

1 **DAFi: A Directed Recursive Filtering and Clustering Approach to Data-Driven**
2 **Identification of Cell Populations from Polychromatic Flow Cytometry Data**

3 Alexandra J. Lee^{1,+}, Ivan Chang^{1,+}, Julie G. Burel², Cecilia S. Lindestam Arlehamn², Daniela Weiskopf²,
4 Bjoern Peters², Alessandro Sette², Richard H. Scheuermann^{1,3}, Yu Qian^{1,*}

5 ¹J. Craig Venter Institute, La Jolla, CA

6 ²La Jolla Institute for Allergy and Immunology, La Jolla, CA

7 ³Department of Pathology, University of California, San Diego, CA

8 ⁺These authors contributed equally

9 ^{*}Corresponding author:

10 Yu Qian, Ph.D.

11 J. Craig Venter Institute

12 4120 Capricorn Lane

13 La Jolla, CA 92037

14 Phone: 1-858-200-1837

15 Fax: 1-858-200-1879

16 Email: mqian@jcv.org

17 This work has been supported by NIH/NIAID U19AI118626 (HIPC) and R24AI108564; NIH/NCATS
18 U01TR001801 (FlowGate); NSF XSEDE allocation MCB170008; NIH/NIAID HHSN272201200005C
19 (RPRC).

20 The authors have no conflict of interest to declare.

1 **Abstract**

2 Computational methods for identification of cell populations from high-dimensional flow cytometry data
3 are changing the paradigm of cytometry bioinformatics. Data clustering is the most common
4 computational approach to unsupervised identification of cell populations from multidimensional
5 cytometry data. We found that combining recursive filtering and clustering with constraints converted
6 from the user manual gating strategy can effectively identify overlapping and rare cell populations from
7 smeared data that would have been difficult to resolve by either a single run of data clustering or manual
8 segregation. We named this new method DAFi: Directed Automated Filtering and Identification of cell
9 populations. Design of DAFi preserves the data-driven characteristics of unsupervised clustering for
10 identifying novel cell-based biomarkers, but also makes the results interpretable to experimental scientists
11 as in supervised classification through mapping and merging the high-dimensional data clusters into the
12 user-defined 2D gating hierarchy. By recursive data filtering before clustering, DAFi can uncover small
13 local clusters which are otherwise difficult to identify due to the statistical interference of the irrelevant
14 major clusters. Quantitative assessment of cell type specific characteristics demonstrates that the
15 population proportions calculated by DAFi, while being highly consistent with those by expert centralized
16 manual gating, have smaller technical variance than those from individual manual gating analysis. Visual
17 examination of the dot plots showed that the boundaries of the DAFi-identified cell populations followed
18 the natural shapes of the data distributions. To further exemplify the utility of DAFi, we show that DAFi
19 can incorporate the FLOCK clustering method to identify novel cell-based biomarkers. Implementation of
20 DAFi supports options including clustering, bisecting, slope-based gating, and reversed filtering to meet
21 various auto-gating needs from different scientific use cases.

22 **Keywords:** auto-gating, data prefiltering, recursive clustering, cell population identification, constrained
23 clustering

1 **1. Introduction**

2 The success of flow cytometry (FCM) is dependent on being able to accurately identify discriminant cell
3 populations. Currently, the most common existing approach is manual gating analysis. In a typical manual
4 gating procedure, an experiment operator would start by inspecting the distribution of cellular events on a
5 selected pair of measured characteristics (scatter parameters or protein markers) on the 2D plot, visually
6 recognize the clusters of the cellular events, draw a 2D polygon to extract a population of interest, inspect
7 the population on another pair of markers to identify its subpopulations, and repeat this procedure to
8 further partition each subpopulation until all cell subsets of interest are identified. Through this 2D by 2D
9 recursive segregation, cell populations are identified and managed in a user-defined hierarchy with
10 phenotypes defined with the markers used at each gating step. Desirable features of manual gating
11 analysis include the flexibility in the analysis procedure and the interpretability of the analysis results.

12
13 However, the manual gating procedure is also subjective, time-consuming, and difficult to reproduce.
14 Technical variance is usually found in independent manual gating analysis conducted across experiments,
15 studies, and labs [12, 31]. Except in the ideal case where each cell population is highly cohesive and
16 segregated from others, there is notable bias in using a sequential 2D by 2D analysis to identify cell
17 populations defined in high-dimensional space since cell populations can be difficult to separate on 2D
18 dot plots. Manual gating analysis typically bisects the overlapped populations with manually drawn lines,
19 resulting in inaccurate identification and calculation of the population characteristics. Additionally, the
20 design of the manual gating approach is not suitable for exploratory data analysis. Gating steps are
21 predefined on user-selected markers and constrained by the operator's knowledge of cell population
22 phenotypes. Recent advance in cytometry instrumentation and reagent technology made these issues more
23 severe (*e.g.*, Becton Dickinson's FACSymphonyTM is claimed to be able to measure up to 50 different
24 characteristics). With these many parameters, it becomes almost impossible to explore the enormous

1 high-dimensional data space exhaustively and accurately using a manual gating approach, considering the
2 time, effort, and human bias involved in the analysis.

3

4 During the last decade, many computational methods have been developed for the identification of cell
5 populations from polychromatic FCM data. State-of-the-art computational approaches are shown to be
6 superior to manual gating analysis in terms of efficiency, reproducibility and reduction of human bias [2,
7 3, 8, 9, 10, 13, 21, 22, 25, 26, 38]. Based on whether user inputs are required, these approaches can be
8 broadly categorized into unsupervised [1, 11, 14, 15, 18, 24, 28, 32, 33, 44, 46] and supervised/semi-
9 supervised [16, 23, 27] approaches. Unsupervised methods are usually based on data clustering methods,
10 making them useful for comprehensive immunophenotyping and identification of novel cell subsets.
11 Because this unsupervised analysis is completely data-driven, there is no direct connection between the
12 identified data clusters and existing knowledge about the cell populations. Each data cluster needs to be
13 annotated and validated, usually manually. The number of clusters identified in different input files can
14 also be different. Therefore, it can be non-trivial to map and interpret all cell populations identified by the
15 unsupervised clustering methods across samples. In contrast, supervised identification methods require
16 prior data analysis results (usually from manual gating analysis) as training data, and thereby guarantee
17 the interpretability of the identified data clusters. The trade-off is that the supervised methods, primarily
18 focused on predefined cell populations, usually do not support exploratory discovery of novel cell subsets.
19 While identifying novel cell subsets is one of the most important features expected by translational
20 researchers when they use computational methods for FCM data analysis, supervised identification
21 methods are preferred for clinical diagnostics.

22

23 Because there is no single model that fits all data, the biggest challenge for adoption of a computational
24 method by experimental scientist is how to select the best method when the dataset changes. Parameters
25 in a computational algorithm often need to be adjusted when being applied to a new dataset, and the

1 adjusting is usually difficult without sufficient understanding about the method and the data. This
2 challenge may be partially addressed with infrastructure efforts such as developing a parallel testing
3 environment (*e.g.*, FlowGate [35]) to assess the performance of multiple applicable computational
4 methods on each specific dataset. Another solution is to support the incorporation of user knowledge to
5 guide the clustering analysis. One example is constrained clustering, in which user-provided constraints
6 about cluster membership of data objects are involved [45]. Constrained clustering is regarded as a special
7 class of semi-supervised learning, which has proven highly effective for solving domain-specific
8 problems.

9 In this paper, we propose a constraint-based recursive filtering and clustering approach – DAFi (directed
10 automated filtering and identification of cell populations) – to address the problem of utilizing
11 computational methods for identification of cell populations from FCM data. Our goal is not to propose a
12 new data clustering or classification method. Instead, we will demonstrate that designing a recursive
13 filtering and clustering approach and combining it with user gating strategy can effectively and reliably
14 accomplish the task of auto-gating for identifying not only the major but also rare and novel cell
15 populations from a variety of FCM datasets.

16

17 **2. Results**

18 Figure 1 illustrates the design of DAFi. At each gating step DAFi supports the use of different clustering
19 methods as well as bisecting (for identifying outlier cells), slope-based (*e.g.*, identifying singlets using
20 FSC/SSC-A vs FSC/SSC-H), and reverse-gating (events inside the hyper-polygon will be filtered out).
21 We experimented with two clustering methods: *K*-means and FLOCK clustering [33], for benchmarking
22 the performance of DAFi. User input is used to identify both predefined and novel populations that can be
23 organized within an easily interpreted gating hierarchy (Figure 1B). We refer to this type of approach as
24 *directed unsupervised clustering*. Figure 1C shows DAFi-identified major (CD4+ T and CD8+ T cells)

1 and rare (CD3+CD56+ T and CD3hiCD56+ T cells) cell populations. The CD3+CD56+ T and
2 CD3hiCD56+ T cell populations are difficult to separate by either manual gating or traditional
3 unsupervised data clustering methods because the two clusters are both relatively rare and close to each
4 other in CD3 expression distributions. However, they were well segregated with natural boundaries
5 (unimodal distribution on each dimension) using DAFi, which applied recursive clustering with manual
6 gating polygons as *constraints* rather than absolute boundaries. DAFi also identified the difficult-to-
7 resolve CD4+CD25+ regulatory T cells (Tregs), yielding cell populations with natural distributions
8 (Figure 1D). In contrast, manual gating analysis using polygon partitions did not capture natural
9 boundaries of cell populations (*i.e.*, an abrupt lower boundary in the CD25 dimension); *K*-means
10 clustering failed to identify the rare Treg cell population at all, even with $K=500$.

11
12 We evaluated the performance of DAFi using FCM data from both the public ImmPort database
13 (Immunology Database and Analysis Portal, <http://www.immport.org>) and our HIPC (Human
14 Immunology Project consortium, <https://www.immuneprofile.org>) studies. Results of DAFi were
15 assessed both quantitatively and by visual examination of the identified cell populations on dot plots. For
16 the quantitative assessment, instead of using the “bulk assessment” such as the sample-level F-measure
17 which is dominated by contributions from the abundant cell populations, we focus on cell type specific
18 statistics for each individual cell population.

19 **2.1 Cell Type Specific Assessment in Comparison with Individual and Centralized Manual Gating** 20 **Analysis**

21 The first assessment focused on the identification of different T cell subsets using a representative 10-
22 color reagent panel on multiple repeat runs of cryopreserved PBMC (peripheral blood mononuclear cells)
23 from one sample donation of a healthy donor [12]. Repeated FCM experiments were performed on
24 various days throughout a 7-month period by three different operators on four different cytometers.

1 Technical variability associated with each cell population across the 24 runs can therefore be estimated
2 and compared between the results of DAFi and those from the manual gating analysis.

3 Two manual gating analysis results were available for comparison - individual manual gating analysis
4 (INDI) performed by different operators when the FCM data were acquired, and centralized manual
5 gating analysis (CENT) performed by one analyst after data from all 24 samples had been acquired. The
6 variety of cell subsets and their relationship specified by expert manual gating in a predefined hierarchy
7 are shown in Figure 2A (also in Supplementary File 1). Among the 22 cell populations identified by
8 manual gating analysis, 17 of them are of special user interest. They were divided into two categories
9 based on the technical variance – “clearly defined” and “poorly resolved” - defined by coefficient of
10 variability (CV) in cell population proportions across the 24 samples from individual manual gating
11 analysis:

- 12 • Clearly defined (low CV): 5: Monocytes; 7: B-cells; 8: NK cells; 11: T-cells; 12: CD4+ T cells;
13 13: CD8+ T cells; 15: Naïve CD4+ T cells; 19: Naïve CD8+ T cells.
- 14 • Poorly resolved (high CV): 9: CD3+CD56+ T cells; 10: CD3highCD56+ T cells; 14: Tregs
15 (regulatory CD4+ T cells); 16: Tcm CD4+ T cells (central memory CD4+ T cells); 17: Tem
16 CD4+ T cells (effector memory CD4+ T cells); 18: Temra CD4+ T cells (effector memory CD4+
17 T cells that express CD45RA); 20: Tcm CD8+ T cells; 21: Tem CD8+ T cells; 22: Temra CD8+
18 T cells.

19 Previously [12], we found that both INDI and CENT could achieve a high degree of concordance for
20 identifying clearly defined cell populations; but for the poorly resolved ones, CENT significantly
21 outperformed INDI.

22 Figure 2B shows the results of DAFi. Visual examination shows the main difference between DAFi and
23 the manual gating analysis is that DAFi identified cell populations with natural boundaries, while manual

1 gating analysis resulted in abrupt bisecting on some of the 2D plots. Based on the dot plots, the three most
2 difficult-to-resolve gating boundaries seem to be: a) between CD3+CD56+ T (Pop#9) and
3 CD3^{high}CD56+ T cells (Pop#10) (Figure 2C); b) among Naïve CD4+ T and three memory CD4+ T cells
4 (Figure 2D); c) between Tregs (Pop#14) and CD4+ helper T cells (Figure 2E). Dot plots of these cell
5 populations across all 24 samples can be found in Supplementary File 2. Visual examination showed that
6 DAFi successfully identified these difficult-to-resolve cell populations.

7 Linear regression analysis of the cell population percentages identified by DAFi and CENT for clearly-
8 defined (Figure 2F) and poorly-resolved (Figure 2G) cell populations are highly consistent, with the
9 degree of concordance on the clearly-defined cell populations being higher than that of the poorly-
10 resolved ones. Both *K*-means and DAFi generated highly consistent population percentages with those of
11 CENT (all *p*-values smaller than 0.0001, ranged from 10^{-4} to 10^{-14}) for clearly-resolved populations
12 (Figure 2H), indicating that both DAFi and a naive application of *K*-means can identify clearly-defined
13 cell populations successfully. However, *K*-means failed to identify many of the poorly-resolved
14 populations in a consistent way with CENT (Figure 2I), including Tregs, Tcm CD8+ T cells,
15 CD3^{hi}CD56+ T cells, and Temra CD4+ T cells. In contrast, populations percentages of the poorly-
16 resolved populations identified by DAFi are consistent with those derived by CENT (with all *p*-values
17 smaller than 0.01, ranged from 10^{-2} to 10^{-9}), indicating the necessity of the recursive filtering and
18 clustering design.

19 Figures 2J and 2K compare CV (coefficient variability) of population percentages across the 24 samples
20 identified by four different approaches: INDI, CENT, *K*-means, and DAFi. For clearly-defined
21 populations, the four approaches generated highly consistent CV values. However, for poorly-resolved
22 populations, *K*-means and INDI generated very large CV values (Figure 2K), especially for rare cell
23 populations (*e.g.*, Tregs and Temra CD4+ T cells), while DAFi and CENT generated similarly smaller
24 CV values.

1 **2.2 Correlation Analysis between Cell Population Proportions with Subject Age and Gender**

2 We extended the single-donor analysis to assess PBMC samples from 132 human subject participants,
3 stained with the same 10-color panel used in Section 2.1. Other details about the FCM experiment can be
4 found in the published study [12]. The goal of the assessment was to determine if T cell population
5 frequency determined by DAFi correlated with subject demographics data, including gender and age.

6 No difference was observed in age distributions among gender groups (Figure 3A) allowing us to mix the
7 subjects from both genders to increase the statistical power in the age-based correlation analysis. We
8 focused on 12 predefined T-cell populations: 11: T-cells; 12: CD4+ T cells; 13: CD8+ T cells; 14: Tregs;
9 15: Naïve CD4+ T cells; 16: Tcm CD4+ T cells; 17: Tem CD4+ T cells; 18: Temra CD4+ T cells; 19:
10 Naïve CD8+ T cells; 20: Tcm CD8+ T cells; 21: Tem CD8+ T cells; 22: Temra CD8+ T cells. Complete
11 set of percentages of these T cell populations identified by DAFi across the 132 samples can be found in
12 Supplementary File 3. Consistently with our previous analysis using CENT [12], we found that the
13 proportion of Naïve CD4+ T and Naïve CD8+ T cells decreased with subject age (Figure 3B, corrected
14 linear regression p -value 3.740E-06 and 9.678E-11, respectively). Figure 3C shows Pearson correlation
15 scores and linear regression p -values across all 12 cell populations identified by DAFi and manual gating
16 analysis. Again, the output of DAFi is highly consistent with that of the centralized manual gating
17 analysis.

18 In the gender-based correlation analysis, DAFi identified that the proportion of the CD4+ T cell
19 population seems to be significantly different between the female and male (corrected p -value 0.023688,
20 Figure 3D). In our previous analysis [12], we were not able to identify this correlation with a significant
21 p -value using manual gating, although the average CD4+ T cell proportion was higher in the female
22 group. A number of previous studies have reported increases in CD4+ T cells in females [4, 17, 20, 37, 39,
23 40, 41, 42, 43]. Most recently, the 10k Immunomes Project based on a meta-analysis of 578 subjects in

1 the ImmPort Database reported the percentages of CD4+ T cells are significantly elevated in women as
2 compared to men (<http://www.biorxiv.org/content/early/2017/08/25/180489>).

3 We also studied the pairwise correlation between population proportions of different T cell subsets.
4 Figure 3E shows the distributions of the population proportion values from the two most significant
5 Pearson correlation scores, both of which are negative. One is between Naïve CD4+ T cells and Tem
6 CD4+ T cells ($r = -0.8601$) and the other is between Naïve CD8+ T cells and Tem CD8+ T cells ($r = -$
7 0.8638). This finding is consistent with the age-based analytics showing that the number of memory T
8 cells increased with age while the number of naïve T cells decreased with age. For both Tem CD4+ and
9 Tem CD8+ T cell populations DAFi identified a stronger association of their proportion increases with
10 age than manual gating analysis (Figure 3C).

11 **2.3 Identification of Known and Novel Cell-Based Biomarkers for Latent Tuberculosis Infection**

12 We also assessed the capability of DAFi for identifying cell populations that have not been defined in the
13 manual gating strategy using a dataset consisting of 12 PBMC samples from 6 latently tuberculosis
14 infected (LTBI) human subjects and 6 *Mycobacterium tuberculosis* (*Mtb*) uninfected control (healthy
15 control; HC) subjects used in a previous study [6]. We divided the sequence of manual gating steps into
16 two stages: prefiltering to identify the CD4+ T cell population (the first row of Figure 4A), and
17 unsupervised clustering (FLOCK [33]) to identify cell subsets within the CD4+ T cell population (Figure
18 4B) associated with subject phenotypes. We noticed that the manual gating strategy after the CD4+ T
19 cells gate was focused on poorly-resolved cell populations with relatively arbitrary gating boundaries on
20 smeared data dimensions including CD25, CCR6, CXCR3, and CCR4 (Figure 4B). Applying a data
21 clustering method that can utilize multiple data dimensions simultaneously to substitute the manual gating
22 strategy starting from CD4+ T cells can be expected to generate more accurate results or identify novel
23 cell-based biomarkers. Complete set of dot plots for data prefiltering by DAFi across all the 12 samples
24 can be found in Supplementary File 4.

