

1 *The curses of performing differential expression analysis using single-cell
2 data*

3
4 Chih-Hsuan Wu, Xiang Zhou, Mengjie Chen

5 **Abstract**

6
7 Differential expression analysis is pivotal in single-cell transcriptomics for unraveling cell-type–
8 specific responses to stimuli. While numerous methods are available to identify differentially
9 expressed genes in single-cell data, recent evaluations of both single-cell–specific methods and
10 methods adapted from bulk studies have revealed significant shortcomings in performance. In
11 this paper, we dissect the four major challenges in single-cell DE analysis: normalization,
12 excessive zeros, donor effects, and cumulative biases. These "curses" underscore the limitations
13 and conceptual pitfalls in existing workflows. In response, we introduce a novel paradigm
14 addressing several of these issues.
15

16
17
18
19 Differential expression (DE) analysis in single-cell transcriptomics provides essential insights into
20 cell-type–specific responses to internal and external stimuli^{1–4}. While many methods are available
21 to identify differentially expressed genes from single-cell transcriptomics, recent studies raise
22 important concerns about the performance of state-of-the-art methods, including both methods
23 tailored to single cell data and techniques that work well in bulk^{5–7}. As population-level single-cell
24 studies rapidly become more feasible, powerful and accurate analytical methods will be essential
25 for obtaining meaningful results. In this context, we discuss the four "curses" that currently plague
26 the differential expression analysis of single-cell data: normalization, zeros, donor effects, and
27 cumulative biases, highlighting the various limitations and conceptual flaws in the current
28 workflows. We demonstrate these limitations using real data from 10X single-cell RNA-seq
29 (scRNA-seq) data from post-menopausal fallopian tubes⁸. Finally, we present a new paradigm
30 that offers a potential solution to some of these issues and illustrate its performance using two
31 case studies.
32

33
34 *The curse of normalization*

35
36 The term 'normalization' has been used to denote multiple distinct approaches in genomics^{9, 10}.
37 For example, it can refer to the process of correcting PCR amplification biases introduced during
38 sequencing library preparation (library size normalization)^{11–13}, the process of harmonizing data
39 across different experimental batches (batch normalization)^{14–18}, or to the process of transforming
40 the data to adhere to a normal distribution (data distribution normalization)¹⁹. All three have been
41 introduced to handle both bulk and single cell RNA-seq data, aiming to minimize unwanted
42 technical variations. Choosing appropriate normalization techniques for DE analysis of scRNA-
43 seq data is clearly important to maintain the integrity of the data, but the field has yet to establish
44 a definitive gold standard outlining the circumstances for which different normalizations should be
45 performed.
46

47 Library size normalization is critical in bulk RNA-seq analysis, as it is impossible to track the
48 absolute abundance of RNA molecules in typical bulk RNA-seq protocols due to an unknown fold
49 of amplification introduced by PCR during library construction. Normalization, in this instance,

50 focuses on estimating and subsequently correcting for a sample-specific size factor. This process
51 allows bulk RNA-seq to estimate relative RNA abundances. Post-normalization, samples are
52 calibrated against a common reference, resulting in most genes displaying similar expression
53 levels across samples. When performing differential expression analysis with bulk RNA-seq data,
54 genes are classified as either up-regulated or down-regulated, based on the assumption that the
55 majority remain unchanged across groups. While this size-factor based normalization technique
56 is suitable for bulk RNA-seq, it does not translate effectively to scRNA-seq. Protocols in scRNA-
57 seq, such as the 10X, employ unique molecular identifiers (UMIs) which discern between genuine
58 RNA molecules and those generated via PCR. This enables the absolute quantification of RNA
59 levels. Unfortunately, size-factor-based normalization methods, like counts per million reads
60 mapped (CPM) convert data into relative abundances erasing useful data provided by the UMIs.
61 Furthermore, CPM-normalized data does not account for competition among genes for cellular
62 resources because the uniform number of molecules found in CPM-normalized data does not
63 accurately represent true expression levels, which ultimately leads to suboptimal DE analysis
64 results.

65
66 In batch effect normalization, dimension reduction methods pinpoint genes with consistent
67 expression patterns across various batches; these genes act as anchors, guiding the alignment
68 and integration of data²⁰. However, in scRNA-seq analysis, only highly expressed or highly
69 variable genes are retained for estimating batch effects and subsequent integration. As a result,
70 gene numbers in integrated scRNA-seq datasets are noticeably reduced compared to the raw
71 UMI data.

