

## **AUTOMATED OPTIMAL PARAMETERS FOR T-DISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING IMPROVE VISUALIZATION AND ALLOW ANALYSIS OF LARGE DATASETS**

ANNA C. BELKINA<sup>1,2,\*</sup>, CHRISTOPHER O. CICCOLELLA<sup>4</sup>, RINA ANNO<sup>5</sup>, RICHARD HALPERT<sup>6</sup>, JOSEF SPIDLEN<sup>6</sup>, JENNIFER E. SNYDER-CAPPIONE<sup>2,3</sup>

<sup>1</sup>Department of Pathology, <sup>2</sup>Flow Cytometry Core Facility and <sup>3</sup>Department of Microbiology, Boston University School of Medicine, Boston, MA; <sup>4</sup>Omiq, Inc, Santa Clara, CA, <sup>5</sup>Department of Mathematics, Kansas State University, Manhattan, KS, <sup>6</sup>FlowJo, LLC, Ashland, OR  
\*CORRESPONDING AUTHOR: ANNA C. BELKINA, M.D., PH.D. BELKINA@BU.EDU

Accurate and comprehensive extraction of information from high-dimensional single cell datasets necessitates faithful visualizations to assess biological populations. A state-of-the-art algorithm for non-linear dimension reduction, t-SNE, requires multiple heuristics and fails to produce clear representations of datasets when millions of cells are projected. We developed opt-SNE, an automated toolkit for optimal t-SNE parameter selection that utilizes Kullback-Liebler divergence evaluation in real time to tailor the early exaggeration and overall number of gradient descent iterations in a dataset-specific manner. The precise calibration of early exaggeration together with opt-SNE adjustment of gradient descent learning rate dramatically improves computation time and enables high-quality visualization of large cytometry and transcriptomics datasets, overcoming limitations of analysis tools with hard-coded parameters that often produce poorly resolved or misleading maps of fluorescent and mass cytometry data. In summary, opt-SNE enables optimal data resolution in t-SNE space and more precise data interpretation.

Keywords: dimensionality reduction, data visualization, t-SNE, mass cytometry, flow cytometry, scRNA-seq, cytometry analysis, machine learning, viSNE

---

Visual exploration of high-dimensional data is imperative for the comprehensive analysis of single cell datasets. Fluorescence, mass and sequencing-based cytometric data analysis requires tools that are able to reveal the combinations of proteomic and/or transcriptomic markers that define complex and diverse cell phenotypes in a mixed population. While traditional biaxial data presentation via expert-driven gating is still the standard analysis method for cytometry data to date, with the advent of the modern multi-parameter era an analysis tool that can accurately and comprehensively visualize multi-dimensional data is direly needed to relieve the current cytometry data-processing bottleneck.

To date, multiple dimensionality reduction techniques have been applied to cytometry data with variable success. Linear methods, such as principal component analysis (PCA) which generates a low-dimensional representation of data with a linear mapping matrix, are mostly unsuitable for cytometry data visualization as such techniques cannot faithfully present the non-linear relationships in the data. t-Distributed Stochastic Neighbor Embedding (t-SNE) is a state-of-the-art dimensionality reduction algorithm for non-linear data representation that creates a low-dimensional distribution, or a ‘map’, of high-dimensional data<sup>1,2</sup>. Conspicuous groupings of datapoints, or ‘islands’, correspond to observations that are similar in the original high-dimensional

space and help to visualize the general structure and heterogeneity of a dataset. t-SNE was developed as a machine learning technique for a broad range of data types and has been adopted<sup>2</sup> for single-cell applications. When t-SNE embeds single cell data, the islands represent cells with similar phenotypes, as defined by a cytometric or genomic signature, thereby allowing to reveal biological data structure and to surface important differences between samples and/or subject groups<sup>3</sup>.

In addition, t-SNE maps are used to categorize single cell data into relevant biological populations for downstream quantification, achievable through expert-guided filtering (‘gating’)<sup>4</sup> or unsupervised clustering of the map<sup>5-7</sup>. Cytometry clustering algorithms that directly interrogate high-dimensional data, such as FlowSOM<sup>8</sup> and PhenoGraph<sup>9</sup>, are often used in conjunction with t-SNE maps to present annotated clusters to the viewer.

A limitation of t-SNE in its current form is its inability to scale to datasets with large numbers of observations<sup>7,10</sup>. This restrains t-SNE’s utility for cytometry datasets that often include millions of observations (‘events’) routinely collected for phenotypic analysis. Unlike PCA, t-SNE learns the embedding non-parametrically, and hence new pieces of data cannot validly be added to an existing analysis, necessitating the whole dataset to be analyzed within one computation. When the full dataset is comprised of multiple samples, each representing a subject in a large cohort or an independent experimental condition, retaining statistically significant representation of small subpopulations in each sample requires inflating the dataset size<sup>11</sup>. However, even within a single measurement, subsampling the data risks preventing rare subsets from being identified. These limitations cannot be overcome via application of the currently understood best practices for t-SNE use. Not only are large datasets computationally expensive to analyze, but also the resulting t-SNE maps provide poor visualization and incomprehensive representation of high-dimensional data. As a result, upon finding initial poor t-SNE visualizations from large datasets, researchers often resort to either subsampling their data to the very limit of detection of rare populations<sup>12</sup> or to exporting specific populations from their dataset, thus compromising the ‘unbiased’ data analysis approach<sup>13,14</sup>.

Although t-SNE has been widely adopted by the scientific community, to our knowledge no rigorous theoretical or empirical testing of t-SNE for cytometry applications has been performed. In 2013, Amir et al first reported the use of t-SNE (or viSNE, as it was renamed<sup>3</sup>) on mass cytometry data; since then, t-SNE has been implemented in the majority of commercial and open-source platforms for cytometry analysis (FlowJo, Cytobank, FCSExpress, cytofkit, etc). In most implementations, few or no adjustments were made to the Barnes-Hut t-SNE algorithm for the requirements of cytometry datasets; the default and hard-coded parameter settings that were originally tested and optimized with non-cytometry datasets like CIFAR (image dataset) or MNIST (handwritten digits) are retained in these cytometry programs. Developing rigorous methodology to release the full potential of t-SNE for single cell data comprehension is the primary motivation for this work.

In this study, we first assessed the behavior of t-SNE computation with routinely used settings that match common best practices, and then iteratively modified parameters of embedding to identify conditions that ensure optimal visualization. As a result of this work, we propose a new method to automatically find optimal t-SNE parameters via fine-tuning of the early exaggeration stage of t-SNE embedding in real time. We call our approach *opt-SNE*, for *optimal* t-SNE. We find that our adjustments can tremendously shorten the number of iterations required to obtain

visualizations of large cytometry datasets with superior quality. Our approach also eliminates the need for trial-and-error runs intended to empirically find the optimal selection of t-SNE parameters, potentially saving hundreds of hours of computation time per research project.

## Materials and methods

**Datasets.** All datasets used in the study are summarized in Table 1.

**Table 1. Datasets used in this manuscript**

Dataset	Data type	Details	References
Mass41parameter	Mass cytometry	41 parameter dataset (14 parameters used for embedding) of 1 million datapoints concatenated from 5 samples of human bone marrow cells	<sup>15</sup>
Flow18parameter	Flow cytometry	18 parameter dataset (11 parameters used for embedding) of 1 million datapoints concatenated from 2 samples of human PBMC	<sup>16</sup>
Flow20M	Flow cytometry	18 parameter dataset (15 parameters used for embedding) of 20 million datapoints concatenated from 27 samples of human PBMC	Panel based on <sup>16</sup>
10X Genomics	scRNA-seq	Single cell gene expression data of E18 mouse brain pre-processed into 20 PCA projections used for t-SNE embedding	<a href="https://support.10xgenomics.com/single-cell-gene-expression/datasets">https://support.10xgenomics.com/single-cell-gene-expression/datasets</a> and <sup>17</sup>
van Unen et al	Mass cytometry	18 parameter dataset (14 parameters used for embedding) of 5.22 million datapoints concatenated from 104 samples of human peripheral blood mononuclear cells (PBMC) and gut biopsy cells	<sup>7</sup>

**Primary samples.** Flow18parameter and Flow20M data were collected as described <sup>18</sup> with minor modifications of the flow cytometry reagent list. Study protocols were approved by the institutional review boards and all subjects provided written informed consent. The committee that approved the research protocol was the Boston University Institutional Review Board, IRB# H-33095.

**Data pre-processing.** Singlet events from several data recordings were digitally concatenated and a randomly subsampled file of 1,000,000 mass <sup>15</sup> or flow cytometry <sup>16</sup> events was created and used for analyses of the mass41parameter and flow18parameter datasets. All observations from 27 recordings of flow cytometry data were concatenated to generate the flow20M dataset. All observations from 104 recordings of mass cytometry data <sup>7</sup> were concatenated to generate the van Unen et al dataset.

All flow cytometry data were compensated with acquisition-defined compensation matrices. Prior to t-SNE analysis, all cytometry data were transformed using asinh (all mass cytometry data and flow14parameter data) or biexponential (flow20M) transformation. Light scatter parameters were log-transformed.

