

# FLASH-MM: fast and scalable single-cell differential expression analysis using linear mixed-effects models

---

Received: 3 April 2025

---

Accepted: 16 January 2026

---

Cite this article as: Xu, C., Pouyabahar, D., Voisin, V. *et al.* FLASH-MM: fast and scalable single-cell differential expression analysis using linear mixed-effects models. *Nat Commun* (2026). <https://doi.org/10.1038/s41467-026-69063-2>

Changjiang Xu, Delaram Pouyabahar, Veronique Voisin, Hamed Heydari & Gary D. Bader

---

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

## FLASH-MM: fast and scalable single-cell differential expression analysis using linear mixed-effects models

Changjiang Xu\* (1), Delaram Pouyabahar\* (1,2), Veronique Voisin (3), Hamed Heydari (1,2), Gary D. Bader (1,2,3,4,5,6,7)

1 - The Donnelly Centre, University of Toronto, Toronto, Ontario, Canada

2 - Department of Molecular Genetics, University of Toronto, Toronto, Ontario, Canada

3 - Princess Margaret Research Institute, University Health Network, Toronto, Ontario, Canada

4 - Department of Computer Science, University of Toronto, Toronto, Ontario, Canada

5 - Lunenfeld-Tanenbaum Research Institute, Toronto, Ontario, Canada

6 - 6 - CIFAR Multiscale Human Program, CIFAR, Toronto, Ontario, Canada

7 - Corresponding author

\*Equal contribution

Corresponding author's email: gary.bader@utoronto.ca

### Abstract

Single-cell RNA sequencing (scRNA-seq) enables detailed comparisons of gene expression across cells and conditions. Single-cell differential expression analysis faces challenges like sample correlation, individual variation, and scalability. We develop a fast and scalable linear mixed-effects model (LMM) estimation algorithm, FLASH-MM, to address these issues. We reformulate aspects of the linear mixed model estimation procedure to make it faster, by reducing computational complexity and memory usage. Simulation studies with scRNA-seq data show that FLASH-MM is accurate, computationally efficient, effectively controls false positive rates, and maintains high statistical power in differential expression analysis. Tests on tuberculosis immune and kidney single cell data demonstrate FLASH-MM's utility in accelerating single-cell differential expression analysis across diverse biological contexts.

### Introduction

Differential expression (DE) analysis is a cornerstone of transcriptomics research. Single-cell RNA sequencing (scRNA-seq) technology enables researchers to profile the transcriptomes of individual cells, uncovering transcriptional similarities and differences across various biological conditions for specific cell types. Advancements in cost efficiency and throughput have facilitated the generation of large-scale datasets comprising hundreds of subjects and millions of cells, opening new avenues for exploring cellular heterogeneity and dynamics across diverse conditions and samples. Cells from the same individual share common genetic and environmental backgrounds, resulting in a hierarchical structure and statistical dependencies among individual cells in scRNA-seq data<sup>1</sup>. This introduces significant challenges for differential expression analysis in single-cell studies, particularly due to the correlation within cell populations of each subject (intra-subject correlation<sup>1</sup>) and the high variability across cell populations from different subjects (inter-subject variability<sup>2</sup>). Ignoring these correlations and variabilities can inflate false positive rates in statistical tests<sup>1,2</sup>. Furthermore, the large scale of single-cell data, often encompassing hundreds of thousands to millions of cells, adds computational complexity, requiring efficient methods to manage and analyze these vast datasets effectively.

Linear mixed-effects models (LMMs) provide a framework to address the challenges of intra-subject correlation and inter-subject variability in single-cell differential expression analysis by

incorporating fixed effects, which capture systematic differences across experimental conditions, and random effects, which model the correlations within subjects and the variations between subjects<sup>1,3–5</sup>. In multi-subject single-cell studies, cells are nested within subjects, and subjects are often nested within experimental conditions, meaning that cells from the same subject are correlated, and subjects within the same condition share common sources of variation. To account for this hierarchical structure, instead of modeling subjects as fixed effects, mixed models treat them as random effects, efficiently capturing both within-subject correlation and between-subject variability. Many methods and software packages have been developed to fit mixed-effects models<sup>5–10</sup>. The most widely used package is *lme4*<sup>6</sup>, which uses maximum likelihood<sup>11</sup> or restricted maximum likelihood<sup>12</sup> methods for model fitting. However, fitting the LMM is computationally demanding, particularly in large-scale single-cell datasets, where standard implementations struggle with memory usage and runtime constraints<sup>5–10</sup>. As a result, the performance of mixed-effects models for single-cell DE analyses has mostly been examined through simulation studies involving small numbers of subjects and cells or pseudobulk methods<sup>2,13</sup>.

To address the challenges of large-scale scRNA-seq data analysis using linear mixed-effects models, we developed FLASH-MM, a fast and scalable LMM estimation algorithm for single-cell differential expression analysis. By leveraging summary statistics, which are precomputed aggregate representations that capture essential information from the data without storing measurements for each individual cell, and by transferring the computation of high-dimension matrices (number of cells) to a lower dimension (numbers of covariates and random effects) in the model estimation step, FLASH-MM achieves both computational efficiency and significantly lower memory usage. Compared to the standard LMM estimation method *lmer* in the *lme4* package<sup>6</sup>, the FLASH-MM algorithm requires orders of magnitude less compute time and memory use while maintaining accuracy. We verified the accuracy, efficiency, and DE analysis performance of FLASH-MM using simulation studies. We simulated multi-subject multi-cell-type scRNA-seq datasets using real reference data based on a negative binomial (NB) distribution. We further demonstrate the application of FLASH-MM for case-control comparisons in a tuberculosis immune atlas and for cell-type-specific sex comparisons in healthy kidney data. In summary, FLASH-MM accelerates accurate single-cell differential expression analysis across diverse biological contexts, supporting the use of mixed models in large-scale, multi-subject single-cell studies.

## Results

### Overview of FLASH-MM

Single-cell RNA sequencing datasets typically consist of gene expression measurements for thousands of genes (approximately 20,000 in humans) across tens of thousands to millions of cells, often collected from multiple subjects or experimental conditions. A LMM can identify differentially expressed genes, correcting for fixed effects, modeled as covariates such as batch, sex, or treatment conditions, and subjects modeled as random effects.

