1 **Title**

2 Scalable Multi-Sample Single-Cell Data Analysis by Partition-Assisted Clustering and Multiple
3 Alignments of Networks

4

5 **Authors**

6 Ye Henry Li[a,*], Dangna Li[b,*], Nikolay Samusik[c], Xiaowei Wang[d], Leying Guan[d], Garry P. Nolan[c], Wing
7 Hung Wong[d,e,1]

8 *Equal Contribution: Y.H.L and D.L.

9 [1]To whom correspondence should be addressed. Email: whwong@stanford.edu

10

11 **Author Affiliations**

12

13 Ye Henry Li

14 [a]Structural Biology Department and Public Policy Program, Stanford University, Stanford, USA.

15

16 Dangna Li

17 [b]Institute for Computational and Mathematical Engineering, Stanford University, Stanford, USA.

18

19 Nikolay Samusik

20 [c]Department of Microbiology and Immunology, Baxter Laboratory, Stanford University, Stanford, USA

21

22 Garry P. Nolan

23 [c]Department of Microbiology and Immunology, Baxter Laboratory, Stanford University, Stanford, USA

24

25 Xiaowei Wang

26 [d]Statistics Department, Stanford University, Stanford, USA.

27

28    Leying Guan

29    [d]Statistics Department, Stanford University, Stanford, USA.

30

31    Wing Hung Wong

32    [d]Statistics Department, Stanford University, Stanford, USA.

33    [e]Department of Biomedical Data Science, Stanford University, Stanford, USA

34

**Contributions**

36    Y.H.L. and W.H.W. conceived the project. Y.H.L., D.L., L.G. and W.H.W. designed the data analysis
37    pipeline; Y.H.L. and D.L. implemented the data analysis pipeline. N.S. generated the hand-gated CyTOF
38    data. Y.H.L., D.L., N.S., and X.W. analyzed the data. Y.H.L., D.L., N.S. and W.H.W. wrote the
39    manuscript and developed the figures. G.P.N and W.H.W. supervised the study.

40

**Conflict of Interest**

42    The authors declare no conflict of interest.

43

44

45

46

47

48

49

50

51

52

53

54

**Abstract**

Mass cytometry (CyTOF) has greatly expanded the capability of cytometry. It is now easy to generate multiple CyTOF samples in a single study, with each sample containing single-cell measurement on 50 markers for more than hundreds of thousands of cells. Current methods do not adequately address the issues concerning combining multiple samples for subpopulation discovery, and these issues can be quickly and dramatically amplified with increasing number of samples. To overcome this limitation, we developed Partition-Assisted Clustering and Multiple Alignments of Networks (PAC-MAN) for the fast automatic identification of cell populations in CyTOF data closely matching that of expert manual-discovery, and for alignments between subpopulations across samples to define dataset-level cellular states. PAC-MAN is computationally efficient, allowing the management of very large CyTOF datasets, which are increasingly common in clinical studies and cancer studies that monitor various tissue samples for each subject.

**Introduction**

Analyses of CyTOF data rely on many of the tools and ideas from flow cytometry (FC) data analysis, as CyTOF datasets are essentially higher dimensional versions of flow cytometry datasets. Currently, the most widely used method in FC is still human hand-gating, as other methods often fail to extract meaningful subpopulations of cells automatically. In hand-gating, we draw polygons or other enclosures around pockets of cell events on a two-dimensional scatterplot to define subpopulations and cellular states that are observed in the data. This process is painfully time-consuming and requires advance knowledge of the marker panel design, the quality of the staining reagents, and, most importantly, *a priori* what cell subpopulations to expect to occur in the data. When presented with a new set of marker panels and biological system, the researcher would find it difficult to delineate the cell events, especially in high-dimensional and multiple-sample datasets.

The inefficient nature of hand-gating in flow cytometry motivated algorithmic development in automatic gating. Perhaps the most popular is flowMeans(1), which is optimized for FC and can learn subpopulations in FC data(2) in an automated manner; however, it has not been successfully applied to CyTOF data analysis. Currently, most data analysis tools created for flow cytometry data analyses are not easily applicable for high-dimensional datasets(3). An exception is SPADE, which was developed and optimized specifically for the analysis of CyTOF datasets(3). flowMeans and SPADE constitute the leading computational methods in cytometry, but as shown later in this work, their performance may become sub-optimal when challenged with large and high-dimensional datasets. There are also other recent clustering-based tools that utilize dimensionality reduction and projections of high-dimensional data, however, these tools do not directly learn the subpopulations for all the cell events, and may be too slow to complete data analysis for an increasing amount of samples.

In this study, we address the data analysis challenges in two major steps. First, we propose the partition-assisted clustering (PAC) approach, which produces a partition of the k-dimensional space (k=number of markers) that captures the essential characteristic of the data distribution. This partitioning methodology is grounded in a strong mathematical framework of partition-based high-dimensional density estimation(4–8). The mathematical framework offers the guarantee that these partitions approximate the

95  underlying empirical data distribution; this step is faster than the recent k-nearest neighbor-based method
96  (9) and is essential to the scalability of our clustering approach to analyze datasets with many samples.
97  The clustering of cells based on recursive partitioning is then refined by a small number of k-means style
98  iterations before a merging step to produce the final clustering.

99   Secondly, the subpopulations learned separately in multiple different but related datasets can be aligned
100  by marker network structures (multiple alignments of networks, or MAN), making it possible to
101  characterize the relationships of subpopulations across different samples automatically. The ability to do
102  so is critical for monitoring changes in a subpopulation across different conditions. Importantly, in every
103  study, batch effect is present; batch effects shift subpopulation signals so that the means can be different
104  from experiment to experiment. PAC-MAN naturally addresses batch effects in finding the alignments of
105  the same or closely related subpopulations from different samples.

106  PAC-MAN finds homogeneous clusters efficiently with all data points in a scalable fashion and enables
107  the matching of these clusters across different samples to discover cluster relationships in the form of
108  clades.

109

110  **Results and Discussion**

111

112  *PAC*

113  PAC has two parts: partitioning and post-processing. In the partitioning part of PAC, the data space is
114  recursively divided into smaller hyper-rectangles based on the number of data points in the locality
115  (Figure 1a). The partitioning is accomplished by either Bayesian Sequential Partition (BSP) with limited
116  look-ahead (Figure 1a and 1b) or Discrepancy Sequential Partition (DSP) (Figure 1a); these are two fast
117  variants of partition-based density estimation methods previously developed by our group (4–8), with
118  DSP being the fastest. BSP and DSP divide the sample space into hyper-rectangles with uniform density
119  value in each of them. The subsetting of cells according to the partitioning provides a principled way of
120  clustering the cells that reflects the characteristics of the underlying distribution. In particular, each
121  significant mode is captured by a number of closely located rectangles with high-density values (Figure
122  1c). Although this method allows a fast and unbiased localization of the high-density regions of the data
123  space, we should not use the hyper-rectangles directly to define the final cluster boundaries for two
124  reasons. First, real clusters are likely to be shaped elliptically, therefore, the data points in the corners of a
125  hyper-rectangle are likely to be incorrectly clustered. Second, a real cluster is often split into more than
126  one closely located high-density rectangles. We designed post-processing steps to overcome these
127  limitations: 1) a small number of k-means iterations is used to round out the corners of the hyper-
128  rectangles, 2) a merging process is implemented to ameliorate the splitting problem, which is inspired by
129  the flowMeans algorithm. The details of post-processing are given in Materials and Methods. The
130  resulting method is named b-PAC or d-PAC depending on whether the partition is produced by BSP or
131  DSP.