**Data analysis.** A desktop C++ Barnes-Hut implementation of t-SNE for Mac OS was used for most t-SNE analyses<sup>2</sup>. All datasets were embedded in 2D space. Original code was edited to allow user input for early exaggeration stop iteration, perplexity, total number of iterations, early exaggeration factor value, and learning rate value. KLD value and t-SNE coordinates were reported during each iteration or as frequently as requested. To allow generation of visually comparable t-SNE maps, the same random seed value was used and all experiments were repeated with several values of random seed. We did not observe noticeable differences in reported results between runs initiated with different seed values. For cross-validation and to benchmark against standard platforms, we utilized cloud-based Cytobank<sup>19</sup>, FlowJo V10.3-10.5 and FlowJo V9.9.6. Cytobank and FlowJo platforms were used to generate FCS files and graphical outputs from tabular data. FIT-SNE (Fast Fourier Transform-accelerated Interpolation-based t-SNE<sup>20</sup> plugin for FlowJo was used as indicated. Logs of t-SNE runs were batch-processed with VBA scripts and analyzed with GraphPad Prism 7. Expert-guided (manual) analysis of cell populations was performed in FlowJo 10.3-10.5 as described previously for specific datasets<sup>15,18</sup> or as explained below and used for map annotations as cluster classifiers.

For scRNAseq analysis, we used SeqGeq 1.3 package. We used PCA projections provided in 10X Genomics dataset to calculate the t-SNE embedding and annotated it using marker genes for major cell types. Louvain cluster classification was adopted from the SCANPY data analysis study<sup>17</sup>

The quality of the embeddings was assessed by a nearest neighbor (NN) classifier strategy as described in previous reports on t-SNE accuracy evaluation<sup>1,2,10</sup>. Briefly, for each observation, the k nearest neighbors (by Euclidean distance) were calculated using the 2D coordinates of the t-SNE map and the class assigned by expert gating was compared to the most common class of its k neighborhood. The rate of correct matches was tallied and represented as the overall nearest neighbor accuracy. The accuracy was also calculated on a per-class basis; different values for k (1,10,20,30,40,50) were reported.

**The standard t-SNE configuration for cytometry applications.** As described in detail in van der Maaten 2014, t-SNE computes low-dimensional coordinates of high-dimensional data resulting in similar and dissimilar data points in the raw data space placed proximally and at a distance, respectively, in the dimensionally reduced map. This map placement is achieved via t-SNE modeling the probabilities as a Gaussian distribution around each data point in the high-dimensional space and modeling the target distribution of pairwise similarities in the lower-dimensional space using Cauchy distribution (Student t-distribution with 1 degree of freedom). Then, the Kullback-Liebler Divergence (KLD) between the distributions is iteratively minimized via gradient descent. The gradient computation is essentially an N-body simulation problem with attractive forces (approximated to nearest neighbors using vantage-point trees) pulling similar points together and repulsive forces (approximated at each iteration using the Barnes-Hut algorithm) pushing dissimilar points apart.

An important part of t-SNE gradient descent computation is the “early exaggeration” (EE) that was proposed by van der Maaten and Hinton (2008) to battle the “overcrowding” artifact of embedding. With EE, all probabilities modeling distances in high-dimensional space are multiplied by a factor (early exaggeration factor, EEf, or  $\alpha$ ) for the duration of the first (typically 250, or 25% of the total number of) iterations. EE coerces data to form tight and widely separated clusters in the map and is considered to enable the map to find a better global structure.

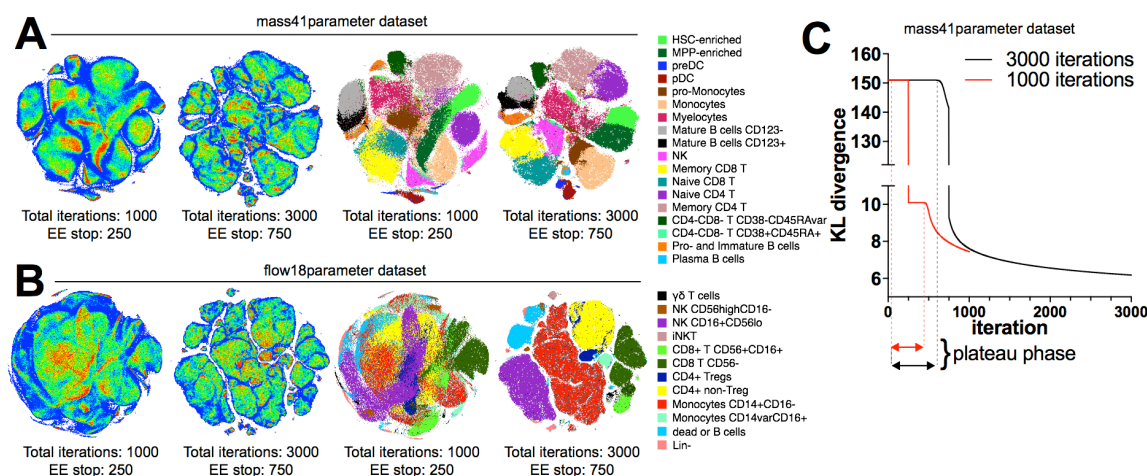
Multiple software platforms incorporate t-SNE algorithm specifically for analysis of cytometry data, including commonly used cytometry analysis desktop packages (FlowJo, FCS Express), and the cloud-based analysis platform Cytobank. Also, implementations of t-SNE are available as open-source packages in popular programming languages such as R (rtsne) and Python (sci-kit learn). Most of implementations wrap or re-write original C++ code of Barnes-Hut t-SNE (van der Maaten, 2014) and produce comparable analysis results upon direct comparison. Here, we customized the standard t-SNE C++ code to implement the parameter adjustments described in this work and published this customization as an open source solution to enable the research community to use this optimized t-SNE algorithm. Also, the equivalent adjustments of t-SNE have been made available as a cloud application from Omiq and integrated into FlowJo and SeqGeq programs.

## Results.

**The standard t-SNE configuration fails to visualize large datasets.** The t-SNE algorithm can be guided by a set of parameters that finely adjust multiple aspects of the t-SNE run<sup>21</sup>. However, cytometry data analysis software often locks or severely restrains the tunability of those parameters, likely to provide a simplified, ‘one-size-fits-all’ solution for t-SNE use in the software packages. Although each software platform has a unique combination of possible adjustments, most allow changes to both the number of iterations and to the perplexity (a soft measure for the number of nearest neighbors considered for each data point).

The datasets used throughout this study include at least 1 million datapoints of fluorescent or mass cytometry data and are therefore considerably larger than the smaller ( $< 5 \times 10^5$ ) datasets previously reported in benchmark studies of cytometry algorithmic tools<sup>22</sup>. Cytometry datasets larger than approximately  $5 \times 10^5$  events are generally observed to produce poor quality t-SNE maps and are therefore usually subsampled prior to analysis.

Empirically, cytometrists have observed that increasing the number of iterations of t-SNE computation results in better quality maps. We hypothesized that the resolution of t-SNE maps created from higher event counts could be dramatically improved via fine-tuning of t-SNE parameters. We first directly tested the relationship between iteration number and map quality by running two datasets (mass41parameter and flow18parameter, as described in Table 1) at the default 1000 iterations per run and with an “extended” 3000 iteration computation (Fig. 1A, B). To