72
73 For data distribution normalization, the field offers both straightforward (e.g., log-transformation)
74 and advanced strategies (e.g., variance stabilizing transformation, or VST). A notable
75 implementation for scRNA-seq of VST is sctransform²¹, which employs a regularized negative
76 binomial regression model, preserving the Pearson residuals for future analytical steps, including
77 DE analysis²². However, if the underlying data distribution deviates significantly from the
78 assumed model, the application of VST may introduce bias into the analysis.

79
80 To demonstrate the effects of various normalization methods on single-cell data, we compared
81 the raw UMI counts of 10x scRNA-seq data obtained from post-menopausal fallopian tubes (see
82 Methods) with data normalized using one of three methods: 1) CPM; 2) integrated log-normalized
83 counts after removing batch effects using the Seurat CCA model²³; and 3) VST data using
84 sctransform²¹. As a result, we see the total UMI counts revealed substantial variations in library
85 sizes across different cell types; notably macrophages (MP) and secretory epithelial (SE) cells
86 exhibited significantly higher RNA content than other cell types (Fig. 1a). Furthermore, SE cells
87 exhibited larger mean library sizes than mast (MA) cells across all donors. These findings align
88 with the understanding that the main active cell types in post-menopausal fallopian tubes are MP
89 and SE cells, with other cell types remaining dormant post-menopause. However, in the
90 integrated data, the disparities in library size distribution were mitigated, even within cell types
91 (Fig. 1a). While integration reduced differences across donors, it came at the cost of diminishing
92 variation across cell types. It is worth mentioning that CPM normalization equalizes library sizes
93 across all cell types; such normalizations may potentially obscure differences between cell types
94 that are vital for understanding their unique biological functions.

95

96

97 *The curse of zeros*

98

99 Bulk RNA-seq provides the average transcriptional output of each gene expressed within a
100 population of heterogenous cell types^{24, 25}. Even a moderate sequencing depth can yield
101 information about many thousands of different genes. In comparison, scRNA-seq data is much
102 sparser in comparison, with fewer genes expressed per sample and a high proportion of genes
103 with zero UMI counts. Zeros in UMI counts for a gene can arise from three scenarios: a genuine
104 zero, indicating that the gene is not expressed, or a sampled zero, indicating that the gene is
105 expressed at a low level, or a technical zero, indicating that the gene is expressed at a high level
106 but not captured by the assay. Despite an increasing body of evidence suggesting that cell-type
107 heterogeneity is the major driver of zeros observed in 10X UMI data²⁶⁻²⁸, the prevailing notion
108 within the single-cell community is that zeros are largely uninformative technical artifacts caused
109 by “drop-out” genes (i.e., technical zeros).
110

111 Accordingly, many single-cell DE studies include pre-processing steps aimed at removing so-
112 called zero inflation. Several popular pre-processing methods include: 1) performing feature
113 selection by aggressively removing genes based on their zero detection rates, such as requiring
114 non-zero values in at least 10% of total cells and restricting DE analysis to a smaller gene set; 2)
115 imputing zeros and performing DE on imputed values²⁹⁻³²; or 3) modeling zeros explicitly as an
116 extra component and essentially performing DE on non-zero values only^{33, 34}.
117

118 However, if zeros are in fact biological zeros due to no expression or very low expression,
119 dismissing or correcting for zeros in scRNA-seq is equivalent to discarding a significant portion of
120 information in the dataset before any analysis. By failing to account for cell-type heterogeneity,
121 zero-inflation pre-processing steps such as normalization and imputation become inappropriate
122 and can introduce unwanted noise into downstream analyses, including DE. Ironically, the most
123 desired markers in single-cell DE analysis—e.g., genes that are exclusively expressed in a rare
124 cell type that accounts for less than 5% of the total population—may be obscured by current pre-
125 processing steps for handling zeros.
126