We developed FLASH-MM to address the computational challenge of fitting the LMM for large-scale scRNA-seq data by efficiently estimating LMM parameters, using maximum likelihood (ML)<sup>11</sup> and restricted maximum likelihood (REML)<sup>12</sup> with a gradient descent approach (see Supplementary Information). Instead of directly processing large cell-level data matrices, our algorithm computes and operates on summary statistics, which are compact data representations,

without storing information for each individual cell. Specifically, FLASH-MM operates the matrix computation by transferring the high-dimension  $n \times n$  matrices (number of cells) to the low-dimension  $p \times p$  and  $q \times q$  matrices (numbers of fixed and random effects). This reformulation substantially reduces computational complexity from  $O(mn^3)$  to  $O(mn(p^2 + q^2))$ , and memory complexity from  $O(mn)$  to  $O(m^2 \max(p, q))$ , where  $m$  represents the number of genes,  $n$  is the number of cells and  $p$  and  $q$  denote the number of fixed and random effects, respectively ( $n > p$  and  $n > q$ ). By precomputing and directly using the summary statistics as inputs, FLASH-MM further reduces the computational complexity to  $O(m(p^3 + q^3))$ , which makes LMM estimation independent of the number of cells and achieves both speed and memory efficiency (see Methods and Supplementary Information).

With usual LMM fitting methods, variance components are constrained to be non-negative. As a result, the asymptotic normality of the maximum likelihood estimation of variance components at the null hypothesis is invalid due to the zero variance components being on the boundary of the parameter space. Also, the asymptotic distribution of the likelihood ratio test (LRT) statistics is a mixture of Chi-squared distributions, making it more difficult to model<sup>14,15</sup>. Our algorithm allows variance component parameters to take negative values such that the zero variance components are no longer on the boundary of the parameter space. Thus, the usual asymptotic properties of maximum likelihood estimation at the null hypothesis remain valid under regularity conditions, which enables the use of t-statistics or z-statistics for hypothesis testing of both fixed effects and variance components, and the LRT statistics asymptotically follow Chi-squared distribution. When the variance component parameter is positive, it suggests that the mixed-effects model is appropriately specified; otherwise, the random effect term may not be needed and should be excluded from the model.

The FLASH-MM workflow for single-cell differential expression analysis is illustrated in Figure 1. Figure 1A shows the input structure of scRNA-seq data, represented as a log-transformed gene-by-cell count matrix, where each row corresponds to a gene's expression profile. Fixed effects can include variables such as log-library size, batch effects, biological conditions of interest, and interactions between conditions and cell types. Random effects can capture variations across subjects and correlations within subjects. Figure 1B outlines the linear mixed-effects model framework, constructed for each gene using design matrices of fixed and random effects, informed by prior knowledge of covariates and the biological question. Figure 1C illustrates the model fitting process, comprising LMM parameter estimation and hypothesis testing. Parameter estimation is performed using a gradient descent algorithm applied to summary statistics. Hypothesis testing evaluates fixed effects and their contrasts using t-statistics (see Methods and Supplementary Information).

### Simulation studies

We validated FLASH-MM's accuracy and computational efficiency by comparing it to the standard linear mixed model method lmer, implemented in the lme4 package<sup>6</sup>, using simulated scRNA-seq data (See Methods and Supplemental Information). We also evaluated the performance of FLASH-MM in single-cell DE analysis through simulations using two key criteria: (1) control of the Type I error rate (false positive rate, FPR) and (2) statistical power (true positive rate, TPR). These

metrics were compared against NEBULA<sup>5</sup>, a generalized linear mixed-effects model (GLMM) designed for DE analysis.

### Simulating scRNA-seq data

Using PBMC 10X droplet-based scRNA-seq data from lupus patients<sup>16</sup> as a reference, we simulated six multi-subject multi-cell-type scRNA-seq datasets with 6,000 genes and sample sizes ranging from 20,000 to 120,000 cells increasing by 20,000 at each step, based on a negative binomial (NB) distribution. Genes were randomly selected from the reference dataset, and cells were simulated from 25 subjects across 12 cell types under two treatment conditions. Treatments, cell types, and subjects were assigned randomly with equal probability. A total of 480 differentially expressed genes were designated, specific to a cell type (Figure S1).

### FLASH-MM has the same accuracy as lmer, but is much faster

We first validated FLASH-MM's accuracy and computational efficiency using simulated data by comparing it to the standard method, lmer, from the lme4 package<sup>6</sup> with its default settings. The linear mixed model (LMM) was fit to the log-transformed counts using FLASH-MM and lmer. Because lmer in the lme4 package doesn't provide p-values for hypothesis tests of coefficients, we computed p-values in this case by refitting the LMM by the lmer using the lmerTest package<sup>17</sup>. The differences in estimated variance components, coefficients, and p-values between FLASH-MM and lmer are shown in Figure 2A. The model parameters (coefficients and variance components) estimated by FLASH-MM and lmer are identical up to the sixth decimal place, demonstrating the high accuracy of the FLASH-MM implementation. FLASH-MM demonstrates substantially greater computational efficiency compared to lmer with default optimizer settings, achieving a speedup of approximately 50-fold to 140-fold as the sample size increases from 20,000 to 120,000 cells (Figure 2B, Table S1, measured on a 2.8 GHz Quad-Core Intel Core i7 processor with 16 GB DDR3 RAM). These results suggest that FLASH-MM achieves computational efficiency, accuracy, and reliable inference in practice.

### FLASH-MM has a similar statistical performance to NEBULA, but is much faster

We compared the performance of FLASH-MM to NEBULA, a state-of-the-art method for differential expression analysis of multi-subject single-cell data based on Negative Binomial and Poisson mixed models. We ran NEBULA on our simulated data using the arguments method='LN' and model='NBLMM', which specify the negative binomial lognormal mixed model- the same model we used to simulate scRNA-seq data above. We computed type I error using the simulated non-differentially expressed (non-DE) and DE genes and compared the power between the two methods at a sample size of 120,000 cells.

FLASH-MM effectively controls Type I error, with p-values remaining within the expected range of a uniform distribution, as shown in the quantile-quantile (QQ) plot of Figure 2C. This expected range is represented by 95% confidence intervals, which indicate the natural variation we would expect if there were no true differences in gene expression. FLASH-MM p-values fall within these confidence intervals, suggesting that the method does not produce an excess of false positives (Figure 2C, Figure S2). FLASH-MM demonstrates a power comparable to NEBULA, with a Receiver Operating Characteristic (ROC) curve achieving an Area Under the Curve (AUC) of 0.97

for both FLASH-MM and NEBULA (Figure 2D, Figure S3). The t-values and p-values calculated by FLASH-MM and NEBULA demonstrate a strong correlation (Figure S4). However, both NEBULA and *lmer* required significantly longer runtimes compared to FLASH-MM (319 and 143 times, respectively, at n=120,000 cells, Table S1).