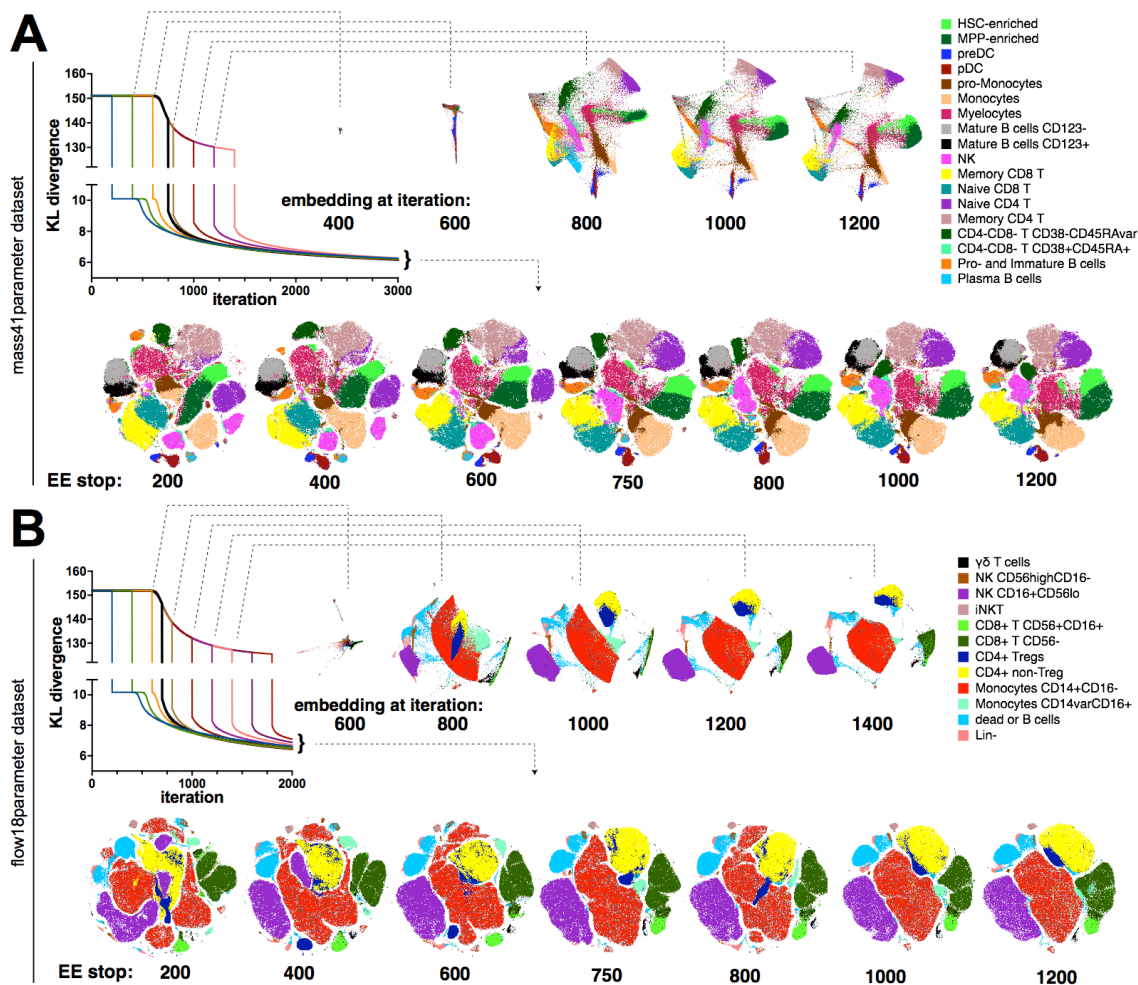
**Figure 1.** Performance of Barnes-Hut t-SNE implementation for cytometry data visualization showing comparison of standard (1000 iterations) and extended (3000 iterations) embeddings of mass cytometry (A) or flow cytometry (B) data presented as heatmap density plots (left) or color-coded population overlays based on ground-truth classification of single cell in the datasets. C. KLD change over iteration time of gradient descent for standard 1000 iterations (red line) or extended 3000 iterations (black line) embeddings of mass41parameter dataset. Representative examples of multiple runs with varying seed values are shown.

aid visualization in Fig 1 and the following figures maps are overlaid with color-coded populations derived from expert-driven (‘manual’) gating of the same dataset which serves as a ground-truth basis for data classification. As expected, 1000-iteration runs produced maps with poor visualization (overall 1-NN accuracy of embedding was 65% and as low as 18% for certain populations). Specifically, massive overlaps and random fragmentation of populations were observed. In contrast, the 3000-iteration runs resulted in maps with defined “islands” comprised of clearly isolated populations and no random fragmentations (overall 1-NN accuracy of embedding 96%; see Suppl. Table 1 for the detailed results of accuracy evaluations). Therefore, these findings are in agreement with the concept of a higher number of iterations resulting in higher quality t-SNE maps.

**KLD plateau phase resolves global cluster structure in t-SNE visualization.** In order to determine the cause of the difference in cluster resolution between the “default” and “extended” t-SNE runs, we examined the behavior of KLD over the duration of t-SNE embeddings (Fig 1C). As expected, the KLD value was inflated during the EE since EE is factored into gradient and KLD value calculation<sup>1</sup>. The EE is applied over 250 iterations in the “default” (red line, Fig. 1C, 1000 iterations) t-SNE configuration and 750 iterations in the “extended” run with 3000 iterations (black line Fig. 1C, 3000 iterations) since most platforms have EE scaled to 25% of total iteration number.

Notably, the KLD did not immediately minimize at the start of the EE in both the “default” and “extended” t-SNE runs; instead, the graph of KLD over time is a plateau that is followed by a curve that captures the incremental decrease of KLD, indicating the gradient descent. In the “default” run, the plateau was interrupted when the EE was stopped and KLD dropped, then continued with a non-exaggerated value of KLD.

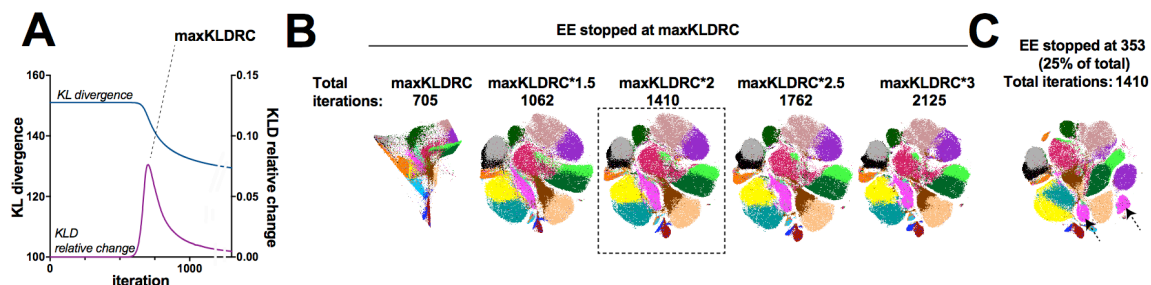
According to van der Maaten and Hinton, EE was introduced as a “trick” to improve resolution of the global structure of the data visualization that would not otherwise converge to separated clusters. As the suboptimal quality of the 1000-iteration t-SNE maps shown in Fig. 1A, B demonstrates poor global structure resolution, we hypothesized that by increasing the total



**Figure 2.** Effect of EE plateau phase on t-SNE visualization. EE was stopped after varying number of iterations and embedding output was examined at several intermediate timepoints and in the end of embedding for mass cytometry (total of 3000 iterations, A) and fluorescent cytometry (total of 2000 iterations, B) data visualization. Graphs showing KLD change over iteration time are color-labeled to distinguish curves corresponding to experiment perturbations, with black line indicating the run with the shortest EE but uninterrupted plateau. t-SNE maps are annotated with color-coded population overlays based on ground-truth classification of single cell in the datasets. Representative examples of multiple runs with varying seed values are shown.

iterations an analyst who uses t-SNE via conventional cytometry analysis platforms may inadvertently increase the number of EE iterations and this specific alteration may be the cause of the improvement in map visualization. To test this hypothesis, we compared multiple 2000- or 3000-iteration runs that differed in timing of the EE stop (Fig. 2A) by plotting the embedding at the iteration both when the EE stops and at later iterations, thus assessing the effects of EE and our perturbations on both mass cytometry (Fig. 2A) and flow cytometry (Fig. 2B) data visualization.

We found notable differences in map quality between the shorter and longer EE runs. Although the map after EE200/total3000 iterations appears visually superior than EE250/total1000 (Fig 1A, B) and could be considered a successful visualization, ground-truth labeling indicated that cluster fragmentation is present in both maps. When cluster fragments were plotted on a biaxial plot against parameters that were used in the t-SNE dimension reduction, we were not able to identify parameters that immediately contributed to their fragmentation (Suppl. Fig. 1).



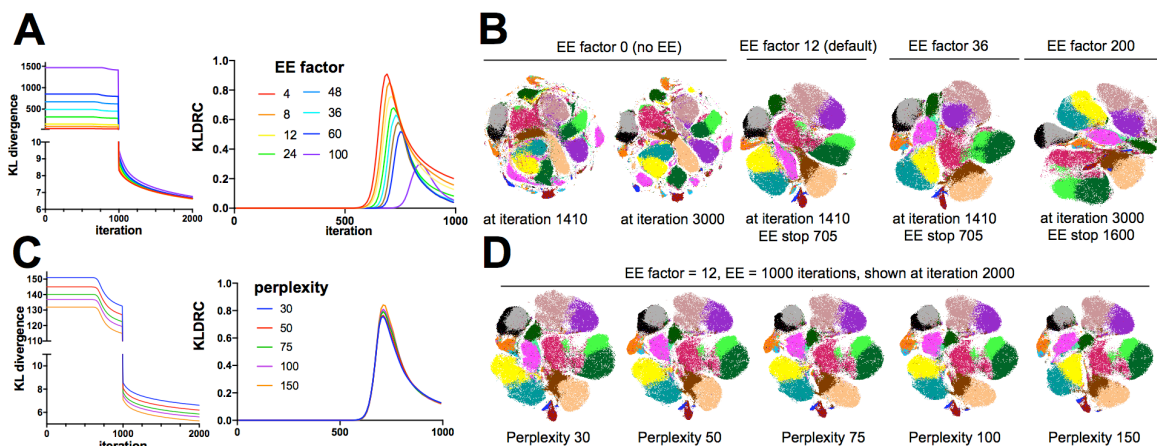
**Figure 3.** Early exaggeration plateau ensures optimal quality of visualization. A, KLD and KLD relative change plotted against iteration time. B, Mass cytometry data visualizations generated with varying duration of post-EE iteration time. Dashed box indicates visualizations used for comparison with default setup. C, Mass cytometry data visualization generated with default (25% of total) duration of EE. Arrows mark cluster fragmentation. All color overlays correspond to cell type classes labeled as in Fig. 1-2. Representative examples of multiple runs with varying seed values are shown.