127 In the fallopian tube dataset, we observed that distinct cell types display varied gene expression
128 patterns in UMI counts. However, these differences become less apparent in imputed or certain
129 transformed datasets (Fig. 1a). Gene expression frequency differs among cell types (Fig. 1b).
130 However, normalization processes can substantially alter the distribution of both non-zero UMI
131 (Fig 1c) and zero UMI counts (Fig. 1d) counts. For example, while the frequency of genes
132 exponentially decline as raw UMI counts increase, VST data forms a more bell-shaped curve with
133 a mode around 1.5 for non-zero raw UMI counts. Non-zero CPM-normalized data, (scaled by
134 1000) peaks near 0.2 and is more right-skewed than the VST data. Following batch integration,
135 UMI counts primarily fall below 5 and are not as strongly right-skewed. It is noteworthy that zero
136 UMI counts can be given non-zero values via normalization (except with CPM normalization); for
137 example, zeros in VST data are adjusted to negative values and are left-skewed (Fig. 1d).
138 Conversely, the integration process transforms original zeros to values clustered closely around
139 zero. We further examined the distributions of gene expression from one gene. Using the gene
140 RUNX3 as an example (Fig. 1e), the distributions in raw UMI counts and CPM data remain right-
141 skewed. In contrast, the VST and integrated data showcase broader, bell-shaped distributions.
142 The handling of zeros in these latter datasets (VST and integrated) intrinsically sets them apart
143 from the former. This variability, combined with shifts in distribution skewness, may raise concerns
144 when performing DE analysis with normalized values.
145

146 *The curse of donor effects*

147
148 Recent reviews have highlighted that many single-cell DE analysis methods are susceptible to
149 generating false discoveries⁵. This is mainly due to failing to account for variations between

150 biological replicates, commonly referred to as "donor effects". In single-cell studies, donor effects
151 are always confounded with batch effects since cells from one biological sample are typically
152 processed in the same experimental batch. While single-cell studies that contain multiple samples
153 will perform batch correction as pre-processing, they usually do not correct for donor effects when
154 performing DE tests in the downstream analysis.
155

156 One question that arises is whether batch effect correction alone suffices to eliminate donor-
157 related effects. To address this, we investigated the contributions of variations from different
158 sources before and after batch correction. Using the same fallopian tube dataset, we further
159 separated 4553 T/NK cells into 20 subtypes using HIPPO³⁵ (Fig. 2a, S1). With the aid of canonical
160 markers, we identified specific subtypes, including NK, CD4+ T, CD8+ T and mature naive T cells.
161 We then focused on subtypes that were observed in all donors (Fig. 2bc).
162

163 To quantify the proportion of variation originating from different sources, we fit a linear model,
164 using cell types and donors as covariates, for each gene in several subtype pairs. Through all
165 pairs, the integration led to a reduction in donor variation (Fig. 2d, S2). However, in comparisons
166 of two subtypes of the same cell type (12 vs. 13) and two subtypes of different cell types (13 vs.
167 19), we observed a decrease in the proportion of cell-type-related variation. This underscores
168 that integration not only mitigates batch effects but also impacts the phenotypes of interest.
169 Importantly, our analysis indicated that even after implementing batch correction, a notable
170 percentage of genes still exhibited donor-related effects (Fig. 2e). As batch effects are often
171 estimated from leading principal components, representing a consensus from a subset of genes,
172 it is quite possible that residual donor effects persist on some, if not all, genes. Therefore, it is
173 crucial to account for donor effects when performing DE tests to avoid false discoveries and obtain
174 accurate results, even after removing batch effects.
175

176 One popular solution to address the issue of donor effects in single-cell studies is the use of
177 pseudo-bulk analysis. This approach involves merging cells from the same donor and treating the
178 resulting data as bulk RNA-seq. DE analysis is then performed using tools such as DESeq2³⁶ or
179 edgeR³⁷. However, this framework ignores within-sample heterogeneity by treating donor effects
180 as a fixed effect and assumes that each cell from the same donor is equally affected. As a result,
181 this type of analysis can be overly conservative and potentially lead to missed discoveries⁵.
182 Moreover, bulk RNA-seq DE tools typically perform normalization by default, which may have the
183 same drawbacks mentioned earlier in the context of single-cell studies. Thus, caution is advised
184 when using pseudo-bulk analysis as it may not always provide an accurate solution to the problem
185 of donor effects in single-cell studies.
186

187 *The curse of cumulative biases*

188 In scRNA-seq analysis, it is common to follow a hierarchical, sequential workflow for clustering
189 and DE analysis. This approach can carry forward biases from one step to the next, from batch
190 correction through to normalization, imputation, and feature selection. Such cumulative biases
191 can ultimately diminish the power to detect differentially expressed genes.
192

