

Supplementary information for **Differential cell-state**
abundance testing using KNN graphs with *Milo*

Emma Dann, Neil C. Henderson, Sarah A. Teichmann, Michael D. Morgan,
John C. Marioni

21 March, 2021

Contents

1	Supplementary notes	2
1.1	Description of workflow for <i>Milo</i> analysis	2
1.1.1	Preprocessing and dimensionality reduction	2
1.1.2	Minimizing batch effects	2
1.1.3	Building the KNN graph	3
1.1.4	Definition of cell neighbourhoods and index sampling algorithm	4
1.1.5	Testing for differential abundance in neighbourhoods	4
1.2	Guidelines on parameter choice	7
1.3	Notes on experimental design	8
	References	10

1 Supplementary notes

1.1 Description of workflow for *Milo* analysis

Given a single-cell dataset of gene expression profiles of M cells collected from S experimental samples, *Milo* aims to quantify systematic changes in the abundance of cells between biological conditions compared to within-condition variability. Here we provide a step-by-step description of the workflow for differential abundance analysis. Of note, we focus on the application to single-cell gene expression profiles, and we provide guidelines for pre-processing on this type of data. However, the core of the *Milo* framework, from KNN graph construction to differential abundance testing, is applicable to any kind of single-cell dataset that can be embedded in a low-dimensional space.

1.1.1 Preprocessing and dimensionality reduction

For pre-processing of scRNA-seq profiles we recommend following standard practices in single-cell analysis [1,2]: we normalize UMI counts by the total number of counts per cell, apply log-transformation and identify highly variable genes (HVGs). Then we project the $H \times M$ gene expression matrix, where M is the number of cells and H is the number of HVGs, to the first d principal components (PCs). While downstream analysis is generally robust to the exact choice of the number of HVGs [1], an optimal value for d can be selected by detecting the “elbow” in the variance explained by PCs or using the “jackstraw” method [3].

1.1.2 Minimizing batch effects

Comparing biological conditions often requires acquiring single-cell data from multiple samples, that can be generated with different experimental conditions or protocols. This commonly introduces batch effects, which can have a substantial impact on the data composition and subsequently the topology of any KNN graph computed across the single-cell data. Consequently, this will have an impact on the ability of *Milo* to resolve genuine differential abundance of cells between experimental conditions of interest. In addition, other biological nuisance covariates can impact DA analysis i.e. biological factors that are not of interest for the analyst, such as donor of origin or sex of the donor. We recommend mitigating the impact of technical or other nuisance covariates *before* building the KNN graph, by using one of the many *in silico* integration tools designed for this task in single-cell datasets. Defining the best tool for this task is beyond the scope

of this work; we refer the reader to a large number of integration methods that have been reviewed and benchmarked in [4–6]. However, users should consider the type of output produced by their integration method of choice, typically one of (A) a corrected feature space, (B) a joint embedding or (C) an integrated graph. The refined neighbourhood search procedure in *Milo* relies on finding neighbors in reduced dimension space. Therefore using an batch-correction method that produces an integrated graph (e.g. BBKNN [7], Conos [8]) may lead to sub-optimal results in DA testing with *Milo*, as the refined neighbourhood search procedure would still be affected by the batch effect.

In addition, the effect of nuisance covariates should be modelled in the generalized linear model used for DA testing in *Milo* to minimize the emergence of false positives in case of imperfect batch correction (see Section 1.1.5) (Fig.2D).

We wish to emphasize that, in the presence of confounding factors, an appropriate experimental design is crucial to obtain reliable results from differential abundance analysis: if nuisance factors are 100% confounded with the biological condition used for differential abundance (e.g. if the samples from diseased and healthy donors are processed in separate sequencing batches), there is no way to disentangle the abundance differences that are truly driven by the biology of interest. In a similar case applying a batch integration strategy before graph construction could lead to a loss of biological signal.

