

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

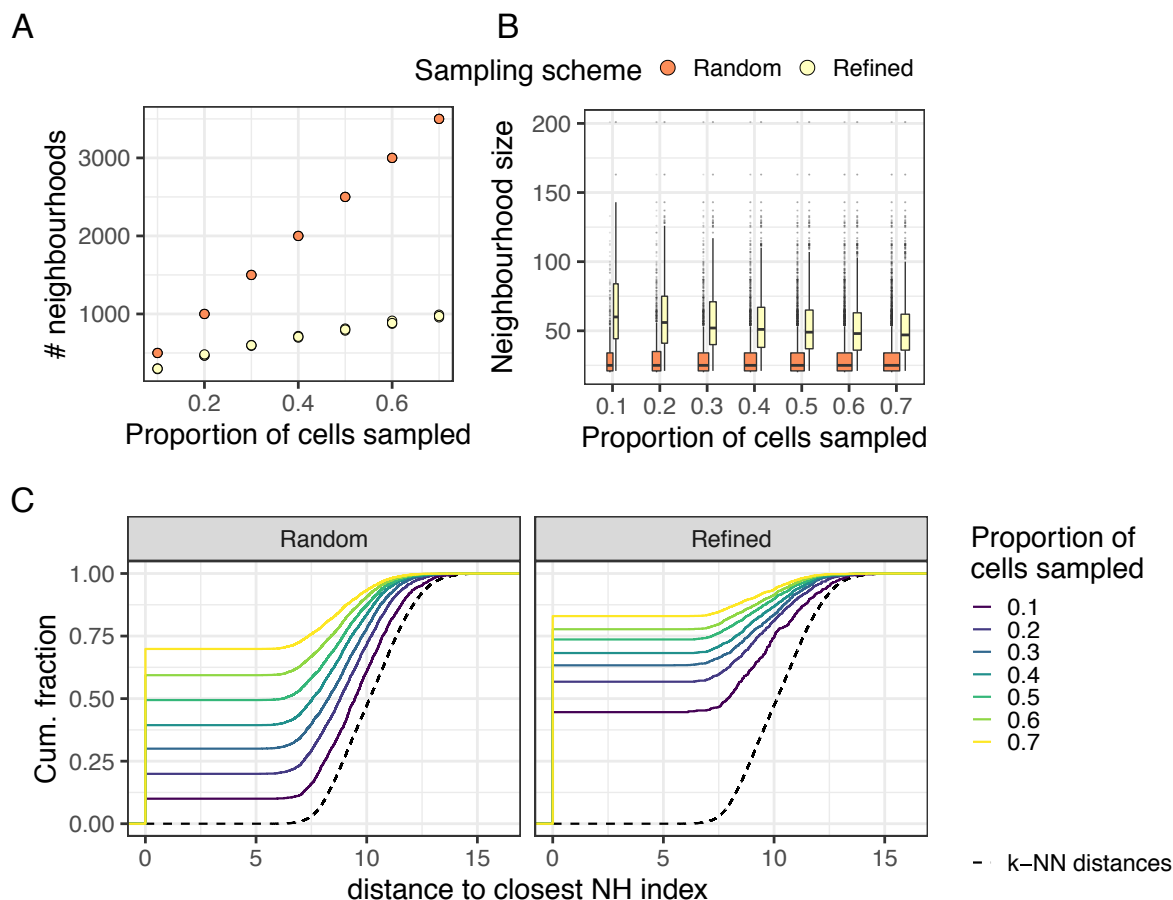
12 July, 2021

Contents

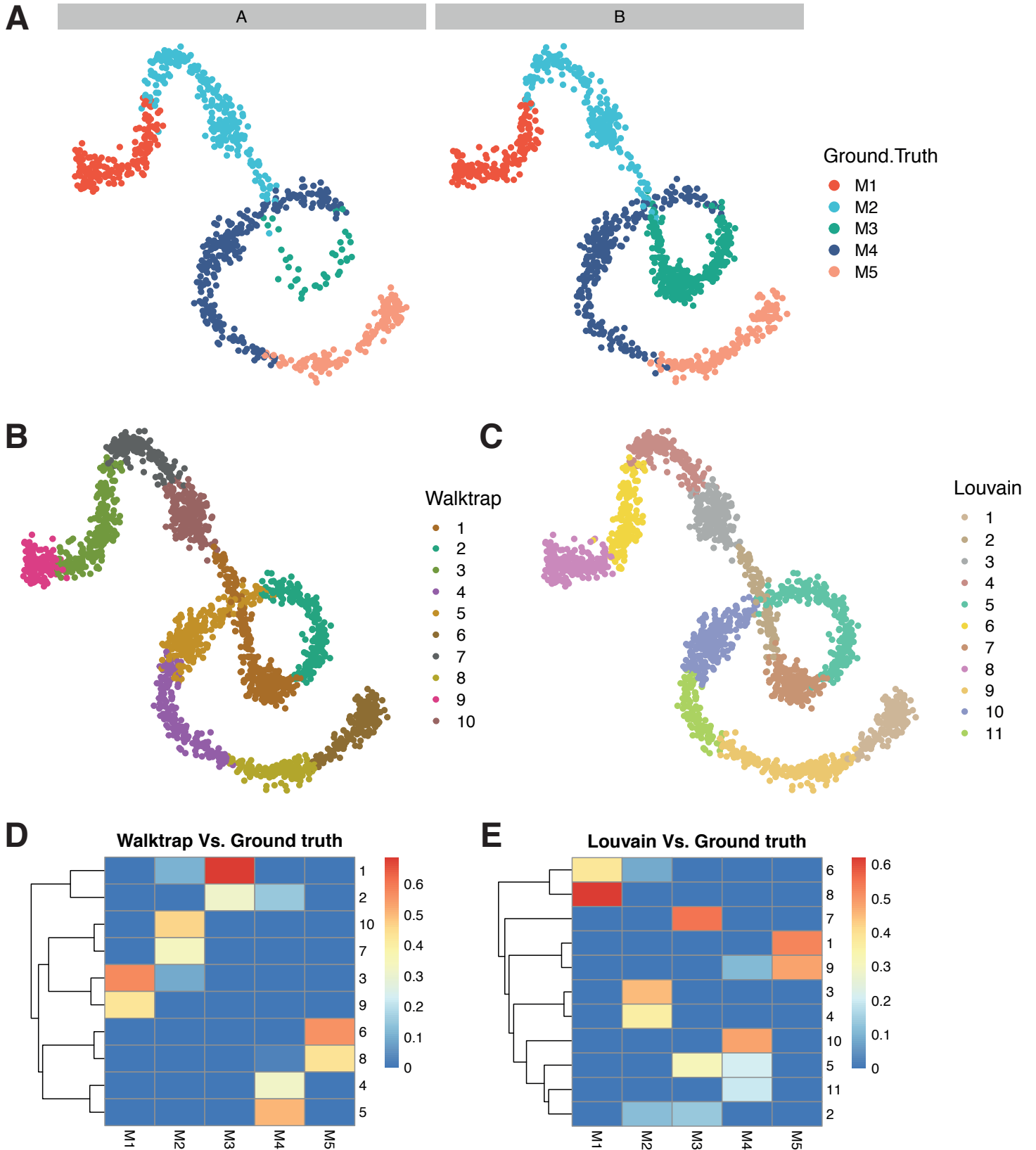
1	Supplementary Figures	2
2	Supplementary Tables	9
3	Supplementary notes	12
3.1	Description of workflow for <i>Milo</i> analysis	12
3.1.1	Preprocessing and dimensionality reduction	12
3.1.2	Minimizing batch effects	12
3.1.3	Building the KNN graph	13
3.1.4	Definition of cell neighbourhoods and index sampling algorithm	13
3.1.5	Testing for differential abundance in neighbourhoods	13
3.2	Guidelines on parameter choice	15
3.3	Notes on experimental design	16
4	Supplementary Note Figures	