132

133    *MAN*

134    An approach to analyze multiple related samples of CyTOF data is to pool all samples into a combined
135    sample before detection of subpopulations. This is a natural approach under the assumptions that there are
136    no significant batch effects or systematic shifts in cell subpopulations across the different samples.
137    However, such assumptions may not hold due to one or more of the following reasons:

138    1)    Dataset size and instruments used. Large number of samples usually means the samples were
139          collected on different days with different experimental preparations. Many steps can introduce
140          significant shifts in measurement levels.
141    2)    Staining reagents. Reagents such as antibodies, purchased from different vendors and batch
142          preparations can affect the overall signal. While saturation of reagents in the protocol could help
143          eliminate the batch effects in the staining procedure, this approach is costly and might not work
144          for all antibodies, especially those with poor specificity.
145    3)    Normalization beads stock. While normalization beads(10) help to control for the signal level,
146          especially within one experiment, the age of the beads stock and their preparation could lead to
147          significant batch effects. In addition, there are different types of normalization beads and
148          normalization calculations.
149    4)    Human work variation. While many researchers are studying the same system (e.g., immune
150          system), different protocols and implementation by different researchers, who sometimes perform
151          experimental steps slightly differently, can lead to batch effects.
152    5)    Subpopulation dynamics. The subpopulation centers can move from sample to sample due to
153          treatments on the cells in treatment-control studies or perturbation studies. General practice is to
154          cluster by phenotypic markers.
155    6)    Sample background. If the data came from different cell lines or individuals in a clinical study,
156          the measurement levels and proportions of cell subpopulations would be expected to change from
157          sample to sample. Without expert scrutiny, it would be difficult to make sense of the data with
158          current data analysis tools.

159    Could we extract shared information that allows us to interpret cross-sample similarities and differences?
160    To ameliorate these difficulties, we have designed an alternative approach that is effective in the presence
161    of substantial systematic between-sample variation. In this approach, each sample is analyzed separately
162    (by PAC) to discover within-sample subpopulations. Over-partitioning in this step is allowed in order not
163    to miss small subpopulations. The subpopulations from all samples are then compared to each other based
164    on a pairwise dissimilarity measure designed to capture the differences in within-sample distributions
165    (among the markers) across two subpopulations. Using this dissimilarity, we perform bottom-up
166    hierarchical clustering of the subpopulations to represent the relationship among the subpopulations. The
167    resulting tree of subpopulations is then used to guide the merging of subpopulations from the same
168    sample, and to establish linkage of related subpopulations from different samples. We note that the design
169    of a dissimilarity measure (Materials and Method) that is not sensitive to systematic sample-to-sample
170    variation is a novel aspect of our approach. The merging of subpopulations from the same sample is also
171    important, as it offers a way to correct any over-partitioning that may have occurred during the initial
172    PAC analysis of each sample. We emphasize that, as with the usage of all statistical methods, the user
173    must utilize samples or datasets that are considered as good as possible; interpretation of the analysis
174    results rely on the researchers to collect data with validated reagents for all samples.

5

175

*Rational initialization for PAC increases clustering effectiveness*

177  Appropriate initialization of clustering is very important for eventually finding the optimal clustering
178  labels; PAC works well because the implicit density estimation procedure yields rational centers to learn
179  the modes of sample subpopulations. When tested on the hand-gated CyTOF data on the bone marrow
180  sample in (14), compared to k-means alone, PAC gives lower total sums of squares and higher F-
181  measures in the subpopulations (Figure 1d and 1e). This process also helps PAC to converge in 50
182  iterations (Figure 1f) in post-processing, whereas k-means performs very poorly even after 5000 iterations
183  (Figure 1g). Through the lens of t-sne plots (Figure 1g), the PAC results are more similar to the hand-
184  gating results, while the k-means, flowMeans, and SPADE clustering results perform poorly. In
185  flowMeans, several large subpopulations are merged. SPADE's separation of points is inconsistent and
186  highly heterogeneous, probably due to its down-sampling nature. On the other hand, by inspection, PAC
187  obtains similar separation for both the major and minor subpopulations as the hand-gating results.

188

189  *PAC is consistently better than flowMeans and SPADE for simulated datasets and hand-gated cytometry*
190  *datasets*

191  In the systematic simulation study, we challenged the methods with different datasets with varying
192  number of dimensions, number of subpopulations, and separation between the subpopulations. The F-
193  measure and p-measures for the PAC methods are consistently equal or higher than that of flowMeans
194  and SPADE (Table 1 and Supporting Figure 2a). In addition, we observe that flowMeans gives
195  inconsistent F-measures for similar datasets (Table 1), which may be due to the convergence of k-means
196  to a local minimum without a rational initialization.

197  Next, we tested the methods based on published hand-gated cytometry datasets to see how similar the
198  estimated subpopulations are to those obtained by human experts. We applied the methods on the
199  hematopoietic stem cell transplant and Normal Donors datasets from the FlowCAP challenges(2) and on
200  the subset of gated mouse bone marrow CyTOF dataset (Dataset 5) recently published(11). The gating
201  strategy of the CyTOF dataset is provided in Supporting Figure 1. The dataset and expert gating strategy
202  are the same as described earlier(12). Note that in the flow cytometry data, the computed F-measures are
203  slightly lower than that reported in FlowCAP; this is due to the difference in the definition of F-measures.
204  Overall, the PAC outperforms flowMeans and SPADE by consistently obtaining higher F-measures
205  (Table 1). In particular, in the CyTOF data example, PAC generated significantly higher F-measures
206  (greater than 0.82) than flowMeans and SPADE (0.59 and 0.53, respectively). In addition, PAC gives
207  higher overall subpopulation-specific purities (Supporting Figure 2b and Supporting Table 1). These
208  results indicate that PAC gives consistently good results for both low and high-dimensional datasets.
209  Furthermore, PAC results match human hand-gating results very well. The consistency between PAC-
210  MAN results and hand-gating results in this large data set confirms the practical utility of the
211  methodology.

212

213  *Separate-then-combine outperforms Pool Approach when Batch Effect is present*

214   It is natural to analyze samples separately then combine the subpopulation features for downstream
215   analysis in the multiple samples setting. However, we need to resolve the batch effects.
216   Two distinct subpopulations could overlap in the combined/pooled sample, such as in the case when the
217   data came from two generations of CyTOF instruments (newer instrument elevates the signals). On the
218   other hand, in cases with changing means, two subpopulations can evolve together such that their means
219   change slightly, but enough to shadow each other when samples are merged prior to clustering.

220   First, we consider the overlapping scenario (Figure 2b). When viewed together in the merged sample, the
221   right subpopulation from sample 1 overlaps with the left subpopulation in sample 2 (Figure 2c). There is
222   no way to use expression level alone to delineate the two overlapping subpopulations (Figure 2d). By
223   learning more subpopulations using PAC, there are some hints that multiple subpopulations are present
224   (Figure 2e). Despite these hints, it would not be possible to say whether the shadowed subpopulations
225   relate in any way to other distinct subpopulations.

226   PAC-MAN resolves the overlapping issue by analyzing the samples separately (Figure 2f). Considering
227   the case in which we do not know *a priori* the number of true subpopulations, we learn three
228   subpopulations per sample. The network structures of the subpopulations discovered are presented in
229   Figure 2g and we see that the third subpopulations from the two samples share the same network
230   structures, while the first subpopulations of the two samples differ by only one edge; these respective
231   networks are clustered together in the dendrogram (Figure 2h, bottom panel). By utilizing the networks,
232   the clades that represent the same and/or similar subpopulations of cells can be established. Clustering by
233   network structures alone resolves the majority of points in the data (Figure 2h, top panel). Furthermore, as
234   discussed next, by incorporating marker levels into the alignment process, all the points can be resolved
235   (Figure 2i).

236   Next we consider the case with dynamic evolution of subpopulations that models the treatment-control
237   and perturbation studies. The interesting information is in tracking how subpopulations change over the
238   course of the experiment. In the simulation, we have generated two subpopulations that nearly converge
239   in mean expression profile over the time course (Figure 3a). The researcher could lose the dynamic
240   information if they were to combine the samples for clustering analysis. As in the previous case, we could
241   use PAC to learn several subpopulations per sample (Figure 3b). Then, with the assumption that there are
242   two evolving clusters from data exploration, we align the subpopulations to construct clades of same
243   and/or similar subpopulations (Figure 3c) based on the network structural information (Supporting Figure
244   3). With network and expression level information in the alignment process, the two subpopulations or
245   clades can be resolved naturally (Figure 3c).