1.1.3 Building the KNN graph

Milo uses a KNN graph computed based on similarities in gene expression space as a representation of the phenotypic manifold in which cells lie. While *Milo* can be used on graphs built with different similarity kernels, here we compute the graph as follows: given the reduced dimension matrix X_{PC} of dimensions $M \times d$, for each cell j , the Euclidean distances to its K nearest neighbors in X_{PC} are computed and stored in a $M \times M$ adjacency matrix D . Then, D is made symmetrical, such that cells c_i and c_j are nearest neighbors (i.e. connected by an edge) if either c_i is a nearest neighbor of c_j or c_j is a nearest neighbor of c_i . The KNN graph is encoded by the undirected symmetric version \tilde{D} of D , where each cell has at least K nearest neighbors.

1.1.4 Definition of cell neighbourhoods and index sampling algorithm

Next, we identify a set of representative cell neighbourhoods on the KNN graph. We define the neighbourhood n_i of cell c_i as the group of cells that are connected to c_i by an edge in the graph. We refer to c_i with $i = 1, 2, \dots, N$ as the index cell of the neighbourhood, so that $N \leq M$. Formally, a cell c_j belongs to neighbourhood n_i if $\tilde{D}_{i,j} > 0$.

In order to define neighbourhoods that span the whole KNN graph, we sample index cells by using an algorithm previously adopted for waypoint sampling for trajectory inference [9,10]. Briefly, we start by randomly sampling $p \cdot M$ cells from the dataset, where $p \in [0, 1]$ (we use $p = 0.1$ by default). Given the reduced dimension matrix used for graph construction X_{PC} , for each sampled cell c_j we consider its K nearest neighbors with PC profiles x_1, x_2, \dots, x_k and compute the mean position of the neighbors in PC space \bar{x} :

$$\bar{x}_j = \frac{\sum_k x_k}{K}$$

Then, we search for the cell c_i such that the Euclidean distance between x_i and \bar{x} is minimized. Because the algorithm might converge to the same index cell from multiple initial samplings, this procedure yields a set of $N \leq p \cdot M$ index cells that are used to define neighbourhoods.

Having defined a set of N neighbourhoods from the sampled index cells, we construct a count matrix of dimensions $N \times S$ which reports, for each sample, the number of cells that are present in each neighbourhood.

1.1.5 Testing for differential abundance in neighbourhoods

To test for differential abundance between biological conditions, *Milo* models the cell counts in neighbourhoods, estimating variability across biological replicates using a generalized linear model (GLM). We build upon the framework for differential abundance testing implemented by *Cydar* [11]. In this section, we briefly describe the statistical model and adaptations to the KNN graph setting.

Quasi-likelihood negative binomial generalized linear models We consider a neighbourhood n with cell counts y_{ns} for each experimental sample s . The counts are modelled by the negative binomial (NB) distribution, as it is supported over all non-negative integers and can accurately model both small and large cell counts. For such non-Normally distributed data we use generalized-linear models (GLMs) as an extension

of classic linear models that can accomodate complex experimental designs. We therefore assume that

$$y_{ns} \sim NB(\mu_{ns}, \phi_n),$$

where μ_{ns} is the mean number of cells from sample s in neighbourhood n and ϕ_n is the NB dispersion parameter.

The expected count value μ_{ns} is given by

$$\mu_{ns} = \lambda_{ns} M_s$$

where λ_{ns} is the proportion of cells belonging to experimental sample s in n and M_s is the sum of counts of cells of s over all the neighbourhoods. In practice, λ_{ns} represents the biological variability that can be affected by treatment condition, age or any biological covariate of interest.

We use a log-linear model to model the influence of a biological condition on the expected counts in the neighbourhood:

$$\log \mu_{ns} = \sum_{g=1}^G x_{sg} \beta_{ng} + \log M_s \quad (1)$$

Here, for each possible value g taken by the biological condition of interest, x_{sg} is the binary vector indicating the condition value applied to sample s . β_{ng} is the regression coefficient by which the covariate effects are mediated for neighbourhood n , that represents the log fold-change between number of cells in condition g and all other conditions. If the biological condition of interest is ordinal (such as age or disease-severity) β_{ng} is interpreted as the per-unit linear change in neighbourhood abundance.

Estimation of β_{ng} for each n and g is performed by fitting the GLM to the count data for each neighbourhood, i.e. by estimating the dispersion ϕ_n that models the variability of cell counts for replicate samples for each neighbourhood. Dispersion estimation is performed using the quasi-likelihood method in `edgeR`[12], where the dispersion is modelled from the GLM deviance and thereby stabilized with empirical Bayes shrinkage, to stabilize the estimates in the presence of limited replication.