References

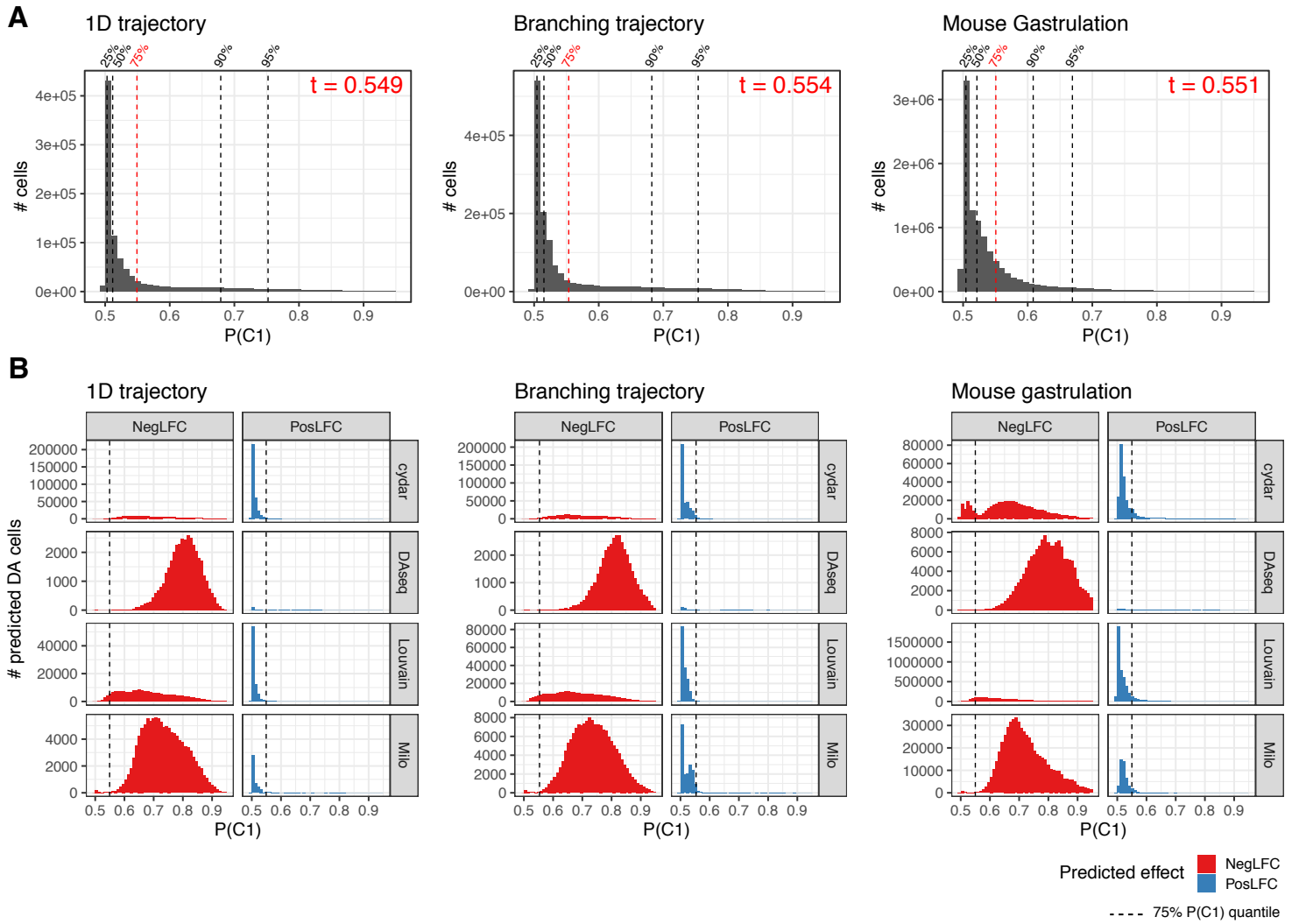
1 Supplementary Figures



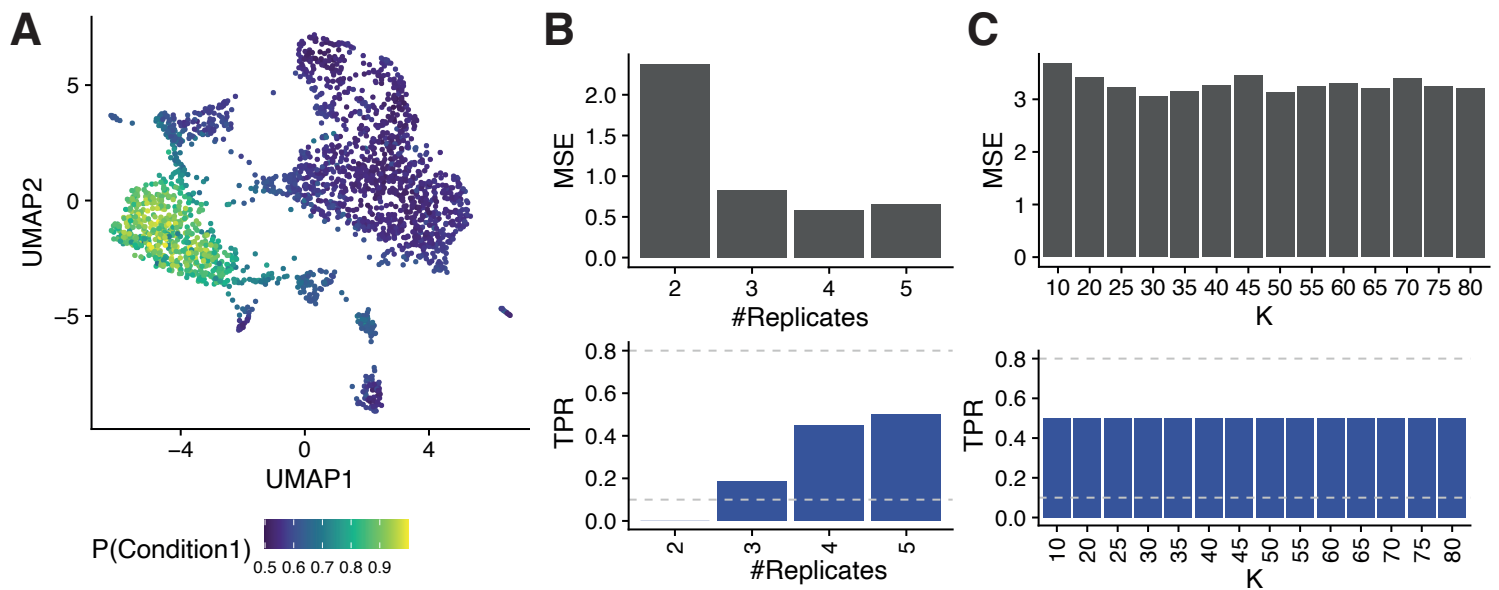
Supplementary Figure 1: **Random sampling of KNN graph vertices is suboptimal compared to sampling with refinement.** (A) Sampling with refinement leads to selection of fewer neighbourhoods (B) Sampling with refinement leads to selection of bigger neighbourhoods for DA testing, independently of the initial proportion of cells sampled. Box plots show the median with interquartile ranges (25–75%); whiskers extend to the largest value no further than 1.5x the interquartile range from the distance from the box, with outlier data points shown beyond this range. (C) Sampling with refinement generates robust neighbourhoods across initializations: for each index cell we calculate the distance from the closest index in a sampling with different initialization. The cumulative distribution of distances to the closest index is shown. The black dotted line denotes the distribution of distances between K nearest neighbours in the dataset (K=30) (NH: neighbourhood). Neighbourhood statistics were calculated using a simulated trajectory dataset of 5000 cells. All plots show results from three sampling initializations for each proportion.