246

247   *Network and expression alignment is better than network or expression alignment alone*

248   With networks in hand, we could further characterize the relationships between subpopulations across
249   samples. However, the alignment process needs to work well for true linkage to be established. We could
250   align by network alone, by expression (or marker) means, or both. Figures 2h, 2i, and 3c present these
251   alternatives in comparison. By using all the subpopulation networks, the results still contain subsets of
252   misplaced cells (Figures 2h top panel and 3c left panel). This is because small clusters of cells have noisy
253   underlying covariance structure; therefore, the networks cannot be accurately inferred. These structural

254    inaccuracies negatively impact the network clustering. The (mean) marker level approach also does not
255    work well (Figure 3c center panel) due to the subpopulation mean shifts across samples. On the other
256    hand, the sequential approach works well (Figures 2i and 3c right panel). In the sequential approach,
257    larger (>1500 in batch effect case; >1000 in dynamic case) subpopulations' networks are utilized for the
258    initial alignment process. Next, the smaller subpopulations, which have noisy covariance, are merged
259    with the closest larger, aligned subpopulations. Thus, more subpopulations could be discovered upstream
260    (in PAC), and the network alignment would work similarly as the smaller subpopulations, which could be
261    fragments of a distribution, do not impact the alignment process (Supporting Figure 4a and b). Moreover,
262    in the network inference step, unimportant edges can negatively impact the alignment process (Supporting
263    Figure 4c) in the network-alone case. Biologically, this means that edges that do not constrain or define
264    the cellular state should not be utilized in the alignment of cellular states. Effectively, the threshold placed
265    on the number of edges in the network inference controls for the importance of the edges. Thus, the
266    combined alignment approach works well and allows moderate over-saturation of cellular states to be
267    discovered in the PAC step so that no advance knowledge of the exact number of subpopulations is
268    necessary.

269

270    *PAC-MAN efficiently outputs meaningful data-level subpopulations for mouse tissue dataset*

271

272    We use the recently published mouse tissue dataset(11) to illustrate the multi-sample data analysis
273    pipeline. The processed dataset contains a total of more than 13 million cell events in 10 different tissue
274    samples, and 39 markers per event (Supporting Table 2). The original research results centered on
275    subpopulations discovered from hand-gating the bone marrow tissue data to find 'landmark'
276    subpopulations; the rest of the data points were clustered to the most similar landmark subpopulations.
277    While this enables the exploration of the overall landscape from the perspective of bone marrow cell
278    types, a significant amount of useful information from the data remains hidden.

279    In contrast, using d-PAC-MAN, the fastest approach by our comparison results, we can perform
280    subpopulation discovery for each sample automatically and then align the subpopulations across samples
281    to establish dataset-level cellular states. On a standard Core i7-44880 3.40GHz PC computer, the single-
282    thread data analysis process with all data points takes about one hour to complete, which is much faster
283    than alternative methods. With multi-threading and parallel processing, the data analysis procedure can be
284    completed very quickly. As mentioned earlier, PAC results for the bone marrow subsetted data from this
285    dataset matches closely to that of the hand-gated results. This accuracy provides confidence for applying
286    PAC to the rest of the dataset.

287    Figure 4 shows the t-sne plots for subpopulation discovered (top panel of each sample) and the
288    representative subpopulation established (bottom panel of each sample) for the entire dataset. In the PAC
289    discovery step, we learn 35 subpopulations per sample without advance knowledge of how many
290    subpopulations are present. This moderate over-partitioning of the data samples leads to a moderate
291    heterogeneity in the t-sne plots. Next, the networks are inferred for the larger subpopulations (with
292    number of cell events greater than 1000), and the networks are aligned for all the tissue samples. We
293    output 80 representative subpopulations or clades for the entire dataset to account for the traditional

294    immunological cellular states and sample-specific cellular states present. Within samples, the
295    subpopulations that cluster together by network structure are aggregated. The smaller subpopulations (not
296    involved in network alignment) are either merged to the closest larger subpopulation or establish their
297    own sample-specific subpopulation by expression alignment; small subpopulations were clamped with
298    larger clades by grouping the subpopulations into 5 clusters per sample based on the means (of marker
299    signal). The representative subpopulations (90 total) follow the approximate distribution of the cell events
300    on the t-sne plots and the aggregating effect cleans up the heterogeneities due to over-partitioning in the
301    PAC step.

302    The cell type clades are the representative subpopulations for the entire dataset, and they could either be
303    present across samples or in one sample alone. Their distribution is visualized by a heatmap (Figure 5).
304    While the bone marrow sample contains many cell types, only a subset of them are directly aligned to cell
305    types in other samples, which means using the bone marrow data as the reference point leaves much
306    information unlocked in the dataset. The cell types in the blood and spleen samples have more alignments
307    with cell types in other samples. The lymph node samples share many clades; the small intestine and
308    colon samples also share many clades, probably due to closeness in biological function. The thymus
309    sample has few clades shared with other samples, which may be due to its functional specificity.

310    PAC-MAN style analysis can be applied to align the tissue subpopulations by their means instead of
311    network similarities (Supporting Figure 5). As done previously, representative clades (88 total) were
312    outputted. The same aggregating effect is observed (Supporting Figure 5a), and this is due to the
313    organization from dataset-level variation in the means. Comparing to the network alignment, the means
314    linkage approach has slightly more subpopulations per sample; the subpopulation proportion heatmap
315    (Supporting Figure 5b) shows more linking. Although the bone marrow sample subpopulations co-occur
316    in the same clades slightly more with other sample subpopulations, this sample does not co-occur with
317    many clades in the dataset. Thus, a PAC-MAN style analysis with means linkage also harvests additional
318    information from the entire dataset.

319    To compare the network and means approaches with PAC-MAN, we study the F-measure and p-measure
320    results with 88 total clades from each approach. The overall F-measure with all cell events is 0.7969 and
321    the overall F-measure with clades assignments of PAC-discovered subpopulations is 0.3143. The two F-
322    measure values suggest that the assignment of PAC-discovered subpopulations is more consistent for
323    larger subpopulations.

324    To illustrate the assignment purities, the p-measures are computed for the following two cases. 1)
325    Network clade assignment is the basis (network-justified), similar to the ground truth in the clustering
326    comparisons previously; or 2) means clade assignment is the basis (means-justified) (Supporting Table 4).
327    P-measure cutoff is set at 0.3 (to remove unreliable comparisons) to obtain purer clade assignments. In the
328    network-justified case, PAC subpopulations with more than 0.3 in p-measure constitute 93.44 % of all
329    cell events. In the means-justified case, PAC subpopulations with more than 0.3 in p-measure constitute
330    92.67 % of all cell events. Furthermore, if the p-measure cutoff were to increase to 0.5, the percentages of
331    cells left for the network-justified and mean-justified cases are 6.25% and 75.16%, respectively. The
332    network-justified case yields drastically lower numbers of cell events in the purer PAC subpopulations
333    because the means approach has more heterogeneity in the linkages (defined as PAC-subpopulation
334    participants in each shared clade with size of at least 2). In fact, the network approach has 100 linkages

9

335    while the means approach has 209 linkages. Therefore, the extra linkages in the means approach would
336    yield greater impurities in the network-justified case. The linkage plot (Supporting Figure 6a) shows that
337    the low linkages occur slightly more frequently for the network approach. One consequence is that the
338    network approach aggregates PAC subpopulations within sample more frequently; for instance, in the
339    thymus sample, the network approach yields 14 clades while the means approach yields 21 clades.

340    After aggregating, the clade sizes (with unique participants per sample) are plotted (Supporting Figure
341    6b). The network approach tends to find fewer linkages, as more clades have sizes of less than 4, while
342    the means approach has more clades than the network approach with clade sizes greater than 4. The
343    network approach is more conservative due to the additional constraints from network structures.
344    Conventionally, in the cytometry field, only the means are considered in the definition of cellular states.
345    Assuming the absence of batch and dynamic effects, the researcher could view the purer shared clade
346    assignments in the network-justified case (general agreement between constrained network approach and
347    means approach) as more reliable candidates of cross-sample relationships to investigate in future
348    experiments (Supporting Figure 6c).