1 DAFi using FLOCK identified 101 cell populations from the 12 samples (percentages of populations can
2 be found in Supplementary File 5). The population percentages were then associated with the subject
3 phenotype using non-parametric Wilcoxon rank sum test with a null hypothesis that there is no difference
4 between the LTBI and HC group. Figure 4C-D shows the top 5 cell populations with most significant p -
5 values against the null hypothesis. Distributions of the percentages of three relatively abundant
6 populations (Pop#23, 27, 28) are shown in Figure 4C, with the two rare ones (Pop#18, 65) in Figure 4D.
7 Due to the limited number of subjects, the best possible p -value in the rank sum test is 0.003948 when
8 there is no overlapping between the two groups in the ranks of their data objects (Figure 4D).

9 2D dot plots of the top 5 significant cell populations are shown in Figures 4E-F. In the previous
10 publication [6], a single cell population was identified by manual gating analysis in the CD25-
11 CCR6+CCR4-CXCR3+ region that significantly differed between LTBI and HC in frequency (Figure 4B,
12 corrected p -value < 0.01). In contrast, FLOCK identified three subsets within the same region: Pop#23,
13 27, and 65 (Figure 4E), which differ from each other based on CCR7 and peptide-MHC tetramer staining.
14 DAFi not only identified the known cell-based biomarker but also elucidated the composition of the
15 CD25-CCR6+CCR4-CXCR3+ cell population containing the vast majority of *Mtb*-specific cells. Two
16 DAFi-identified cell populations that were ignored in the original manual gating analysis are: Pop#18 and
17 #28 (Figure 4F), which differ in CCR4, CCR6, and tetramer staining. Their corresponding positions in the
18 predefined cell type hierarchy are indicated with the red arrows “novel” in Figure 4B. Figure 4G shows a
19 tSNE map [5] of the CD4+ T cells of the same LTBI sample used in Figure 4E-F, color-coded based on
20 the tetramer staining levels of the DAFi-identified cell populations. Separation of the tetramer+ cells from
21 the other cells on the tSNE map indicates that two very rare Pops 18 and 65 are indeed distinct (Figure
22 4H).

23 **2.4 Quantification of Human Immune Response to Influenza and Pneumococcal Vaccination**

1 Finally, we applied DAFi to identify plasmablasts/plasma cells to measure human immune responses to
2 vaccination. The FCM dataset used, SDY180 [29], was downloaded from the ImmPort database
3 (www.immport.org). 36 human subjects were enrolled into three immunization arms for FCM
4 experiments: Fluzone (2009-2010 seasonal influenza vaccine, N=12), Pneumovax23 (23-valent
5 pneumococcal vaccine, N=12), and Saline (N=12). PBMC samples were collected at 10 different time
6 points: Day-7, 0 (vaccination day), 0.5, 1, 3, 7, 10, 14, 21, and 28. FCM data files used in our data
7 analysis (306 FCS files in total) were acquired using an 8-color reagent panel focused on identification of
8 the plasmablasts/plasma cells and other types of B-cells: FSC-A, SSC-A, FITC-A_IgD, Pacific-Orange-
9 A_CD45, APC-A_CD138, APC-Cy7-A_CD27, PE-A_CD24, PE-Texas-Red-A_CD19, PE-Cy5-A_CD20,
10 and PE-Cy7-A_CD38. Manual gating analysis was used to identify the cellular composition at the 10
11 different time points before and after vaccination, which revealed a peak in plasmablast frequencies at
12 Day 7 post vaccination for both vaccines [29].

13 We reanalyzed the B-cell phenotyping FCM data of SDY180 using DAFi. Dot plots of CD19+ B cells
14 (Figure 5A, blue) and Plasmablasts (Figure 5B, magenta) identified by DAFi are shown with their
15 defining rectangle boundaries. Note that events outside the 2D rectangles may still belong to the cell
16 population as long as the centroid of their data cluster is within the hyper-rectangle. Similarly, an event
17 inside the 2D rectangle may not be assigned to the cell population if its cluster centroid is outside the
18 hyper-rectangle. In Day 7 samples only, the IgD-CD27^{high} plasmablasts can be clearly seen. The clear
19 peaks on Day 7 post both Fluzone and Pneumovax23 vaccinations (Figure 5C) confirmed the finding
20 reported previously [29]. We applied 0-1 min-max normalization to both results (medians of the
21 population percentages across different samples on the same day were used) so that the time-series
22 patterns identified by both approaches could be compared. The time-series pattern identified by DAFi is a
23 close match with that by manual gating analysis (Figure 5D), with a peak on Day 7 for both Fluzone and
24 Pneumovax 23 groups. The second peak post vaccination in the Fluzone group is on Day 14, which also
25 seems a close match between the two approaches. The baseline identified by DAFi seems smoother than

1 that of manual gating analysis in all three groups. 11-fold and 47-fold increase were reported in the
2 previous publication in the absolute numbers of plasmablasts following vaccinations of Fluzone and
3 Pneumovax23, respectively [29]. When comparing the median percentage values of DAFi-identified
4 plasmablasts on Day 7 with the baseline (Day -7), we achieved 16-fold and 43-fold increase post Fluzone
5 and Pneumovax23 vaccinations, respectively, a close match to the manual gating analysis result reported
6 previously.

7 3. **Methods**

8 Design of DAFi consists of four major steps (Figure 1A): unsupervised data clustering, encoding
9 predefined gating boundaries, merging of data clusters based on predefined gating strategy and
10 boundaries, and output for recursive filtering and population statistics. At each gating step, DAFi supports
11 four different gating options/modes to identify the individual cell populations: clustering, bisecting, slope-
12 based bisecting, and reversed filtering. Clustering is the default mode in DAFi. Because of the lack of
13 gold standard in defining and assessing cell populations, we keep the bisecting as an option in DAFi to
14 support the exact recapitulation of manual gating analysis. In some cases, the user prefers to do the gating
15 in a reversed way, *i.e.*, keeping the cells outside of the gate and removing those inside (*e.g.*, the second
16 step in Figure 4B, the identification of memory T cells based on CCR7 vs CD45RA can be achieved by
17 drawing a reversed gate around the naïve T cells in the double positive region), which is supported by the
18 reversed filtering mode in DAFi.

19 We have implemented two existing data clustering methods: *K*-means and FLOCK to be used with DAFi.
20 We named them DAFi-filtering and DAFi-gating:

21 DAFi-filtering (for data prefiltering and identification of predefined cell populations):

- 22 • Step 1 FCS file preprocessing: each FCS file converted to a data matrix for clustering analysis.
- 23 • Step 2 Generation of DAFi configuration file based on manual gating strategy.

- 1 • Step 3 Apply *K*-means data clustering to identify data clusters in each input file.
- 2 • Step 4 Merge data clusters whose centroids are within the hyper-rectangle formed by gating
- 3 boundaries; output the merged data as the input file for next-run clustering analysis.
- 4 • Step 5 Repeat Steps 3 and 4 until all predefined cell populations of interest are identified.
- 5 • Step 6 Output the dot plots and statistics of the identified cell populations together with their
- 6 names and phenotypes.

7 DAFi-gating replaces the *K*-means clustering in DAFi-filtering with the FLOCK clustering method that
8 can identify the undefined cell populations:

- 9 • Step 1 Run DAFi-filtering with the FLOCK clustering method to identify the predefined cell
- 10 population that needs to be explored for undefined cell subsets.
- 11 • Step 2 Normalize the filtered data across samples and merge the normalized events across
- 12 samples into a single data file.
- 13 • Step 3 Apply FLOCK to the data file to identify data clusters in a fully unsupervised way, and
- 14 map the clustering results back to individual samples.
- 15 • Step 4 Output the dot plots and statistics of the identified cell populations for manual annotation
- 16 of their names and phenotypes.

17 **3.1 FCS file preprocessing**

18 FCSTrans [34] was used to convert and transform (logicle transformation [30]) the binary FCS files
19 generated in all the FCM experiments used in this paper into data matrices for computational processing
20 and analysis.

1 **3.2 Converting manual gating strategies into configuration file for DAFi-filtering**

2 Different gating software (FlowJo, FCSExpress, FACS Diva *etc.*) and their different versions use
3 different ways to record the gating boundaries in different formats. Coordinate values of gating
4 boundaries on FCM data from one data transformation cannot be directly applied to gate the data
5 generated by a different transformation or with a different set of transformation parameters. Based on the
6 design of DAFi, we only need to simulate the gating boundaries by drawing rectangles. The Results
7 Section shows that even without using the exact gating boundaries DAFi still achieved highly consistent
8 results with expert manual gating analysis across a variety of experiments and cell populations of interest.

9 Table 1 illustrates an example configuration file used in DAFi-filtering. The *RecursiveParent* specifies
10 whether this cell population will be used in downstream recursive filtering and clustering. By default, a
11 cell population will be identified from its direct parent in clustering mode, while it can also be identified
12 from its grandparent population or even from the input FCM sample when the user chooses to skip the
13 intermediate gating steps. The former way is recursive clustering while the latter way is not recursive with
14 a hyper-polygon as the constraint, depending on user preference.

15 **3.3 Impact of K of K -means and size of gating boundaries on DAFi-filtering**

16 While K -means is easy to implement and use, one challenge is to set the value of K . We have
17 experimented different values of K from 100-600 using the LTBI dataset in Section 2.3. Supplementary
18 File 6 shows the F -measure values for each of the 5 cell populations comparing between the bisecting (*i.e.*,
19 manual gating analysis) and the clustering mode of DAFi. The box plot shows that the average F1 scores
20 across the 12 samples are larger than 0.95 for all of $K=100$ to 600. The variation of the F1 scores across
21 samples is also small. The larger number of K , the closer the result is to the bisecting (the F1 scores seems
22 the largest for $K=600$). By default, we set $K=500$.

1 We used the same LTBI dataset in Section 2.3 and tested three sizes of gating boundaries for identifying
2 the 5 cell populations across the 12 samples: a) normal: the same size as the bisecting boundaries; b)
3 small: 10% smaller on each dimension than the bisecting boundaries, and c) large: 10% larger on each
4 dimension than the bisecting boundaries. We calculated the precision, recall, and F1 scores of comparing
5 results of DAFi-filtering using these three sets of rectangles against the bisecting results (Supplementary
6 File 7). All the F1 scores as well as precisions and recalls are very high, while using a small rectangle
7 seems increasing the precision but reducing the recall, compared with using a large rectangle. The
8 variation of the F1 scores across the 12 samples is also small without being affected by the slight change
9 of the size of the rectangle gate. For the downstream cell populations (e.g., Pop5: CD3+CD4+ live T
10 lymphocytes), the accumulated change in F1 scores from the previous DAFi-filtering steps is not obvious
11 either.

12 **3.4 DAFi-gating: prefiltering and identification of undefined cell populations using FLOCK**

13 For generating the results in Section 2.3, after prefiltering, the remaining events across the 12 samples
14 were first normalized and then merged together. The cross-sample normalization (*a.k.a.*, sample
15 alignment) was done using the GaussianNorm approach [19]. Supplementary File 8 illustrates the
16 application of the GaussianNorm method to normalize the individual data dimensions CCR6 and
17 CD45RA. Only data dimensions needed in the unsupervised FLOCK analysis are kept, resulting in a 7-
18 dimensional data matrix (CD25, CXCR3, CCR4, CCR6, CCR7, CD45RA, and Tetramer), from which
19 FLOCK was applied to identify the 101 cell subsets (*number of bins* = 12 and *density threshold* = 3). The
20 population membership of the events in the merged file is then mapped back to the individual samples for
21 cross-sample statistics and comparisons. Phenotypes of the FLOCK-identified populations were visually
22 examined and manually annotated under the predefined cell type hierarchy (*i.e.*, CD3+CD4+ live T
23 lymphocytes).

1 **4. Discussions**

2 The most significant design of DAFi is that it implements the recursive data filtering and clustering along
3 the user-defined manual gating hierarchy, which improves the interpretability of the generated data
4 clusters, including identifying consistent numbers of cell populations across samples and generating user-
5 familiar phenotype definitions of the cell populations. DAFi does not aim to recapitulate the manual
6 gating analysis. Although the manual gating strategy is used, the results of DAFi are data-driven based on
7 the results by unsupervised clustering methods. Both the predefined and novel cell populations identified
8 by DAFi are managed under the same cell type hierarchy for knowledge integration. DAFi can work with
9 different data clustering methods for generating comparable and interpretable results. Because of these
10 characteristics, DAFi will help accelerate the adoption of computational methods by experimental
11 scientists.

12 The idea of incorporating user inputs into FCM data analysis is not new. Existing approaches such as
13 SPADE [36] and SWIFT [28] require manual operation at the end of the data clustering to group or
14 partition the data clusters into cell populations. Approaches like viSNE [5] and SPADE plot the single
15 cell data in a graph or a transformed space for providing a 2D overview of the high-dimensional data,
16 which can be difficult to interpret or operate on (*e.g.*, grouping the nodes in a SPADE tree into a cell
17 population can be error-prone without checking the events of the nodes on the original 2D plots; a viSNE
18 map is on the tSNE-transformed data space whose dimensions have no biological meaning). In contrast,
19 results of DAFi based on manual gating strategy are much easier to validate and interpret. The use of the
20 manual gating strategy guarantees the consistency of DAFi with manual gating analysis. Supplementary
21 File 9 shows that DAFi outperforms the best unsupervised method on identifying the 4 major cell
22 populations in the FlowCAP-I GvHD dataset in F-measure.

23 Sample quality control (QC) and cross-sample normalization are important components in any
24 computational pipeline of FCM data analysis. They may impact the results of DAFi. Slight data shifting is

1 not a problem for DAFi. When a data cluster is slightly shifted outside the gating boundaries, its centroid
2 remains within and its events outside of the boundaries will not be lost. However, if there is huge cross-
3 sample variance, currently we need to manually adjust the gates used in DAFi for each group of samples.
4 One solution is to integrate DAFi into a pipeline with components of QC and cross-sample normalization.
5 Though DAFi was shown to be able to address the existing challenges faced by computational methods,
6 there continue to be improvements needed in the future including eliminating the requirement for a user-
7 provided gating example, which in some cases may be unavailable. For example, one idea is to use
8 flowDensity [27] with DAFi to estimate the boundary coordinates based on 2D data distributions instead
9 of relying on predefined gating boundaries. There are also computational methods being developed to
10 identify the optimal gating path for a given set of cell population phenotypes. The Cell Ontology (CL) [7]
11 can provide a standardized cell type hierarchy to support meta-analysis across different FCM experiments.
12 Development of a graphical user interface for allowing the data analyst to create different gating
13 sequences, connect with FlowJo workspace files, compare result statistics, and integrate with other data
14 filtering and clustering methods will help improve the usability of DAFi.

15 We implemented and benchmarked the performance of DAFi on the Comet cluster at the San Diego
16 Supercomputer Center. Through parallel computing, DAFi processing and analysis of the 306 files in the
17 ImmPort SDY 180 was completed in about 30 minutes using a single compute node with 24 CPU cores.
18 We are collaborating with FlowJo to release DAFi as a plug-in tool. We are also integrating DAFi into the
19 FlowGate cyberinfrastructure [35], and implementing it on a JupyterHub server for interactive auto-gating
20 analytics of FCM data.

21 **5. Conclusions**

22 The advancement of FCM data with increased dimensionality brings in challenges in data analytics, but
23 also provides possibilities for the identification of novel cell-based biomarkers based on measurements on

1 additional combinations of markers. The large number of measured characteristics also provide
2 information to accurately define a cell population, which is facilitated by the development and use of
3 computational methods for automated identification of cell populations. How to integrate human
4 intelligence on pattern recognition with the power of computation to identify cell populations from high-
5 dimensional FCM data robustly and interpretably is a challenge that has not been sufficiently addressed.
6 In this paper, we propose a new computational method and framework - DAFi. Datasets from four
7 different study settings were used to evaluate the performance of DAFi, demonstrating DAFi's
8 characteristics in

- 9 • Generation of consistent cell type specific statistical measurements with expert centralized
10 manual gating analysis;
- 11 • Identification of natural shapes of both major and rare cell populations;
- 12 • Identification of both clearly-defined and poorly-resolved cell populations, and
- 13 • Easy interpretation and management of the identified cell populations using user-defined manual
14 gating strategy.

15

16 **Acknowledgements**

17 La Jolla Institute for Allergy and Immunology: Veronique Schulten, Jason Greenbaum; Human Longevity
18 Inc.: Rick Stanton; University of Rochester: David Topham, David Roumanes, Edward Walsh, Gloria
19 Pryhuber, Nathan Laniewski, Kristin Scheible, Jeanne Holden-Wiltse; FlowJo LLC.: Josef Spidlen; San
20 Diego Supercomputer Center: Robert Sinkovits.

21 **Author Contributions:**

22 DAFi method design: YQ and RHS; DAFi method implementation: YQ, IC, and AJL; Computational
23 processing and analysis of FCM data: AJL, IC, YQ, and JB; Flow cytometry data acquisition and manual

1 gating analysis: JB, CLA, and DW; Immunology use cases and result interpretation: RHS, BP, and AS;
2 Manuscript preparation: YQ, AJL, and RHS. All authors helped with revision of the manuscript. All
3 authors read and approved the final manuscript.

4 **Availability of Data and Software:**

5 The FCM datasets used in this manuscript have been submitted to FlowRepository under accessions:

6 FR-FCM-ZYBS, FR-FCM-ZYBT, FR-FCM-ZYBU

7 The HIPC study data used in this manuscript is also publicly available at the ImmPort database (SDY820);
8 the ImmPort SDY 180 dataset is publicly accessible at ImmPort (<https://www.immport.org>) and can be
9 downloaded at: <https://aspera-immport.niaid.nih.gov:9443/browser?path=SDY180>

10 **Literature Cited:**

11 1. Aghaeepour N, Nikolic R, Hoos HH, Brinkman RR. Rapid cell population identification in flow
12 cytometry data. *Cytometry A* 2011;79: 6–13.

13 2. Aghaeepour N, Finak G; FlowCAP Consortium; DREAM Consortium, Hoos H, Mosmann TR,
14 Brinkman R, Gottardo R, Scheuermann RH. Critical assessment of automated flow cytometry data
15 analysis techniques. *Nat Methods*. 2013 Mar;10(3):228-38. doi: 10.1038/nmeth.2365. Epub 2013 Feb 10.
16 Erratum in: *Nat Methods*. 2013 May;10(5):445. PubMed PMID: 23396282; PubMed Central PMCID:
17 PMC3906045.