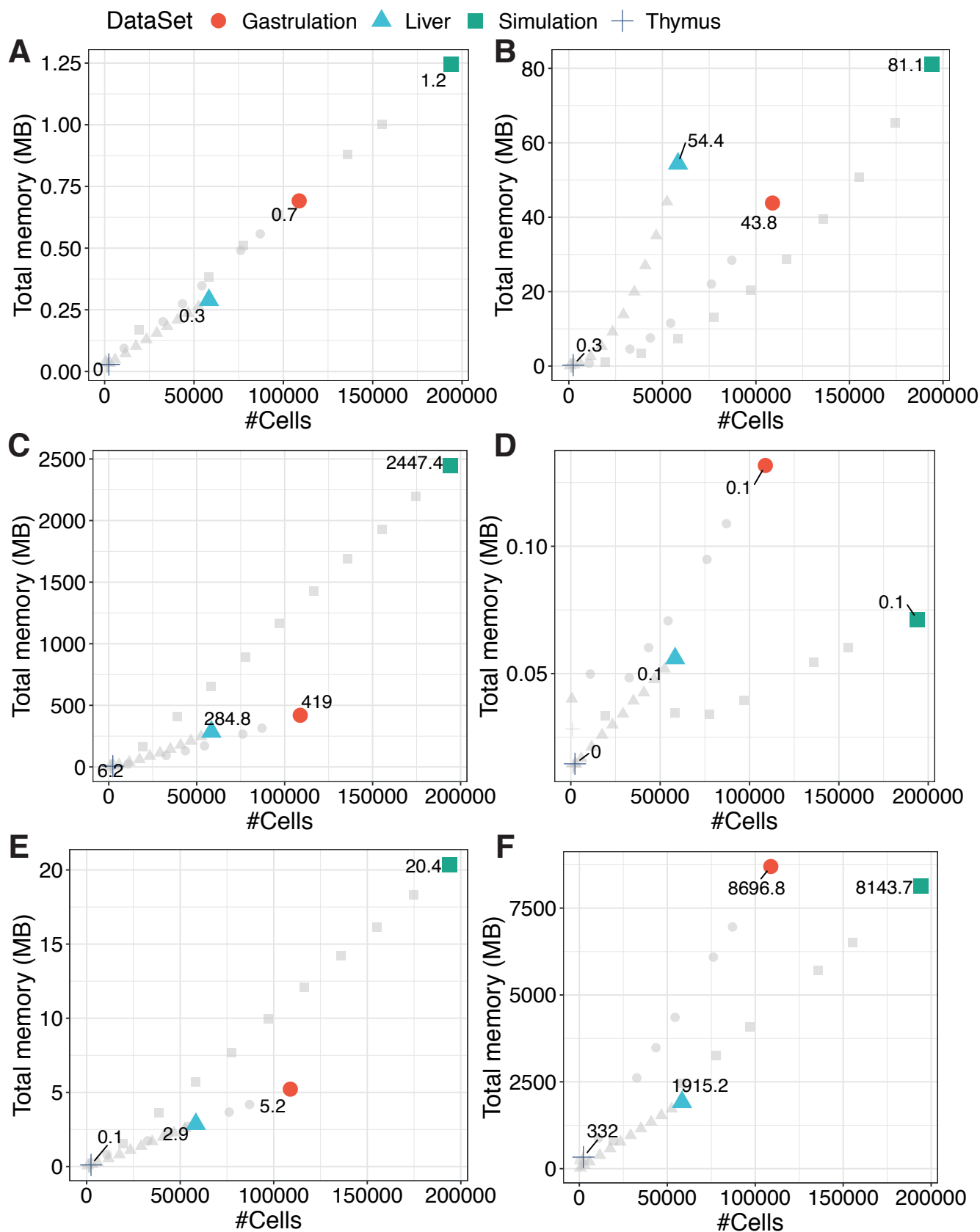
Supplementary Figure 2: **Graph-clustering does not faithfully capture simulated groups and differentially abundant subpopulations in a simulated continuous trajectory.** (A) A simulated linear trajectory of 2000 single-cells generated from 5 different groups, with cells assigned to either condition ‘A’ (left) or condition ‘B’ (right). (B) A Walktrap clustering of the data in (A) using the same KNN graph. Cells are coloured by Walktrap cluster identity. (C) A Louvain clustering of the data in (A) using the same KNN graph. Cells are coloured by the Louvain clustering identity. (D-E) Heatmaps comparing the numbers of cells in each cluster with respect to the ground truth groups in (A). Each entry in the heatmap is coloured by the proportion of cells from the column groups (ground truth) that are assigned to the respective cluster.



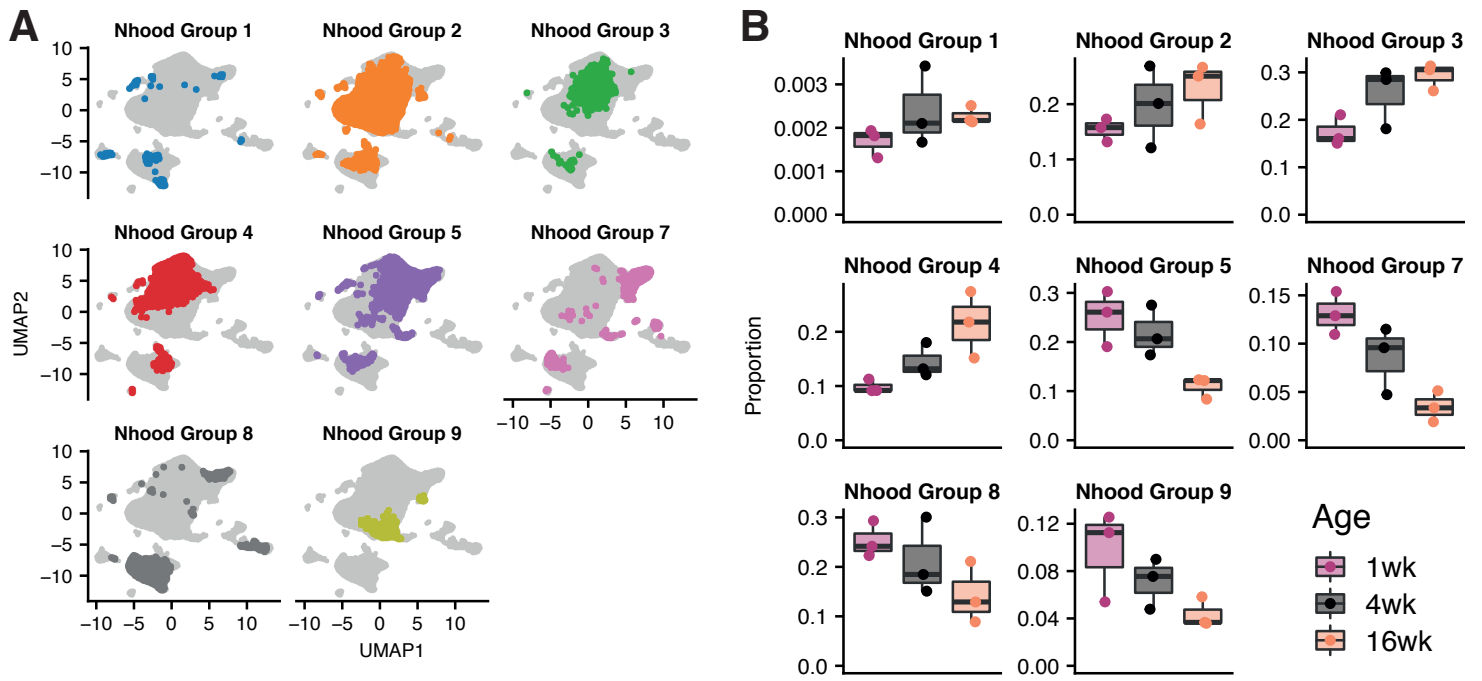
Supplementary Figure 3: **Selection of probability threshold to define ground-truth DA regions** (A) Histograms of $P(C1)$ for cells from *all* the simulations for each continuous dataset topology. The dotted lines indicate 25%, 50%, 75%, 90% and 95% quantiles. The red dotted line indicates the 75% quantile, that was chosen as the threshold t to define the DA region. The value of the threshold t is indicated in red. (B) Histograms of $P(C1)$ for cells predicted to be in DA regions by different methods in *all* the simulations on datasets of different topologies. Cells are split by the predicted direction of the effect inferred by DA methods, where NegLFC indicates predicted negative log-fold change (true positives) and PosLFC indicates predicted positive log-Fold Change (false positives). The dotted line indicates the threshold used to define the true DA region, corresponding to the 75% quantile of all $P(C1)$ values.