Conversely, tight clusters that form at the end of the plateau remain mostly unchanged as long as the EE is being applied to the computation (Fig 2 A, B). The KLD minimization in that case could be explained by the gradual shrinking of the 2D space (data not shown). Once EE is removed, the attractive forces within each cluster are weakened and the local structure of the data is fully resolved within each cluster. Overall, these observations suggest that the EE stage of the gradient descent is essential for data clustering while the non-exaggerated descent results in resolution of local structures.

**Real-time monitoring of the KLD plateau results in optimal quality of t-SNE maps.** We have demonstrated that when the EE is too short, cell clusters continue to be resolved simultaneously with the local structure of each cluster being unfolded, leading to fragmented, overlapped or deformed “islands” in the resulting map (Fig. 2). Due to these results, we constructed an equation to find optimal EE timing. Specifically, we tracked the relative rate of KLD change ( $KLDRC_N = 100\% * ((KLD_{N-1} - KLD_N) / KLD_{N-1})$  where N is the iteration number) and identified the local maximum of KLDRC (maxKLDRC) (Fig. 3A). Since KLD is computed at each iteration, the maxKLDRC ‘sensor’ can be easily added to the algorithm programmatically and would stop EE at the next iteration past maxKLDRC. For the mass41parameter dataset of 1M datapoints, the maxKLDRC was detected at iteration 705 (Fig. 3A). Next, we ran t-SNE with an EE stop at iteration 706 and sampled map development at 706, 1.5x706, 2x706, 2.5x706 and 3x706 iterations (Fig 3B). As expected, at the maxKLDRC iteration the map contained the primordial clusters only; it was well shaped at 2 x max KLDRC and there was no visible improvement in map quality past that step and the visualization was very similar to the EE750/3000 map at Fig. 1A. When compared to the map created with ‘default’ settings of EE taking 25% of the run yet computed with the same number of iterations (1410), the maxKLDRC-triggered t-SNE produced visually and KLD-wise superior results within the similar computation time, and it also eliminated extensive trial-and-error calibration of t-SNE parameters (Fig. 3C). We propose a conservative approach to finalize the embedding automatically when  $(KLD_{N-1} - KLD_N) < KLD_N / 10,000$ . Alternatively, t-SNE projection output can be evaluated in real time to justify the termination of embedding.

**Moderate adjustments of EE factor and perplexity do not impact visualization.** Once we found EE to be crucial for map optimization, we next examined if the value of the EE factor  $\alpha$  can also be tuned to improve the results of t-SNE. We made  $\alpha$  user-accessible in our C++ t-SNE



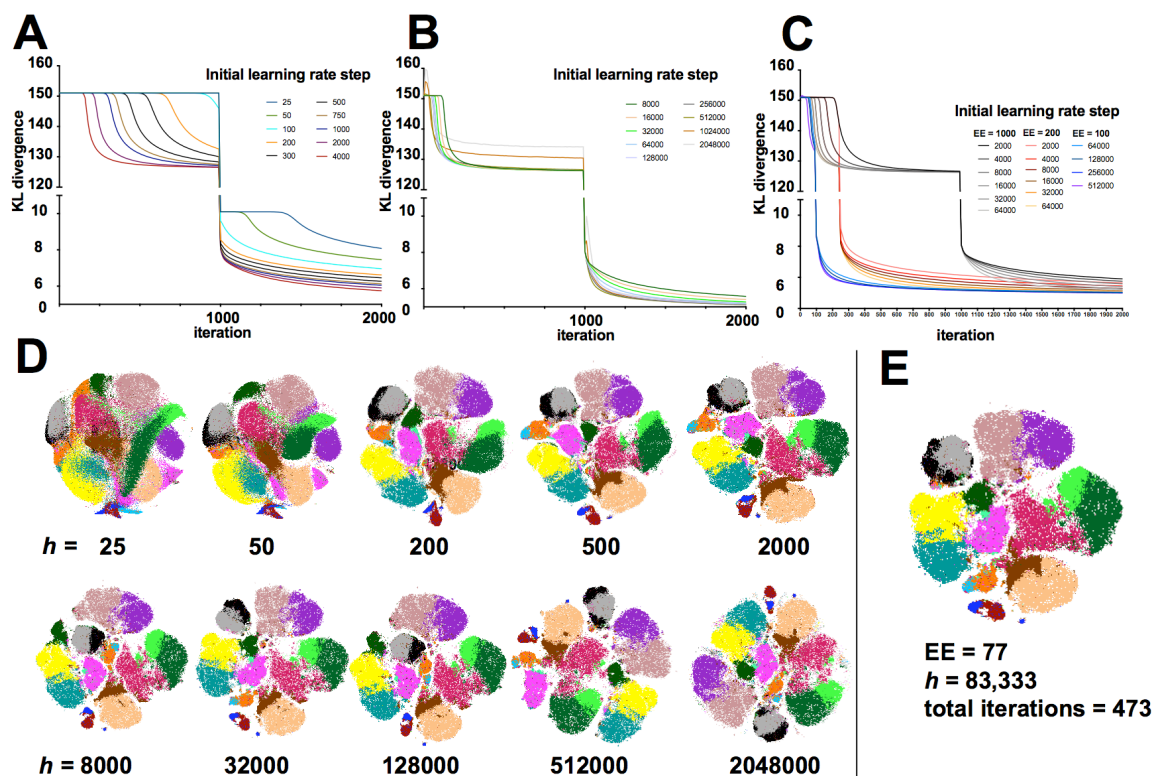


**Figure 4. Effects of perplexity and EE factor adjustments on t-SNE visualization of cytometry data.** A, B. KLD, KLDRC and t-SNE biaxial plots generated with varying EE factor values. C, D. KLD, KLDRC and t-SNE biaxial plots generated with varying perplexity. Graphs showing KLD and KLDRC change over iteration time are color-labeled to distinguish curves corresponding to experiment perturbations. All color overlays on t-SNE plots correspond to cell type classes labeled as in Fig. 1-2. Representative examples of multiple runs with varying seed values are shown.

code since it is hard-coded in the original Barnes-Hut C++ t-SNE implementation and all aforementioned results were obtained with default value of  $\alpha = 12$ . We chose the parameters defined above ( $\alpha = 12$ , EE = 706 iterations, 1410 iterations total) as our baseline for comparison since they provided optimal balance of map quality versus computation time. First, we tested how the optimization would proceed without EE ( $\alpha = 1$ ). We expected the run to fail or produce extremely crowded results as explained in the original t-SNE report<sup>1</sup>; however, we did not see much overlap in cluster positioning, probably due to the substantial number of map iterations run (Fig 4B). Nevertheless, the resulting map showed a lot of fragmentation proving to be an extreme case of an interrupted plateau phase. Even when run for as many as 3000 iterations, the fragmentation could not be remedied, again demonstrating the necessity of EE.

As expected, higher values of  $\alpha$  lead to much higher KLD during EE, however, the KLD values were similar at 2000 iterations when  $\alpha$  varied between 4 and 60 (Fig. 4A). Larger  $\alpha$  prolonged the plateau phase and became detrimental for KLD values when over 100. Visually  $\alpha = 200$  results in a distorted map with smaller populations lost. Therefore, we suggest that for cytometry applications the  $\alpha$  parameter may remain unchanged and set to 12, as suggested in van der Maaten 2014, or reverted to  $\alpha = 4$ , as originally proposed in van der Maaten and Hinton (2008) since per our results any value between 4 and 20 leads to comparable outcomes.

Increased perplexity has been proposed to be an intuitively beneficial method for visualization improvement since it translates to a larger number of considered nearest neighbors and hence a more accurate approximation of attractive forces, while decreased perplexity can completely fail the visualization<sup>21</sup>. KLD values for runs with varying perplexity cannot be directly compared since the KLD value is related to perplexity; however, KLD records over time do not show that increased



**Figure 5.** Learning step size optimization for t-SNE visualization of large datasets. A, B, C KLD change over iterations for embeddings with varying values of initial learning rate step size, color coded as indicated. A, EE = 1000 iterations, learning rate step = 25-4000; B, EE = 1000 iterations, learning rate step = 8K-2048K; C, EE = 100-1000 iterations, learning rate step = 2K - 512K. D, representative t-SNE plots of embeddings graphed on A. E, t-SNE plot of an optimal embedding. All color overlays on t-SNE plots correspond to cell type classes labeled as in Fig. 1-2. Representative examples of multiple runs with varying seed values are shown.

perplexity results in faster resolution of clusters (Fig 4C) or cleaner data visualization (Fig 4D). However, while changing  $\alpha$  does not affect t-SNE computation time, perplexity is linearly related to the time and memory required to create the embedding (data not shown). Although we and others have found some benefits of perplexity increases to map quality in otherwise suboptimal t-SNE runs, optimizing the EE step as described above and further in this work does not leave much space for improvement with perplexity tuning (Fig 4E).