Count model normalisation and compositional biases In equation (1) above the $\log M_s$ term is provided as an offset to the NB GLM which effectively normalises the cell counts in each neighbourhood by the total number of cells in each sample S , thus accounting for variation in cell numbers across samples. If there is a single strong region of differential abundance then the counts for these samples will increase, which can negatively bias the model log fold-change estimates. This results in an underestimate of the

117 true log fold-changes and the appearance of false discoveries in the opposite direction to the true DA effect
 118 direction. This compositional effect arises due to the finite sampling of cells without saturation, such that
 119 as the counts for one region increase it drives down the counts for all other neighbourhoods generating
 120 false-positive differential abundance. To address this issue we turn to the RNA-seq literature, specifically
 121 the trimmed mean of M-values (TMM) method for estimating normalisation factors that are robust to such
 122 compositional differences across samples [13]. Under the assumption that the majority of neighbourhoods
 123 are not differentially abundant, the TMM approach first computes the per-neighbourhood log count ratios
 124 for a pair of samples s and s' (M values):

$$M_n = \log \frac{Y_{ns}/M_s}{Y_{ns'}/M_{s'}}$$

125 And the absolute neighbourhood abundance (A values):

$$A_n = \frac{1}{2} \log_2(Y_{ns}/M_s \cdot Y_{ns'}/M_{s'}), \text{ for } Y_n \neq 0$$

126 Both the M and A distribution tails are trimmed (30% for M, 5% for A by default) before taking a weighted
 127 average over neighbourhoods using precision weights, computed as the inverse variance of the neighbourhood
 128 counts, to account for the fact that more abundant neighbourhoods have a lower variance on a log scale.
 129 Thus, the normalisation factors are computed, with respect to a reference sample, r :

$$\log_2(TMM_s^{(r)}) = \frac{\sum_{n \in N} w_{ns}^r M_{ns}^r}{\sum_{n \in N} w_{ns}^r}$$

130 where, M_{ns}^r is computed as above for samples s and r , and:

$$w_{ns}^r = \frac{M_s - Y_{ns}}{M_s Y_{ns}} + \frac{M_r - Y_{nr}}{M_r Y_{nr}}$$

131 In practice, M_r and Y_{ns} are computed from the sample with the counts per million upper quartile that is
 132 closest to the mean upper quartile across samples.

133 **Adaptation of Spatial FDR to neighbourhoods** To control for multiple testing, we need to account for
 134 the overlap between neighbourhoods, that makes the differential abundance tests non-independent. We apply

a weighted version of the Benjamini-Hochberg (BH) method, where p-values are weighted by the reciprocal of the neighbourhood connectivity, as an adaptation to graphs of the Spatial FDR method introduced by *Cydar* [11]. Formally, to control for FDR at a selected threshold α we reject null hypothesis i where the associated p-value is less than the threshold:

$$\max_i p_{(i)} : p_{(i)} \leq \alpha \frac{\sum_{l=1}^i w_{(l)}}{\sum_{l=1}^n w_{(l)}}$$

Where the weight $w_{(i)}$ is the reciprocal of the neighbourhood connectivity c_i . As a measure of neighbourhood connectivity, we use the Euclidean distance between the neighbourhood index cell c_i and its k th nearest neighbour in PC space.

1.2 Guidelines on parameter choice

In this section we provide practical guidelines to select default parameters for KNN graph and neighbourhood construction for DA analysis with Milo. We recognize that DA analysis will also be impacted by choices made during feature selection and dimensionality reduction. However these depend strongly on the nature of the single-cell dataset used as input. For example feature selection strategies suitable for UMI-based scRNA-seq data might be suboptimal for data generated with non-UMI protocols, or dimensionality reduction methods alternative to PCA might be used for single-cell epigenomics data. We point the reader to existing resources and heuristics for the application to scRNA-seq in section 1.1.1.