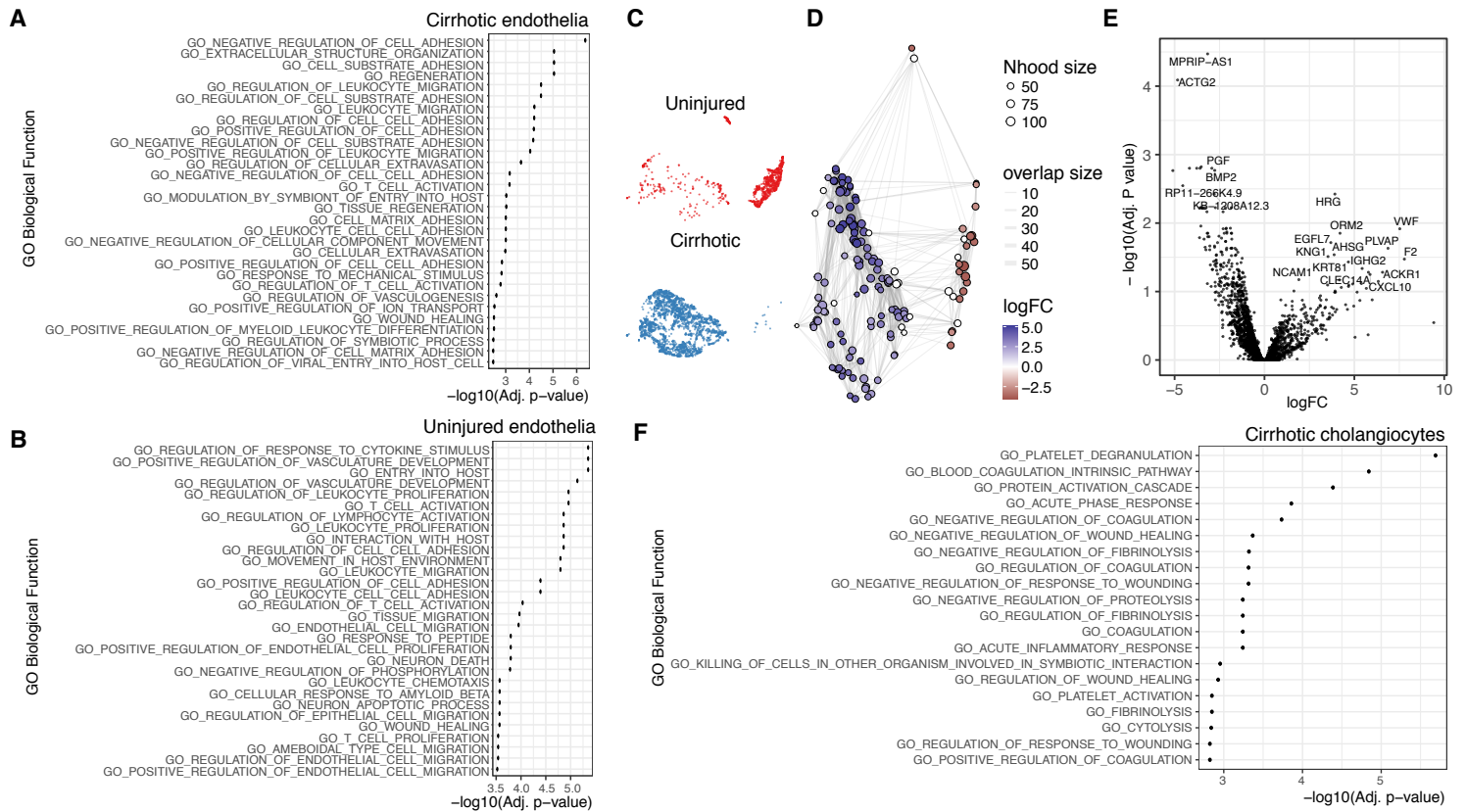
Supplementary Figure 4: **The impact of replication and k selection on effect size estimation variance.** (A) A UMAP of the mouse thymus data with a single simulated DA. Points are single cells coloured by the $P(\text{Condition1})$. (B) Increasing the number of replicates reduces the difference between the true simulated and estimated effect sizes, using the mean squared error (MSE; top panel) and increases the testing true positive rate (TPR; bottom panel). (C) Increasing k marginally reduces the estimation variance (top panel), and has less of an impact on power (bottom panel) compared to increased replication.