**Learning step size is a key parameter to ensure t-SNE visualization of large datasets.** The step size in t-SNE gradient descent is updated at each iteration per Jacobs adaptive learning rate scheme<sup>23</sup>. This method increases the learning rate in directions in which the gradient is stable. A conservative initial value of 200 is hard-coded into most platforms. We hypothesized that larger datasets may stay longer on KLD plateau due to the number of iterations it takes to build up a sufficient learning rate step size. To evaluate this possibility, we titrated the step size  $h$  while observing the KLD with fixed EE=1000 iterations in Mass41parameter dataset. In agreement with our hypothesis,  $h = 25$  and  $h = 50$  runs failed to resolve from KLD plateau within 1000 iterations of EE (Fig. 5A) and  $h = 200$  finished the plateau in  $\sim 700$  iterations as previously shown. With further increases in  $h$ , we found that not only are progressively fewer iterations required to complete the plateau, but also that the final KLD of the maps scored at lower values. KLD is directly related

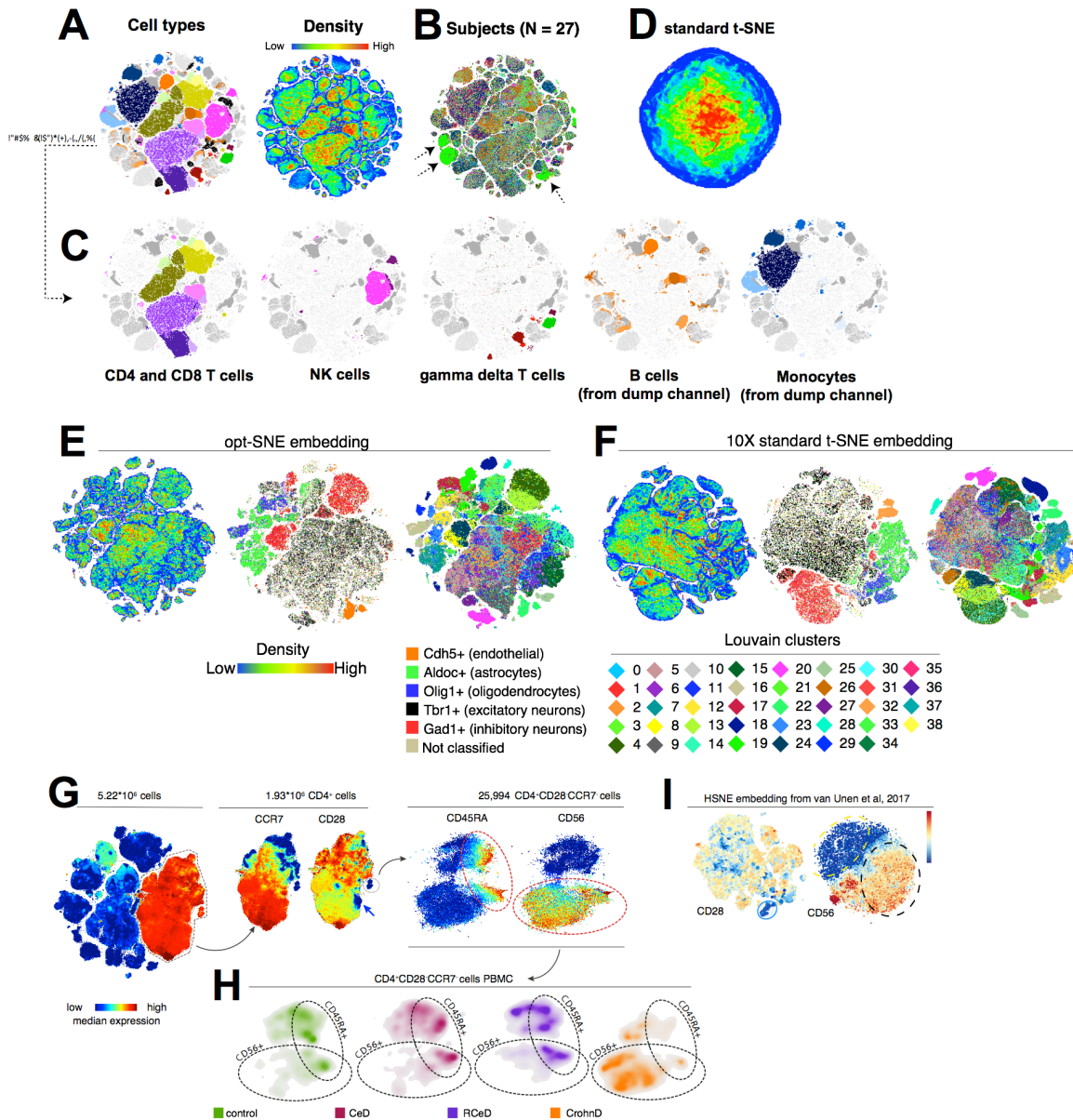
to the quality of visualization since it reflects the faithfulness of representation of high-dimensional data in t-SNE space; therefore, lower KLD values indicate superior visualization quality.

We continued to see improvement in plateau duration and KLD values with higher  $h$  values up to  $h \sim 64,000$ , a value that is drastically far from the “default”  $h = 200$  setting (note that in most platforms,  $h$  is restricted to ranges below 3,000) (Fig. 5B). At  $h \sim 256,000$  we observed irregular peaks in KLD graph indicating that the prescribed step size rendered gradient descent ineffective. However, using lower values of  $h$  we were able to converge the map with lowest KLD values at a fraction of time when limiting the EE step to 200 and even 100 iterations despite the  $10^6$  size of the dataset (Fig. 5C). Visual inspection of the embedded maps over the range of  $h$  values agrees with KLD values (Fig 5D).

In a recent publication, Linderman and Steinerberger<sup>24</sup> prove that in general t-SNE embedding will not converge if a product of EE factor  $\alpha$  and of fixed learning rate step size  $h$  is larger than the number of datapoints  $n$  (i.e. if  $\alpha h > n$ ). Since we employ an adaptive learning rate, our selection of initial  $h$  value is more forgiving; however, in our experiments we found the optimal settings of  $h$  to be close to  $h = n / 12$  for computations where  $\alpha = 12$ . Therefore, we propose to initiate the gradient descent with  $h = n / \alpha$  to create optimal t-SNE visualization (Fig 5E).

**Opt-SNE allows optimal embedding of massive datasets.** We implemented all proposed techniques including: (1) a dataset-specific automated early exaggeration step controlled by the KLDR sensor, (2) an optimal learning rate step size, and (3) a KLDR-driven embedding termination into a single workflow labeled ‘opt-SNE’, for ‘optimized t-Stochastic neighbor embedding’. To test the performance of opt-SNE, we used a  $20.1 \times 10^6$  event fluorescent cytometry dataset concatenated from two independent batches of PBMC samples ( $N=27$ ) stained with a variation of the OMIP-037 fluorescent cytometry panel<sup>18</sup> that allows detailed assessment of naïve and memory CD4+ and CD8+ T cells, NK cells and  $\gamma\delta$  T cells (Fig. 6A). The embedding completed in 770 iterations with only 73 iterations required to pass the EE step (Suppl. Fig. 3) and resulted in clear separation of cell clusters as evaluated by cell type annotation (Fig 6C). The majority of clusters appear to be populated by cells from all subjects with the exceptions of several populations that contained sample-unique debris features (Fig 6B, dashed arrows), confirming an absence of batch effects. A detailed breakdown of identified populations is presented in Fig 6C that shows subsets of CD4+ and CD8+ T cells, NK cells,  $\gamma\delta$  T cells, B cells, and monocytes. Importantly, B cell and monocyte lineage markers were detected together with a viability dye in a dump channel in this panel and therefore cannot be gated accurately via traditional biaxial plot analysis. However, opt-SNE identified them in high-dimensional space through the combination of several surface labels and light scatter characteristics, and each was successfully clustered into populations that were minimally mixed with dead cells. Use of the standard t-SNE algorithm completely failed to reveal the structure of the multi-million event flow cytometry dataset, even with several thousands of iterations (Fig 6D).

In order to test the suitability of opt-SNE for applications beyond flow and mass cytometry, we analyzed a  $1.3 \times 10^6$  cell single-cell RNA-seq dataset of mouse embryonic brain cells published by 10X Genomics. We used pre-calculated PCA projections included in the dataset to generate opt-SNE maps that we compared with 10X standard t-SNE embedding (Fig 6E). 10X used EE = 1000/total 4000 iterations of standard t-SNE while we used opt-SNE settings with  $h = 97,959$ , EE = 66/total 885 iterations (Suppl. Fig. 2). Non-immune single cell transcriptomics data are more



**Figure 6.** opt-SNE allows high-quality visualization of large cytometry and transcriptomics datasets. **A-D:** 20 million datapoints from fluorescent cytometry dataset concatenated from 27 subjects visualized in 2D space. **A** and **C**, cell type classes and density overlaid on 2D opt-SNE embedding. **B**, subject identifier overlaid on 2D opt-SNE embedding. Dashed arrows indicate clusters represented by datapoints from a single subject; **D**, standard t-SNE visualization (4000 iterations). **E-F**, 10X Genomics mouse brain scRNA-seq dataset (1.3 million datapoints) visualized in 2D space with opt-SNE (**E**) or standard t-SNE (**F**). From left to right: density features, single gene classes, and Louvain clusters (0-38) overlays. **G**, 5.22 million datapoints from mass cytometry dataset used in van Unen et al (2017) visualized in 2D space with opt-SNE. From left to right: CD4 expression overlaid on opt-SNE embedding; CCR7 and CD28 expression overlaid on CD4<sup>+</sup> opt-SNE cluster; CD45RA and CD56 expression intensity overlaid on CD4<sup>+</sup>CD28<sup>+</sup>CCR7<sup>-</sup> cluster. **H**, CD4<sup>+</sup>CD28<sup>+</sup>CCR7<sup>-</sup> cells from control, celiac disease (CeD), refractory celiac disease (RCeD) and Crohn's disease (CrohnD) subjects presented on density plots. Dashed encirclements indicate CD45RA<sup>+</sup> and CD56<sup>+</sup> areas of the cluster as defined in **G**. **I**, hierarchical t-SNE (HSNE) embedding of the CD4 (left) and CD4<sup>+</sup>CD28<sup>+</sup> cluster (right) from van Unen et al (2017). Color indicates marker expression intensity.