349    Hence, the network alignment approach is in agreement that of the means approach, with network
350    alignment being more stringent in the establishment of linkages. The network PAC-MAN approach
351    defines cellular states with the additional information from network structures, and it has the effect of
352    constraining the number of linkages between samples while finding linkages for subpopulations that are
353    distant in their means.

354

355    *Network hubs provide natural annotations*

356    To further characterize the cell types, we annotate the clades within each sample using the top network
357    hub markers, which constrain the cellular states. The full annotation, along with mean average expression
358    profiles, is presented in Supporting Table 3. The clade information is presented in the ClusterID column.
359    The annotations for cells across different samples but within the same clades share hub markers. For
360    example, in clade 1 for the blood and bone marrow samples, the cells share the hub markers Ly6C and
361    CD11b. In the bone marrow sample, one important set of subpopulations is the hematopoietic stem cell
362    subpopulations.  One such subpopulation is present as clade 18 with the annotation CD34-CD27-cKit-
363    Sca1 and is about 1.87 percent in the bone marrow sample. Clade 18 is only present in the bone marrow
364    sample, indicating that the PAC-MAN pipeline defines this as a sample-specific and coherent
365    subpopulation using dataset-level variation. The thymus contains a large subpopulation (84.07 percent)
366    that is characterized as CD5-CD4-CD43-CD3, suggesting it to be the maturing T-cell subpopulation.

367

**Conclusion**

368

369    We have presented the PAC-MAN data analysis pipeline. This pipeline was designed to remove major
370    roadblocks in the utilization of existing and future CyTOF datasets. First, we established a quick and
371    accurate clustering method that closely matches expert gating results; second, we demonstrated the
372    management of multiple samples by handling mean shifts and batch effects across samples. The
373    alignment allows researchers to find relationships between cells across samples without resorting to

10

374     pooling of all data points. This pipeline can be efficiently utilized to analyze large datasets of high-
375     dimension. PAC-MAN allows the cytometry field to harvest information from the increasing amount of
376     CyTOF data available.

377

378     **Materials and Methods**

379

380     *Partition-assisted clustering has two parts*

381     1) Partitioning: a partition method (BSP(5) or DSP(7)) is used to learn N initial cluster centers from the
382     original data.

383     2) Post-processing: A small number (m) of k-mean iterations is applied to the rectangle-based clusters
384     from the partitioning, where m is a user-specified number. We used m=50 in our examples. After this k-
385     means refinement, we merge the N clusters hierarchically until the desired number of clusters (this
386     number is user-specified) is reached. The merging is based on a given distance metric for clusters. In the
387     current implementation, we use the same distant metric as in flowMeans(1). That is, for two clusters X
388     and Y, their distance $D(X, Y)$ is defined as:

$$D(X, Y) = \min \{(\bar{x} - \bar{y})^T S_x^{-1}(\bar{x} - \bar{y}), (\bar{x} - \bar{y})^T S_y^{-1}(\bar{x} - \bar{y})\}$$

389     where $\bar{x}, \bar{y}$ are the sample mean of cluster X and Y, respectively. $S_x^{-1}$ is the inverse of the sample
390     covariance matrix of cluster X. $S_y^{-1}$ is defined similarly. In each step of the merging process, the two
391     clusters having the smallest pairwise distance will be merged together into one cluster.

392

393     *Partition Methods*

394     There are two partition methods implemented in the comparison study: d-PAC and b-PAC. The results
395     are similar, with d-PAC being the faster algorithm. Figure 1a illustrates this recursive process.

396     d-PAC is based on the discrepancy density estimation (DSP)(7). Discrepancy, which is widely used in the
397     analysis of Quasi-Monte Carlo methods, is a metric for the uniformity of points within a rectangle. DSP
398     partitions the density space recursively until the uniformity of points within each rectangle is higher than
399     some pre-specified threshold. The dimension and the cut point of each partition are chosen to
400     approximately maximize the gap in uniformity of two adjacent rectangles.

401     BSP + LL is an approximation inference algorithm for Bayesian sequential partitioning density estimation
402     (BSP)(5). It borrows ideas from Limited-Look-ahead Optional Pólya Tree (LL-OPT), an approximate
403     inference algorithm for Optional Pólya Tree(8). The original inference algorithm for BSP looks at one
404     level ahead (i.e. looking at the possible cut points one level deeper) when computing the sampling
405     probability for the next partition. It then uses resampling to prune away bad samples. Instead of looking at
406     one level ahead, BSP + LL looks at h levels ahead (h > 1) when computing the sampling probabilities for
407     the next partition and does not do resampling (Figure 1b). In other words, it compensates the loss from

408   not performing resampling with more accurate sampling probabilities. For simplicity, 'BSP + LL' is
409   shortened to 'BSP' in the rest of the article.

410

411   *F-measure*

412   We use the F-measure for comparison of clustering results to ground truth (known in simulated data, or
413   provided by hand-gating in real data). This measure is computed by regarding a clustering result as a
414   series of decisions, one for each pair of data points. A true positive decision assigns two points that are in
415   the same class (i.e. same class according to ground truth) to the same cluster, while a true negative
416   decision assigns two points in different classes to different clusters. The F-measure is defined as the
417   harmonic mean of the precision and recall. Precision P and recall R are defined as:

418   $P = \frac{TP}{TP+FP}$ , $R = \frac{TP}{TP+FN}$, where TP is the total number of true positives, FP is the total number of false
419   positives and FN is the total number of false negatives.

420   F-measure ranges from 0 to 1. The higher the measure, the more similar the estimated cluster result is to
421   the ground truth. This definition of F-measure is different than that of FlowCAP challenge(2). The use of
422   co-assignment of labels in this definition is a more accurate way to compute the true positives and
423   negatives.

424

425   *Purity-measure (p-measure)*

426   Most of the existing measurements for clustering accuracy aim at measuring the overall accuracy of the
427   entire datasets, i.e. comparing with the ground truth over all clusters. However, we are also interested in
428   analyzing how well a clustering result matches the ground truth within a certain class. Specifically,
429   consider a dataset D with K classes: $\{C_1, C_2, ..., C_K\}$ and a given ground truth cluster labels g, we construct
430   an index called the purity measure, or p-measure for short, to measure how well our clustering result
431   matches g for each class $C_i$. This index is computed as follows:

432   1) For each class $C_k$, look for the cluster that has the maximum number of overlapping points with this
433   class, denoted by $L_{i_k}$.

434   2) Define $S_1 = \frac{|C_k \cap L_{i_k}|}{|L_{i_k}|}$, $S_2 = \frac{|C_k \cap L_{i_k}|}{|C_k|}$, where $|\cdot|$ denotes the number of points in a set.

435   3) The final P-index for class $C_k$ is given by: $P = \frac{2S_1S_2}{S_1+S_2}$.

436   If we were to match a big cluster with a small class, even though the overlapping may be large, $S_1$ would
437   still be low since we have divided the score by the size of the cluster in $S_1$. In addition, we are interested
438   in knowing how many points in $C_k$ are clustered together by $L_{i_k}$, which is measured by $S_2$.

439

12

440    *Network construction and comparison*

441    After PAC, the discovered subpopulations typically have enough cells for the estimation of mutual
442    information. This enables the construction of networks as the basis for cell type characterization.
443    Computationally, it is not good to directly use the mutual information networks constructed this way to
444    organize the subpopulations downstream. The distance measure used to characterize the networks could
445    potentially give the same score for different network structures. Thus, it is necessary to threshold the
446    network edges based on the strength of mutual information to filter out the noisy and miscellaneous
447    edges. In this work, these subpopulation-specific networks are constructed using the MRNET network
448    inference algorithm in the Parmigene (13) R package. The algorithm is based on mutual information
449    ranking, and outputs significant edges connecting the markers. The top *d* edges (*d* is set to be 1x the
450    number of markers in all examples) are used to define a network for the subpopulation. This process
451    enables a careful calculation of the distance measure.

