Using a jupyter notebook, we have previously shown the ability of Scanpy's Leiden function  
334 (*scanpy.tl.leiden*) to analyze a 50-color human PBMC dataset<sup>15</sup>. CAFE acts as a wrapper  
335 combining packages within Streamlit to provide a web interface, offering more accessible and  
336 extensive functionality than Jupyter notebooks' CLI. We demonstrated analysis of a 35-color  
337 human PBMC dataset using CAFE in this study with 350,000 cells and 12.25M data points. Major  
338 steps including data processing, dimension reduction, batch correction and Leiden clustering  
339 were completed in under 12 minutes using an Apple M3 Pro laptop. In our analysis, we observed  
340 COVID-19 patients with altered immune cell distributions and marker expression profiles,  
341 consistent with prior findings, as well as MAIT cells expressing high CD161 and B cells expressing  
342 high CD20<sup>20</sup>.

343 While developing CAFE, we have balanced compatibility and performance and included many  
344 options for customization of how the code processes and analyzes data while integrating default  
345 options and tooltips to help guide users. Our implementation of Pandas with PyArrow significantly  
346 improves processing speed over Pandas alone. However, transitioning to Polars' lazy evaluation  
347 framework could further speed up processing once compatibility issues between ARM and x86  
348 machines are resolved. CAFE's web app design and functionality also revolved around simplicity  
349 as we de-emphasized features that are not commonly used in order to streamline user-

350 experience. For data dimension reduction, we adhered to a convention of using PCA as the  
351 primary method and using PCA-reduced representations for constructing UMAP neighbor graphs,  
352 as opposed to utilizing UMAP directly for primary dimensionality reduction. Although PCA is  
353 designed for linear data, it effectively reduces noise and enhances clustering performance. Users  
354 have the choice to skip PCA to perform Leiden clustering on the raw data or use UMAP  
355 embeddings (i.e., `X_umap`) to use UMAP reduced data for clustering. One viable alternative to  
356 PCA for non-linear data structure is Kernel PCA, but we have skipped adding the kernel PCA  
357 option in the CAFE workflow because it may not be practical since it is computationally taxing.

358 UMAP parameters and Leiden resolution largely influence the number of clusters for community  
359 detection. Leiden clustering is performed on the graph structure, so evaluation of clustering quality  
360 solely based on methods such as elbow or silhouette score is not ideal. Rather a combined  
361 approach with prior biological knowledge can inform the most correct clustering resolution. We  
362 recommend using CAFE's advanced clustering evaluation tab to generate plots with a range of  
363 varied UMAP parameters and Leiden resolutions for visual inspection. Using this approach, users  
364 can select the most appropriate clustering resolution for each dataset. Since this is an  
365 unsupervised algorithm, setting up an incorrect resolution can heavily skew the interpretation of  
366 data.

367 Manual gating continues to be the gold standard in flow cytometry analysis, but it is limited by  
368 sequentially drilling down into subsets of cells with 2-dimensional bi-axial gating. Thus, our goal  
369 was to complement this hypothesis-driven approach with the unsupervised computational  
370 algorithms. In this way, users can perform hypothesis-driven analysis with manual gating and  
371 hypothesis-generating analysis with unsupervised clustering. Compared to other open-source  
372 tools, CAFE provides a wide range of publication-ready visualization options. Due to its underlying  
373 code in python, CAFE is highly scalable to datasets of millions of cells and takes advantage of  
374 multi-threading to obtain higher performance. Previously, python-based Pytometry and CRUSTY  
375 integrated unbiased clustering algorithms within their tools<sup>21,22</sup>. Pytometry incorporates the Leiden  
376 algorithm<sup>10</sup>, which has been shown to be an improvement over the predecessor Louvain  
377 algorithm<sup>23</sup>; however, Pytometry requires coding using Python. CRUSTY incorporates an easy-  
378 to-use GUI but it does not offer the Leiden algorithm and relies on FlowSOM and Phenograph.  
379 Another limitation of CRUSTY is that the cloud-based service limits users to analyzing 100,000  
380 total events. There are also limited visualization and analysis options in CRUSTY and they rely  
381 on most of the Phenograph and FlowSOM default settings which cannot be customized. Cloud-  
382 based solutions may face limitations in availability, scalability, and data security. Users may be

383 prohibited from uploading data to cloud-based systems that have protected health information  
384 due to HIPAA. Among a few other GUI based tools, Cytoflow<sup>24</sup>, Floreada (floreada.io), EasyFlow<sup>25</sup>  
385 allow for flow cytometry data analysis but do not offer clustering. FlowPy (flowpy.wikidot.com)  
386 allows for clustering but uses k-mean clustering rather than the most advanced algorithms  
387 currently in use (i.e. Leiden). Additional tools like terraFlow<sup>26</sup> and CellEngine (CellCarta, Montreal,  
388 Canada) are for-profit spectral flow cytometry analysis softwares and the price of these may be  
389 restrictive for many users. Finally, FlowJo is a staple for many immunologists and has some native  
390 clustering capabilities. It also supports plugins for additional clustering algorithms, but these add-  
391 ons do not offer much customization in the clustering parameters.