18 3. Aghaeepour, N., Chattopadhyay, P., Chikina, M., Dhaene, T., Gassen, S.V., Kursu, M., Lambrecht,
19 B.N., Malek, M., Qian, Y., Qiu, P., Saeys, Y., Stanton, R., Tong, D., Vens, C., Walkowiak, S., Wang, S.,
20 Finak, G., Gottardo, R., Mosmann, T., Nolan, G., Scheuermann, R.H., Brinkman, R.R. (2016), A

- 1 Benchmark for Evaluation of Algorithms for Identification of Cellular Correlates of Clinical Outcomes,
2 *Cytometry Part A*, 89: 16–21. doi: 10.1002/cyto.a.22732. PMID: PMC4874734.
- 3 4. Aina O, Dadik J, Charurat M, Amangaman P, Gurumdi S, Mang E, Guyit R, Lar N, Datong P,
4 Daniyam C, Kanki P, Abimiku A, Institute of Human Virology/Plateau State Specialist Hospital AIDS
5 Prevention in Nigeria Study Team. Reference Values of CD4 T Lymphocytes in Human
6 Immunodeficiency Virus-Negative Adult Nigerians. *Clin Diagn Lab Immunol* 2005;12: 525–530.
- 7 5. Amir ED, Davis KL, Tadmor MD, Simonds EF, Levine JH, Bendall SC, Shenfeld DK, Krishnaswamy
8 S, Nolan GP, Pe'er D. viSNE enables visualization of high dimensional single-cell data and reveals
9 phenotypic heterogeneity of leukemia. *Nat. Biotechnol.* 2013;31: 545-552. PMID: PMC4076922.
- 10 6. Arlehamn CL, Seumois G, Gerasimova A, Huang C, Fu Z, Yue X, Sette A, Vijayanand P, Peters B.
11 Transcriptional profile of tuberculosis antigen-specific T cells reveals novel multifunctional features. *J.*
12 *Immunol.* 2014;193: 2931–2940.
- 13 7. Bard J, Rhee SY, Ashburner M. An ontology for cell types. *Genome Biol* 2005;6: 2, R21. PMID:
14 PMC551541.
- 15 8. Bashashati A, Brinkman RR. A Survey of Flow Cytometry Data Analysis Methods. *Advances in*
16 *Bioinformatics* 2009. PMID: PMC2798157.
- 17 9. Brinkman RR, Aghaeepour N, Finak G, Gottardo R, Mosmann T, Scheuermann RH. State-of-the-Art in
18 the Computational Analysis of Cytometry Data. *Cytometry A.* 2015 Jul;87(7):591-3. doi:
19 10.1002/cyto.a.22707. PubMed PMID: 26111230.
- 20 10. Brinkman RR, Aghaeepour N, Finak G, Gottardo R, Mosmann T, Scheuermann RH. Automated
21 analysis of flow cytometry data comes of age. *Cytometry A.* 2016 Jan;89(1):13-5. doi:
22 10.1002/cyto.a.22810. PubMed PMID: 26812230.

- 1 11. Bruggner RV, Bodenmiller B, Dill DL, Tibshirani RJ, Nolan GP. Automated identification of
2 stratifying signatures in cellular subpopulations. *Proc. Natl. Acad. Sci. U.S.A.* 2014;111: E2770-2777.
- 3 12. Burel JG, Qian Y, Lindestam Arlehamn C, Weiskopf D, Zapardiel-Gonzalo J, Taplitz R, Gilman RH,
4 Saito M, de Silva AD, Vijayanand P, Scheuermann RH, Sette A, Peters B. An Integrated Workflow To
5 Assess Technical and Biological Variability of Cell Population Frequencies in Human Peripheral Blood
6 by Flow Cytometry. *J. Immunol.* 2017;198: 1748–1758.
- 7 13. Chester C, Maecker HT. Algorithmic Tools for Mining High-Dimensional Cytometry Data. *J*
8 *Immunol* 2015;195: 773–779.
- 9 14. Cron A, Gouttefangeas C, Frelinger J, Lin L, Singh SK, Britten CM, Welters MJP, van der Burg SH,
10 West M, Chan C. Hierarchical modeling for rare event detection and cell subset alignment across flow
11 cytometry samples. *PLoS Comput. Biol.* 2013;9: e1003130.
- 12 15. Finak G, Bashashati A, Brinkman R, Gottardo R. Merging mixture components for cell population
13 identification in flow cytometry. *Adv Bioinformatics* 2009:247646.
- 14 16. Finak G, Frelinger J, Jiang W, Newell EW, Ramey J, Davis MM, Kalams SA, De Rosa SC, Gottardo
15 R. OpenCyto: an open source infrastructure for scalable, robust, reproducible, and automated, end-to-end
16 flow cytometry data analysis. *PLoS Comput. Biol.* 2014;10: e1003806.
- 17 17. García-Dabrio MC, Pujol-Moix N, Martinez-Perez A, Fontcuberta J, Souto JC, Soria JM, Nomdedéu
18 JF. 2012. Influence of age, gender and lifestyle in lymphocyte subsets: report from the Spanish Gait-2
19 Study. *Acta Haematol.* 27:244–249.
- 20 18. Ge Y, Sealfon SC. flowPeaks: a fast unsupervised clustering for flow cytometry data via K-means and
21 density peak finding. *Bioinformatics* 2012;28: 2052–2058.

- 1 19. Hahne F, Khodabakhshi AH, Bashashati A, Wong C-J, Gascoyne RD, Weng AP, Seyfert-Margolis V,
2 Bourcier K, Asare A, Lumley T, Gentleman R, Brinkman RR. Per-channel basis normalization methods
3 for flow cytometry data. *Cytometry A* 2010;77: 121–131.
- 4 20. Kam KM, Leung WL, Kwok MY, Hung MY, Lee SS, Mak WP. 1996. Lymphocyte subpopulation
5 reference ranges for monitoring human immunodeficiency virus-infected Chinese adults. *Clin. Diagn. Lab.*
6 *Immunol.* 3:326.
- 7 21. Kidd BA, Peters LA, Schadt EE, Dudley JT. Unifying immunology with informatics and multiscale
8 biology. *Nature Immunology* 2014;15: 118–127.
- 9 22. Kvistborg P, Gouttefangeas C, Aghaeepour N, Cazaly A, Chattopadhyay PK, Chan C, Eckl J, Finak G,
10 Hadrup SR, Maecker HT, Maurer D, Mosmann T, Qiu P, Scheuermann RH, Welters MJ, Ferrari G,
11 Brinkman RR, Britten CM. Thinking outside the gate: single-cell assessments in multiple dimensions.
12 *Immunity.* 2015 Apr 21;42(4):591-2. doi: 10.1016/j.immuni.2015.04.006. PubMed PMID: 25902473.
- 13 23. Lee H-C, Kosoy R, Becker CE, Dudley JT, Kidd BA. Automated cell type discovery and
14 classification through knowledge transfer. *Bioinformatics* 2017;33: 1689–1695.
- 15 24. Lo K, Hahne F, Brinkman RR, Gottardo R. flowClust: a Bioconductor package for automated gating
16 of flow cytometry data. *BMC Bioinformatics* 2009;10: 145.
- 17 25. Lugli E, Roederer M, Cossarizza A. Data analysis in flow cytometry: the future just started.
18 *Cytometry A* 2010;77: 705–713.
- 19 26. Mair F, Hartmann FJ, Mrdjen D, Tosevski V, Krieg C, Becher B. The end of gating? An introduction
20 to automated analysis of high dimensional cytometry data. *Eur. J. Immunol.* 2016;46: 34–43.

- 1 27. Malek M, Taghiyar MJ, Chong L, Finak G, Gottardo R, Brinkman RR. flowDensity: reproducing
2 manual gating of flow cytometry data by automated density-based cell population identification.
3 *Bioinformatics* 2015;31: 606–607.
- 4 28. Naim I, Datta S, Rebhahn J, Cavanaugh JS, Mosmann TR, Sharma G. SWIFT-scalable clustering for
5 automated identification of rare cell populations in large, high-dimensional flow cytometry datasets, part
6 1: algorithm design. *Cytometry A* 2014;85: 408–421.
- 7 29. Obermoser G, Presnell S, Domico K, Xu H, Wang Y, Anguiano E, Thompson-Snipes L, Ranganathan
8 R, Zeitner B, Bjork A, Anderson D, Speake C, Ruchaud E, Skinner J, Alsina L, Sharma M, Dutartre H,
9 Cepika A, Israelsson E, Nguyen P, Nguyen Q-A, Harrod AC, Zurawski SM, Pascual V, Ueno H, Nepom
10 GT, Quinn C, Blankenship D, Palucka K, Banchereau J, Chaussabel D. Systems scale interactive
11 exploration reveals quantitative and qualitative differences in response to influenza and pneumococcal
12 vaccines. *Immunity* 2013; 38:831–844; PMID: PMC3681204.
- 13 30. Parks DR, Roederer M, Moore WA. A new “Logicle” display method avoids deceptive effects of
14 logarithmic scaling for low signals and compensated data. *Cytometry A* 2006;69: 541–551.
- 15 31. Pedersen NW, Chandran PA, Qian Y, Rebhahn J, Petersen NV, Hoff MD, White S, Lee AJ, Stanton R,
16 Halgreen C, Jakobsen K, Mosmann T, Gouttefangeas C, Chan C, Scheuermann RH, Hadrup SR.
17 Automated Analysis of Flow Cytometry Data to Reduce Inter-Lab Variation in the Detection of Major
18 Histocompatibility Complex Multimer-Binding T Cells. *Front Immunol* 2017;8: 858.
- 19 32. Pyne S, Hu X, Wang K, Rossin E, Lin T, Maier LM, Baecher-Allan C, McLachlan GJ, Tamayo P,
20 Hafler DA, De Jager PL, Mesirov JP. Automated high-dimensional flow cytometric data analysis. *PNAS*
21 2009; 106(21):8519-8524. PMID: PMC2682540

- 1 33. Qian Y, Wei C, Eun-Hyung Lee F, Campbell J, Halliley J, Lee JA, Cai J, Kong YM, Sadat E,
2 Thomson E, Dunn P, Seegmiller AC, Karandikar NJ, Tipton CM, Mosmann T, Sanz I, Scheuermann RH.
3 Elucidation of seventeen human peripheral blood B-cell subsets and quantification of the tetanus response
4 using a density-based method for the automated identification of cell populations in multidimensional
5 flow cytometry data. *Cytometry B Clin Cytom.* 2010;78 Suppl 1: S69-82. doi: 10.1002/cyto.b.20554.
6 PubMed PMID: 20839340; PubMed Central PMCID: PMC3084630.
- 7 34. Qian Y, Liu Y, Campbell J, Thomson E, Kong YM, Scheuermann RH. FCSTrans: an open source
8 software system for FCS file conversion and data transformation. *Cytometry A.* 2012 May;81(5):353-6.
9 PMCID: PMC3932304.
- 10 35. Qian Y, Kim H, Purawat S, Wang J, Stanton R, Lee A, Xu W, Altintas I, Sinkovits R, Scheuermann
11 RH. FlowGate: Towards Extensible and Scalable Web-Based Flow Cytometry Data Analysis, In
12 Proceedings of the 4th annual XSEDE (Extreme Science and Engineering Discovery Environment)
13 conference. 2015. St. Louis, ACM Press.
- 14 36. Qiu P, Simonds EF, Bendall SC, Gibbs KD Jr, Bruggner RV, Linderman MD, Sachs K, Nolan GP,
15 Plevritis SK. Extracting a cellular hierarchy from high-dimensional cytometry data with SPADE. *Nat.*
16 *Biotechnol.* 2011;29: 886-891. PMCID: PMC3196363.
- 17 37. Rudy BJ, Wilson CM, Durako S, Moscicki A-B, Muenz L, Douglas SD. Peripheral Blood
18 Lymphocyte Subsets in Adolescents: a Longitudinal Analysis from the REACH Project. *Clin. Diagn. Lab.*
19 *Immunol.* 2002;9: 959–965.
- 20 38. Saeys Y, Gassen SV, Lambrecht BN. Computational flow cytometry: helping to make sense of high-
21 dimensional immunology data. *Nature Reviews Immunology* 2016;16: 449–462.

- 1 39. Thakar MR, Abraham PR, Arora S, Balakrishnan P, Bandyopadhyay B, Joshi AA, Devi R,
2 Vasanthapuram R, Vajpayee M, Desai A, Mohanakrishnan J, Narain K, Ray K, Patil SS, Singh R, Singla
3 A, Paranjape RS. 2011. Establishment of reference CD4+ T cell values for adult Indian population. *AIDS*
4 *Res. Ther.* 8:35.
- 5 40. Tollerud, DJ, Clark JW, Brown LM, Neuland CY, Pankiw-Trost LK, Blattner WA, and Hoover RN.
6 1989. The influence of age, race, and gender on peripheral blood mononuclear-cell subsets in healthy
7 nonsmokers. *J. Clin. Immunol.* 9:214-222.
- 8 41. Tollerud, DJ, Ildstad ST, Brown LM, Clark JW, Blattner WA, Mann DL, Neuland CY, Pankiw-Trost
9 L, and Hoover RN. 1990. T-cell subsets in healthy teenagers: transition to the adult phenotype. *Clin.*
10 *Immunol. Immunopathol.* 56:88-96.
- 11 42. Tugume, SB, Piwowar EM, Lutalo T, Mugenyi PN, Grant RM, Mangeni FW, Pattishall K, and
12 Katongole-Mbidde E. 1995. Hematological reference ranges among healthy Ugandans. *Clin. Diagn. Lab.*
13 *Immunol.* 2:233-235.
- 14 43. Uppal SS, Verma S, and Dhot PS. 2003. Normal values of CD4 and CD8 lymphocyte subsets in
15 healthy Indian adults and the effects of sex, age, ethnicity and smoking. *Cytometry* 52B:32-36.
- 16 44. Van Gassen S, Callebaut B, Van Helden MJ, Lambrecht BN, Demeester P, Dhaene T, Saeys Y.
17 *FlowSOM: Using self-organizing maps for visualization and interpretation of cytometry data. Cytometry*
18 *A*, 2015;87:636–645.
- 19 45. Wagstaff K, Cardie C, Rogers S, Schrödl S, Constrained K-means Clustering with Background
20 Knowledge, *Proceedings of the Eighteenth International Conference on Machine Learning*, p.577-584,
21 June 28-July 01, 2001.

1 46. Zare H, Shoostari P, Gupta A, Brinkman RR. Data reduction for spectral clustering to analyze high
2 throughput flow cytometry data. BMC Bioinformatics 2010;11: 403.

3 **Figure Legends:**

4 **Figure 1. Design features of DAFi.** A) Steps in the DAFi workflow. In Step 1, putative cell populations
5 are identified by data clustering in high dimensional space, with cell events colored by population
6 membership. In Step 2, a hyper-polygon is provided from combining 2D manual gating boundaries to
7 identify the dataspace region of interest. Cell clusters are selected if their centroids are located within the
8 hyper-polygon (two clusters shown, in light blue and magenta). In Step 3, all cell events associated with
9 the centroids are selected and retained as the filtered population (in red), which is used as the input to the
10 next iteration in Step 4. B) An example gating hierarchy in which the DAFi framework can be used to
11 identify both predefined (solid lines) and novel (dotted lines) cell populations, and organize them within a
12 user-provided gating hierarchy for simplified annotation and interpretation. C) DAFi identification of
13 CD4+T, CD8+T, CD3+CD56+ T and CD3hiCD56+ T cells. CD4+T and CD8+T cells are shown on CD4
14 vs CD8 dot plots, while CD3+CD56+ T and CD3hiCD56+ T cells are on CD3 vs CD56 plots. Cell
15 populations identified by DAFi are colored in red. D) Comparison of DAFi with other clustering and
16 filtering methods. Putative CD4+CD25+ regulatory T cells (Tregs) were identified using *K*-means
17 clustering, manual gating, and DAFi. The identified Treg cells are colored in red and the remaining cells
18 colored in white.

19 **Figure 2. Performance evaluation of DAFi in comparison with individual and centralized manual**
20 **gating analysis.** A) Illustration of the manual gating hierarchy for identifying the 22 predefined cell
21 populations from the 10-color T cell panel, with gating boundaries shown on each 2D dot plot. Along the
22 direction of the red arrows is the sequence of the gates with their parent populations. The cell populations
23 are numbered. Names of the cell types are listed to the right. B) DAFi results for identifying the
24 corresponding 22 predefined cell populations. Events from the whole sample are colored in white. The

1 black colored dots are events of the parent population, with events identified by DAFi highlighted in red,
2 yellow, green, and blue. C) T cells (Blue), NK cells (Red), CD3+CD56+ T (yellow), and CD3hiCD56+
3 T (green) from four data files are shown. D) Naïve CD4+ T (red), effector memory CD4+ T (Tem CD4,
4 green), central memory CD4+ T (Tcm CD4, yellow), and effector memory CD4+ T expressing CD45RA
5 (Temra CD4, blue) identified by DAFi from four data files are shown. E) Regulatory T cells (Tregs,
6 CD4+CD25+, red) identified by DAFi from four data files are shown. F) Linear regression and
7 correlation analysis of cell population percentages of clearly-defined cell populations identified by DAFi
8 (y-axis) compared with centralized manual gating (x-axis). G) Linear regression and correlation analysis
9 of percentages of poorly-resolved cell populations identified by DAFi (y-axis) compared with centralized
10 manual gating analysis (x-axis). H) P-values (-log₁₀ transformed) of x-variable in linear regression
11 analysis between percentages generated by data clustering methods (DAFi and K-means) and centralized
12 manual gating analysis for clearly-defined cell populations. The cell populations were sorted based on
13 their average percentage of their parents from the largest to the smallest, as shown on the x-axis. I) P-
14 values (-log₁₀ transformed) of linear regression analysis between percentages generated by data
15 clustering methods (DAFi and K-means) and centralized manual gating analysis for poorly-resolved cell
16 populations. J) Coefficient variability (CV) of population percentages across the 24 samples for clearly-
17 defined cell populations. K) CV of population percentages across the 24 samples for poorly-resolved cell
18 populations.

19 **Figure 3. Correlation analysis of DAFi-defined cell population proportions with subject age and**
20 **gender.** A) Age distribution of participants separated by gender. B) Proportions of naïve CD4 T and naïve
21 CD8 T cells (with CD4+ and CD8+ T cells as parents, respectively) versus age with linear regression p-
22 value reported. C) Pearson correlation and linear regression analysis of proportions of T cell subsets with
23 subject age. Parent population definitions of the T-cell subsets can be found in Figure 2A. P-values of x-
24 variable in linear regression analysis were -log₁₀ transformed and multiple comparison corrected by

1 Bonferroni correction. D) Proportion of CD4+ T cells in female and male participants. E) Correlation
2 between the proportions of effector memory T cells versus naïve T cells.