452    For each pair of subpopulation networks, we calculate a network distance, which is defined as follows. If
453    $G_1$ and $G_2$ are two networks, let S be the set of shared edges and A be union of the of the edges in the two
454    networks, then we define

455    $$\text{Similarity}(G_1, G_2) = \frac{|S|}{|A|},\text{ where } |\cdot| \text{ denotes the size of a set.}$$

456    This is known as the Jaccard coefficient of the two graphs. The Jaccard distance, or 1- Jaccard coefficient,
457    is then obtained. This is a representation of the dissimilarity between each pair of networks; the Jaccard
458    dissimilarity is the measure used for the downstream hierarchical clustering.

459

460    *Cross-sample linkage of subpopulations*

461    We perform agglomerative clustering of the pool of subpopulations from all samples. This clustering
462    procedure greedily links networks that are the closest in Jaccard dissimilarity, and yields a dendrogram
463    describing the distance relationship between all the subpopulations. We cut the dendrogram to obtain the
464    *k* clades of subpopulations. Subpopulations from the same sample and falling into the same clade are then
465    merged into a single subpopulation (Figure 2a). This merging step has the effect of consolidating the
466    over-partitioning in the PAC step. No merging is performed for subpopulations from different samples
467    sharing the same clade. In this way, we obtain *k* clades of subpopulations, with each clade containing no
468    more than one subpopulation from each sample. We regard the subpopulations within each clade as being
469    linked across samples.

470    In the above computation, only subpopulations with enough cells to define a stable covariance are used
471    for network alignment via the Jaccard distance; the rest of the cell events from very small subpopulations
472    are then merged with the closet clade by marker profile via distance of mean marker signals. If the small
473    subpopulations are distant from the defined clades, then a new sample-specific clade is created for these
474    small subpopulations.

475

476

13

477    *Annotation of Subpopulations*

478    To annotate the cellular states, we first apply PAC-MAN to learn the dataset-level subpopulation/clade
479    labels. Next, these labels are used to learn the representative/clade networks. The top hubs (i.e. the most
480    connected nodes) in these networks are used for annotation. This approach has biological significance in
481    that important markers in a cellular state are often central to the underlying marker network, which is
482    analogous to important genes in gene regulatory networks; these important markers have many
483    connections with other markers. If the connections were broken, the cell would be perturbed and
484    potentially driven to other states.

485

486    *Running Published Methods*

487    To run t-SNE (14) a dimensionality reduction visualization tool, we utilized the scripts published here
488    (https://lvdmaaten.github.io/tsne/). Default settings were used.

489    To run SPADE, we first converted the simulated data to fcs format using Broad Institute's free
490    CSVtoFCS online tool in GenePattern(15)  (http://www.broadinstitute.org/cancer/software/genepattern#).

491    Next,    we    carried    out    the    tests    using    the    SPADE    package    in    Bioconductor    R(16)
492    (https://bioconductor.org/packages/release/bioc/html/spade.html).

493    To run flowMeans, we carried out the tests using the flowMeans package in Bioconductor R(1)
494    (https://bioconductor.org/packages/release/bioc/html/flowMeans.html).

495    In the comparisons, we selected only cases that work for all methods to make the tests as fair as possible.

496    To calculate the mutual information of the subpopulations, we use the infotheo R package (https://cran.r-
497    project.org/web/packages/infotheo/index.html).

498    To run network inference, we use the mrnet algorithm in the parmigne R package (13). (https://cran.r-
499    project.org/web/packages/parmigene/index.html).

500

501    *Code Availability*

502    The PAC R package can be accessed at:

503    https://cran.r-project.org/web/packages/PAC/index.html

504

505    *Simulated Data for Clustering Analysis*

506    To compare the clustering methods, we generated simulated data from Gaussian Mixture Model varying
507    dimension, the number of mixture components, mean, and covariance. The dimensions range from 5 to
508    39. The number of mixture components is varied along each dimension. The mean of each component

14

509   was generated uniformly from a d-dimensional hypercube; we generated datasets using hypercube of
510   different sizes, but kept all the other attributes the same. The covariance matrices were generated as $AA^T$,
511   where $A$ is a random matrix whose elements were independently drawn from the standard normal
512   distribution. The sizes of the simulated dataset range from 100k to 200k.

513   The simulated data are provided as (Datasets 1-6). Datasets 1-4 are for the PAC part. Dataset 1 contains
514   data with 5 dimensions; Dataset 2 contains data with 10 dimensions; Datasets 3a and 3b contain data with
515   20 dimensions; and Datasets 4a and 4b contain data with 35 dimensions. The ground truth labels are
516   included as separate sheets in each dataset.

517   When applying flowMeans, SPADE, and the PAC to the data, we preset the desired number of
518   subpopulations to that in the data to allow for direct comparisons.

519

520   *Gated Flow Cytometry Data*

521   Two data files were downloaded from the FlowCAP challenges(2). One data file is from the
522   Hematopoietic stem cell transplant (HSCT) data set; it has 9,936 cell events with 6 markers, and human
523   gating found 5 subpopulations. Another data file is from the Normal Donors (ND) data set; it has 60,418
524   cell events with 12 markers, and human gating found 8 subpopulations. The files are the first ('001') of
525   each dataset. These data files were all 1) compensated, meaning that the spectral overlap is accounted for,
526   2) transformed into linear space, and 3) pre-gated to remove irrelevant events. We used the data files
527   without any further transformation and filtering. When applying flowMeans, SPADE, and the PAC to the
528   data, we preset the desired number of subpopulations to that in the data to allow for direct comparisons.

529

530   *Gated Mass Cytometry Data*

531   Human gated mass cytometry data was obtained by gating for the conventional immunology cell types
532   using the mouse bone marrow data recently published(11). The expert gating strategy is provided as
533   Supporting Figure 1. The gated sample subset contains 64,639 cell events with 39 markers and 24
534   subpopulations and it is provided as Dataset 7.

535   To test the performance of different analysis methods, the data was first transformed using the asinh(x/5)
536   function, which is the transformation used prior to hand-gating analysis; For SPADE analysis, we utilize
537   the asinh(x/5) option in the SPADE commands. The post-clustering results from flowMeans, SPADE, b-
538   PAC, and d-PAC were then subsetted using the indexes of gated cell events. These subsetted results are
539   compared to the hand-gated results.

540

541   *Simulated Data for MAN Analysis*

542   To test the linking of subpopulations, we generated simulated data from multivariate Gaussian with preset
543   signal levels and randomly generated positive definite covariance matrices. There are two cases, batch
544   effect and dynamic. Each simulated sample file has five dimensions, with two of these varying in levels;

545    these are the dimensions that are visualized. Dataset 5 contains the data for general batch effects case and
546    Dataset 6 contains the data for dynamic effects case. The ground truth labels are included as separate
547    sheets in each dataset.

548

549    *General batch scenario*. Sample 1 represents data from an old instrument (instrument 1) while sample 2
550    represents data from a new instrument (instrument 2). There are two subpopulations per sample. These
551    two subpopulations are the same, but their mean marker levels shifted higher up in sample 2 due to higher
552    sensitivity of instrument 2 (Figure 2b). The subpopulations have different underlying relationships
553    between the markers. In this simulated experiment, five markers were measured. Out of the five markers,
554    two markers show significant shift, and we focus on these two dimensions by 2-dimensional scatterplots.
555    In Figure 2b, the left subpopulation in sample 1 is the same as the left subpopulation in sample 2; the
556    same with the right subpopulation. The same subpopulations were generated from multivariate Gaussian
557    distributions with changing means with fixed covariance structure.

558    *Dynamic scenario*. Dynamic scenario models the treatment-control and perturbation studies. In the
559    simulation, we have generated two subpopulations that nearly converge over the time course (Figure 3a).
560    The researcher could lose the dynamic information if they were to combine the samples for clustering
561    analysis. The related subpopulations were generated from multivariate Gaussian distributions with
562    changing means with fixed covariance structure.