Selecting the number of nearest neighbors K For construction of the KNN graph and neighbourhoods, the user has to select the number of nearest neighbors K to use for graph construction. The choice of K influences the distribution of cell counts within neighbourhoods, as K represents the lower limit in the number of cells in each neighbourhood ($\sum(y_{n,s})$). Hence, if K is too small the neighbourhoods might not contain enough cells to detect differential abundance. In order to perform DA testing with sufficient statistical power, the analyst should consider the number of experimental samples S (that will correspond to the columns in the count matrix for DA testing) and the desired minimum number of cells per neighbourhood and experimental sample. For example, having on average 5 cells per sample in each neighbourhood allows to detect n fold changes between 2 experimental conditions. The median number of cells per sample in each neighbourhood $\hat{y}_{n,s}$ increases with the total neighbourhood size (Suppl.Fig. . . .A), with:

$$\hat{y}_{ns} \sim \frac{\sum_s y_{ns}}{S}$$

Therefore a conservative approach to minimize false positives is to select $K \geq S \times 3-5$.

On the other hand increasing K increases power, but can come at the cost of FDR control, as we illustrate by testing for DA with increasing values for K in the mouse gastrulation dataset with synthetic condition labels on 4 different populations (Suppl.Fig. . . .A). We recommend users to inspect the histogram of neighbourhood sizes after sampling of neighbourhoods (Suppl.Fig. . . .B) and to consider the number of cells that would be considered a “neighbourhood” in the dataset at hand. As a heuristic for selecting a lower bound on K to increase the resolution of neighbourhoods for capturing rare sub-populations or states, the user can select K such that the mean neighbourhood size is no more than 10% of the expected size of the rare population. We provide the utility function `plotNhhoodSizeHist` to visualize the neighbourhood size distribution as part of our R package.

Selecting the proportions of cells sampled as neighbourhood indices p The proportion of cells sampled for search of neighbourhood indices can affect the total number of neighbourhoods used for analysis, but this number will converge for high proportions thanks to the sampling refinement step described in section 1.1.4 (Suppl.Fig. 1A). In practice, we recommend initiating neighbourhood search with $p = 0.05$ for datasets with more than 100k cells and $p = 0.1$ otherwise, which we have found to give appropriate coverage across the KNN graph while reducing the computational and multiple-testing burden. We recommend increasing $p > 0.1$ only if the dataset appears to contain rare disconnected subpopulations.

1.3 Notes on experimental design

Of key consideration when designing any single-cell experiment is how the sample collection relates to the biological variables of interest, and how these samples are processed and experiments are performed. Moreover, the experimenter (and analyst together), should design their experiment to minimise the impact of confounding effects on differential abundance testing, and incorporate appropriate replication to achieve enough power to detect the expected effect size for their experiment.

Statistical power considerations Increases in statistical power can be achieved by several means: (1) Increased cell numbers in neighbourhoods and (2) higher signal-to-noise ratio. The first can be achieved

by collecting more cells for each sample, increasing K during graph building such that neighbourhoods are on average larger, and by increasing the number of replicate samples. Collecting more cells gives a greater coverage of the cell-to-cell heterogeneity and different cell states/types, including increased detection for rarer sub-populations. Increasing K increases power by constructing larger neighbourhoods, however, this increase in power comes at a cost of reduced sensitivity for rarer sub-populations and an increased false discovery rate (Suppl Fig...). Designing an experiment with more replicate samples has multiple benefits in terms of increasing statistical testing power, increasing the signal-to-noise ratio, and increasing the accuracy of effect size estimates (Suppl Fig...). Therefore, in order of their impact on power and differential abundance testing, we would recommend: (1) collecting more replicate samples, with a minimum of $n=3$, (2) collecting more cells per sample, (3) increasing K to generate larger neighbourhoods. We would like to stress that choosing (3) should not be used to justify not considering either (1) or (2), and is way to increase power post sample collection, at a price.

Batch effects and experimental design Proper experimental design is crucial for answering scientific questions, particularly in the presence of confounding effects. In single-cell experiments these can range from batch effects introduced between samples processed on different days, owing to logistical constraint or sample availability, to biological sample collections from a heterogenous population; the latter being particularly apparent for non-model organisms.