3 **Figure 4. Identification of known and novel cell biomarkers for LTBI using FLOCK-based**
4 **clustering of DAFi filtered populations.** A) Upper: manual gating strategy for identifying CD4+ T cells.
5 The gating path sequentially identifies lymphocytes (FSC-A vs. SSC-A), singlet lymphocytes based on
6 FSC-A/W, singlet lymphocytes based on SSC-A/W, live CD8- T lymphocytes (the DUMP channel
7 includes CD8/CD14/CD19/LiveDead), and CD3+CD4+ T lymphocytes. Lower: CD4+ T cell population
8 identified by DAFi highlighted on corresponding 2D dot plots with density contour lines showing the
9 natural data distribution. A hyper-polygon based on gating boundaries on FSC-A/W, SSC-A/W, DUMP,
10 CD3, and CD4 was used to identify FLOCK data clusters within the CD4+ T cell population. B) Manual
11 gating strategy for identifying subset populations from the CD4+ T cells, based on CD25, CCR7,
12 CD45RA, CCR4, CCR6, and CXCR3 expression. Note that tetramer staining was not used in the manual
13 gating analysis, and different memory T cell regions were not separated based on CCR7 vs. CD45RA. C-
14 D) Percentages of the five most significant DAFi-identified cell subsets (CD4+ T cell population as
15 parent) that differed between LTBI and HC and their corresponding p-values in Wilcoxon rank sum test.
16 Mean and standard deviation of the percentage values of each cell population are shown with the
17 individual values. When N=12, 0.003948 is the best possible *p*-value with the rank sum test when there is
18 no overlap between the ranking of the two groups. E) The three DAFi-identified cell subsets that differ
19 between LTBI and HC in the known CD25-CCR6+CCR4-CXCR3+ region. Events are highlighted in red
20 and shown on different 2D dot plots. The three subsets (Pop#23: Tet-CCR7-CD45RA-, Pop#27: Tet-
21 CCR7+CD45RA-, and Pop#65: Tet+CCR7+CD45RA-) differ from each other based on tetramer and
22 CCR7. F) The two DAFi-identified cell subsets that were not reported in the previous publication with
23 their events highlighted in red on 2D plots of different markers. Both are very rare (average < 0.1% of
24 CD4+ T cells). Pop#18 is Tet+CD25-CCR6+CCR4dimCXCR3+ while Pop#28 is Tet-CD25-CCR6-
25 CCR4-CXCR3+. G) tSNE map of the filtered data. FLOCK clusters of CD4+ T cells (generated without

1 scatter/DUMP/CD3/CD4 values as input) are color-coded based on expression level of tetramer to
2 highlight the tetramer+ population in the mid-upper left region. H) Zoomed-in tSNE map shows that the
3 “island” of the tetramer+ population consists of two separated regions, corresponding to the Pop#18
4 (highlighted in yellow) and the Pop#65 (highlighted in blue).

5 **Figure 5. Quantification of human immune response to influenza and pneumococcal vaccination**
6 **using DAFi.** From left to right under each vaccine/saline treatment are three selected time points from
7 one individual in each treatment group: 7 days before the treatment (Day -7), and Day 7 and Day 28 after
8 treatment. A) CD19+ B cells were identified by DAFi using the 2D rectangular gates in FSC/SSC-A and
9 CD19/SSC-A plots illustrated in the first two rows. The two following rows show the B-cell events
10 (colored in blue) on IgD vs CD27 and CD20 vs CD138 dot pots. B) Plasmablast cells identified by DAFi
11 from the CD19+ B cell population. The plasmablasts, defined as IgD-CD27high, are shown in the red box.
12 C) Percentage of plasmablast cells (with CD19+B cell as parent) identified across times and treatment
13 groups by DAFi in box plots. D) Normalized proportions of the plasmablast population (with CD19+ B
14 cell as parent) identified by DAFi and manual gating analysis across times and treatment groups.

15 **Table 1. An Example Configuration Table Used in DAFi to Specify Gating Boundary Coordinates**
16 **and Hierarchical Relationships.** Each row corresponds to a cell population defined by the user with
17 Population ID (*PopID*) and Population Name (*PopName*). The two markers used to define the gate were
18 *Xname* and *Yname*. Different data ranges across instruments and experiments are 0-1 min-max normalized
19 into 0-200. DAFi configuration simplified the shapes of the gate by using rectangle gates only. A
20 rectangle gate on a 2D plot is defined by four values: *Xmin*, *Xmax*, *Ymin*, and *Ymax*. For one-dimensional
21 gate on X-axis only, the *Yname* can be an arbitrary dimension and *Ymin/Ymax* can be voided by putting
22 0/200. The parent of each cell population is specified by the ID of its parent population (*ParentID*), with
23 the *Mode* specifying the way of identifying the cell population from its parent (cluster, slope, bisecting,
24 and reversed).

1

2 **Supplementary Materials:**

3 Supplementary File 1: Definitions of 22 cell populations using a 10-color T-cell FCM reagent panel.

4 Supplementary File 2: Different T cell subsets identified by DAFi across 24 repeat runs of PBMC from a
5 control donation, including clearly defined and poorly resolved ones.

6 Supplementary File 3: Proportions of the 22 predefined cell populations identified by DAFi across 132
7 subjects for correlation with age and gender

8 Supplementary File 4: Results of DAFi prefiltering by applying FLOCK clustering method across 12
9 samples in LTBI study. From left to right: lymphocytes, singlet lymphocytes, live singlet lymphocytes,
10 and CD3+CD4+ live T lymphocytes.

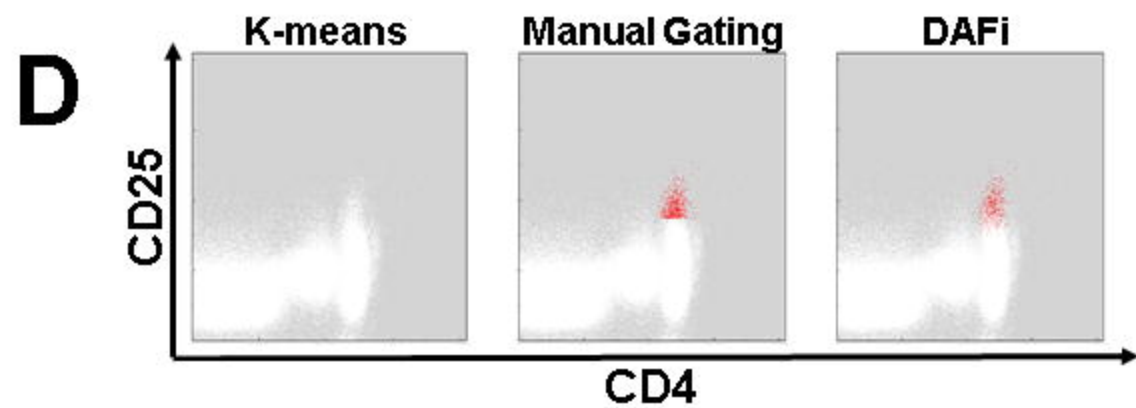
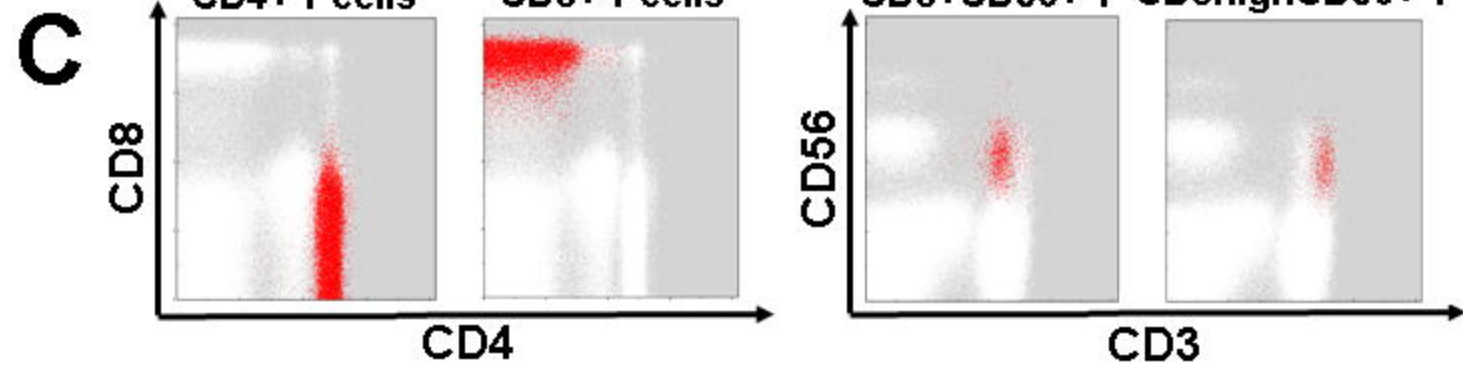
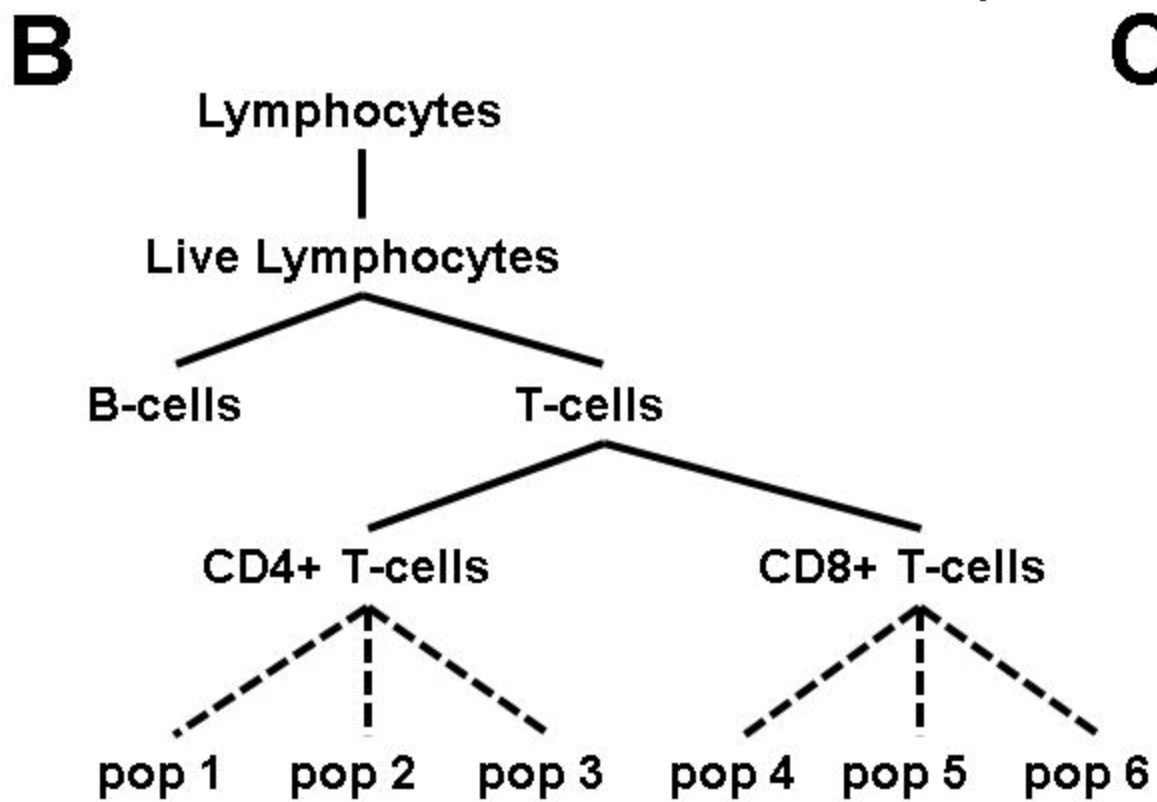
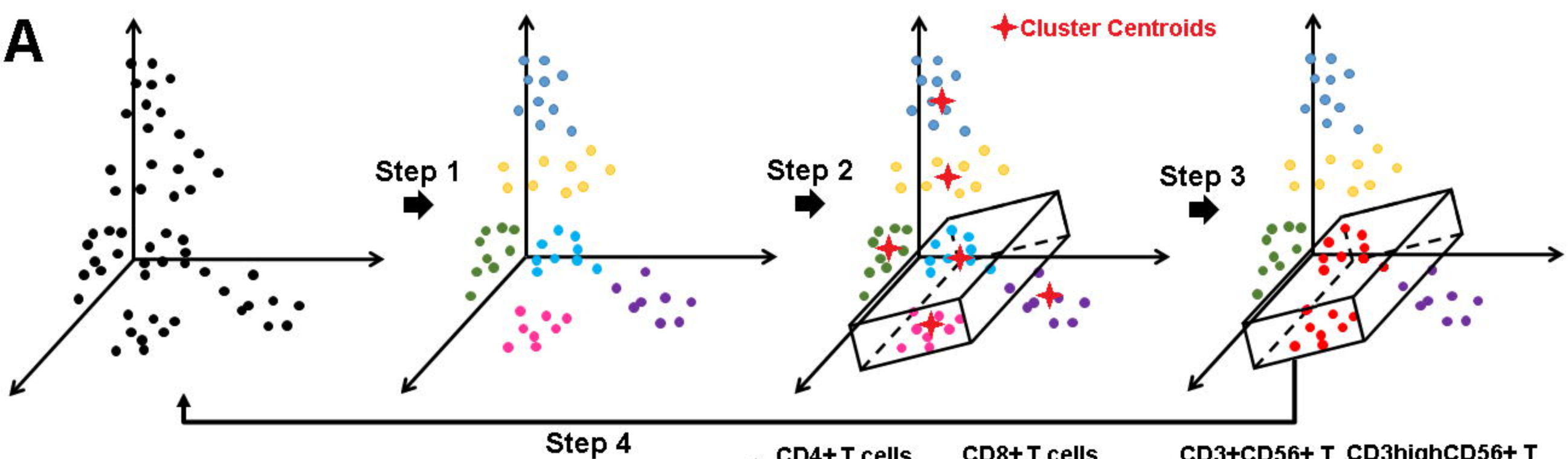
11 Supplementary File 5: Population percentages of the 101 cell populations identified by FLOCK (with
12 CD3+CD4+live T lymphocytes as parent) and their rank sum test p -values between LTBI and HC groups.