Supplementary Figure 5: **Memory usage across the Milo analysis workflow.** Total memory usage across the steps of the Milo analysis workflow in 4 datasets containing different numbers of cells (Gastrulation: circles, Liver: triangles, Thymus: crosses, Simulation: squares). Grey points denote down-sampled datasets of the corresponding type. Coloured points denote the total number of cells for the respective dataset. Total memory usage (y-axis) is shown in megabytes (MB). (A) KNN graph building, (B) neighbourhood sampling and construction, (C) within-neighbourhood distance calculation, (D) cell counting in neighbourhoods according to the input experimental design, (E) differential abundance testing, (F) total in memory R object size. A fixed value was used in all datasets for graph building and neighbourhood construction ($K=30$).



Supplementary Figure 6: **Label transferred neighbourhood groups onto droplet scRNA-seq cells.** (A) Joint UMAP embedding for SMART-seq and droplet scRNA-seq datasets, points are coloured by label-transferred neighbourhood groups for the droplet scRNA-seq cells. (B) Proportions of label-transferred neighbourhood groups across mouse ages (n=3 replicates per age), corresponding to (A). Boxplots show the median with interquartile ranges (25–75%); whiskers extend to the largest value no further than 1.5x the interquartile range from the distance from the box, with outlier data points shown beyond this range.



Supplementary Figure 7: **Downstream analysis of disease-specific subpopulations in liver cirrhosis.** (A) GO term enrichment analysis on marker genes of cirrhosis-enriched endothelia. (B) GO term enrichment analysis on marker genes of healthy-enriched endothelia. The top 30 significant terms are shown. (C-D) UMAP embedding (C) and graph representation (D) of neighbourhoods of 3369 cells from cholangiocyte lineage. (E) Volcano plot for DGE test on cholangiocytes DA subpopulations: the x-axis shows the log-fold change between expression in cirrhotic and healthy cholangiocytes. The y-axis shows the $-\log_{10}(\text{adjusted p-value})$. (F) GO term enrichment analysis on marker genes of cirrhosis-enriched cholangiocytes. The top 20 significant terms are shown.

2 Supplementary Tables

Parameter name	Description	Range of tested values			
KNN graph	Underlying k-NN graph, generated from simulated or real scRNA-seq datasets	Clusters [dyntoy]	1D trajectory [dyntoy]	Branching trajectory [dyntoy]	Mouse gastrulation atlas [Pijuan-Sala et al. 2019]
DA population	cell population selected as centroid for differential abundance region	M1	M1	M1	Caudal_neurectoderm
		M2	M2	M2	Erythroid2
		M3	M3	M3	Gut
			M4	M4	Somitic_mesoderm
			M5	M5	Pharyngeal_mesoderm
			M6	M6	Erythroid1
			M7	M7	Mesenchyme
				M8	ExE_endoderm
				M9	
				M10	
Logit parameter	Coefficient of logit transformation	0.5	0.5	0.5	0.5
Max C1 probability	maximum probability of Condition 1	0.75 - 0.95	0.75-0.95	0.75-0.95	0.75-0.95
Seed for label sampling	Random seed for sampling of condition labels and assignment to replicates	43, 44, 45	43, 44, 45	43, 44, 45	43, 44, 45
Batch effect magnitude	Standard deviation of gaussian vector added to all cells in the same batch	0	0	0	0, 0.25, 0.5, 0.75, 1
	Total # simulations	54	126	180	810

Supplementary Table 1: **Summary of parameters used for DA simulations**

	Milo	MELD	DAseq	Cydar	Louvain + GLM
Framework	R	python	R	R	R
Unit for DA estimate	KNN graph neighbourhood	Cell	Cell	PC hypersphere	Cluster
Hyperparameters	d	d	d	d	d
	K	K	minimum K	hypersphere radius	K
	prop		maximum K	downsampling fraction	(resolution)
			K step		
Clustering-free	yes	yes	yes	yes	no
Representative sampling across dataset	yes	yes	yes	no	yes
Output					
Effect size estimate	log-Fold change	condition likelihood	DA score	log-Fold change	log-Fold change
Statistical testing	yes	no	yes	yes	yes
Modelling variation between replicates	yes	no	no	yes	yes
Spatial FDR control	yes	no	no	yes	NA
DA testing experimental design					
two condition	yes	yes	yes	yes	yes
multi-condition	yes	yes	no	yes	yes
continuous condition	yes	no	no	yes	yes
nuisance covariate control	yes	no	no	yes	yes
interaction	yes	no	no	yes	yes

Supplementary Table 2: **Qualitative comparison of evaluated methods for DA analysis**

Method	Parameters	Values - clusters	Values - 1D trajectory	Values - branching trajectory	Values - mouse gastrulation	Significance threshold
Milo	K	15	20	20	50	10% FDR
	d	30	30	30	30	
	prop	0.1	0.1	0.1	0.1	
MELD	K	NA	NA	NA	50	NA
	d	NA	NA	NA	30	
Cydar	tol	0.8	6	6	1	10% FDR
	downsample	3	3	3	3	
	d	30	30	30	30	
DAseq	k.vector	15-500, steps of 50	20-500, steps of 50	20-500, steps of 50	50-500, steps of 50	DA score > permutation threshold (pred.thres = NULL)
	d	30	30	30	30	
Louvain + GLM	K	15	20	20	50	10% FDR
	d	30	30	30	30	

Supplementary Table 3: Summary of parameters used for benchmarking of DA methods

3 Supplementary notes

3.1 Description of workflow for *Milo* analysis

Given a single-cell dataset of gene expression profiles of L cells collected from S experimental samples, *Milo* aims to quantify systematic changes in the abundance of cells between biological conditions. 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.

3.1.1 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 $H \times L$ gene expression matrix, where L 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].

3.1.2 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 a 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 3.1.5) (Fig 2D, Extended Data Fig 4C).