In the context of differential abundance testing with Milo, we recommend designing experimental procedures and sample processing such that samples from different conditions are randomised across batches. One example is to pair samples between conditions, such that during batch effect removal the variability between these pairs of samples is minimally removed. This will help to facilitate removal of technical batch effects, whilst retaining the relevant biological variability.

As described above, the exact choice of batch integration method should be carefully considered before applying Milo, with a preference for methods that generate a batch-integrated space (either reduced dimensions or gene expression). The keypoint is that sample processing and experimental batches are not perfectly confounded with the biological variable of interest. We expect *some* technical variability to remain (no batch integration is perfect), which can be handled in Milo's GLM framework by including the batch identity as a blocking factor in the design model. Examples of this correction are shown in the benchmarking in Figure ... and Suppl Fig ...

References

1. Luecken, M.D., and Theis, F.J. (2019). Current best practices in single-cell RNA-seq analysis: A tutorial. *Molecular Systems Biology* *15*, e8746.
2. Amezquita, R.A., Lun, A.T.L., Becht, E., Carey, V.J., Carpp, L.N., Geistlinger, L., Marini, F., Rue-Albrecht, K., Risso, D., and Soneson, C. *et al.* (2020). Orchestrating single-cell analysis with Bioconductor. *Nature Methods* *17*, 137–145.
3. Chung, N.C., and Storey, J.D. (2015). Statistical significance of variables driving systematic variation in high-dimensional data. *Bioinformatics* *31*, 545–554.
4. Luecken, M.D., Büttner, M., Chaichoompu, K., Danese, A., Interlandi, M., Mueller, M.F., Strobl, D.C., Zappia, L., Dugas, M., and Colomé-Tatché, M. *et al.* (2020). Benchmarking atlas-level data integration in single-cell genomics. *bioRxiv*, 2020.05.22.111161.
5. Chazarra-Gil, R., Dongen, S. van, Kiselev, V.Y., and Hemberg, M. (2020). Flexible comparison of batch correction methods for single-cell RNA-seq using BatchBench. *bioRxiv*, 2020.05.22.111211.
6. Tran, H.T.N., Ang, K.S., Chevrier, M., Zhang, X., Lee, N.Y.S., Goh, M., and Chen, J. (2020). A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biology* *21*, 12.
7. Polański, K., Young, M.D., Miao, Z., Meyer, K.B., Teichmann, S.A., and Park, J.-E. BBKNN: Fast batch alignment of single cell transcriptomes. *Bioinformatics*.
8. Barkas, N., Petukhov, V., Nikolaeva, D., Lozinsky, Y., Demharer, S., Khodosevich, K., and Kharchenko, P.V. (2019). Joint analysis of heterogeneous single-cell RNA-seq dataset collections. *Nat Methods* *16*, 695–698.
9. Gut, G., Tadmor, M.D., Pe’er, D., Pelkmans, L., and Liberali, P. (2015). Trajectories of cell-cycle progression from fixed cell populations. *Nature Methods* *12*, 951–954.
10. Setty, M., Tadmor, M.D., Reich-Zeliger, S., Angel, O., Salame, T.M., Kathail, P., Choi, K., Bendall, S., Friedman, N., and Pe’er, D. (2016). Wishbone identifies bifurcating developmental trajectories from single-cell data. *Nature Biotechnology* *34*, 637–645.
11. Lun, A.T.L., Richard, A.C., and Marioni, J.C. (2017). Testing for differential abundance in mass cytometry data. *Nature Methods* *14*, 707–709.

- 241 12. Robinson, M.D., McCarthy, D.J., and Smyth, G.K. (2010). edgeR: A Bioconductor package for differ-
242 ential expression analysis of digital gene expression data. *Bioinformatics* *26*, 139–140.
- 243 13. Robinson, M.D., and Oshlack, A. (2010). A scaling normalization method for differential expression
244 analysis of RNA-seq data. *Genome Biology* *11*, R25. Available at: [http://genomebiology.biomedcentral.](http://genomebiology.biomedcentral.com/articles/10.1186/gb-2010-11-3-r25)
245 [com/articles/10.1186/gb-2010-11-3-r25](http://genomebiology.biomedcentral.com/articles/10.1186/gb-2010-11-3-r25) [Accessed March 18, 2021].