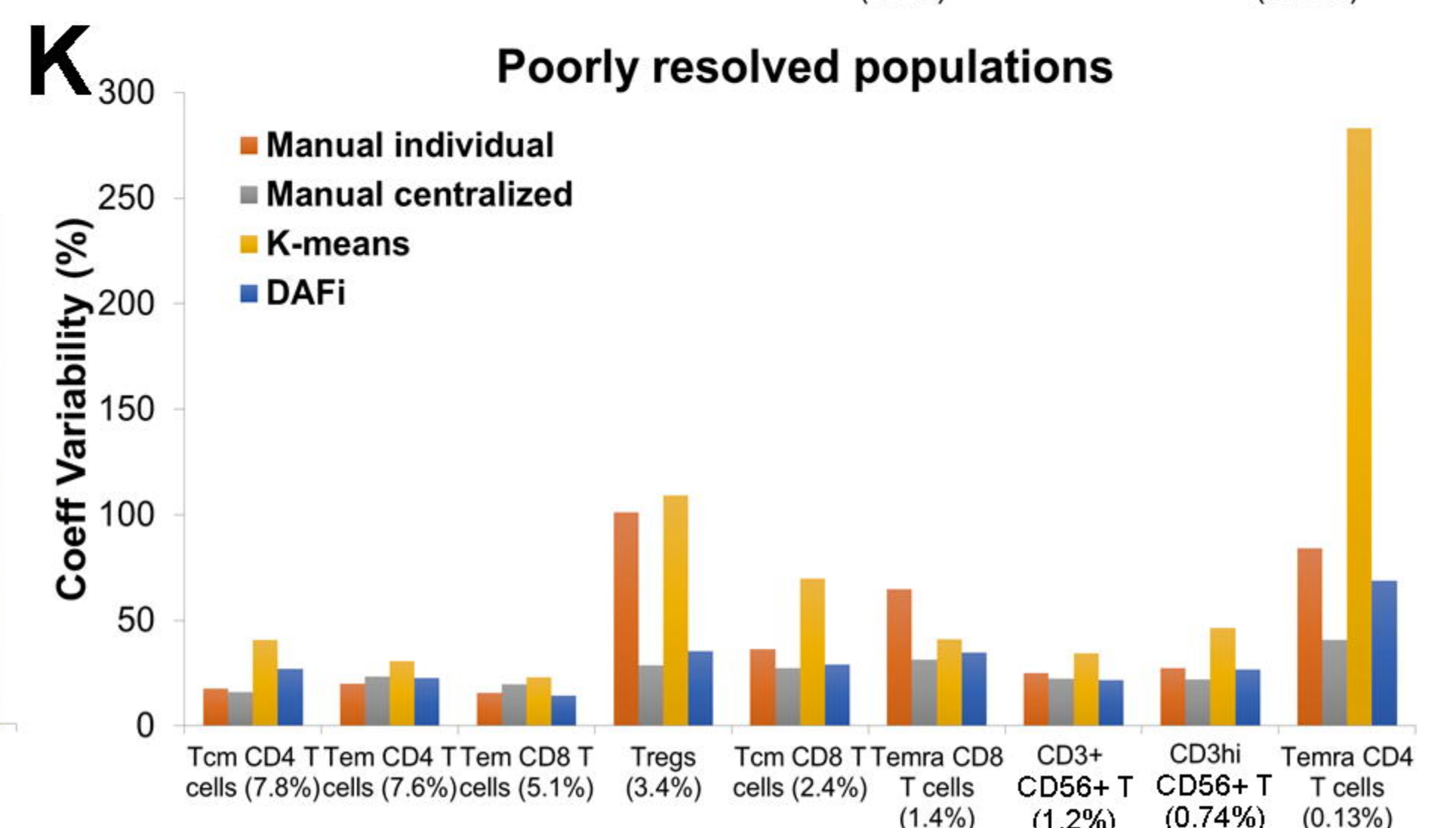
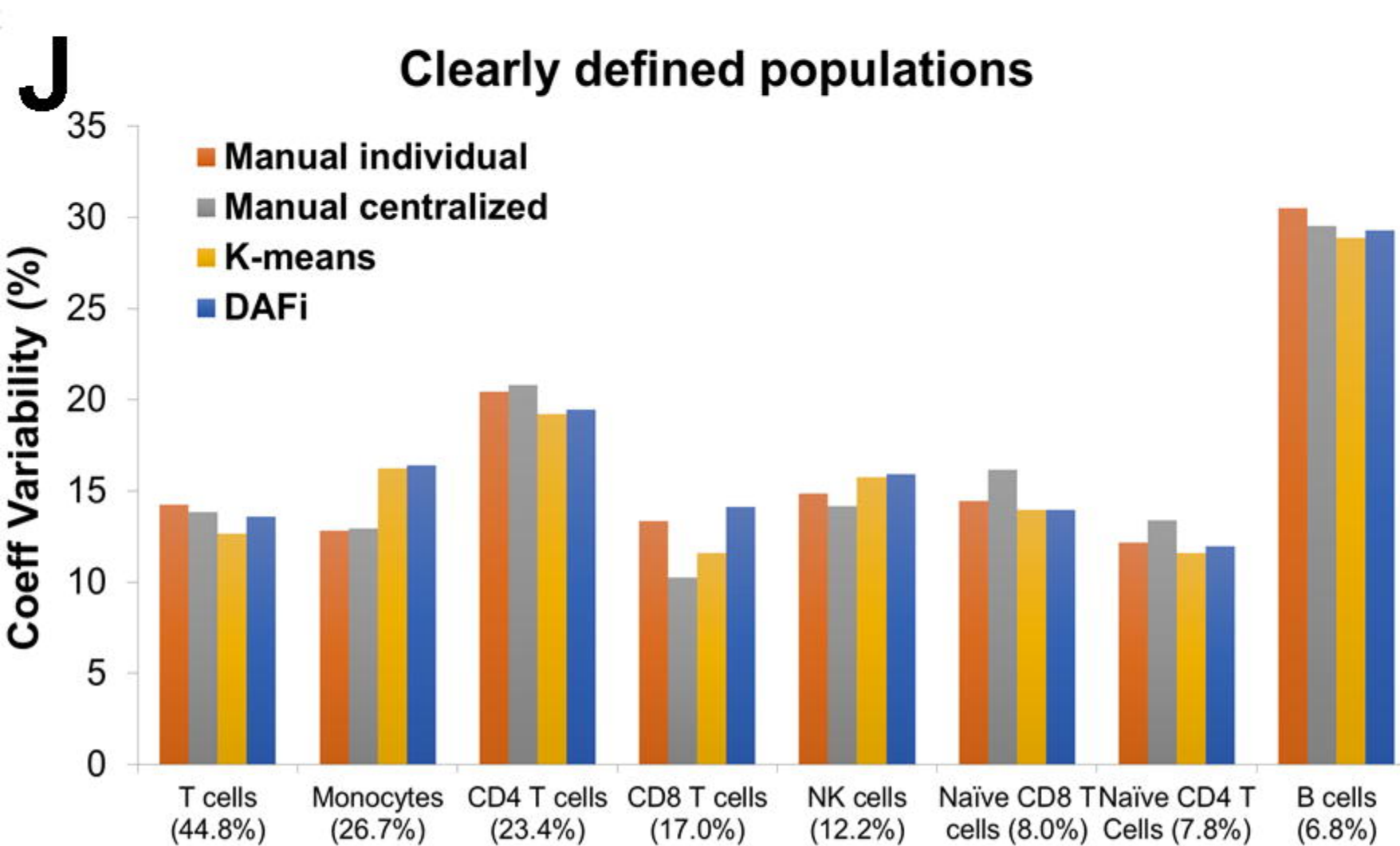
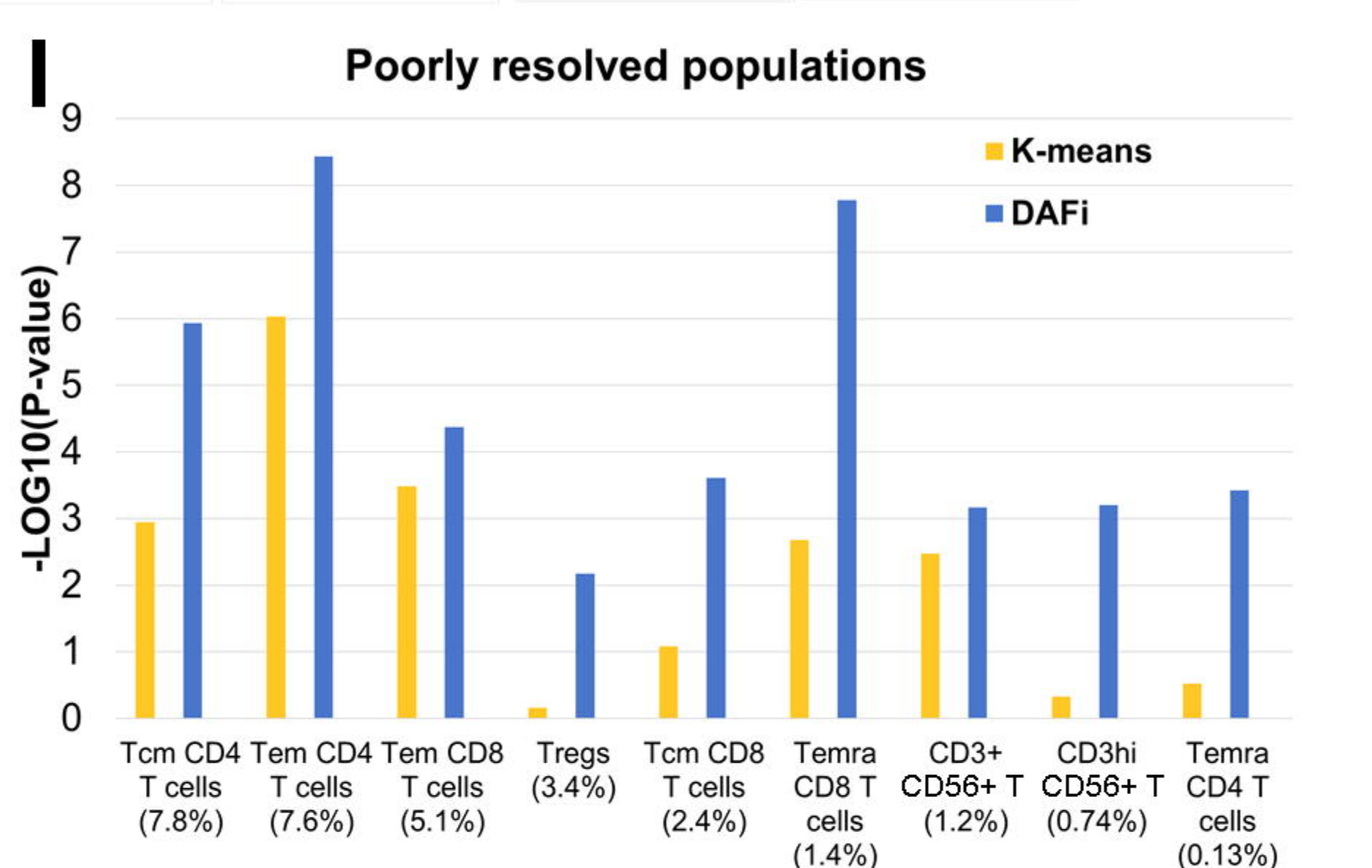
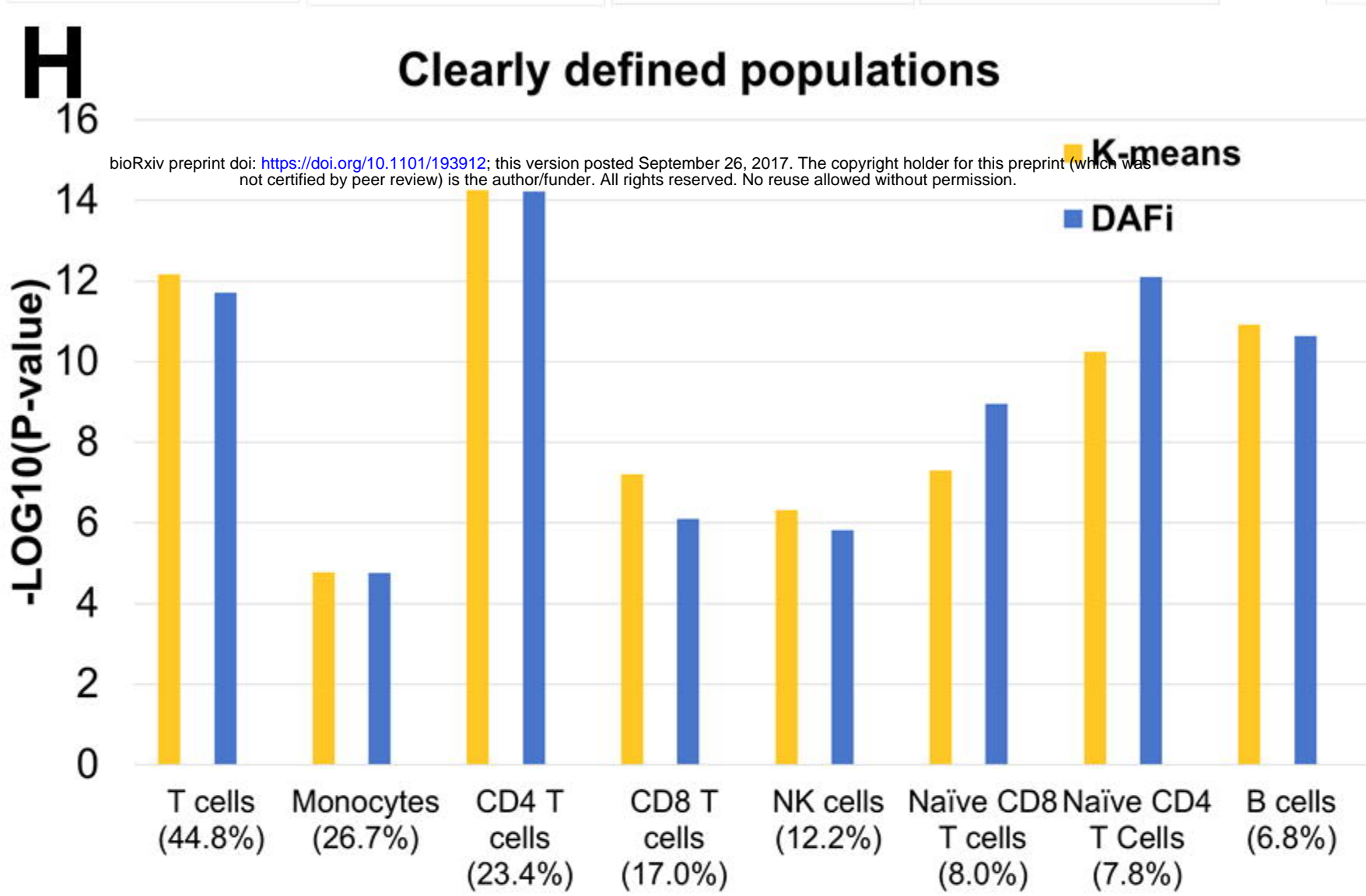
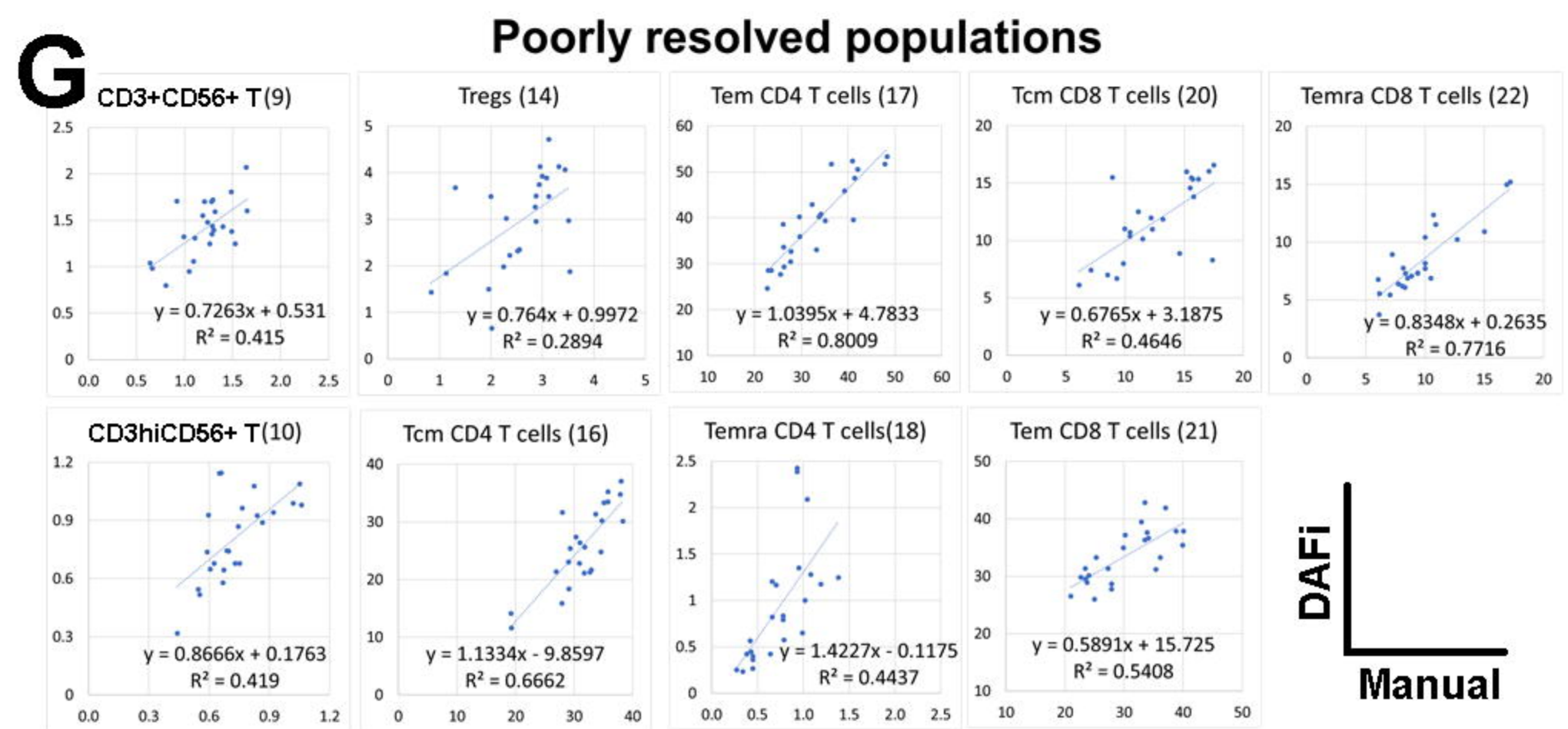
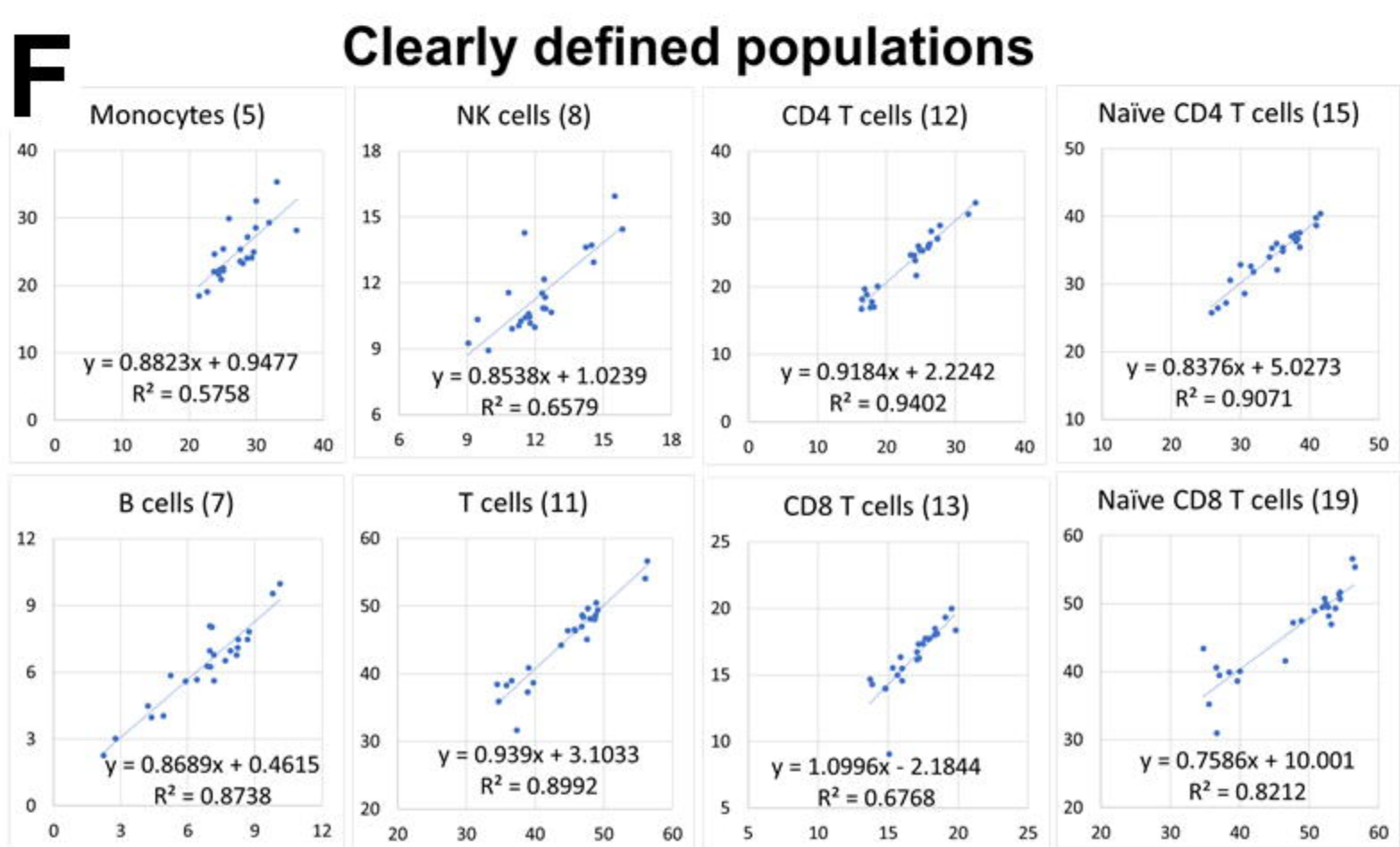
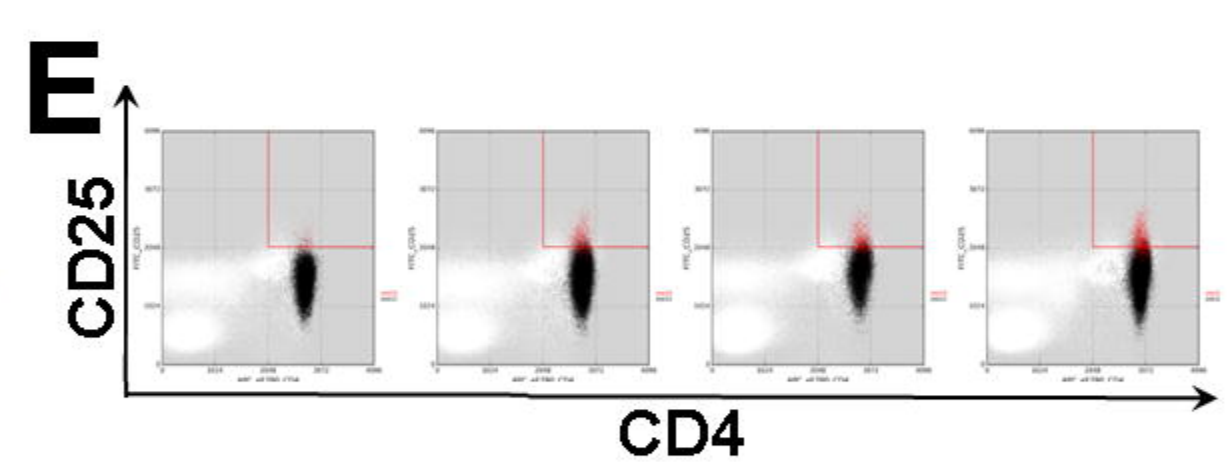
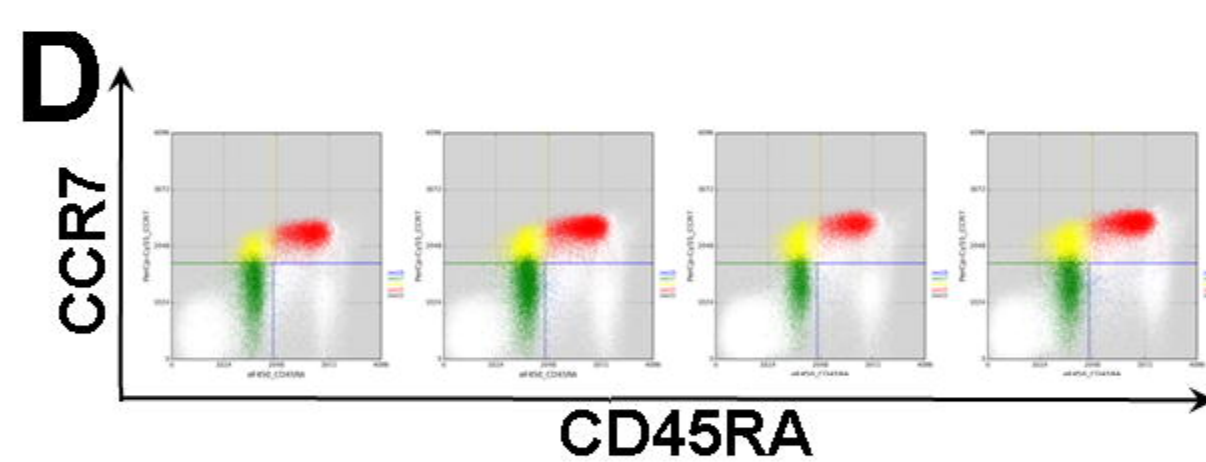
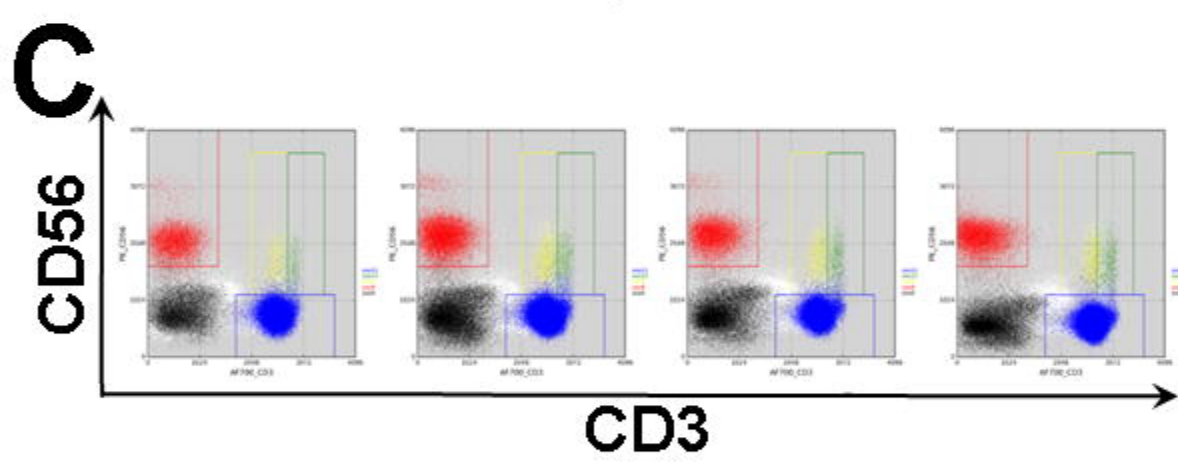
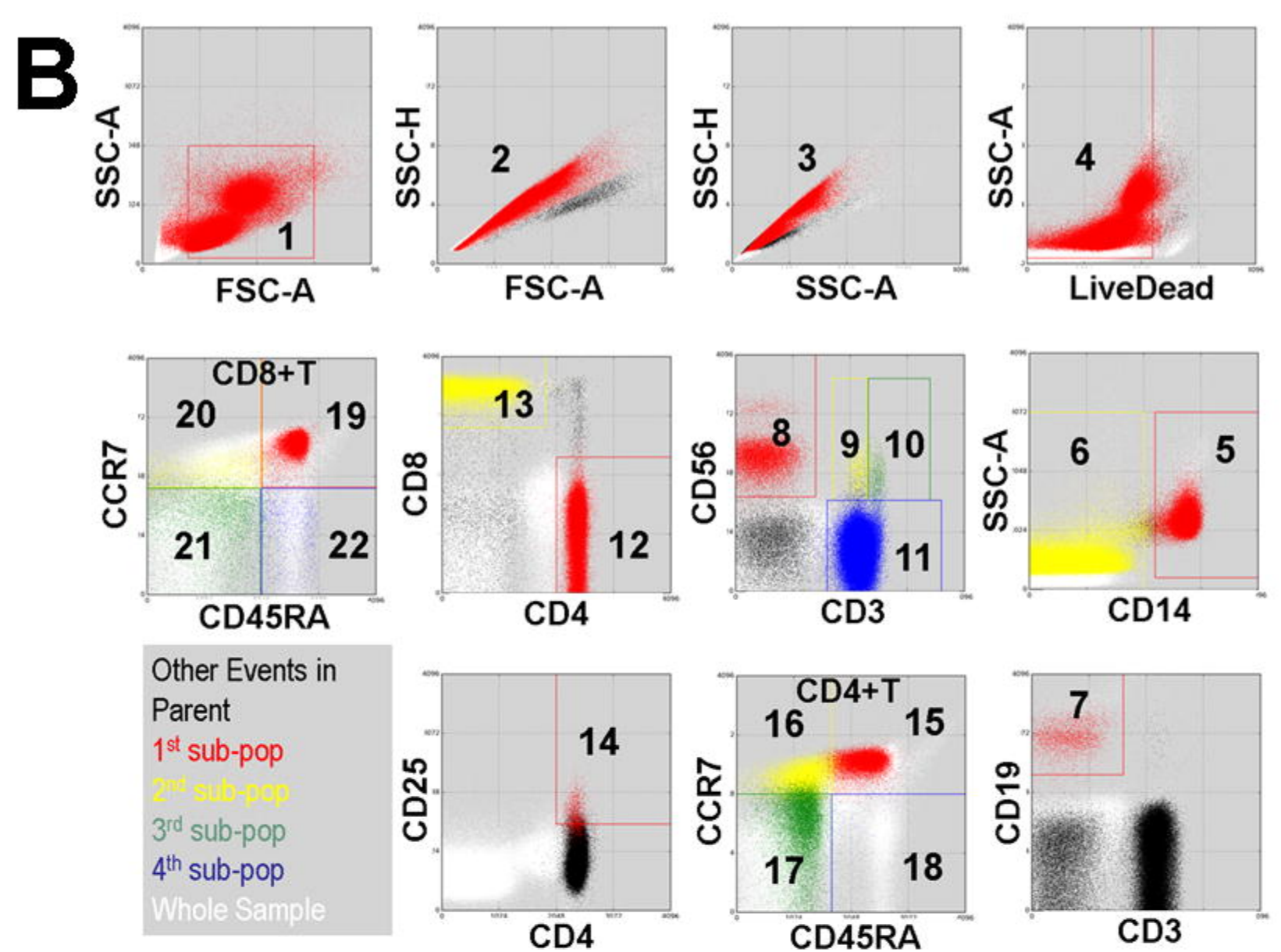
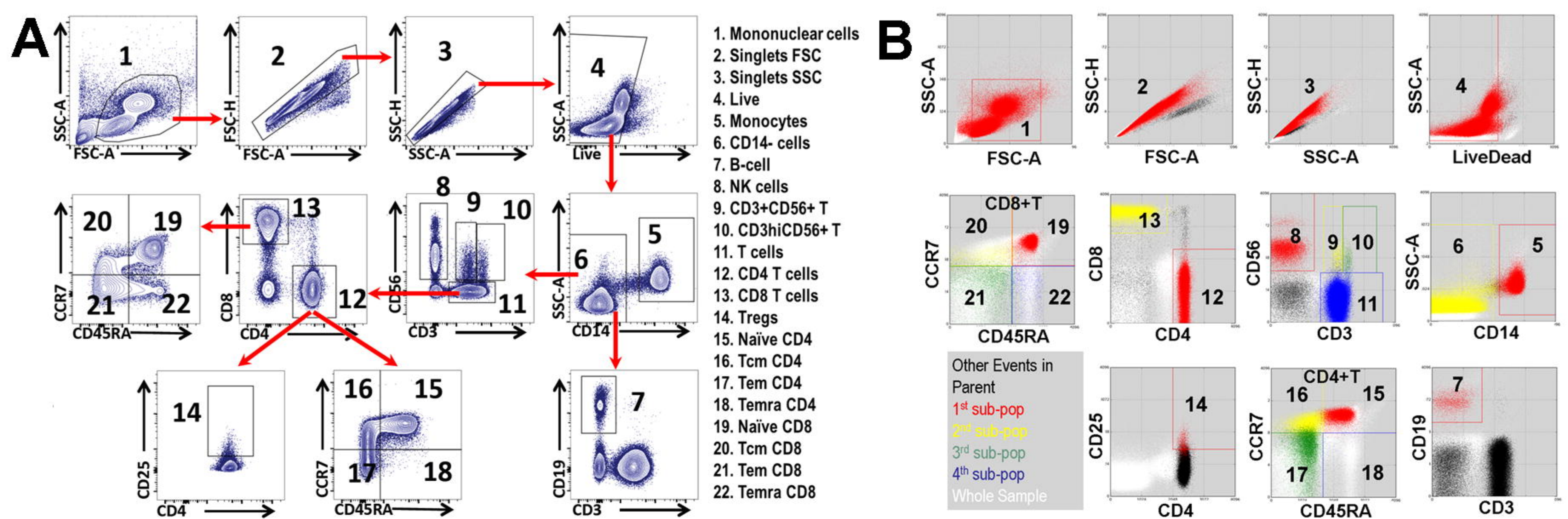
13 Supplementary File 6: F1 scores of using different K in K -means clustering for DAFi-filtering across the
14 12 samples in LTBI study for each cell population compared with bisecting analysis.