392  
393 CAFE, while addressing many of these limitations as an open-source alternative, has its own  
394 practical considerations. CAFE is intended to be run locally which requires installing a package  
395 manager such as Pixi or Anaconda/Miniconda3 through terminal. Performance is also dependent  
396 on the user's machine. For larger datasets, we recommend utilizing our provided scripts in Github  
397 to run data processing step through an HPC cluster by allocating more RAM. Once the user has  
398 Anndata file generated with cluster information, all the downstream analysis and figure generation  
399 steps become significantly less computationally demanding. Ultimately, CAFE's aim is to become  
400 a secure, scalable, and open-source platform accessible to a broad range of researchers to run  
401 complex analyses through a simple intuitive graphical user interface.

402  
403 **Acknowledgements:** This work was supported by the National Institutes of Health/National  
404 Institute on Aging Grant R00 AG068309 (to D.J.T.); Nathan Shock Center, which is supported by  
405 the National Institutes of Health/National Institute on Aging Grant P30 AG050886; and this  
406 research was conducted while (Daniel Tyrrell) was an AFAR Grant for Junior Faculty awardee  
407 A24063 (to D.J.T.).

408  
409 **Data availability:** The raw data are publicly available at FlowRepository: FR-FCM-Z3WR. The  
410 downsampled .csv files are available from Figshare: 27940719. The processed Anndata object  
411 after dimensionality reduction and clustering is also available from Figshare: 27940752.

412  
413 **Code availability:** CAFE is freely available at <https://github.com/mhbsiam/cafe>.

414



415 **Contributions:** M.H.B.S. and D.J.T. conceived the software development. M.H.B.S. created the  
416 software and web application; M.H.B.S. analyzed the data. M.H.B.S. and D.J.T. wrote the paper.  
417 All authors reviewed and edited the paper.

418

419 **Corresponding author:** Correspondence to Daniel J. Tyrrell at [danieltyrrell@uabmc.edu](mailto:danieltyrrell@uabmc.edu).

420

421 **Ethics declarations:** The authors declare they have no competing interests.

422

423 **References:**

424

- 425 1. McKinnon, K. M. Flow Cytometry: An Overview. *Current Protocols in Immunology* **120**,  
426 5.1.1-5.1.11 (2018).
- 427 2. Nolan, J. P. The evolution of spectral flow cytometry. *Cytometry Part A* **101**, 812–817 (2022).
- 428 3. Konecny, A. J., Mage, P. L., Tyznik, A. J., Prlic, M. & Mair, F. OMIP-102: 50-color  
429 phenotyping of the human immune system with in-depth assessment of T cells and dendritic  
430 cells. *Cytometry Part A* **105**, 430–436 (2024).
- 431 4. Tyrrell, D. J. *et al.* Clonally expanded memory CD8+ T cells accumulate in atherosclerotic  
432 plaques and are pro-atherogenic in aged mice. *Nat Aging* **3**, 1576–1590 (2023).
- 433 5. Na, S., Choo, Y., Yoon, T. H. & Paek, E. CyGate Provides a Robust Solution for Automatic  
434 Gating of Single Cell Cytometry Data. *Anal Chem* **95**, 16918–16926 (2023).
- 435 6. Van Gassen, S. *et al.* FlowSOM: Using self-organizing maps for visualization and  
436 interpretation of cytometry data. *Cytometry A* **87**, 636–645 (2015).
- 437 7. Samusik, N., Good, Z., Spitzer, M. H., Davis, K. L. & Nolan, G. P. Automated mapping of  
438 phenotype space with single-cell data. *Nat Methods* **13**, 493–496 (2016).
- 439 8. Qiu, P. *et al.* Extracting a cellular hierarchy from high-dimensional cytometry data with  
440 SPADE. *Nat Biotechnol* **29**, 886–891 (2011).

- 441 9. Levine, J. H. *et al.* Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells  
442 that Correlate with Prognosis. *Cell* **162**, 184–197 (2015).
- 443 10. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing well-  
444 connected communities. *Sci Rep* **9**, 5233 (2019).
- 445 11. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587.e29  
446 (2021).
- 447 12. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data  
448 analysis. *Genome Biology* **19**, 15 (2018).
- 449 13. Rich, J. M. *et al.* The impact of package selection and versioning on single-cell RNA-seq  
450 analysis. *bioRxiv* 2024.04.04.588111 (2024) doi:10.1101/2024.04.04.588111.
- 451 14. McInnes, L., Healy, J., Saul, N. & Großberger, L. UMAP: Uniform Manifold Approximation  
452 and Projection. *Journal of Open Source Software* **3**, 861 (2018).
- 453 15. Vardaman, D. *et al.* Development of a Spectral Flow Cytometry Analysis Pipeline for High-  
454 dimensional Immune Cell Characterization. *J Immunol* **213**, 1713–1724 (2024).
- 455 16. Burton, R. J. *et al.* CytoPy: An autonomous cytometry analysis framework. *PLoS Comput*  
456 *Biol* **17**, e1009071 (2021).
- 457 17. Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data  
458 using empirical Bayes methods. *Biostatistics* **8**, 118–127 (2007).
- 459 18. Csardi, G. & Nepusz, T. The Igraph Software Package for Complex Network Research.  
460 *InterJournal Complex Systems*, 1695 (2005).
- 461 19. Silverman, B. W. *Density Estimation for Statistics and Data Analysis*. (Routledge, New York,  
462 2017). doi:10.1201/9781315140919.
- 463 20. Yu, C. *et al.* Mucosal-associated invariant T cell responses differ by sex in COVID-19. *Med*  
464 **2**, 755-772.e5 (2021).