We performed additional simulation studies to further evaluate hypothesis testing performance, and showed: 1) The FLASH-MM t-test for fixed effects achieves type-I error control comparable to that of the *lmerTest* Satterthwaite approximation (Figure S5); and 2) FLASH-MM's z-test and LRT for variance components maintain proper type-I error control (Figure S6).

### FLASH-MM supports DE analysis of diverse biological scRNA-seq data

#### Kidney scRNA-seq data

We first examined the sex variation within healthy kidney cell types using kidney scRNA-seq data<sup>18</sup> (Figure 3A). The kidney data from 19 subject samples contains 14,175 genes and 27,550 cells consisting of 19 cell types after quality control (Figure 3B). We performed a differential expression analysis using FLASH-MM to identify the DE genes between males and females within each cell type while considering the subjects as a random effect.

Among the various cell populations, connecting tubule (CNT) cells have the highest number of sex-specific differentially expressed genes (200), meeting the criteria of  $FDR < 0.05$  and  $|LogFC| > 0.5$  (Figures 3C). Pathway analysis of these genes highlights distinct enrichments: male CNT cells show enrichment for pathways related to acid secretion, transporter activity, blood pressure regulation, and ion importation, whereas female CNT cells exhibit enrichment for kinase activity and positive regulation of receptor recycling (Figures 3D).

On the Kidney dataset, FLASH-MM required only 1.1 minutes of runtime, whereas *lmer* took 119.4 minutes under identical hardware conditions (2.8 GHz Quad-Core Intel Core i7, 16 GB DDR3 RAM).

#### Tuberculosis (TB) scRNA-seq data

We then applied FLASH-MM to single-cell transcriptomics data from 500K memory T cells from 259 donors in a tuberculosis (TB) progression cohort<sup>19</sup>. After quality control, the large TB dataset contains 11,596 genes and 499,713 cells covering 29 cell states from 46 batches and 259 individual donors. We applied FLASH-MM to identify genes associated with TB status within each cellular state. FLASH-MM identified a varying number of differentially expressed genes associated with TB progression across different cell states (Figure 4A), using a threshold of  $FDR < 0.05$  and a positive effect size (i.e., upregulation in TB samples). The cell types with the highest number of DE genes are activated CD4+ and CD8+ T cell populations, with 1266 and 268 DE genes, respectively (Figures 4A and 4B). To further investigate TB-associated signatures, we identified the top TB-enriched genes within these two cell states (Figure 4C and 4D). Pathway enrichment analysis of these DE genes identifies cell-cycle pathways in activated CD4+ cells, while activated CD8+ cells show enrichment for pathways related to immune response, TCR-mediated T-cell activation, and chemokine signaling (Figure 4E and 4F).

Notably, FLASH-MM completed the 500K T cell dataset analysis in 1.4 hours, compared to 55.6 hours (2 days and 7.6 hours) for *lmer*, measured on a 2.8 GHz Quad-Core Intel Core i7 processor with 16 GB DDR3 RAM. These results demonstrate that FLASH-MM is substantially more computationally efficient and thus a more practical choice for large-scale datasets.

## Discussion

Differential expression analysis to detect changes in gene expression across conditions has long been a fundamental aspect of transcriptomics research. However, single-cell RNA sequencing data introduces unique statistical challenges, such as sample correlation, individual variation, and scalability. We developed FLASH-MM as a fast and scalable linear mixed-effects model (LMM) estimation algorithm to address these issues.

Mixed-effects models are powerful tools in single-cell studies due to their ability to model intra-subject correlations and inter-subject variabilities. Classic LMM estimation methods, like *lmer* in the *lme4* package<sup>6</sup>, face limitations of speed and memory use in the analysis of large-scale scRNA-seq data. These limitations have encouraged researchers to use traditional bulk RNA-seq differential expression analysis methods with pseudobulk counts by summing reads within each cell type for each subject<sup>3,20</sup>. While this simplifies the analysis, it sacrifices the resolution inherent in single-cell data.

To support our simulation studies, we developed a scRNA-seq simulator, implemented in the *simuRNAseq* function in the FLASH-MM software distribution, to generate multi-subject multi-cell-type scRNA-seq data based on a negative binomial distribution (Supplemental information). *simuRNAseq* shares similarities with *muscat*<sup>20</sup> and *GLMsim*<sup>21</sup>. Like *muscat*, the simulator captures key characteristics of real single-cell RNA-seq data by modulating zero-inflation, overdispersion, variance differences, cell-level library size variation, number of clusters or cell populations, and the number of expected differentially expressed genes (Supplemental information). However, both *muscat* and *GLMsim* have some limitations when used for scRNA-seq simulations. *Muscat* estimates the dispersion of the negative binomial distribution using the *edgeR* package, but it relies on only a subset of the reference data. As a result, it cannot scale to large scRNA-seq datasets and captures only partial information from the reference dataset. *GLMsim* estimates the coefficients and dispersion parameters of the NB model for each gene using *glm.nb* from the *MASS* package<sup>22</sup>. While effective, this approach is computationally intensive and limited to generating data of a fixed size that matches the reference data. Our scRNA-seq simulator uses the method-of-moments estimate (MME)<sup>23</sup> to compute dispersion parameters for the NB distribution. This approach is faster, more flexible, and uses the full biological reference dataset. The performance of the scRNA-seq simulator, named *simuRNAseq*, is illustrated in Figure S7.

Constructing design matrices for fixed and random effects is an important step in LMM-based DE analysis, requiring identification of the aspects of the data to be modeled, and a balance between reducing residual variance, avoiding overfitting, and managing collinearity among covariates<sup>24</sup>. A design matrix encodes variables such as sample conditions, batch effects, and cell types, specifying how observations are mapped to model parameters. While dataset integration methods (e.g., *Harmony*<sup>25</sup>, *Seurat*<sup>26</sup>) are often applied prior to modeling, mixed-effect models can directly

account for batch effects by modeling batches as fixed or random effects, depending on the study design and underlying biological question. However, in many scRNA-seq studies, samples or batches are perfectly confounded with experimental conditions. In such cases, including batch as a fixed effect may introduce collinearity and inadvertently remove the biological signal of interest. Mixed-effect models can address this issue by modeling the batch as a random effect.