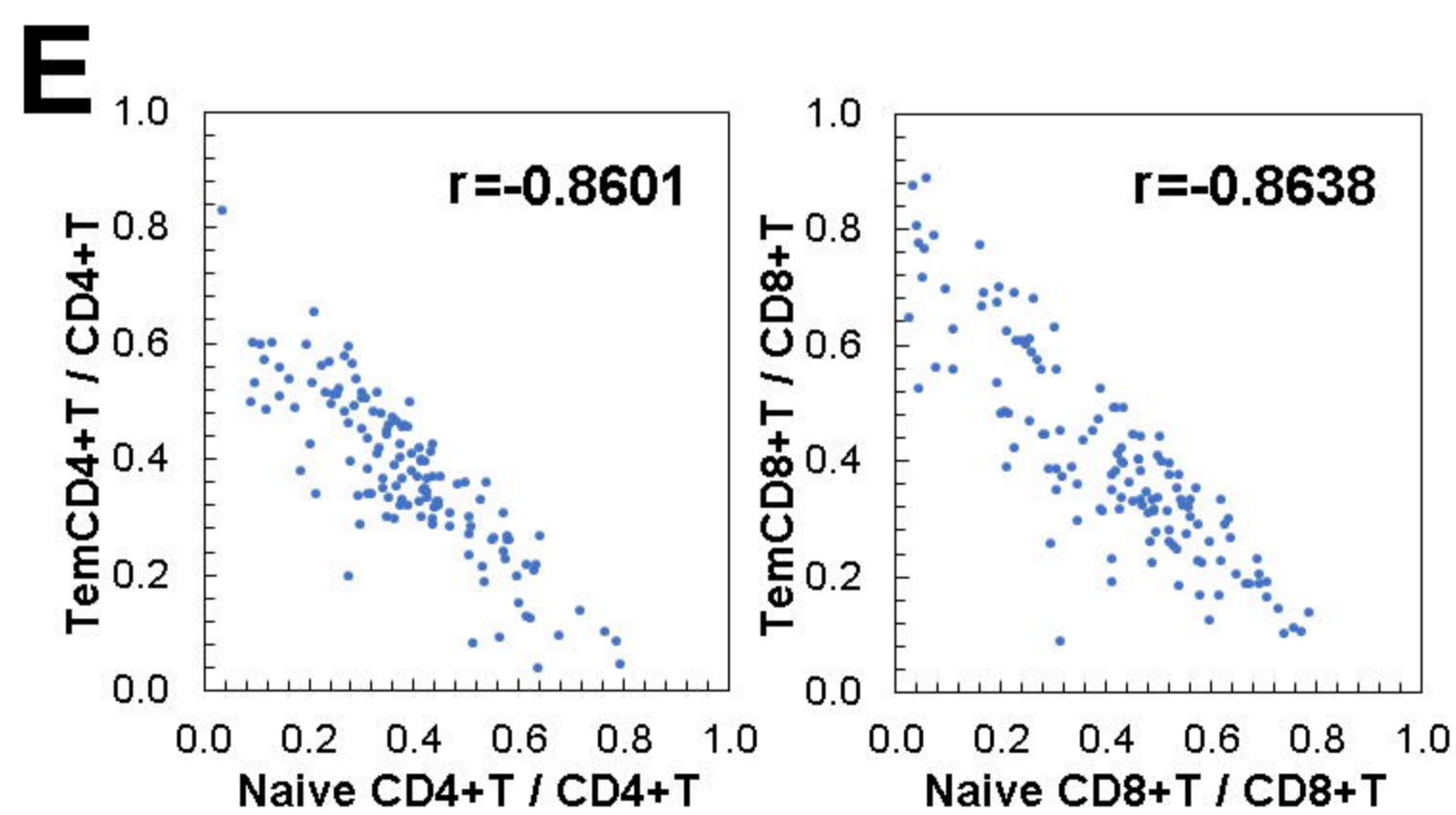
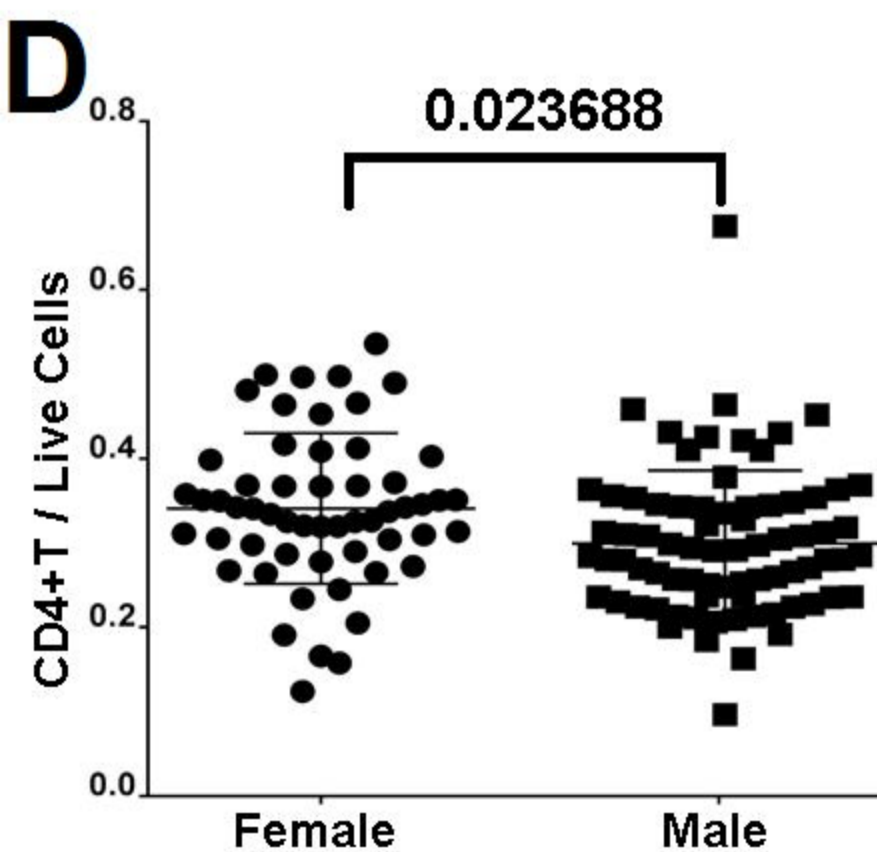
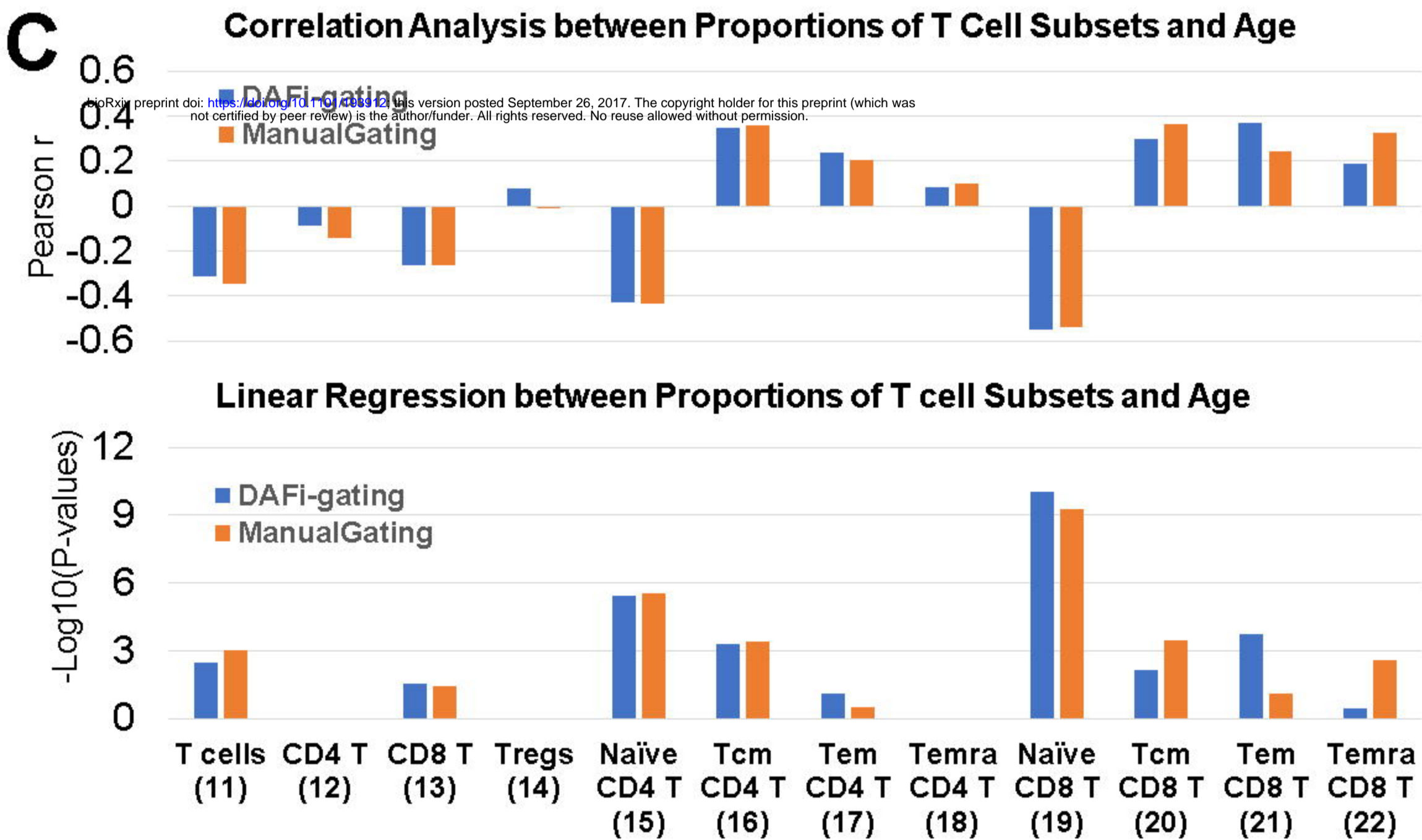
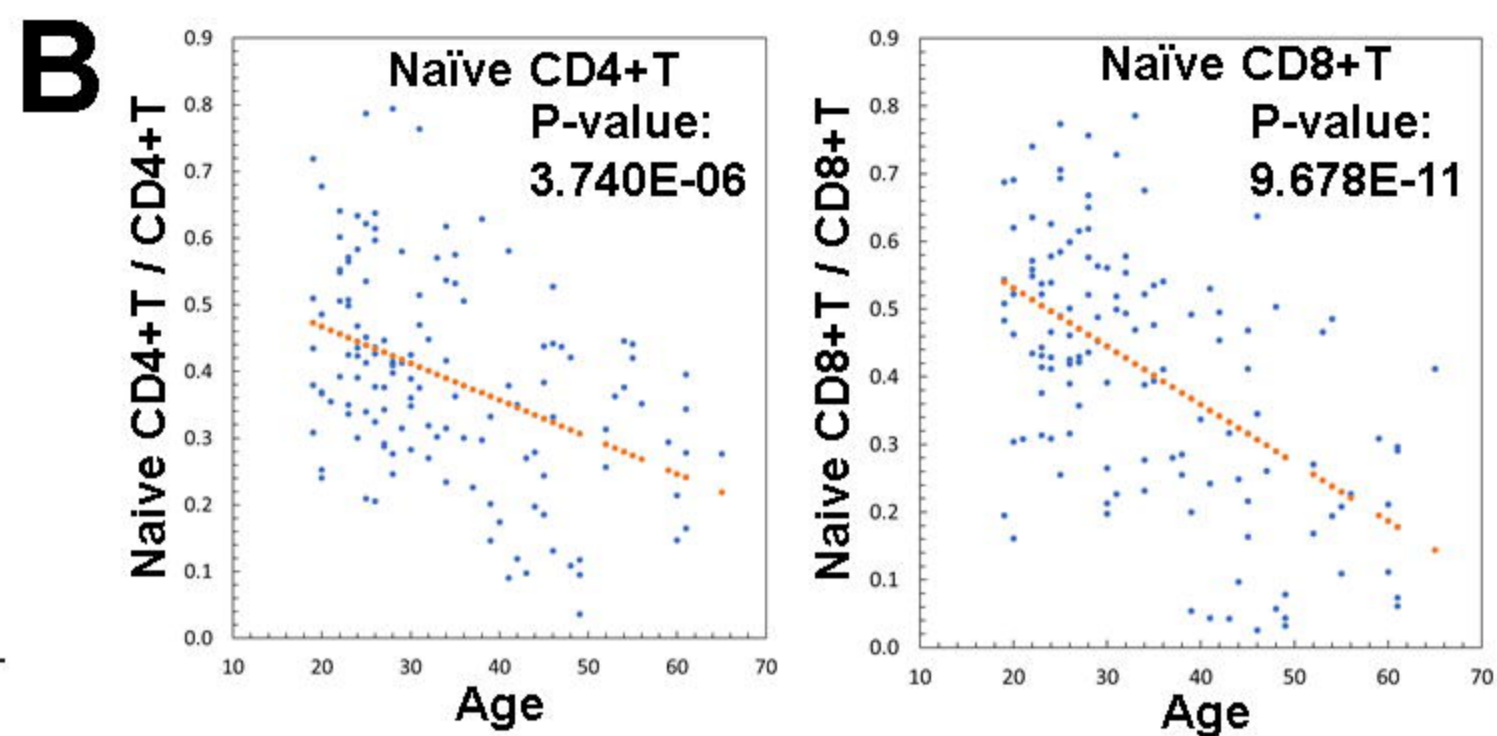
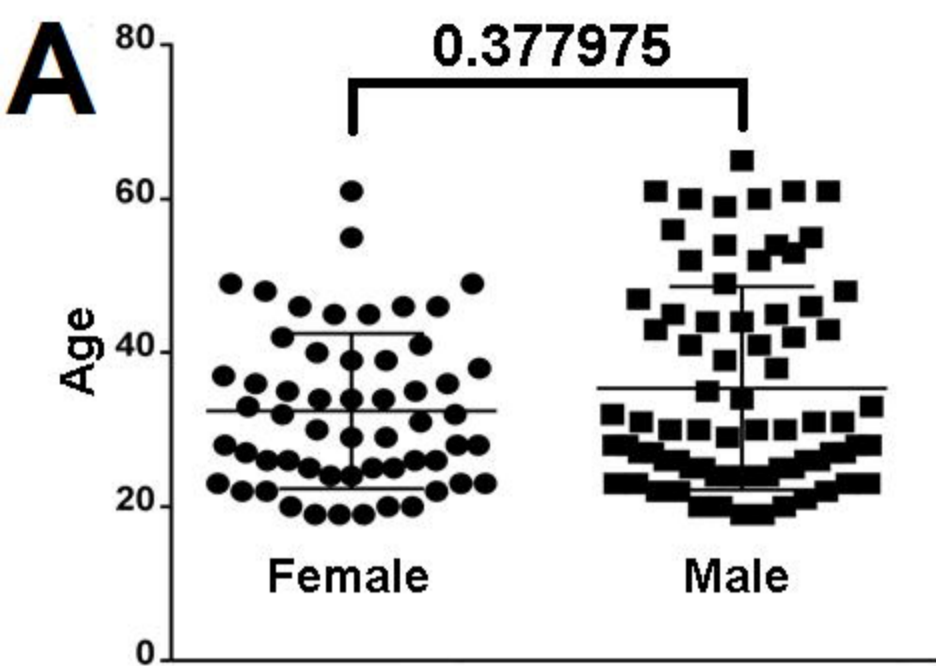
15 Supplementary File 7: Precision, recall, and F1 scores when using different sizes of gating boundaries in
16 DAFi-filtering across the 12 samples in LTBI study for each cell population compared with bisecting
17 analysis.