difficult to interpret with ground-truth classes since much fewer scRNA-seq markers can be easily interpreted for population identification. Therefore, we utilized both single gene classification and classification through the Louvain algorithm clustering using the Scanpy Python package<sup>17</sup> to

annotate the data. While opt-SNE and t-SNE both capture the macro-structure illustrated by gene overlays, several Louvain clusters (3, 18, 19, 25, 28) appear partially or completely overlapped by other clusters in standard t-SNE but not in the opt-SNE embedding. Therefore, opt-SNE allowed equivalent or superior resolution of single cell transcriptomics data as with standard t-SNE but with ~5x smaller iteration time (885 iterations of opt-SNE vs 4000 iterations of standard t-SNE).

HSNE (hierarchical SNE) is a t-SNE adaptation that was recently reported to facilitate analysis of large datasets by constructing a hierarchy of embeddings that can be explored from the overview of ‘landmark’ populations up to a single-cell level or resolution<sup>7</sup>. We applied opt-SNE to the 5.2 million point dataset that was reported in HSNE analysis of mass cytometry data, and compared opt-SNE visualization to the full resolution level of HSNE embedding (Fig. 6G, Suppl. Fig. 3). In the CD4<sup>+</sup> subset, opt-SNE visualization revealed two groups of CD4<sup>+</sup>CD28<sup>-</sup> cells likely representing terminally differentiated memory CD4<sup>+</sup> T cells<sup>25</sup> with different levels of CCR7 expression that may reflect distinct differentiation states of the two populations. While HSNE allowed identification of CD4<sup>+</sup>CD28<sup>-</sup>CCR7<sup>-</sup> cells, it was unable to visualize CD4<sup>+</sup>CD28<sup>-</sup>CCR7<sup>+</sup> cells as a single cluster and projected these cells sparsely in the CD4<sup>+</sup> island (Fig. 6I, left). Also, both algorithms projected heterogeneous expression of CD56 in CD4<sup>+</sup>CD28<sup>-</sup>CCR7<sup>-</sup> cells, but HSNE embedding did not resolve separate CD56<sup>+</sup> and CD56<sup>-</sup> clusters within CD4<sup>+</sup>CD28<sup>-</sup>CCR7<sup>-</sup> cells, loosely mapping them to the poles of the single round cluster (Fig. 6I, right). On the contrary, opt-SNE embedding visualized both the CD56<sup>+</sup> and CD56<sup>-</sup> clusters and disparate CD45RA expression within each cluster, revealing distinct phenotypes for the control and diseased subject groups (Fig. 6H). Also, opt-SNE demonstrated that the CD56<sup>+</sup>CD45RA<sup>-</sup> cells in the cluster originate from several subjects with Crohn’s disease (Fig. 6H). In summary, these results confirm that opt-SNE embedding provides superior visualization quality for complex cytometry data.

## Discussion

Visual exploration of data drives hypothesis formation and/or serendipitous discoveries; therefore, t-SNE is an extremely valuable tool for data comprehension. It is often used to facilitate data perception when hypothesis generation is automated by robust computational methods<sup>8,26</sup>. Comparison of t-SNE embeddings from multiple experimental conditions, timepoints, or subjects is invaluable to visualize sample-to-sample differences including disease hallmarks and longitudinal observations<sup>27</sup>. t-SNE is also valuable for quality assessment of data, when abnormal clustering could be traced back to sample preparation, data acquisition and preprocessing artifacts<sup>28</sup>. Therefore, batch embedding of multiple experimental points is essential for sample comparison and can only be enabled when t-SNE accommodates large datasets.

t-SNE was first introduced in cytometry research as a tool to visualize CyTOF data, since fluorescence-based high-parameter datasets were less common at that time. With recent advances in instrumentation and reagent availability, flow cytometry datasets with >20 parameters are quickly becoming prevalent and even standard in the field<sup>29-31</sup>, yet the proper data assessment tools are lacking for general use. DNA-barcoded antibodies have been recently utilized to allow simultaneous protein-epitope and transcriptome measurements in single cells<sup>32</sup> thus expanding the repertoire of traditional cytometry methods that could employ t-SNE as a staple method of data visualization and presentation.

One approach to large dataset t-SNE embedding is to model visualization with a subset of datapoints curated<sup>7,33</sup> or randomly selected (as implemented in cytometry data analysis platforms such as FCS Express and Gemstone) from the nearest neighbor graph. Such techniques may fail to project extremely rare datapoints, as demonstrated by CD4+CD28-CCR7+ cells that were identified with opt-SNE, but not HSNE in van Unen et al dataset (Fig. 6). This T cell subset was significantly less abundant in subjects with severe inflammatory conditions including Crohn's disease as compared to controls (Suppl. Fig. 3), marking this population as extremely rare. On the other hand, the CD4+CD28-CCR7-, while also low in frequency in most PBMC samples (<1% of all CD4+ cells), was likely successfully clustered by HSNE because two of the 11 Crohn disease subjects in the dataset showed an unusually high frequency of this population (40.7% and 31.0% of all CD4+ cells), allowing it to be well represented in the kNN graph of the full dataset.

Several attempts to successfully apply t-SNE-like methods to massive datasets have been recently reported including aforementioned HSNE<sup>7,33,34</sup>, LargeVis<sup>10</sup> and net-SNE<sup>35</sup>. However, these improved methods, when applied to large datasets, often require/benefit from considerable computational resources; for instance, the LargeVis study was performed on a 512Gb RAM, 32 core station. However, routine data analysis in a typical research laboratory should be possible with less available resources. We have not explicitly focused on computational efficiency in this work since we benchmarked the algorithm against itself with no specific emphasis on shortening computation time, which occurred due to the fewer number of iterations required to complete the data embedding with our adjustments. However, we have addressed several aspects of computation efficiency. For public use, we released an opt-SNE modification of multicore t-SNE C++ implementation (See Data and Software Availability) since the original Barnes-Hut t-SNE does not employ multi-threading. FIt-SNE, a recently published alternative to Barnes-Hut t-SNE that uses fast Fourier transform for much faster computation of repulsive forces approximation<sup>20</sup>, renders opt-SNE even more feasible on personal computers; when we combined our automated opt-SNE setup with FIt-SNE approximation in the van Unen dataset analysis (Fig. 6G-H), it was completed in about two hours on a 16Gb RAM personal notebook with 2 cores. Notably, we have not observed differences in embedding quality between embeddings generated with fast Fourier transform versus Barnes-Hut approximations when opt-SNE was used to control the visualization (Suppl. Fig. 4 and data not shown). All other analyses performed in this manuscript were performed using Barnes-Hut approximation on personal computers with the exception of the 20M embedding that required ~60Gb RAM at its peak and was run for several days on a multicore workstation using SeqGeq implementation of opt-SNE. Therefore, we expect opt-SNE to be applicable for existing or future adaptations of t-SNE even if alternative methods of computation are utilized<sup>20,36</sup> provided that they retain the core principles of t-SNE embedding. A promising approach that may be integrated with opt-SNE is the smart EE adjustment implemented in A-tSNE (approximated t-SNE)<sup>34</sup> algorithm where EE is removed gradually and on a per-point basis. Cytosplore, a novel software platform that includes HSNE and A-tSNE, allows the analyst to interactively initiate the local refinement of the map, resulting a significant improvement in computation time<sup>7</sup>.

Similar to other types of biological data, the structure of cytometry data is difficult to project due to its mixed nature, often comprised of cluster-like, manifold-like and/or hierarchical components<sup>28,37</sup>. In this paper we propose multiple techniques that are essential for optimal t-SNE data projection and are all germane to the fine-tuning of the early exaggeration stage of t-SNE embedding. EE facilitates cluster formation on a 2D plane<sup>1</sup> and serves as a necessary compromise

that allows clusters to escape the crowding effect. We designed an efficient measure to ensure that the cluster-like global structure of the data is fully revealed during the EE stage by monitoring the KLD output of the embedding in real time. As indicated by the lower KLD values of opt-SNE embedding where EE was limited to fewer iterations (Fig. 5E), prolonged amplification of the attractive forces that drive tight cluster formation in EE may be detrimental for the manifold-like local data structure represented by signal distributions that are continuous with background signal. These ‘continuum expression’ molecules include the proteins used to define classic immune cell subsets, as well as those linked with activation and/or exhaustion and markers indicating disease phenotypes<sup>38</sup>. Conversely, the non-exaggerated stage of t-SNE allows local data structures to be revealed.<sup>1</sup> Therefore, cytometry data analysis would be missing valuable information if we limited t-SNE applicability to finding only well separated clusters, especially since other techniques would perform that task better and faster. However, some workflows call for t-SNE pre-processing to facilitate extraction of cluster features from multidimensional data<sup>5,39</sup>. In those cases, it may be helpful to adapt the opt-SNE toolkit to terminate the embedding calculation immediately at the EE stop iteration and re-assess the high-dimensional structure within each cluster. Alternatively, a ‘late exaggeration approach’<sup>20</sup> can be cautiously applied to create very tight clusters, although in our experience this approach was only marginally beneficial for global structure representation but detrimental for local structure (data not shown).