For single cell transcriptomics data, including library size as a fixed effect helps control p-value inflation and should generally be included in the model to help with normalization. In the analyses of the kidney and tuberculosis data, 98.8% and 99% of genes have a significant covariate of log-library size with Bonferroni correction p-value < 0.05, respectively. Based on this empirical evidence, we recommend considering including the log-library size as a covariate in the mixed-effects model. Other model design decisions should be defined by the user based on the structure of their data and the biological question under study. If samples are modeled as random effects, their number may be high (e.g. hundreds or higher) and the FLASH-MM LMM algorithm may slow down (Figure S8). In such cases, subsampling the data or applying dimensionality reduction techniques, such as PCA, on the random effects can potentially help reduce their number to improve computational efficiency (see Supplemental Information). When the covariates of random effects are highly correlated, the dimensionality reduction technique can substantially reduce the number of random effects. Otherwise, if the covariates of random effects are independent or uncorrelated, the dimensionality reduction technique would not reduce the number of random effects.

FLASH-MM makes mixed models computationally feasible for large single-cell datasets. This scalability enables future benchmark studies to determine the optimal modeling approach in biologically relevant contexts such as subtle perturbations, continuous gradients, and rare cell populations where accounting for cell-level variation is likely important. FLASH-MM's versatile framework can be extended to other data modalities, such as spatial transcriptomics and multiomics. Expanding its application across diverse biological data could provide opportunities for uncovering novel insights and facilitating integrated analyses in a wide range of research contexts.

## Methods

### LMM estimation and inference

Consider the linear mixed-effects model (LMM) as expressed below<sup>27</sup>

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\epsilon} \quad (1)$$

where  $\mathbf{y}$  is a  $n \times 1$  vector of observed response (expression for a gene),  $\mathbf{X}$  is a  $n \times p$  design matrix for fixed effects  $\boldsymbol{\beta}$ ,  $\mathbf{Z}$  is a  $n \times q$  design matrix for random effects  $\mathbf{b}$ , and  $\boldsymbol{\epsilon}$  is a  $n \times 1$  vector of residual errors. The term random effects may be a combination of various random-effect components,  $\mathbf{Z}\mathbf{b} = \mathbf{Z}_1\mathbf{b}_1 + \dots + \mathbf{Z}_K\mathbf{b}_K$ , where  $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_K]$ ,  $\mathbf{b} = [\mathbf{b}_1^T, \dots, \mathbf{b}_K^T]^T$ ,  $K$  is the number of the random-effect components, and  $\mathbf{Z}_k$  is a  $n \times q_k$  design matrix for the  $k$ -th component. The superscript  $T$  denotes a transpose of a vector or matrix. The basic assumptions are as follows: 1) The design matrix  $\mathbf{X}$  is of full rank, satisfying conditions of estimability for the parameters; 2) The random vectors  $\mathbf{b}_k$  and  $\boldsymbol{\epsilon}$  are independent and follow a normal distribution,  $\mathbf{b}_k \sim N(\mathbf{0}, \sigma_k^2 \mathbf{I}_{q_k})$  and  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ . Here  $\sigma_k^2$  and  $\sigma^2$  are unknown parameters, called variance components,  $\mathbf{0}$  is a

vector or matrix of zero elements, and  $\mathbf{I}_n$  is an  $n \times n$  identity matrix. The random effects reflect variations between groups (subjects) and correlations within groups (subjects). Assumption (1) implies  $p < n$ . We also assume  $q_k < n$ . If  $q_k > n$ , we can use principal component analysis (PCA) to obtain an equivalent LMM with the number of random effects less than  $n$  (see Supplementary Information).

Maximum likelihood estimation (MLE) and restricted maximum likelihood (REML) are methods for estimating fixed effects and variance components in LMMs. MLE estimates all parameters of fixed effects and variance components together but can produce biased variance component estimates, whereas REML removes fixed effects from variance estimation, resulting in unbiased estimates. Both methods are asymptotically identical. Estimating variance components using either MLE or REML requires numerical methods, among which the iterative gradient-based methods are the most commonly used.

**Fast and scalable algorithm:** The gradient descent methods usually have a computational complexity of  $O(n^3)$ . We developed the summary statistics-based algorithm, FLASH-MM, to implement the gradient methods for speeding up the LMM estimation and reducing the computer memory usage. Instead of the individual cell-level data:  $\mathbf{X}$ ,  $\mathbf{Z}$  and  $\mathbf{y}$ , FLASH-MM uses the summary statistics:  $\mathbf{X}^T\mathbf{X}$ ,  $\mathbf{X}^T\mathbf{y}$ ,  $\mathbf{Z}^T\mathbf{X}$ ,  $\mathbf{Z}^T\mathbf{y}$  and  $\mathbf{Z}^T\mathbf{Z}$  to estimate the LMM parameters. FLASH-MM achieves a computational complexity of  $O(n(p^2 + q^2))$ , which is fast and linearly scalable with the sample size  $n$ . These summary statistics have a low dimension and require less computer memory. By precomputing and directly using the summary statistics as inputs, the algorithm complexity is reduced to  $O(p^3 + q^3)$ , which makes computations independent of the sample size  $n$  (number of cells) and achieves both speed and memory efficiency (see Supplementary Information section 2.1).

scRNA-seq data typically consist of gene expression measurements for thousands of genes (approximately 20,000 in human) across thousands to millions of cells. In the single-cell differential expression analysis, the complexity of FLASH-MM is  $O(mn(p^2 + q^2))$ ,  $m$  is the number of genes, which is linearly scalable with the number of genes. By using the pre-computed summary statistics as inputs, the algorithm complexity becomes  $O(m(p^3 + q^3))$ .