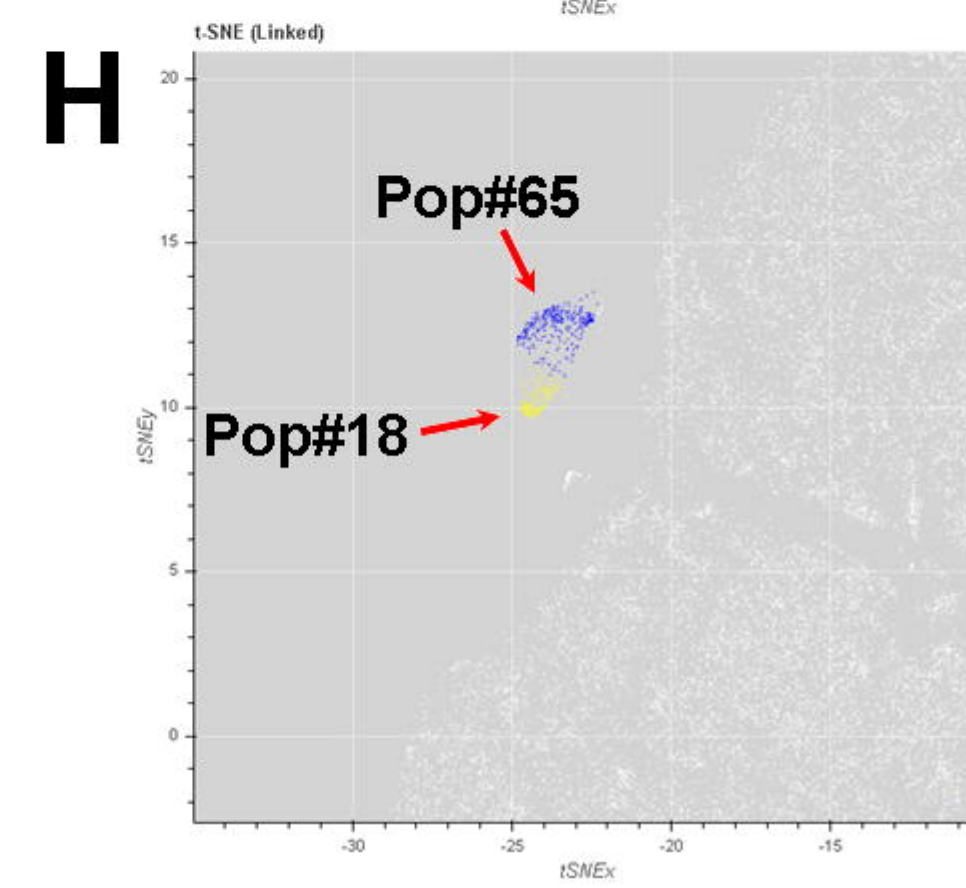
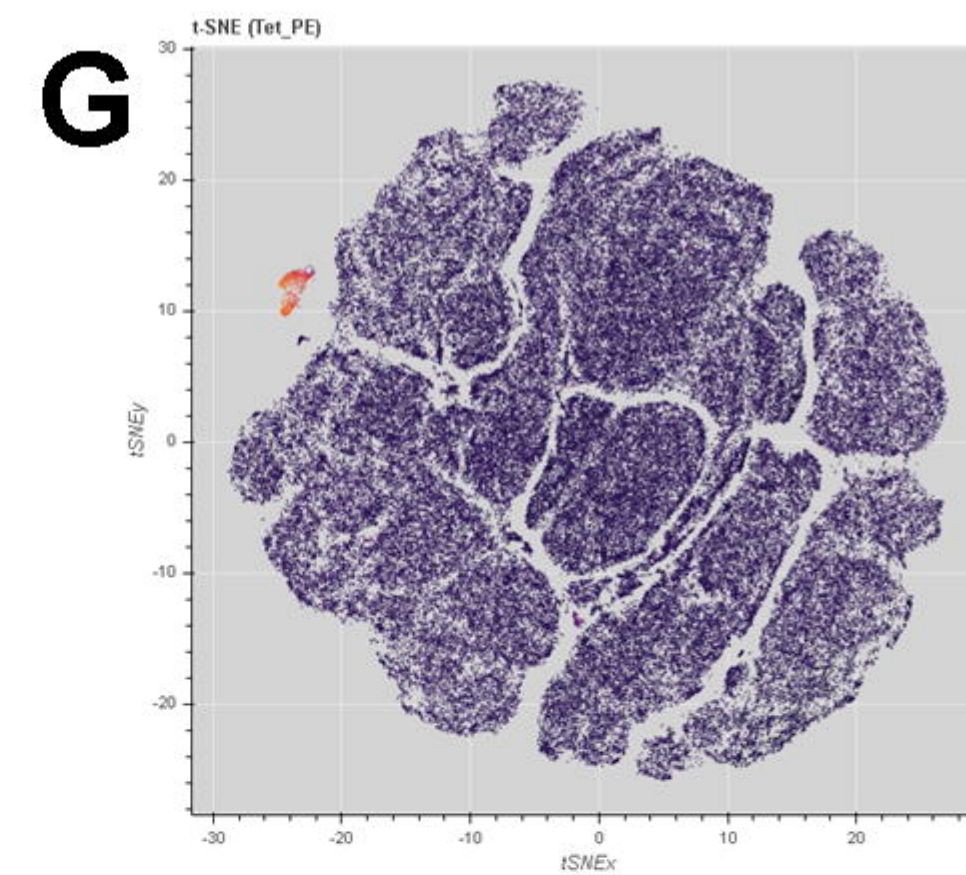
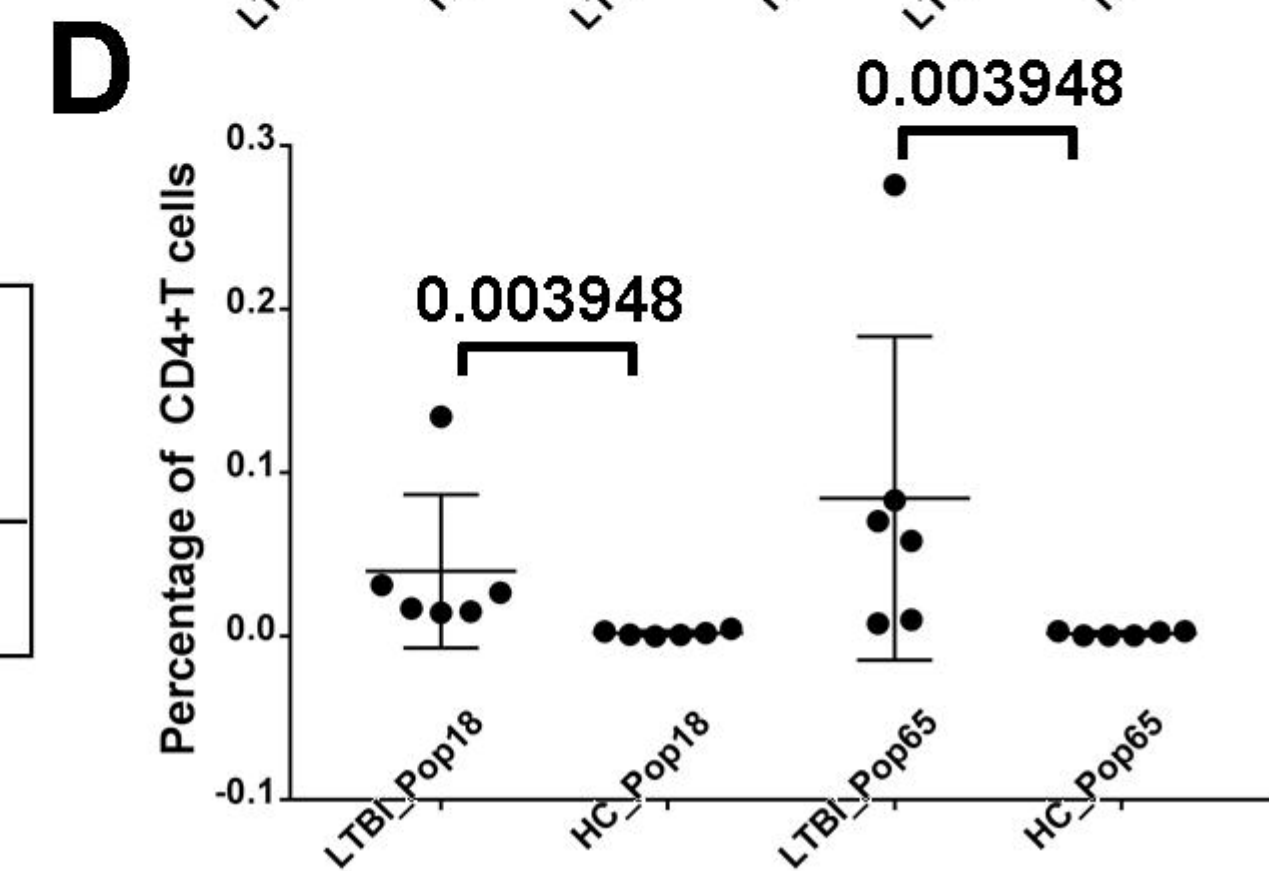
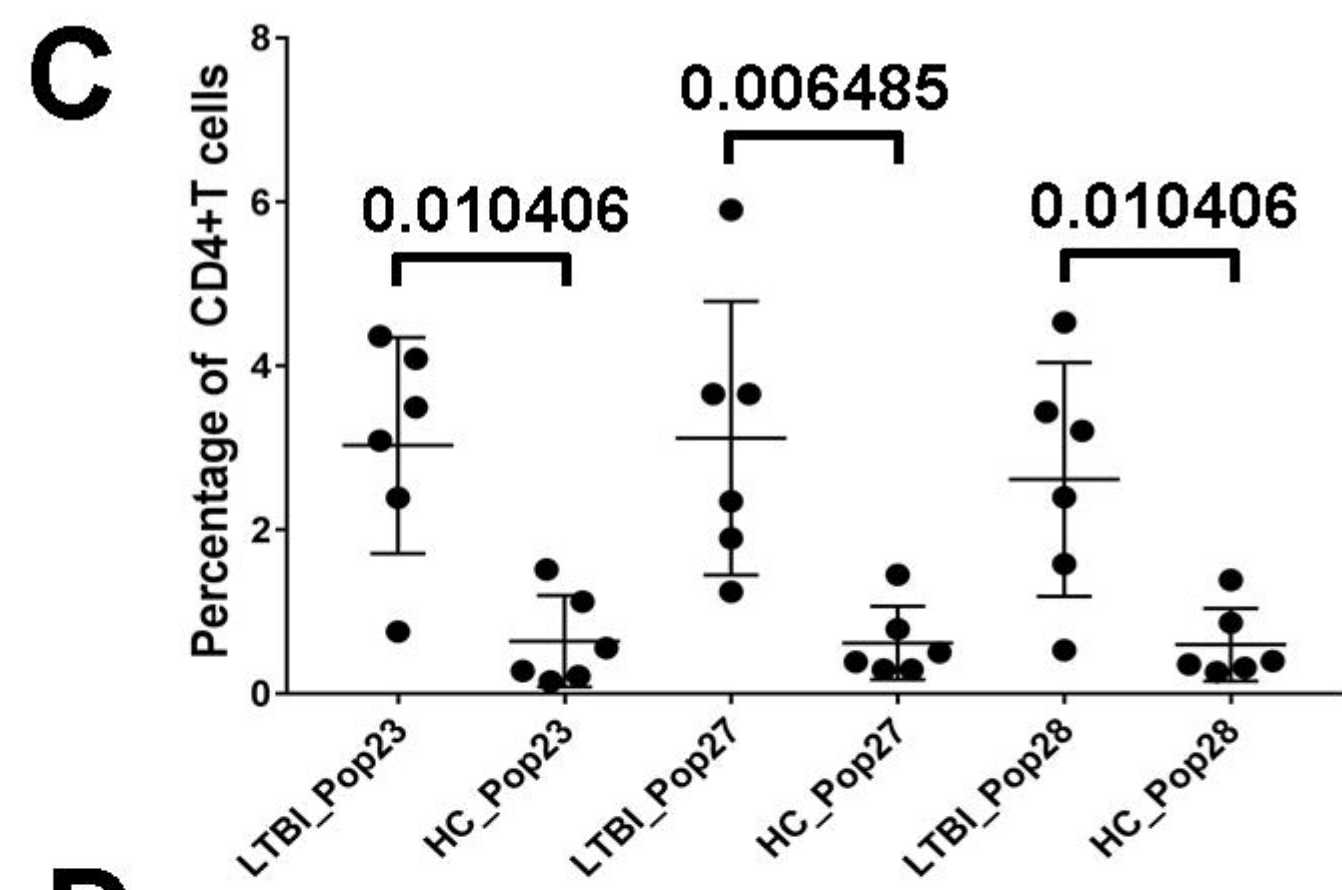
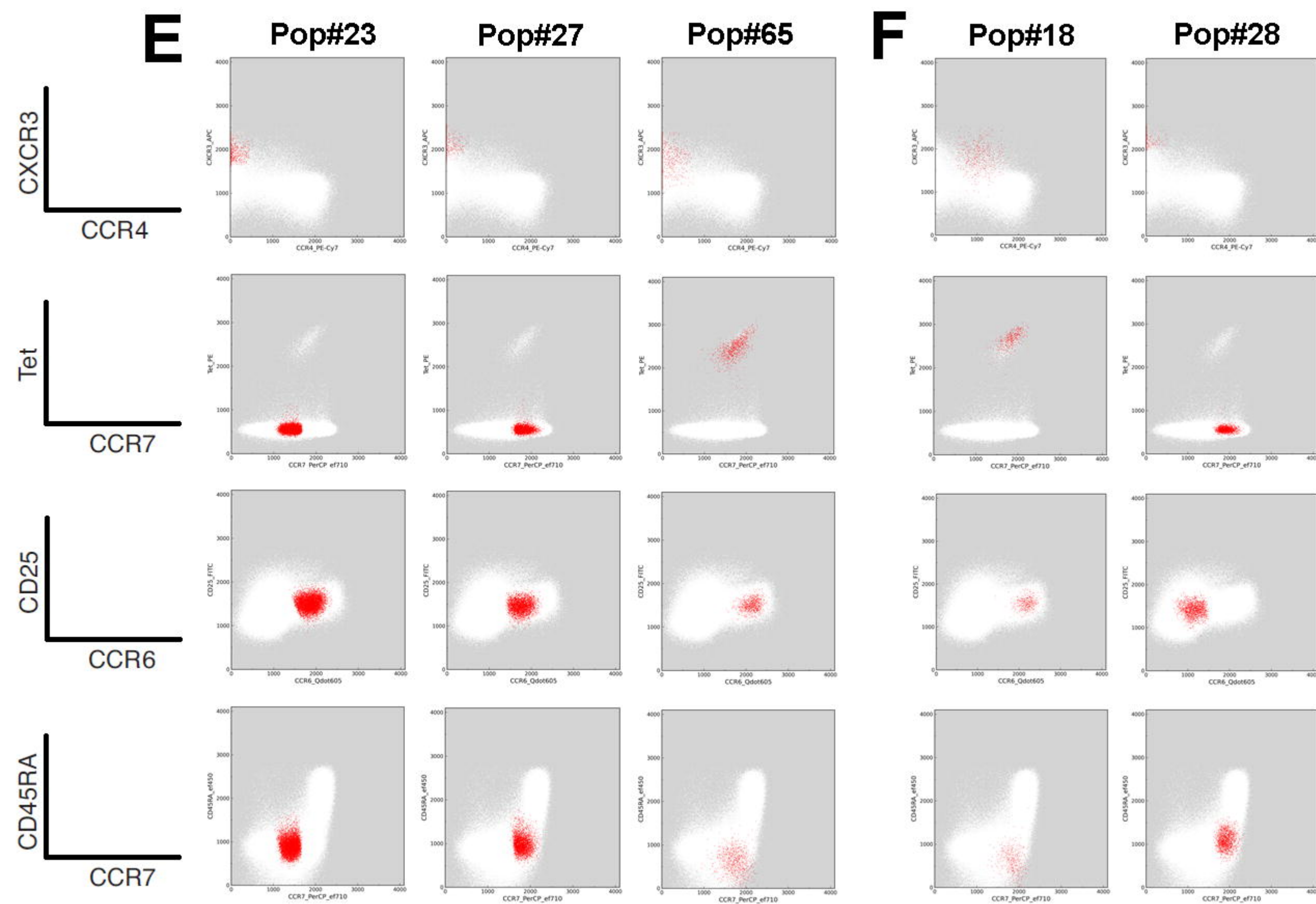
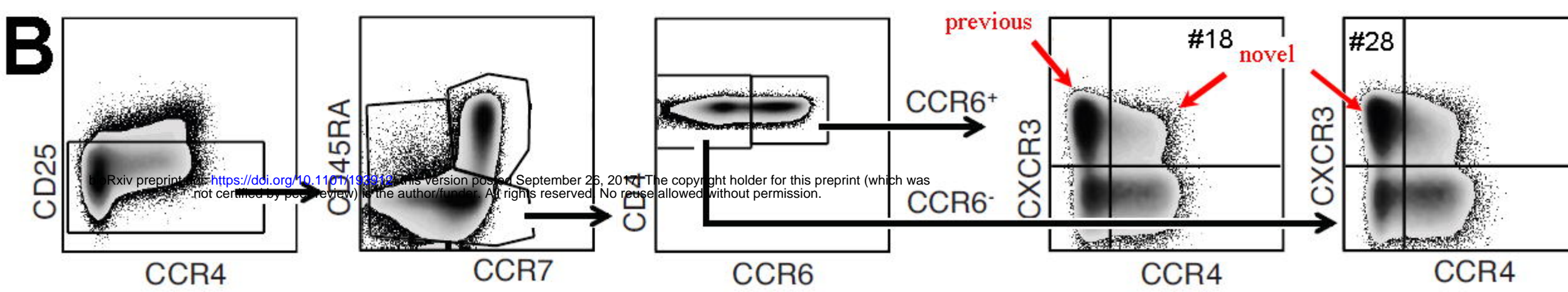
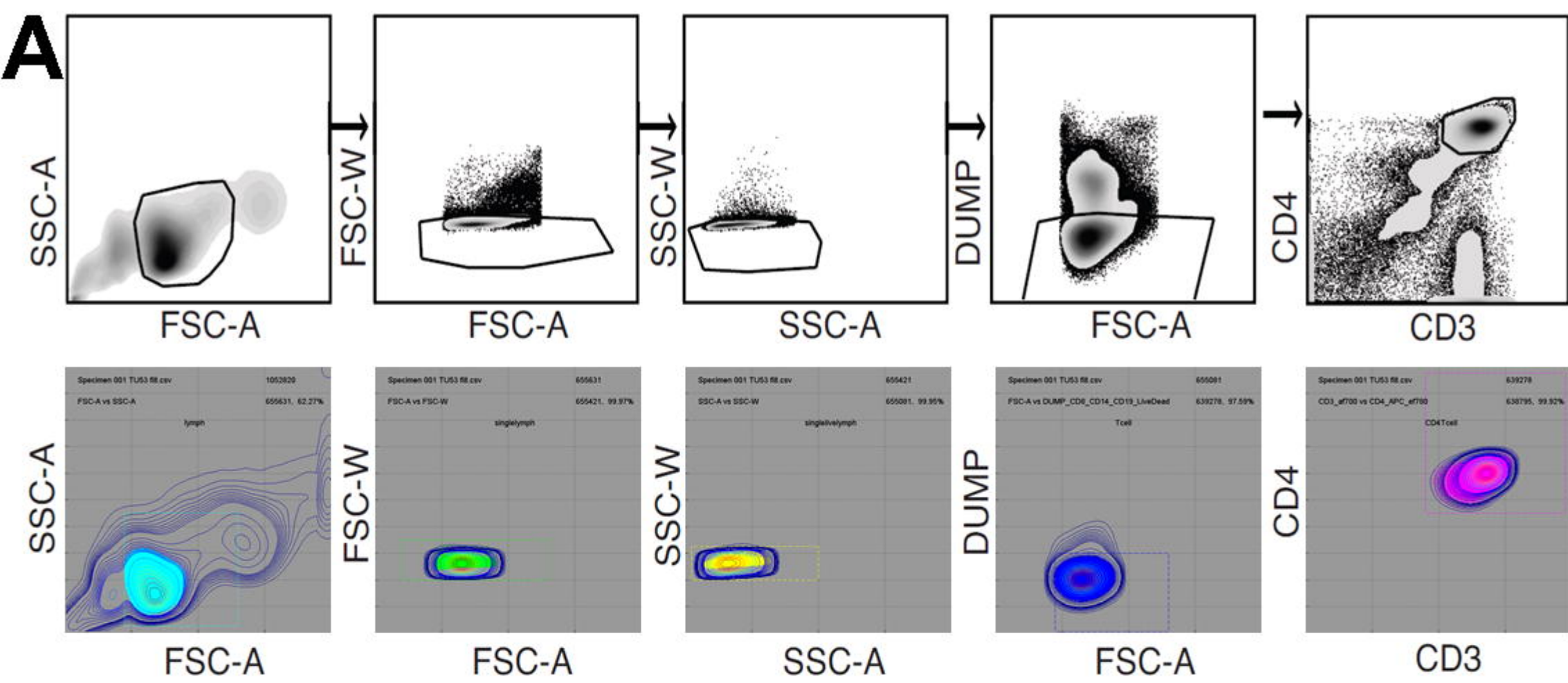
18 Supplementary File 8: Results shown in histograms before and after applying GaussianNorm method to
19 normalize CCR6 and CD45RA channels across the 12 samples in LTBI study.