563

564    *Raw CyTOF Data Processing*

565    The researcher preprocesses the data to 1) normalize the values to normalization bead signals, 2) de-
566    barcode the samples if multiple barcoded samples were stained and ran together, and 3) pre-gate to
567    remove irrelevant cells and debris to clean up the data(10,17). Gene expressions look like log-normal
568    distributions(18); given the lognormal nature of the values, the hyperbolic arcsine transform is applied to
569    the data matrix to bring the measured marker levels (estimation of expression values) close to normality,
570    while preserving all data points. Often, researchers use the asinh(x/5) transformation, and we use the same
571    transformation for the CyTOF datasets analyzed in this study.

572

573    *Mouse Tissue Data*

574    In the Spitzer et al., 2015 dataset(11), three mouse strains were grown, and cells were collected from
575    different tissues: thymus, spleen, small intestine, mesenteric lymph node, lung, liver, inguinal lymph
576    node, colon, bone marrow, and blood. In each experiment, 39 expression markers were monitored. The
577    authors used the C57BL6 mouse strain as the reference(11); the data was downloaded from Cytobank,
578    and we performed our analysis on the reference strain.

579    First, all individual samples were filtered by taking the top 95% of cells based on DNA content and then
580    the top 95% of cells based on cisplatin: DNA content allows the extraction of good-quality cells and
581    cisplatin level (low) allows the extraction of live cells. Overall, the top 90% of cell events were extracted.

16

582 The filtered samples were then transformed by the hyperbolic arcsine (x/5) function, and merged as a
583 single file, which contains 13,236,927 cell events and 39 markers per event (Supporting Table 2).

584 Using PAC-MAN, we obtained 35 subpopulations in each sample then 80 clades for the entire dataset.
585 The 80 clades account for the traditional immune subpopulations and sample-specific subpopulations.
586 Small subpopulations not used in alignment are later merged into the closest clades; this is done by
587 performing hierarchical clustering with the marker signals to obtain 5 "expression" subclades per sample.
588 Subsequently, any clade with less than 100 cells is discarded. Subpopulation proportion heatmap was
589 plotted to visualize the subpopulation-specificities and relationships across the samples. Finally,
590 annotation was performed using the hub markers of each representative subpopulation in each sample.

591

## Acknowledgements

597

## References

599

600  1. Aghaeepour N, Nikolic R, Hoos HH, Brinkman RR. Rapid cell population identification in flow
601     cytometry data. Cytometry A. 2011 Jan 1;79A(1):6–13.

602  2. Aghaeepour N, Finak G, Consortium TF, Consortium TD, Hoos H, Mosmann TR, et al. Critical
603     assessment of automated flow cytometry data analysis techniques. Nat Methods. 2013
604     Mar;10(3):228–38.

605  3. Qiu P, Simonds EF, Bendall SC, Gibbs Jr KD, Bruggner RV, Linderman MD, et al. Extracting a
606     cellular hierarchy from high-dimensional cytometry data with SPADE. Nat Biotechnol. 2011
607     Oct;29(10):886–91.

608  4. Wong WH, Ma L. Optional Pólya tree and Bayesian inference. Ann Stat. 2010 Jun;38(3):1433–59.

609  5. Lu L, Jiang H, Wong WH. Multivariate Density Estimation by Bayesian Sequential Partitioning. J
610     Am Stat Assoc. 2013 Dec 1;108(504):1402–10.

611  6. Yang K, Wong WH. Discovering and Visualizing Hierarchy in the Data. ArXiv14034370 Stat
612     [Internet]. 2014 Mar 18 [cited 2015 Nov 27]; Available from: http://arxiv.org/abs/1403.4370

613  7. Yang K, Wong WH. Density Estimation via Adaptive Partition and Discrepancy Control.
614     ArXiv14041425 Stat [Internet]. 2014 Apr 4 [cited 2015 Nov 27]; Available from:
615     http://arxiv.org/abs/1404.1425

616    8.    Jiang H, Mu JC, Yang K, Du C, Lu L, Wong WH. Computational Aspects of Optional Pólya Tree. J
617          Comput Graph Stat. 2015 Feb 13;0(ja):00–00.

618    9.    Samusik N, Good Z, Spitzer MH, Davis KL, Nolan GP. Automated mapping of phenotype space
619          with single-cell data. Nat Methods. 2016 Jun;13(6):493–6.

620    10.   Finck R, Simonds EF, Jager A, Krishnaswamy S, Sachs K, Fantl W, et al. Normalization of mass
621          cytometry data with bead standards. Cytometry A. 2013 May 1;83A(5):483–94.

622    11.   Spitzer MH, Gherardini PF, Fragiadakis GK, Bhattacharya N, Yuan RT, Hotson AN, et al. An
623          interactive reference framework for modeling a dynamic immune system. Science. 2015 Jul
624          10;349(6244):1259425.

625    12.   Samusik N, Good Z, Spitzer MH, Davis KL, Nolan GP. Automated mapping of phenotype space
626          with single-cell data. Nat Methods. 2016 Jun;13(6):493–6.

627    13.   Sales G, Romualdi C. parmigene—a parallel R package for mutual information estimation and gene
628          network reconstruction. Bioinformatics. 2011 Jul 1;27(13):1876–7.

629    14.   Maaten L van der, Hinton G. Visualizing Data using t-SNE. J Mach Learn Res. 2008;9(Nov):2579–
630          605.

631    15.   Reich M, Liefeld T, Gould J, Lerner J, Tamayo P, Mesirov JP. GenePattern 2.0. Nat Genet. 2006
632          May;38(5):500–1.

633    16.   Linderman MD, Bjornson Z, Simonds EF, Qiu P, Bruggner RV, Sheode K, et al. CytoSPADE:
634          high-performance analysis and visualization of high-dimensional cytometry data. Bioinformatics.
635          2012 Sep 15;28(18):2400–1.

636    17.   Zunder ER, Finck R, Behbehani GK, Amir ED, Krishnaswamy S, Gonzalez VD, et al. Palladium-
637          based mass tag cell barcoding with a doublet-filtering scheme and single-cell deconvolution
638          algorithm. Nat Protoc. 2015 Feb;10(2):316–33.

639    18.   Bengtsson M, Ståhlberg A, Rorsman P, Kubista M. Gene expression profiling in single cells from
640          the pancreatic islets of Langerhans reveals lognormal distribution of mRNA levels. Genome Res.
641          2005 Oct 1;15(10):1388–92.

642

643

644 **Figure and Table Legends**

645

646 **Figure 1: PAC utilizes rational initialization for fast and accurate clustering convergence**

647 (a) Partition-based methods estimate data density by cutting the data space into smaller rectangles.
648 Bayesian Sequential Partition (BSP) divides the data space via binary partition in the middle of the
649 bounded region, while that of Discrepancy Sequential Partition (DSP) occur at the location that balances
650 the data point uniformly on both sides of the cut. The numbers denote sequential order of partitions. Since
651 DSP adapts to the data points, it converges on the estimated density faster than BSP. (b) In the (one-step)
652 look-ahead of version of partition, the algorithm cuts the data space for all potential cuts plus one step
653 more (steps 2 and 3), and it finds the optimal future version (after step 3), which determines the actual cut
654 (step 2). (c) The partitioning of simulated data space containing five subpopulations; the hyper-rectangles
655 surround high-density areas, approximating the underlying distribution. (d-e) The rational initialization
656 step helps PAC to outperform random initialization. The handgated CyTOF data was used. In this case,
657 the overall sum of squares error is lower and the F-measure is higher for PAC. (f) The convergence of
658 PAC toward the hand-gated results, or ground truth, is fast. It takes less than 50 downstream post-
659 processing kmeans iterations for the PAC to achieve a significantly higher F-measure than the alternative
660 methods. In contrast, flowMeans convergence is poor. (g) Visualization of clustering results with t-sne
661 plot. The t-sne plots contain 10,000 cell events of the handgated CyTOF data with different set of labels
662 drawn. Note that the colors are informative only within each panel. These labels are from kmeans,
663 SPADE, flowMeans, b-PAC, and d-PAC. The subpopulation numbers for all methods were set to be the
664 same as that of handgated results.