**Hypothesis testing:** The hypothesis testing for fixed effects and variance components can be respectively defined as:

$$H_{0,i}: \beta_i = 0 \text{ versus } H_{1,i}: \beta_i \neq 0,$$

$$H_{0,k}: \sigma_k^2 = 0 \text{ versus } H_{1,k}: \sigma_k^2 > 0.$$

The variance components under null hypothesis,  $\sigma_k^2 = 0$ , are on the boundary of the parameter space, in which case the MLE asymptotic normality is inappropriate. By reparameterizing the variance components,  $\theta_k = \sigma^2 \gamma_k$ , the covariance matrix,  $\mathbf{V}_\theta = \sigma^2 (\mathbf{I} + \gamma_1 \mathbf{Z}_1 \mathbf{Z}_1^T + \cdots + \gamma_K \mathbf{Z}_K \mathbf{Z}_K^T)$ , becomes positive-definite and well-defined when  $\gamma_k > -1/\lambda_{max}$ , where  $\lambda_{max} > 0$ , is the largest singular value of  $\mathbf{Z}\mathbf{Z}^T$  or  $\mathbf{Z}^T\mathbf{Z}$ . Now the parameters of variance components,  $\theta_k$ ,  $k = 1, \dots, K$ , can be negative. If  $\theta_k > 0$ ,  $\sigma_k^2 = \theta_k$  is definable and the mixed model is well-specified. Otherwise, the

term of random effects is not needed in the model. Then the hypotheses for the variance components are extended as:

$$H_{0,k}: \theta_k \leq 0 \text{ versus } H_{1,k}: \theta_k > 0,$$

in which the zero components,  $\theta_k = 0$ , are no longer on the boundary of the parameter space and the MLE normal asymptotic properties hold. Then we can use z-statistic or t-statistic for hypothesis testing of fixed effects and variance components. We can also test whether there are no random effects, that is, the variance components are equal to zero, using likelihood ratio test (LRT) statistic. See Supplementary Information section 1.2 for details.

### Simulation methods

We generated the multi-subject multi-cell type scRNA-seq dataset by using reference data based on a negative binomial (NB) distribution. The mean of NB distribution is taken as the sample mean for each gene. The dispersion of NB distribution is computed by the method-of-moments estimate (MME)<sup>23</sup>. Compared to the maximum likelihood-based estimates, such as the `glm.nb` function in the MASS package <sup>22</sup>, the MME is computationally simpler and performs reasonably well. See detailed simulation methods and performance in Supplementary Information.

### Data preprocessing and model design

In general, we selected thresholds to be in line with typical single cell genomics projects<sup>28</sup>, as well as to ensure there is enough data per cell and gene for the LMM to operate on. We processed the data by filtering the outliers based on: total numbers of UMI counts per cell and numbers of detected genes (cell filtering), numbers of UMI counts per gene, and UMI counts per cell (cpc) (gene filtering), which are standard in the literature<sup>29</sup>, and described below for each dataset. Our goal was not to optimize thresholds for biological discovery, but rather to apply reasonable filters that balance noise reduction with data retention.

### Healthy human kidney atlas

The healthy human kidney transcriptomic map<sup>18</sup> was generated from 27,677 cells obtained from 19 living donors (10 female and 9 male). Decontaminated raw data, processed using `SoupX`<sup>30</sup>, was provided upon request by the authors. Cells were filtered based on three criteria: the number of detected features (nFeature), library size, and the number of cells within each cell type (cluster).

The number of features per cell, defined as the total number of non-zero genes, was required to meet a minimum threshold of 100. Library size, calculated as the total counts per cell, was restricted to a range of  $2^9$  (512) to  $2^{16}$  (65,536). Cell types with fewer than 20 cells were excluded, and “Podocyte” cells were removed due to their low sample size (16 cells post filtering). Additionally, genes were filtered based on their expression levels, where the counts per cell ratio had to exceed 0.005 (i.e., the total gene count divided by the number of cells had to be greater than 0.5%). Genes were further filtered to retain those expressed in at least 16 cells, with a minimum of 10 cells in each group, total counts between  $2^6$  (64) and  $2^{20}$  (1,048,576), and a counts per cell ratio above the threshold. After these filtering steps, 27,550 cells and a refined set of genes were retained for downstream analyses.

The FLASH-MM model was designed to account for both technical and biological variation. The model formula used was:

$$\sim \log(\text{library.size}) + \text{Cell\_Types\_Broad} + \text{Cell\_Types\_Broad:sex} + (1|\text{sampleID}),$$

where log-transformed library size was included as a fixed effect to normalize for differences in sequencing depth, Cell\_Types\_Broad captured the cell type-specific effects, and the interaction term Cell\_Types\_Broad:sex identified sex-specific differences within each cell type. A random effect (1|sampleID) was added to account for inter-sample variability.

### Tuberculosis (TB) T cell scRNA-seq data

The Tuberculosis (TB) memory T cell dataset was obtained from Nathan et al., 2021<sup>19</sup>, comprising 500,089 cells. The raw count matrix was used for preprocessing. Metadata associated with the cells includes donor identity, sex, batch, cluster annotations, and TB status. The dataset includes cells from 259 unique donors, spanning 46 batches and 29 cell clusters. In the DE analysis, we modeled the donors as a random effect and ignored the batch effect because the majority of donors were sequenced in a single batch.

Pre-processing steps consisted of removing cells and genes to remove extreme values. Library sizes were assessed using a boxplot, and cells with library sizes outside the lower whisker and above an upper threshold of  $2^{15}$  (32,768) were removed, resulting in the retention of 499,973 cells. Genes were filtered in two steps: first, genes expressed in fewer than  $2^9$  cells (512 cells) were excluded; second, genes with a counts-per-cell ratio less than 0.005 (i.e., total gene count divided by the number of cells below 0.5%) were removed. These filtering steps reduced the dataset to 11,596 genes and 499,973 cells, which were used for downstream analyses.

The FLASH-MM model was designed to analyze this dataset, accounting for both technical and biological variations. The model formula used was:

$$\sim \log(\text{library.size}) + \text{cluster\_name} + \text{cluster\_name:TB\_status} + (1|\text{donor}),$$

where log-transformed library size normalized for differences in sequencing depth, cluster\_name captured the cell type-specific effects, and the interaction term cluster\_name: TB\_status identified TB-associated differences within each cell type. A random effect (1|donor) was added to account for inter-donor variability.

Pathway enrichment analysis for both the kidney and TB datasets was performed using *gprofiler2*<sup>31,32</sup> (v0.2.3), and data visualizations were generated using *ggplot2* (v3.5.1).