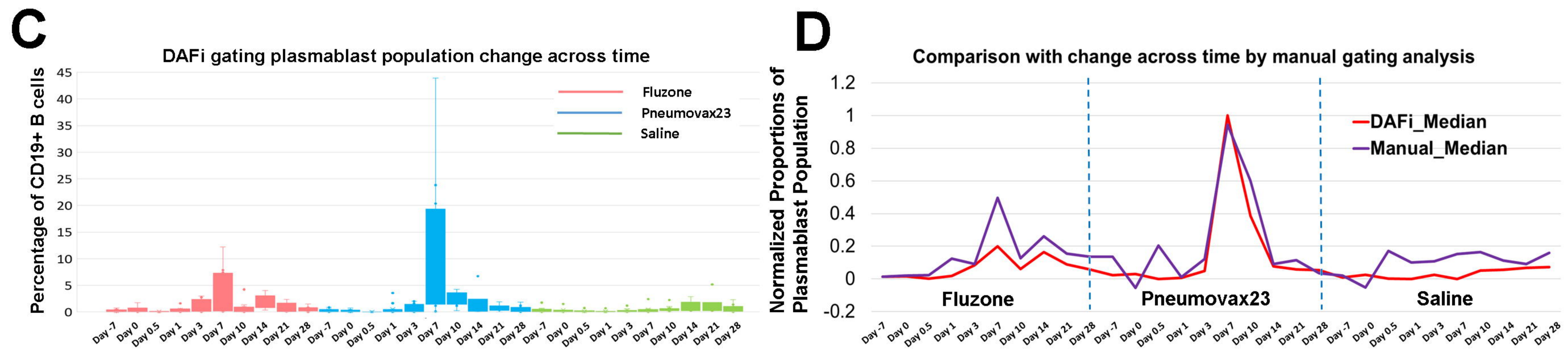
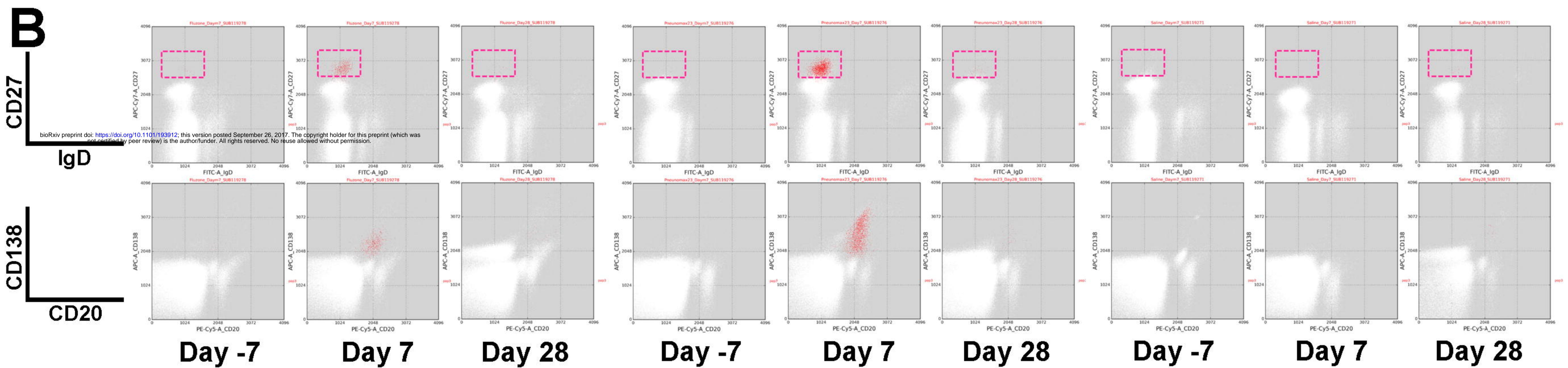
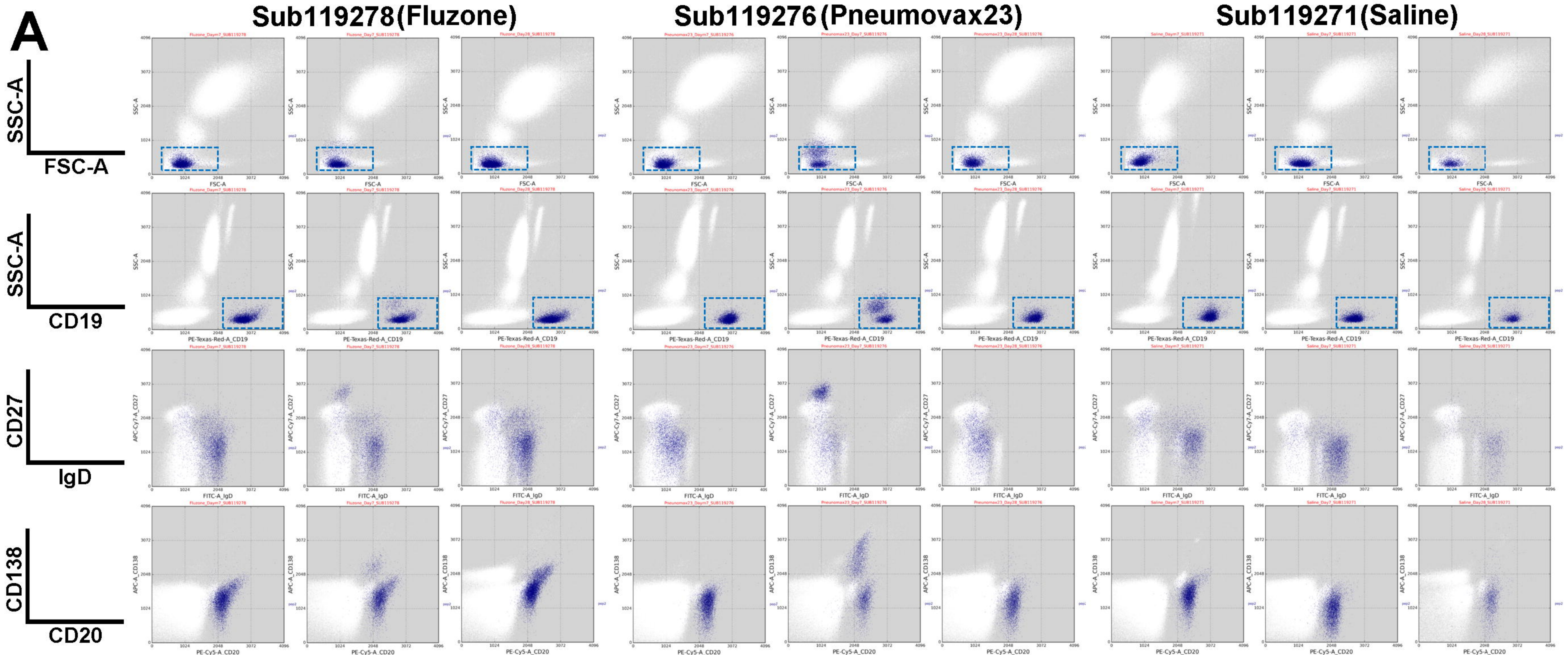
- 1 Supplementary File 9: Comparison of dot plots and F1 score between manual gating analysis in
- 2 FlowCAP-I GvHD dataset and the DAFi result.











PopID	PopName	Xname	Yname	Xmin	Xmax	Ymin	Ymax	ParentID	Mode	RecursiveParent
1	Lymphocyte	FSC-A	SSC-A	30	100	5	70	0	Cluster	Yes
2	Singlet Lymphocytes	FSC-A	FSC-H	200	200	100	200	1	Slope	Yes
3	Singlet Lymphocytes	SSC-A	SSC-H	200	200	110	200	2	Slope	Yes
4	Live CD3T	CD3	LiveDead	100	200	0	100	3	Cluster	Yes
5	CD4T	CD4	CD8	100	200	0	90	4	Cluster	No
6	CD8T	CD4	CD8	0	80	120	200	4	Cluster	No