It is advisable to note that certain data structures, such as cluster hierarchy, cannot be revealed with t-SNE<sup>40,41</sup>. t-SNE accessibility in cytometry analysis software lead to its not infrequent misuse with cytometry data, evident when the cluster-like structure is not prominent in the map. Therefore, the features identified from t-SNE embedding in its current form should be verified with alternative methods when possible for confirmational purposes. Im et al also suggest that if a continuous manifold structure exists in the data, large perplexity values may cause artificial breaks (overclustering)<sup>41</sup>. The perplexity values commonly used in cytometry analysis are on the lower end of the suggested range for efficient clustering, as it is often advised to scale the number of nearest neighbors to the average cluster size<sup>42</sup>; however, if computationally feasible, higher perplexity values might facilitate feature preservation for markers whose expression is not bimodally distributed.

In summary, we believe that opt-SNE is a powerful optimization toolkit that removes major limitations of t-SNE use for cytometric datasets and could thus potentiate novel data-driven findings in single cell research.

**Acknowledgements.** The authors would like to thank Yvan Saeys, Gary Kazantsev, Jonathan Irish, Allison Irvine and Katherine Drake for helpful discussions, and Geoff Kraker for technical assistance. We thank Sean Bendall from Stanford University and Vincent van Unen from Leiden University Medical Center for sharing their data, and Brian Tilton and Riley Pihl from BUSM Flow Cytometry Core Facility for assistance with data collection at BUSM. Anna C. Belkina is an ISAC (International Society for Advancement of Cytometry) SRL Emerging Leader 2015-2019 and thanks the ISAC organization and members for continuous support and encouragement. Josef Spidlen is an ISAC Marylou Ingram Scholar.

## Data and software Availability

All code including the open source multicore t-SNE C++ implementation, usage instructions, flow18parameters and mass41parameters datasets, and the cloud version of opt-SNE are available at <http://www.omic.ai/opt-SNE> (the C++ code with a Python wrapper is available from <https://github.com/omic-ai/Multicore-opt-SNE>).

Van Unen et al dataset is available at <http://flowrepository.org/id/FR-FCM-ZYRM>.

10X Genomics 1.3M scRNA-seq dataset is available at <https://support.10xgenomics.com/single-cell-gene-expression/datasets>

To facilitate availability to flow cytometry and scRNA-seq data analysts, opt-SNE has been incorporated into FlowJo version  $\geq 10.5.2$  and SeqGeq version  $\geq 1.4$ . At the time of this publication, this option was considered experimental and therefore hidden, but users can enable it by adding `<DRPlatform showAutoLearning="1" />` to the FlowJo10.prefs (or SeqGeq.prefs) XML file. It will be fully integrated and available by default in future releases of FlowJo and SeqGeq.

## Author contributions

ACB has conceived the study. ACB and COC developed the opt-SNE method. COC developed implementation in C++ and RH and JS developed implementation in FlowJo and SeqGeq. ACB has analyzed flow and mass cytometry datasets. RH and JS have analyzed the scRNA-seq dataset. RA and JSC provided conceptual input. ACB wrote the manuscript. All authors discussed the results and commented on the manuscript.

## Competing interests:

COC is a founder of Omiq, Inc. RH and JS are employees of FlowJo, LLC.

## References

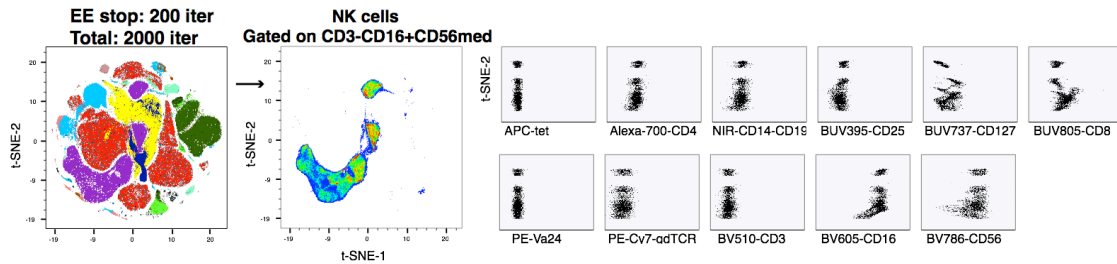
- 1 van der Maaten, L. & Hinton, G. Visualizing data using t-SNE. *The Journal of Machine Learning Research* **9**, 85 (2008).
- 2 Van Der Maaten, L. Accelerating t-SNE using Tree-Based Algorithms. *The Journal of Machine Learning Research* **15**, 3221-3245 (2014).
- 3 Amir el, A. D. *et al.* viSNE enables visualization of high dimensional single-cell data and reveals phenotypic heterogeneity of leukemia. *Nat Biotechnol* **31**, 545-552, doi:10.1038/nbt.2594 (2013).
- 4 Wong, M. T. *et al.* Mapping the Diversity of Follicular Helper T Cells in Human Blood and Tonsils Using High-Dimensional Mass Cytometry Analysis. *Cell reports* **11**, 1822-1833, doi:10.1016/j.celrep.2015.05.022 (2015).
- 5 Becher, B. *et al.* High-dimensional analysis of the murine myeloid cell system. *Nat Immunol* **15**, 1181-1189, doi:10.1038/ni.3006 <http://www.nature.com/ni/journal/v15/n12/abs/ni.3006.html#supplementary-information> (2014).
- 6 Chen, H. *et al.* Cytokit: A Bioconductor Package for an Integrated Mass Cytometry Data Analysis Pipeline. *PLoS Comput Biol* **12**, e1005112, doi:10.1371/journal.pcbi.1005112 (2016).



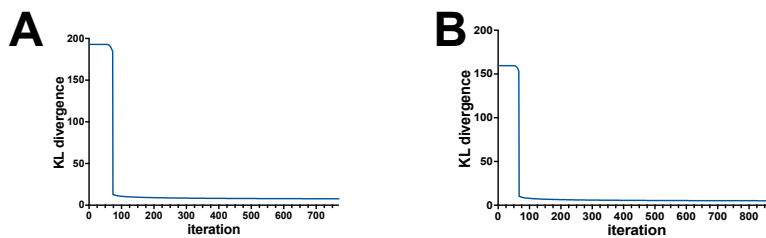
- 7 van Unen, V. *et al.* Visual analysis of mass cytometry data by hierarchical stochastic neighbour embedding reveals rare cell types. *Nat Commun* **8**, 1740, doi:10.1038/s41467-017-01689-9 (2017).
- 8 Van Gassen, S. *et al.* FlowSOM: Using self-organizing maps for visualization and interpretation of cytometry data. *Cytometry A* **87**, 636-645, doi:10.1002/cyto.a.22625 (2015).
- 9 Levine, J. H. *et al.* Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis. *Cell* **162**, 184-197, doi:10.1016/j.cell.2015.05.047 (2015).
- 10 Tang, J., Liu, J., Zhang, M. & Mei, Q. in *Proceedings of the 25th International Conference on World Wide Web %@ 978-1-4503-4143-1* 287-297 (International World Wide Web Conferences Steering Committee, Montr&#233;al, Qu&#233;bec, Canada, 2016).
- 11 Donnenberg, A. D. & Donnenberg, V. S. Rare-event analysis in flow cytometry. *Clin Lab Med* **27**, 627-652, viii, doi:10.1016/j.cll.2007.05.013 (2007).
- 12 DiGiuseppe, J. A., Tadmor, M. D. & Pe'er, D. Detection of minimal residual disease in B lymphoblastic leukemia using viSNE. *Cytometry. Part B, Clinical cytometry* **88**, 294-304, doi:10.1002/cyto.b.21252 (2015).
- 13 Lin, L. *et al.* Identification and visualization of multidimensional antigen-specific T-cell populations in polychromatic cytometry data. *Cytometry A* **87**, 675-682, doi:10.1002/cyto.a.22623 (2015).
- 14 Hirakawa, M. *et al.* Low-dose IL-2 selectively activates subsets of CD4(+) Tregs and NK cells. *JCI insight* **1**, e89278, doi:10.1172/jci.insight.89278 (2016).
- 15 Bendall, S. C. *et al.* Single-cell mass cytometry of differential immune and drug responses across a human hematopoietic continuum. *Science (New York, N.Y.)* **332**, 687-696, doi:10.1126/science.1198704 (2011).
- 16 Belkina, A. *et al.* Multivariate Computational Analysis of Gamma Delta T cell Inhibitory Receptor Signatures Reveals the Divergence of Healthy and ART-Suppressed HIV+ Aging. *bioRxiv* (2018).
- 17 Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol* **19**, 15, doi:10.1186/s13059-017-1382-0 (2018).
- 18 Belkina, A. C. & Snyder-Cappione, J. E. OMIP-037: 16-color panel to measure inhibitory receptor signatures from multiple human immune cell subsets. *Cytometry A* **91**, 175-179, doi:10.1002/cyto.a.22983 (2017).
- 19 Chen, T. J. & Kotecha, N. Cytobank: providing an analytics platform for community cytometry data analysis and collaboration. *Curr Top Microbiol Immunol* **377**, 127-157, doi:10.1007/82\_2014\_364 (2014).
- 20 Linderman, G. C., Rachh, M., Hoskins, J. G., Steinerberger, S. & Kluger, Y. Efficient Algorithms for t-distributed Stochastic Neighborhood Embedding. *arXiv:1712.09005* (2017).
- 21 Wattenberg, M. V., Fernanda; Johnson, Ian. How to Use t-SNE Effectively. *Distill*, doi:10.23915/distill.00002 (2016).
- 22 Weber, L. M. & Robinson, M. D. Comparison of clustering methods for high-dimensional single-cell flow and mass cytometry data. *Cytometry Part A* **89**, 1084-1096, doi:10.1002/cyto.a.23030 (2016).