### Data Availability

The healthy human kidney atlas<sup>18</sup> data files were downloaded from the UCSC Cell Browser at <https://cells.ucsc.edu/?ds=living-donor-kidney>. The Tuberculosis (TB) memory T cell dataset<sup>19</sup> can be accessed from the GEO with accession code GSE158769 [<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE158769>]. The stimulated PBMC data<sup>16</sup> was downloaded from muscData package (Kang18\_8vs8) at <https://github.com/HelenaLC/muscData>. Source data is provided as a Source Data file.

## Code Availability

The FLASH-MM software is openly available in both R and Python implementations, with example case studies, through the following repositories:

R package (CRAN): <https://cran.r-project.org/web/packages/FLASHMM/index.html>

GitHub: <https://github.com/BaderLab/FLASHMM>

Python package (PyPI): <https://pypi.org/project/FLASH-MM/>

The package is distributed under the MIT License. Analysis scripts used for the case studies and data simulations are available at: <https://github.com/BaderLab/FLASH-MM-analysis/>

The code is archived and citable via Zenodo: <https://doi.org/10.5281/zenodo.18222187>

## References

1. Zimmerman, K. D., Espeland, M. A. & Langefeld, C. D. A practical solution to pseudoreplication bias in single-cell studies. *Nat. Commun.* **12**, 738 (2021).
2. Squair, J. W. *et al.* Confronting false discoveries in single-cell differential expression. *Nat. Commun.* **12**, 5692 (2021).
3. Hoffman, G. E. & Roussos, P. Dream: powerful differential expression analysis for repeated measures designs. *Bioinformatics* **37**, 192–201 (2021).
4. Gagnon, J. *et al.* Recommendations of scRNA-seq Differential Gene Expression Analysis Based on Comprehensive Benchmarking. *Life (Basel)* **12**, (2022).
5. He, L. *et al.* NEBULA is a fast negative binomial mixed model for differential or co-expression analysis of large-scale multi-subject single-cell data. *Commun. Biol.* **4**, 629 (2021).
6. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting linear mixed-effects models using lme4. *J. Stat. Softw.* **67**, 1–48 (2015).
7. Finak, G. *et al.* MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol.* **16**, 278 (2015).
8. Brooks, M., E. *et al.* glmmTMB Balances Speed and Flexibility Among Packages for Zero-inflated Generalized Linear Mixed Modeling. *R J.* **9**, 378 (2017).
9. Bürkner, P.-C. brms: an R package for bayesian multilevel models using stan. *J. Stat. Softw.* **80**, (2017).
10. Pinheiro, J. C. & Bates, D. M. *Mixed-Effects Models in S and S-PLUS*. (Springer-Verlag, 2000). doi:10.1007/b98882.
11. Hartley, H. O. & Rao, J. N. K. Maximum-likelihood estimation for the mixed analysis of variance model. *Biometrika* **54**, 93–108 (1967).
12. Patterson, H. D. & Thompson, R. Recovery of inter-block information when block sizes are unequal. *Biometrika* **58**, 545–554 (1971).
13. Nguyen, H. C. T., Baik, B., Yoon, S., Park, T. & Nam, D. Benchmarking integration of single-cell differential expression. *Nat. Commun.* **14**, 1570 (2023).
14. Self, S. G. & Liang, K.-Y. Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *J. Am. Stat. Assoc.* **82**, 605 (1987).
15. Stram, D. O. & Lee, J. W. Variance components testing in the longitudinal mixed effects model. *Biometrics* **50**, 1171–1177 (1994).
16. Kang, H. M. *et al.* Multiplexed droplet single-cell RNA-sequencing using natural genetic

- variation. *Nat. Biotechnol.* **36**, 89–94 (2018).
17. Kuznetsova, A., Brockhoff, P. B. & Christensen, R. H. B. Imertest package: tests in linear mixed effects models. *J. Stat. Softw.* **82**, 1–26 (2017).
  18. McEvoy, C. M. *et al.* Single-cell profiling of healthy human kidney reveals features of sex-based transcriptional programs and tissue-specific immunity. *Nat. Commun.* **13**, 7634 (2022).
  19. Nathan, A. *et al.* Multimodally profiling memory T cells from a tuberculosis cohort identifies cell state associations with demographics, environment and disease. *Nat. Immunol.* **22**, 781–793 (2021).
  20. Crowell, H. L. *et al.* muscat detects subpopulation-specific state transitions from multi-sample multi-condition single-cell transcriptomics data. *Nat. Commun.* **11**, 6077 (2020).
  21. Wang, J., Chen, L., Thijssen, R., Phipson, B. & Speed, T. P. GLMsim: a GLM-based single cell RNA-seq simulator incorporating batch and biological effects. *BioRxiv* (2024) doi:10.1101/2024.03.20.586030.
  22. Venables, W. N. & Ripley, B. D. *Modern applied statistics with S* 4th ed. (Springer, New York, 2002).
  23. Clark, S. J. & Perry, J. N. Estimation of the Negative Binomial Parameter  $\kappa$  by Maximum Quasi-Likelihood. *Biometrics* **45**, 309 (1989).
  24. Gelman, A. & Hill, J. *Data Analysis Using Regression and Multilevel/Hierarchical Models*. (Cambridge University Press, 2006). doi:10.1017/CBO9780511790942.
  25. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat. Methods* **16**, 1289–1296 (2019).
  26. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888–1902.e21 (2019).
  27. Searle, S. R., Casella, G. & McCulloch, C. E. *Variance Components*. 536 (Wiley-Interscience, 2006).
  28. Heumos, L. *et al.* Best practices for single-cell analysis across modalities. *Nat. Rev. Genet.* **24**, 550–572 (2023).
  29. Luecken, M. D. & Theis, F. J. Current best practices in single-cell RNA-seq analysis: a tutorial. *Mol. Syst. Biol.* **15**, e8746 (2019).
  30. Young, M. D. & Behjati, S. SoupX removes ambient RNA contamination from droplet-based single-cell RNA sequencing data. *Gigascience* **9**, giaa151. (2020).
  31. Reimand, J., Kull, M., Peterson, H., Hansen, J. & Vilo, J. g:Profiler--a web-based toolset for functional profiling of gene lists from large-scale experiments. *Nucleic Acids Res.* **35**, W193–200 (2007).
  32. Reimand, J. *et al.* g:Profiler-a web server for functional interpretation of gene lists (2016 update). *Nucleic Acids Res.* **44**, W83–9 (2016).

### Acknowledgments

This research was supported by the Canadian Institutes for Health Research (grant PJT 469829 to GDB) and by the University of Toronto's Data Sciences Institute Doctoral Student fellowship program.