665

666 **Figure 2: Overlapping batch effects can be resolved by PAC-MAN**

667 (a) Schematic of MAN. Consider a deck of networks (in analogy to cards), with each "suit" representing a
668 sample and each "rank" representing a unique network structure. The networks are aligned by similarity
669 and organized on a dendrogram. The tree is cut (red line) at the user-specified level to output the desired $k$
670 clades. Within each clade, the network structures are similar or the same. If the same sample has multiple
671 networks in the same clade, then these networks are merged (black box around same cards). (b) Simulated
672 data samples with two of the same subpopulations. The means shifted due to measurement batch effect.
673 (c) When the samples are combined, as in the case of analyzing all samples together, two different
674 subpopulations overlap. (d) The overlapped subpopulations cannot be distinguished by clustering. (e)
675 PAC could be used to discover more subpopulations, however, the hints of the present of another
676 subpopulation do not help to resolve the batch effect. (f) PAC was used to discover several
677 subpopulations per sample without advanced knowledge of the exact number of subpopulations. (g) The
678 networks of the subpopulations discovered in (f). Networks can be grouped by similarities to organize the
679 subpopulations across samples; the alignment is based on Jaccard dissimilarity network structure
680 characterization matrix; dendrogram of the hierarchical clustering results. (h) Resolution of batch effect
681 by networks of all subpopulations discovered. (i) Resolution of batch effect first by gene networks of
682 larger subpopulations and then by merging smaller subpopulations into the aligned clades.

683

**Figure 3: Dynamic information can be extracted by PAC-MAN**

684

685    (a) Ground truth of simulated samples. Two subpopulations, in blue color, almost converge in time by
686    mean shifts. (b) PAC discovers several subpopulations per sample without advanced knowledge of the
687    number of subpopulations present. (c) Comparison of PAC-MAN results between representative clades
688    (number of clades set to 2). Using gene networks and expression information alone do not resolve the
689    dynamic information. On the other hand, dynamic information is resolved first by gene networks of larger
690    subpopulations and then by merging smaller subpopulations into the aligned clades.

691

**Figure 4: Mouse tissue data analysis results visualized by t-sne plots.**

692

693    Each t-sne plot was generated using 10,000 randomly drawn cell events from each mouse tissue sample.
694    The results from PAC (top panel) and MAN (bottom panel) steps are presented as a pair. Initial PAC
695    discovery was set to 35 subpopulations without advanced knowledge of the number of subpopulations in
696    each sample. In MAN, 80 network clades were outputted, and the cellular states are defined by gene
697    expression, network structure, and dataset-level variation. This composite definition naturally aggregates
698    the initial 35 subpopulations to yield smaller number of subpopulations in less variable samples.

699

**Figure 5: Clade proportions and annotation**

700

701    Heatmap of clade proportions across the samples. Sample-specific clades have a value of 1, while shared
702    clades have proportions spread across different samples. Physiologically similar samples share more
703    clades.

704

**Supporting Figure 1: Gating strategy of CyTOF data for methods comparison**

705

706    Biaxial gating hierarchy for the mouse bone marrow CyTOF dataset. Gating strategy that was used to find
707    24 reference populations in the mouse bone marrow CyTOF data. Pre-gating step involved removal of
708    doublets, dead cells, erythrocytes and neutrophils. Non-neutrophils population was either subject to
709    cluster analysis by computational tools or subsequent gating. Dotted boxes represent 24 terminal gates
710    that were selected as reference populations for the comparison analysis.

711

**Supporting Figure 2: Subpopulation purity of simulated and real CyTOF data**

712

713    (a) Subpopulation-specific purity plot of 35-dimensional simulated data with 10 subpopulations. The blue
714    points denote the differences between the p-measures of the partition-based method (either d-PAC or b-
715    PAC) and flowMeans, while the red points denote the p-measure differences between the partition
716    methods and SPADE. The horizontal line at 0 means no difference between the methods. Most of the blue

717 and red points are above 0, indicating that the PAC generates purer subpopulations compared to the
718 ground truth. The two subplots are very similar, which means that d-PAC and b-PAC give very similar p-
719 measures. More precisely, the sum of differences between d-PAC and flowMeans and d-PAC and
720 SPADE are 0.85 and 1.09, respectively; and the overall difference between b-PAC and flowMeans and b-
721 PAC and SPADE are 0.84 and 1.08, respectively.

722 (b) Subpopulation-specific purity plot of the hand-gated CyTOF data. The same convention is used as in
723 (Supporting Figure 2a). Again, more blue and red points are above 0, indicating that the partition-based
724 methods generate purer subpopulations compared to the ground truth. There is a cluster of points below 0
725 occurring in the middle of the plot, suggesting that flowMeans and SPADE capture the mid-size
726 subpopulations more similar to hand-gating than the partition-based methods. More specifically,
727 flowMeans does better (p-measure difference of 0.1 or better; difference of less 0.1 is considered
728 practically no difference) with finding subpopulations of GMP, CD8 T cells, MEP, CD4 T cells
729 (compared to d-PAC), and Plasma cells, while SPADE does better with CD19+IgM- B cells, NK cells
730 (compared to d-PAC), CD8 T cells, NKT cells, Basophils, Short-Term HSC, and Plasma cells. However,
731 overall, PAC has a much better performance, as the absolute sum of points above 0 is higher than that of
732 points below 0. More precisely, the sum of differences between d-PAC and flowMeans and d-PAC and
733 SPADE are 1.21 and 1.45, respectively; and the overall difference between b-PAC and flowMeans and b-
734 PAC and SPADE are 2.06 and 2.31, respectively. The difference table is provided in Supporting Table 1.

735

736 **Supporting Figure 3: Gene Networks inferred from subpopulations in the dynamic example**
737 **simulated dataset**

738 Figure 3 introduced the dynamic example in which five samples each having 2 true subpopulations
739 capture the almost-convergence of means. Here the underlying gene network structures for the PAC
740 discovered subpopulations (three per sample) are presented.

741

742 **Supporting Figure 4: Comparison between aligning cross-sample subpopulations by gene network,**
743 **expression profile, or both**

744 (a) PAC can be used to discover more subpopulations, with the effect of more partitions from the true
745 clusters. (b) When over-partitioning is present, gene network or gene expression profile alone cannot
746 resolve the dynamic (or batch) effects due to noisy covariance for small fragments of distributions.
747 However, first aligning the larger subpopulations with more stable covariance, and thus network
748 structures, and then merge in the smaller subpopulations by expression profile resolves the effects. (c) If
749 more irrelevant edges were introduced, network alignment would fail due to the negative impact of the
750 miscellaneous edges; however, eliminating small subpopulations from the alignment step alleviates the
751 increased edge count problem.

752

753 **Supporting Figure 5: PAC-MAN style linkage by means**

754    (a) t-sne plots of mouse tissue samples colored by representative subpopulations labels from linkage by
755    means. (b) Subpopulation proportion heatmap of clades of samples from linkage by means.

756

757    **Supporting Figure 6: Comparison between network and means PAC-MAN**

758    (a) PAC-discovered subpopulations are aggregated by MAN into clades; the number of PAC
759    subpopulations/clades for the network and means PAC-MAN approaches are plotted. (b) After
760    aggregating shared clades within samples, the number of shared clades for the entire dataset is plotted for
761    the two PAC-MAN approaches. c) Using the network approach results as basis, the clades with strong
762    agreement (high p-measures) with the means PAC-MAN approach are given. The shared clades (present
763    in more than one sample) are reliable candidates for future experiment to find cross-sample relationships.

764

765    **Table 1: F-measure Comparisons of Methods on Simulated and Hand-gated Cytometry Datasets.**

766    F-measure is calculated using the original hand-gate labels and the estimated labels generated by each
767    analysis method. The true-positives are found if the methods assign the same labels to points belonging to
768    the same subpopulation in the hand-gated data. The more true-positives found, the higher the F-measure,
769    which ranges from 0 to 1, with 1 being the highest. Partition-based methods perform consistently well on
770    data ranging from 5 to 39 dimensions. In the simulations, d-PAC and b-PAC perform just as well or better
771    than flowMeans and SPADE. flowMeans gives drastically different F-measures for the cases
772    20_10_40_100k and 20_20_40_100k :  0.25386 vs. 0.92518; this large difference is likely due to the
773    random initiation of cluster centers. In the hand-gated datasets, SPADE has the worst performance.
774    Ultimately, the performance of flowMeans and SPADE deteriorate for the 39-dimensional real CyTOF
775    data, while d-PAC and b-PAC perform consistently well.