- 23 Jacobs, R. A. Increased rates of convergence through learning rate adaptation. *Neural Networks* **1**, 295-307, doi:[https://doi.org/10.1016/0893-6080\(88\)90003-2](https://doi.org/10.1016/0893-6080(88)90003-2) (1988).
- 24 Linderman, G. C. & Steinerberger, S. Clustering with t-SNE, provably. *arXiv:1706.02582* (2017).
- 25 Mou, D., Espinosa, J., Lo, D. J. & Kirk, A. D. CD28 negative T cells: is their loss our gain? *American journal of transplantation : official journal of the American Society of Transplantation and the American Society of Transplant Surgeons* **14**, 2460-2466, doi:10.1111/ajt.12937 (2014).
- 26 Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nature Biotechnology* **36**, 411, doi:10.1038/nbt.4096 <https://www.nature.com/articles/nbt.4096#supplementary-information> (2018).
- 27 Wogtsland, C. E. *et al.* Mass Cytometry of Follicular Lymphoma Tumors Reveals Intrinsic Heterogeneity in Proteins Including HLA-DR and a Deficit in Nonmalignant Plasmablast and Germinal Center B-Cell Populations. *Cytometry. Part B, Clinical cytometry* **92**, 79-87, doi:10.1002/cyto.b.21498 (2017).
- 28 Mazza, E. M. C. *et al.* Background fluorescence and spreading error are major contributors of variability in high-dimensional flow cytometry data visualization by t-distributed stochastic neighboring embedding. *Cytometry Part A* **93**, 785-792, doi:10.1002/cyto.a.23566 (2018).
- 29 Staser, K. W., Eades, W., Choi, J., Karpova, D. & DiPersio, J. F. OMIP-042: 21-color flow cytometry to comprehensively immunophenotype major lymphocyte and myeloid subsets in human peripheral blood. *Cytometry A* **93**, 186-189, doi:10.1002/cyto.a.23303 (2018).
- 30 Mair, F. & Prlic, M. OMIP-044: 28-color immunophenotyping of the human dendritic cell compartment. *Cytometry A* **93**, 402-405, doi:10.1002/cyto.a.23331 (2018).
- 31 Nettey, L., Giles, A. J. & Chattopadhyay, P. K. OMIP-050: A 28-color/30-parameter Fluorescence Flow Cytometry Panel to Enumerate and Characterize Cells Expressing a Wide Array of Immune Checkpoint Molecules. *Cytometry A*, doi:10.1002/cyto.a.23608 (2018).
- 32 Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells. *Nature Methods* **14**, 865, doi:10.1038/nmeth.4380 <https://www.nature.com/articles/nmeth.4380#supplementary-information> (2017).
- 33 Pezzotti, N., Höllt, T., Lelieveldt, B., Eisemann, E. & Vilanova, A. Hierarchical Stochastic Neighbor Embedding. *Computer Graphics Forum* **35**, 21-30, doi:10.1111/cgf.12878 (2016).
- 34 Pezzotti, N. *et al.* Approximated and User Steerable tSNE for Progressive Visual Analytics. *IEEE Trans Vis Comput Graph* **23**, 1739-1752, doi:10.1109/TVCG.2016.2570755 (2017).
- 35 Cho, H., Berger, B. & Peng, J. Generalizable and Scalable Visualization of Single-Cell Data Using Neural Networks. *Cell Syst* **7**, 185-191 e184, doi:10.1016/j.cels.2018.05.017 (2018).

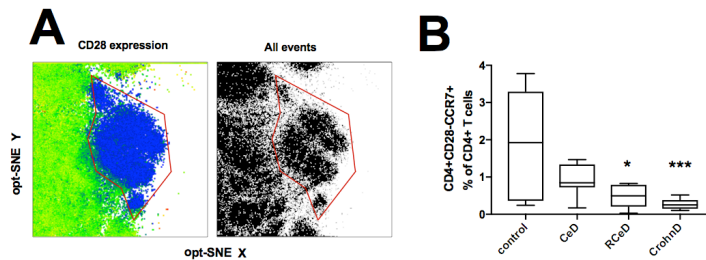
- 36 Chan, D. M., Rao, R., Huang, F. & Canny, J. F. t-SNE-CUDA: GPU-Accelerated t-SNE and its Applications to Modern Data. *arXiv:1807.11824* (2018).
- 37 Finn, W. G., Carter, K. M., Raich, R., Stoolman, L. M. & Hero, A. O. Analysis of clinical flow cytometric immunophenotyping data by clustering on statistical manifolds: Treating flow cytometry data as high-dimensional objects. *Cytometry Part B: Clinical Cytometry* **76B**, 1-7, doi:10.1002/cyto.b.20435 (2008).
- 38 Chattopadhyay, P. K. & Roederer, M. Cytometry: today's technology and tomorrow's horizons. *Methods* **57**, 251-258, doi:10.1016/j.ymeth.2012.02.009 (2012).
- 39 Shekhar, K., Brodin, P., Davis, M. M. & Chakraborty, A. K. Automatic Classification of Cellular Expression by Nonlinear Stochastic Embedding (ACCENSE). *Proceedings of the National Academy of Sciences* **111**, 202 (2014).
- 40 Amid, E. & Warmuth, M. K. A more globally accurate dimensionality reduction method using triplets. *arXiv:1803.00854* (2018).
- 41 Im, D. J., Verma, N. & Branson, K. Stochastic Neighbor Embedding under f-divergences. *arXiv:1811.01247* (2018).
- 42 Cao, Y. & Wang, L. Automatic Selection of t-SNE Perplexity. *arXiv:1708.03229* (2017).



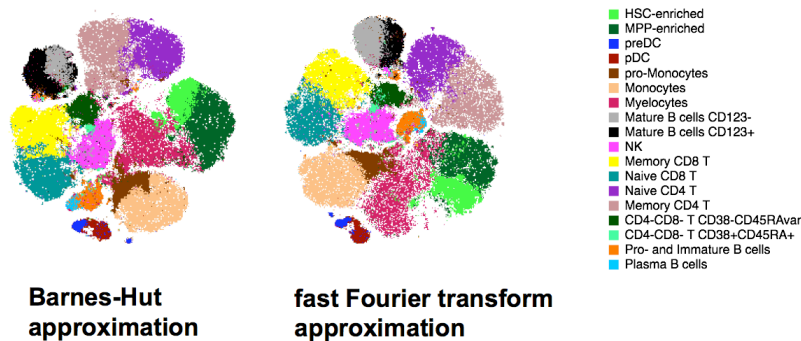
**Supplementary Figure 1. Interrupted EE plateau results in cluster fragmentation.** Clusters corresponding to CD3-CD16+CD56med NK cells were subsetted from the dataset and assessed as biaxial plots of different parameters plotted versus t-SNE-2 axis. Color overlays correspond to cell type classes labeled as in Fig. 2B (left) or density heatmap (middle).



**Supplementary Figure 2.** KLD and KLDRC graphs for opt-SNE embedding of flow20M dataset (A) and 10X Genomics 1.3 million datapoints scRNA-seq dataset (B).



**Supplementary Figure 3.** A. Enlarged fragment of opt-SNE embedding showing CD4+CD28-CCR7+ cell cluster (left: color indicates CD28 expression intensity; right: dot plot). B. Frequencies of CD4+CD28-CCR7+ cells in CD4+ PBMC compared between different cohorts of subjects. CeD, Celiac disease; RCeD, refractory celiac disease; CrohnD, Crohn's disease. Two-tailed T tests were performed for statistical analysis. \*  $p = 0.0258$ ; \*\*\*  $p = 0.0002$ .



**Supplementary Fig. 4.** Opt-SNE embeddings of mass41parameter datasets with Barnes-Hut (left) and fast Fourier transform (right) used to compute approximation. Color overlays correspond to cell type classes.