### Author Contributions Statement

Conceptualization, D.P., V.V, C.X and G.D.B.; Methodology, C.X, H.H; Formal analysis, D.P., C.X; Writing – Original Draft, D.P., C.X; Writing – Review & Editing, D.P., C.X, V.V and G.D.B.; Supervision, G.D.B.; Funding Acquisition, G.D.B.

### Competing Interests Statement

GDB is on the Scientific Advisory Boards of Adela Bio and BioRender. No other competing interests are declared.

### Main Figure Legends

**Figure 1. FLASH-MM workflow for single-cell differential expression analysis.** A. Data: gene expression matrix  $\mathbf{Y} = \log(1 + \text{counts})$ , with each row corresponding to a gene's expression profile and each column corresponding to a cell (gene-by-cell matrix). Metadata includes various variables such as log-library size, batch effects, biological conditions of interest, and interactions between conditions and cell types, which could be modeled as fixed effects, and individual subjects, which could be modeled as random effects. B. Model: the linear mixed-effects model (LMM) for each gene by design matrices  $\mathbf{X}$  and  $\mathbf{Z}$ , which are constructed based on prior knowledge about the covariates and the biological question. C. Model fitting: comprises LMM estimation and tests. LMM estimation is implemented by a gradient descent algorithm over summary statistics. The summary statistics are computed as  $\mathbf{X}^T \mathbf{X}$ ,  $\mathbf{X}^T \mathbf{Y}^T$ ,  $\mathbf{Z}^T \mathbf{X}$ ,  $\mathbf{Z}^T \mathbf{Y}^T$ , and  $\mathbf{Z}^T \mathbf{Z}$ . LMM tests perform hypothesis tests on the fixed effects by t-test and variance components by likelihood ratio test (LRT).

**Figure 2. Computational and statistical performance of FLASH-MM in differential expression analysis of simulated scRNA-seq data.** A) Boxplots of differences of variance components, coefficients, and  $-\log_{10}(p\text{-values})$  between FLASH-MM and lmer fitting for each of 6,000 genes across the six simulated scRNA-seq datasets. The boxplots contain 72,000 ( $=2*6000*6$ ) values for variance components and 900,000 ( $=25*6000*6$ ) values for coefficients and p-values. B) Computation time (in minutes) for FLASH-MM, lmer, and NEBULA across the six datasets with sample sizes from 20,000 to 120,000. C) QQ-plots of non-DE genes (negative controls) p-values for FLASH-MM and NEBULA. The grey area represents the 95% confidence interval, indicating the expected range under the null hypothesis. D) ROC curves for FLASH-MM and NEBULA. Source data are provided as a Source Data file.

**Figure 3. FLASH-MM identifies sex-specific variations in a healthy human kidney map.** A) UMAP projection of the healthy human kidney transcriptomic atlas, highlighting connecting tubule (CNT) cells in purple. Other cell types are shown in lighter shades for contrast. B) Bar plot showing the proportion of male and female cells within each annotated kidney cell type. C) Volcano plot of differential expression between sexes within the CNT cell population. Selected genes with significant male or female bias are labeled. Each point represents a gene, with the x-axis showing the sex-specific log fold change (logFC) and the y-axis showing  $-\log_{10}(p\text{-value})$ . Genes significantly upregulated in males or females (adjusted p-value  $< 0.05$  and  $|\log FC| > 0.5$ ) are colored in blue and pink, respectively. A subset of the top male-biased and female-biased genes (ranked by effect size) are labeled to minimize visual overlap. Differential expression between male and female CNT cells was tested gene-wise using a linear mixed-effects model of the form

expression ~ log(library.size) + Cell\_Types\_Broad + Cell\_Types\_Broad:sex + (1|sampleID); P values correspond to two-sided t-tests on the Cell\_Types\_Broad:sex interaction term and were adjusted for multiple testing across genes within the cell type using the Benjamini–Hochberg false discovery rate. D) Pathway enrichment results for male-biased and female-biased genes within CNT cells. Dot size reflects gene set size, and x-axis position indicates significance,  $-\log_{10}(p\text{-value})$ . Pathway enrichment was performed separately for male-biased and female-biased gene sets using g:Profiler. Source data are provided as a Source Data file.

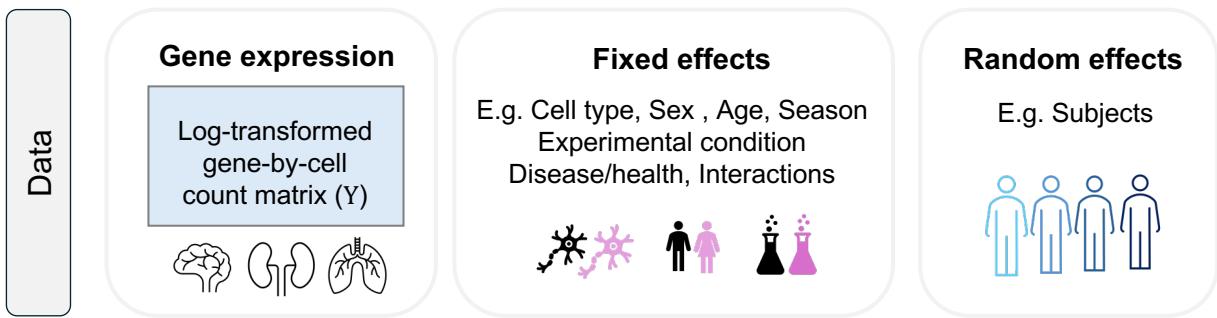
**Figure 4. FLASH-MM identifies TB-enriched signatures within T cell populations while accounting for confounding variables.** A) Bar plot showing the number of TB-associated differentially expressed genes (FDR < 0.05 and positive effect size) identified within each T cell subtype. B) UMAP projection of single-cell RNA-seq data from ~500K memory T cells across 259 individuals in a TB progression cohort. CD4<sup>+</sup> activated and CD8<sup>+</sup> activated T cells are highlighted in blue and dark orange, respectively; all other cell types are shown in lighter tones for context. C–D) Volcano plots of differential expression in CD8<sup>+</sup> activated (C) and CD4<sup>+</sup> activated (D) T cells. Each point represents a gene, with the x-axis showing the log fold change (logFC) and the y-axis showing  $-\log_{10}(p\text{-value})$ . Genes significantly upregulated in TB samples are colored (FDR < 0.05 and logFC > 0.1). Vertical and horizontal dashed lines indicate the logFC and p-value thresholds. A subset of the top 10 TB-upregulated genes, ranked by effect size, are labeled. Differential expression between TB and control samples in each T cell subtype was tested gene-wise using a linear mixed-effects model of the form expression ~ log(library.size) + cluster\_name + cluster\_name:TB\_status + (1|donor); P values correspond to two-sided t-tests on the cluster\_name:TB\_status interaction term and were adjusted for multiple testing across genes within each cell type using the Benjamini–Hochberg false discovery rate. E) and F) Dot plots showing enriched pathways among TB-upregulated genes in CD8<sup>+</sup> (E) and CD4<sup>+</sup> (F) activated T cells. x-axis position reflects enrichment significance,  $-\log_{10}(p\text{-value})$ . Pathway enrichment in panels E and F was performed separately for TB-upregulated gene sets using g:Profiler. Source data are provided as a Source Data file.