776    *Simulated data have the following convention: a_b_c_d, where a denotes the number of
777    dimensions/markers, b denotes the number of subpopulations, c denotes the size of the hypercube for data
778    generation, and d denotes the number of cells.

779    **from rounding up, not originally 1.00

780

781    **Supporting Table 1: Purity (p) Measure Differences in CyTOF Comparison**

782    p-measure differences in gated CyTOF data analysis comparison. The differences are shown for all the
783    annotated cell subpopulations, which are ordered by their sizes. Overall, the PAC methods give more
784    positive p-measures.

785

786    **Supporting Table 2: Sample Sizes in Mouse Tissue CyTOF Dataset**

787    The numbers of cells in the samples of Spitzer et al., 2015 CyTOF dataset. The data is from the C57BL6
788    mouse strain and a total of ten tissue samples are present. The raw column shows the number of cells

22

789 prior to filtering by DNA and cisplatin values. The final cell counts are shown in the filtered file (3$^{rd}$)
790 column.

791

792 **Supporting Table 3: PAC-MAN Subpopulation Characterization Output for Mouse Tissue CyTOF**
793 **Dataset**

794 The full set of annotated results, along with mean expressions, subpopulation proportion and counts, are
795 reported.

796

797 **Supporting Table 4: Network-justified and means-justified p-measures for Alignments of PAC-**
798 **discovered Subpopulations**

799 The PAC-discovered subpopulations were mapped as clades in both the network and means PAC-MAN
800 approaches. The p-measures were calculated for the cases 1) network approach mapping as the basis and
801 2) means approach mapping as the basis. The comparison is the same in principle to the comparison of
802 labels for clustering methods. The results are ordered by p-measures.

803

# Figure 1

a

b

c



BSP

DSP

Partitioned Data Space

d

Comparision of Different Initializations: Total Sum of Squares



e

Comparision of Different Initializations: F-measure



f

Convergence of Clustering Results to Handgate Labels



g



Handgate

kmeans 5000 iterations

SPADE

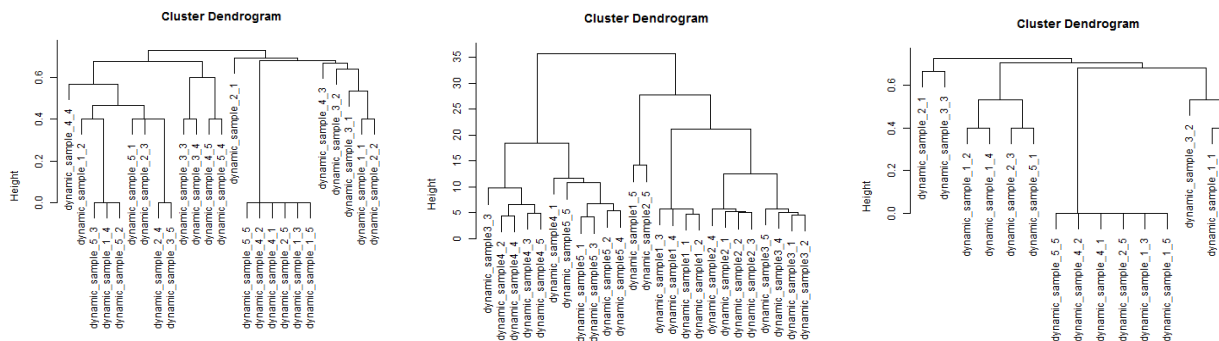flowMeans

b-PAC

d-PAC

# Figure 2

# Figure 3

a

b



c

# Figure 4

# Figure 5

# Supporting Figure 1

# Supporting Figure 2

a

b

# Supporting Figure 3

# Supporting Figure 4

# Supporting Figure 5

a

b

# Supporting Figure 6

## a

Clade Assignments — Frequency vs PAC Subpopulations/Clades (Network, Means)

## b



Frequency of Clade Sizes — Frequency vs Clade Size (Network, Means)

## c

| Clade | Number of PAC Subpop | p-measure | Dataset-level Proportion | Number of Samples Present |
|---|---|---|---|---|
| clade3 | 1 | 1 | 0.016726488 | 4 |
| clade7 | 5 | 1 | 0.826687342 | 5 |
| clade9 | 2 | 1 | 0.076407462 | 2 |
| clade15 | 2 | 1 | 0.031185486 | 2 |
| clade16 | 5 | 1 | 1.283666519 | 4 |
| clade19 | 2 | 1 | 0.277964818 | 1 |
| clade22 | 1 | 1 | 0.063776132 | 1 |
| esc_2_1 | 1 | 1 | 0.000339958 | 1 |
| esc_2_2 | 1 | 1 | 0.000377731 | 1 |
| clade32 | 1 | 1 | 0.010417826 | 1 |
| clade33 | 3 | 1 | 0.216689266 | 2 |
| clade34 | 1 | 1 | 0.065861208 | 1 |
| clade39 | 3 | 1 | 0.027634813 | 1 |
| esc_3_1 | 1 | 1 | 0.001442933 | 1 |
| clade40 | 1 | 1 | 0.007660388 | 1 |
| clade43 | 3 | 1 | 0.034947688 | 1 |
| clade45 | 2 | 1 | 0.043590178 | 2 |
| clade46 | 3 | 1 | 0.023426888 | 1 |
| clade47 | 2 | 1 | 0.291140081 | 2 |
| esc_5_1 | 1 | 1 | 0.004026614 | 1 |
| clade57 | 1 | 1 | 0.043907472 | 1 |
| clade60 | 1 | 1 | 0.015026146 | 1 |
| esc_8_1 | 1 | 1 | 0.000695025 | 1 |
| clade73 | 3 | 1 | 0.614742379 | 1 |
| clade74 | 1 | 1 | 0.090496835 | 1 |
| clade76 | 1 | 1 | 0.008151439 | 1 |
| esc_10_1 | 2 | 1 | 0.002266387 | 1 |
| clade80 | 1 | 1 | 0.013688978 | 1 |
| clade30 | 4 | 0.75 | 0.084717548 | 3 |
| clade11 | 15 | 0.73333 | 0.645202621 | 8 |
| clade8 | 3 | 0.66667 | 0.117451732 | 3 |
| clade28 | 3 | 0.66667 | 0.065249283 | 2 |
| clade38 | 3 | 0.66667 | 0.052648171 | 1 |
| clade67 | 3 | 0.66667 | 0.038407706 | 2 |
| clade1 | 10 | 0.5 | 2.813553327 | 6 |
| clade10 | 8 | 0.5 | 0.297402864 | 4 |
| clade13 | 8 | 0.5 | 0.138785989 | 4 |
| clade17 | 2 | 0.5 | 0.110365495 | 1 |
| clade18 | 3 | 0.5 | 0.276642759 | 1 |
| clade24 | 2 | 0.5 | 0.22283873 | 2 |
| clade41 | 2 | 0.5 | 0.015328331 | 1 |
| clade54 | 2 | 0.5 | 0.022920728 | 1 |
| clade55 | 4 | 0.5 | 0.81990329 | 1 |
| clade62 | 2 | 0.5 | 0.030180721 | 1 |
| clade6 | 26 | 0.42308 | 24.41300764 | 9 |
| clade5 | 29 | 0.41379 | 14.76349458 | 10 |
| clade59 | 5 | 0.4 | 0.030074956 | 1 |
| clade4 | 11 | 0.36364 | 10.7864008 | 4 |
| clade14 | 11 | 0.36364 | 0.371619485 | 4 |
| clade2 | 28 | 0.35714 | 29.79685542 | 8 |
| clade27 | 6 | 0.33333 | 1.545260467 | 3 |
| clade37 | 6 | 0.33333 | 0.567828167 | 2 |
| clade51 | 3 | 0.33333 | 0.174889534 | 1 |