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.

3.1.3 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 $L \times d$, for each cell c_j , the Euclidean distances to its K nearest neighbors in X_{PC} are computed and stored in a $L \times L$ 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.

3.1.4 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 $N \leq L$. 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 $p \cdot L$ 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}_j is minimized. Because the algorithm might converge to the same index cell from multiple initial samplings, this procedure yields a set of $N \leq p \cdot L$ 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.

3.1.5 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

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

where λ_{ns} is the proportion of cells belonging to experimental sample s in n and L_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:

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

Here, for each possible value g taken by the biological condition of interest, x_{sg} is the 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 $\log L_s$ 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 true log fold-changes and the appearance of false discoveries in the opposite direction to the true DA effect direction. To address this issue we turn to the RNA-seq literature, specifically the trimmed mean of M-values (TMM) method for estimating normalisation factors that are robust to such compositional differences across samples [13]. Under the assumption that the majority of neighbourhoods are not differentially abundant, the TMM approach first computes the per-neighbourhood log count ratios for a pair of samples s and s' (M values):

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

And the absolute neighbourhood abundance (A values):

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

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

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

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

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

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

Adaptation of Spatial FDR to neighbourhoods To control for multiple testing, we need to account for 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.

3.2 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 3.1.1.

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. As we illustrate by testing for DA with increasing values for K in the mouse gastrulation dataset with synthetic condition labels (Supp Note Fig 1A-B) increasing K increases power, but can come at the cost of FDR control. 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. The median number of cells per sample in each neighbourhood \hat{y}_{ns} increases with the total neighbourhood size (Supp Note Fig 1C), with:

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

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

We recommend users to inspect the histogram of neighbourhood sizes after sampling of neighbourhoods (Supp Note Fig 1D) 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 `plotNhoodSizeHist` to visualize the neighbourhood size distribution as part of our R package.

To verify how robust the findings from Milo are to the choice of K , we repeated DA analysis on both the mouse thymus and human liver data sets presented in the Results across a range of values of K , from very

small ($K=2$) to very large ($K=100$). We found that for these 2 data sets the DA regions correspond to those identified in Fig 4-5 across a range of values of K (Supp Note Fig 2). We computed the log fold change for the DA neighbourhoods at each value of K and found that the differential abundance results are robust, with loss of power seen only at very small values ($K<10$). As K becomes very large ($K>50$), neighborhoods contain more heterogeneous mixtures of cells, leading to an “over-smoothing” and a loss of resolution for rarer cell states in the thymus dataset (for example in the sTEC cluster).

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 3.1.4 (Supp 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 selecting $p > 0.1$ only if the dataset appears to contain rare disconnected subpopulations.

3.3 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 (Supp Note Fig 1A-B). 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 (Supp Fig 4B). 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.

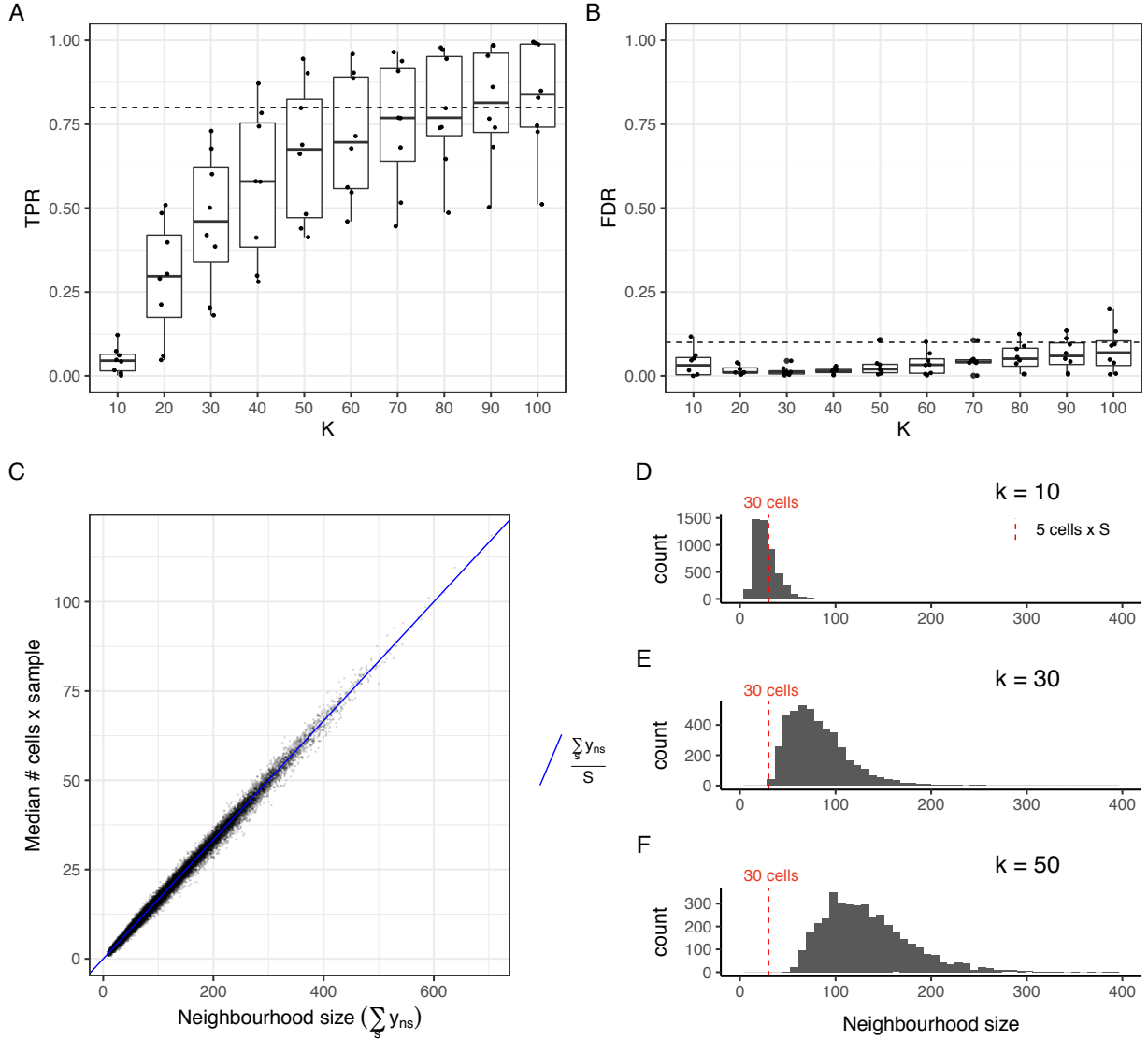
Batch effects and experimental design Proper experimental design is crucial for answering scientific questions, particularly in the presence of confounding effects. In single-cell experiments these can range from batch effects introduced between samples processed on different days, owing to logistical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