### Editor's Summary

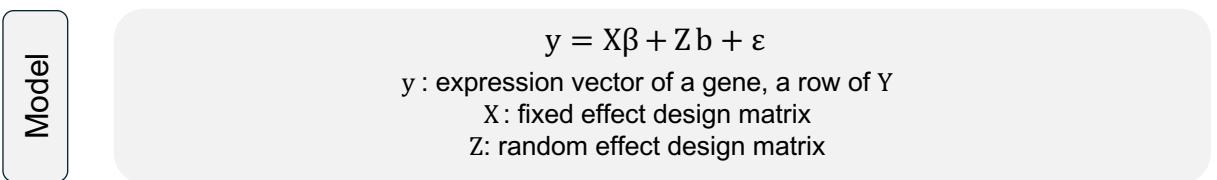
Detecting gene expression changes in single-cell data while accounting for sample structure is vital but computationally demanding. FLASH-MM is a scalable, memory efficient, and statistically robust method that can quickly compute cell-level differential expression across diverse biological contexts.

**Peer Review Information:** *Nature Communications* thanks Matthew Hirschey and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

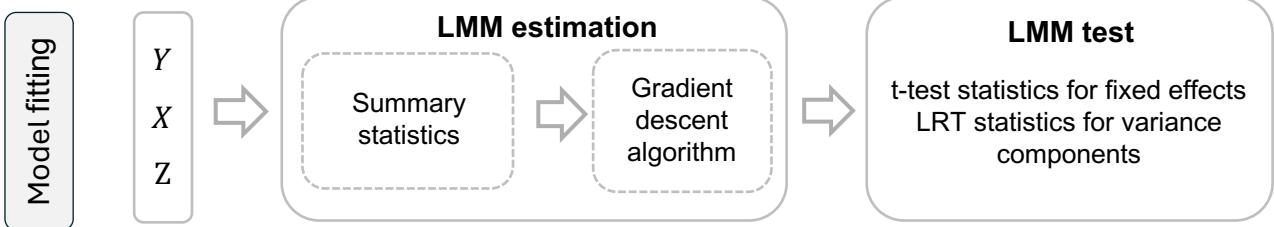
A



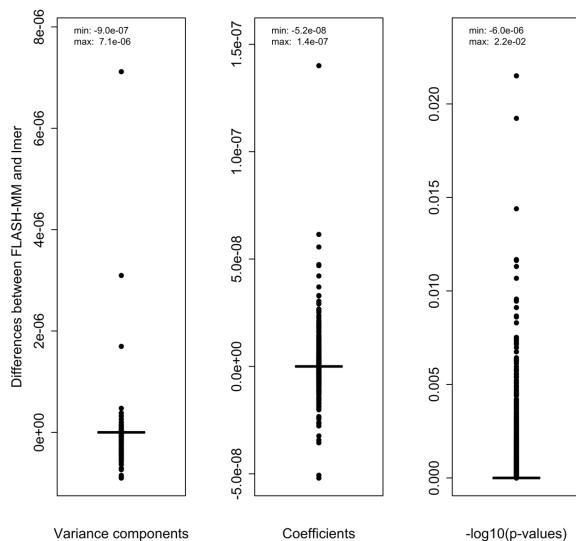
B



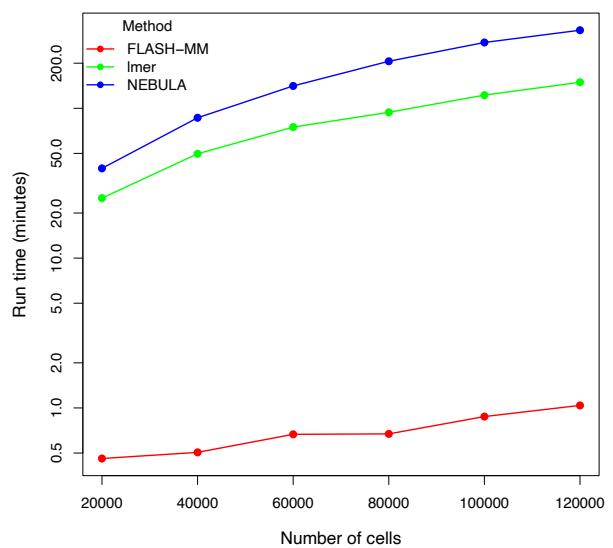
C



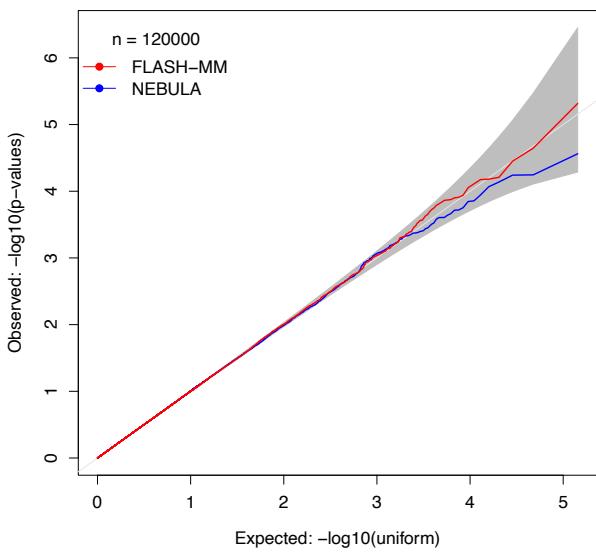
A



B



C



D