- 465 21. Büttner, M., Hempel, F., Ryborz, T., Theis, F. J. & Schultze, J. L. Pytometry: Flow and mass  
466 cytometry analytics in Python. 2022.10.10.511546 Preprint at  
467 <https://doi.org/10.1101/2022.10.10.511546> (2022).
- 468 22. Puccio, S. *et al.* CRUSTY: a versatile web platform for the rapid analysis and visualization of  
469 high-dimensional flow cytometry data. *Nat Commun* **14**, 5102 (2023).
- 470 23. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities  
471 in large networks. *J. Stat. Mech.* **2008**, P10008 (2008).
- 472 24. Teague, B. Cytoflow: A Python Toolbox for Flow Cytometry. 2022.07.22.501078 Preprint at  
473 <https://doi.org/10.1101/2022.07.22.501078> (2022).
- 474 25. Ma, Y., Eizenberg-Magar, I. & Antebi, Y. EasyFlow: An open-source, user-friendly cytometry  
475 analyzer with graphic user interface (GUI). *PLOS ONE* **19**, e0308873 (2024).
- 476 26. Freeman, D. *et al.* terraFlow, a high-parameter analysis tool, reveals T cell exhaustion and  
477 dysfunctional cytokine production in classical Hodgkin's lymphoma. *Cell Reports* **43**, (2024).

478

#### 479 **Figure Legends:**

480 **Figure 1: The flowchart outlines steps and components of CAFE's workflow.** Preprocessing  
481 includes compensation, data scaling/transformation using a standard FCM software and scaled  
482 CSV files are then exported and renamed as Sample\_Group.csv. Data processing performs error  
483 checks and concatenation of CSV files into an AnnData object/H5AD file. Major steps requiring  
484 user input include dimension reduction, batch correction, UMAP (UMAP Uniform Manifold  
485 Approximation and Projection) and community detection. Outputs are downloadable as CSV,  
486 H5AD, PNG, JPG, SVG and PDF files.

487 **Figure 2: Profiling of Human PBMCs Reveals Distinct Immune Subpopulations and Marker**  
488 **Expression Patterns.** (a) UMAP plots showing selected marker expression intensities across all  
489 cells in the UMAP space to highlight lineage-specific marker distribution. (b) Dot plot of all marker  
490 expression across all identified PBMC cell types. Dendrogram highlighted distinct marker-based  
491 groupings. (c) UMAP visualization showing 16 distinct clusters with annotated cell types including  
492 Naive CD4 and CD8 T cells, central memory CD4 and CD8 T cells (Tcm), effector memory CD8  
493 T cells (Tem), terminally differentiated effector memory CD8 T cells (Temra), mucosal-associated  
494 invariant (MAIT) T cells, classical monocytes (cMO), intermediate monocytes (iMO), B cells, NK

495 cells, gamma delta ( $\gamma\delta$ ) T cells (Tgd), conventional dendritic cell (cDC) and plasmacytoid dendritic  
496 cell (pDC).

497 **Figure 3: Comparative Analysis of Immune Cell Subpopulations in Healthy and COVID-19**  
498 **Individuals.** (a) UMAP plots displaying distinct clustering patterns and differential distribution of  
499 cell types in PBMC across healthy and COVID-19 group. (b) Sankey diagram illustrates the  
500 distribution of cells across groups, with thicker flow indicating more cells. (c) Composite bar-strip  
501 plot summarizing cell count distribution across cell subpopulations. Dots represent each individual  
502 samples colored by group. (d) Stacked bar chart showing distribution of cells in percentage across  
503 two groups. (e) Effect size calculated using Cohen's d indicating changes in the number of cells  
504 in COVID-19 compared to reference healthy control. (f) Comparison of individual cell type  
505 frequencies between healthy and COVID-19 groups with p-values for statistical significance.  
506 N=9/group.

507 **Figure 4: Marker expression and distribution differences between COVID-19 and healthy**  
508 **individuals.** (a) Violin plot showed the median expression levels of CD161 across immune cell  
509 subtypes. (b) Sankey diagram illustrated marker expression in CD8+ T cells, with thicker flows  
510 indicating more cells expressing that marker. (c) Bar chart showed median expression of markers  
511 across all cell types between COVID-19 and healthy individuals where positive values indicated  
512 upregulation. Box plots displayed (d) the number of cells expressing IgG, IgM and IgD in B cells,  
513 and (e) the number of CD8+ T cells expressing CD25, HLA-DR, and CD38.

514