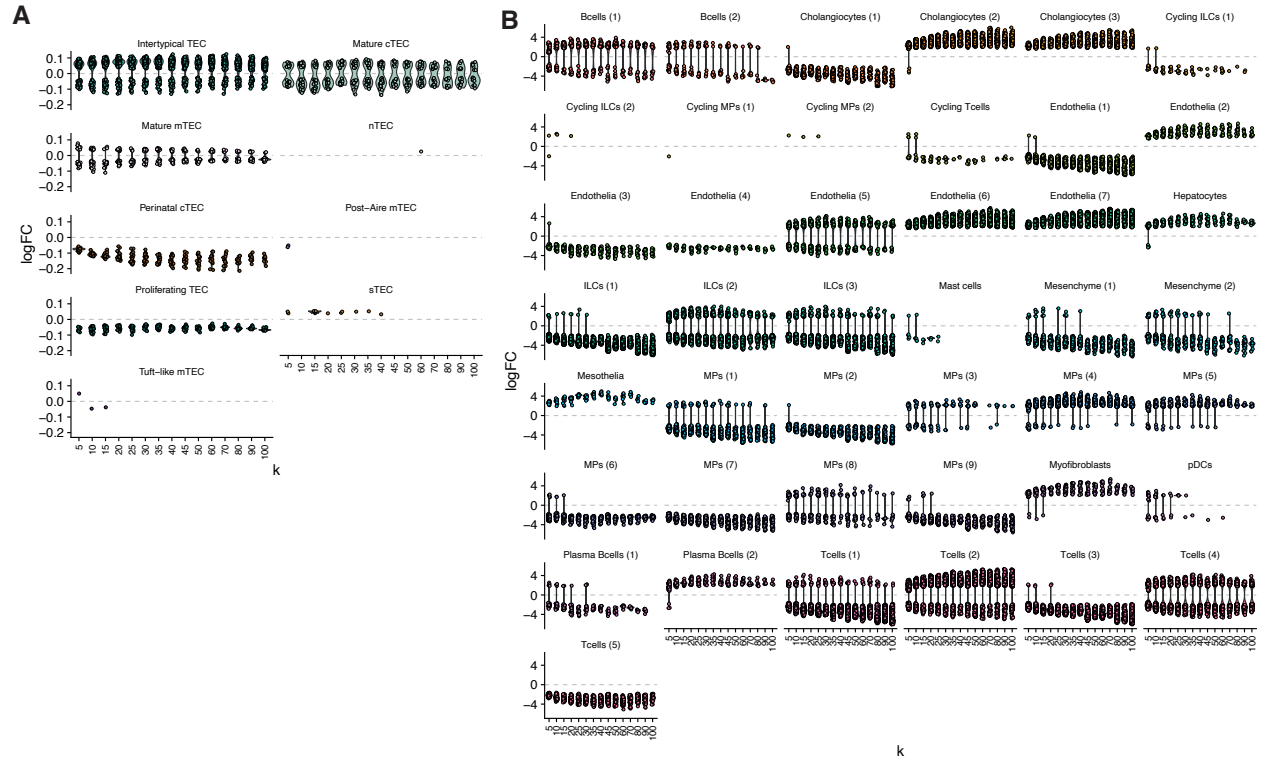
As described above, the exact choice of batch integration method should be carefully considered before applying Milo, with a preference for methods that generate a batch-integrated space (either reduced dimensions or gene expression). The key point is that sample processing and experimental batches are 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 Fig 2E and Extended Data Fig 4C.

4 Supplementary Note Figures



Supplementary Note Figure 1: **Selection of K parameter** (A-B) Example trends for TPR and FDR for increasing values of K used for KNN graph building on simulated DA on 8 regions ($P(C1) = 0.8$). Dotted lines highlight TPR=0.8 and FDR=0.1 thresholds. Results for simulations on $n=8$ populations are shown. Box plots show the median with interquartile ranges (25–75%); whiskers extend to the largest value no further than 1.5x the interquartile range from the distance from the box, with outlier data points shown beyond this range. (C) The median number of cells per experimental sample is a function of the neighbourhood size $\sum_s y_{n,s}$ divided by the total number of samples S . (D-F) Histogram of neighbourhood sizes for different choices of K . The red dotted line denotes the minimum neighbourhood size to obtain 5 cells per sample on average.



Supplementary Note Figure 2: **Robustness of Milo DA testing to varying K**. Distributions of DA neighbourhoods across values of K for the mouse ageing thymus (A) and human cirrhotic liver (B) data sets. Shown are the distributions of log fold-changes (y-axis) for DA (FDR 10%) neighbourhoods using different values of K (x-axis) from 5-100, illustrating that DA testing is robust across a broad range of values of K.

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