193 Unsupervised learning, especially clustering analysis, is essential in single-cell studies. It groups
194 cells based on gene expression patterns, facilitating the cell-type annotation. While clustering is
195 effective with normalized values like CPMs, it essentially reweights gene features based on their
196 relative contributions. As a result, clustering provides a generalized perspective of variation in
197 gene expression across cell types. The reliance on relative expression also makes clustering fairly
198 resilient to errors and biases introduced by the pre-processing steps.
199

200

201 On the other hand, DE analysis operates at the gene level, using group labels from the clustering
202 process. The effects of biases, whether from donors or batch processing, can vary for each gene.
203 Although DE analysis technically follows clustering—given its reliance on group labels—the
204 metrics used do not need to be identical for both. As we show later in the case studies with data
205 that complete clustering and annotation successfully, if DE analysis is performed using processed
206 expression levels, the cumulative biases can still lead to false discoveries or overlook of certain
207 DEs.

208
209 *An alternative paradigm – mixed effects model on UMI counts*
210

211 To minimize the pre-processing biases discussed above, we proposed an approach that conducts
212 DE analysis on raw UMI counts prior to implementing batch correction, normalization, imputation,
213 or feature selection. This approach, which uses a generalized linear mixed model (GLMM)³⁸,
214 preserves sample-specific structures and biological signals in the data. Furthermore, our
215 proposed approach can adjust for any potential confounding factors, such as batch, age, sex, or
216 ancestry, by incorporating them as covariates with fixed effects. This framework enables us to
217 explicitly account for the variation among biological replicates in comparison to other effects (Fig.
218 3). The proposed procedures have been implemented in software LEMUR (<https://github.com/C-HW/LEMUR>).
219

220
221 Unlike existing packages that utilize GLMMs, such as Muscat³⁹, LEMUR treats group-of-interest
222 as a fixed effect while accounting for donor-specific variations as random effects. In contrast,
223 Muscat assigns a random effect term for each combination of donor and group-of-interest.
224 Muscat's approach treats certain aspects of group-of-interest variability as random effects,
225 potentially masking differences between groups. Furthermore, Muscat's GLMMs use library size
226 as an offset to normalize counts, essentially focusing on relative abundance rather than raw
227 counts. Overall, Muscat's GLMMs operate similarly to pseudo-bulk methods, grouping counts
228 within the same donor before performing the normalization, which can result in comparable
229 performance, as demonstrated in the later examples.
230

231 To benchmark the performance of our new paradigm, we implemented eight distinct methods for
232 DE analysis: two new paradigm methods, Poisson-glmm and Binomial-glmm; two traditional
233 pseudo-bulk methods DESeq2 and edgeR; and four existing single-cell–specific methods,
234 MAST³⁴, Wilcox in Seurat, and two Muscat GLMMs (MMvst and MMpoisson).
235

236 Binomial-glmm fits a GLMM model on the zero proportion of each gene, adding donors as random
237 effect. Pseudo-bulk DESeq2 applies both VST and library size normalization. EdgeR applies
238 library size normalization. MAST adopts a zero-inflated negative binomial model, using log-
239 transformed CPM counts and incorporating the cellular detection rate as covariates. The Wilcox
240 test is non-parametric, using integrated normalized counts. The two Muscat models, MMvst with
241 VST counts and MMpoisson with raw UMI counts, account for library size. Both Muscat models
242 consider donor–group combinations as random effects. See 'Methods' for more details.
243

244 Each method was rigorously evaluated in two case studies (across cell types and across cell
245 states) and under different scenarios, such as variations in library size between groups and
246 pronounced heterogeneity within groups.
247

248 ***Case study 1 – DE analysis on different immune cell types in fallopian tube*** 249

250 In this dataset, we examined the efficacy of various methods across three distinct scenarios:
251 homogeneous groups with differing library sizes, homogeneous groups with similar library sizes,

252 and heterogenous groups. For each scenario, we illustrate the overarching gene expression
253 profile, describe the DE results using diagnostic plots, and conduct a gene ontology (GO) analysis
254 to investigate the biological foundations of our DE findings.
255

256 *Contrasting CD8+ T cell subgroups with marked library size differences*

257
258 The first comparison is between groups of CD8+ T cells (clusters 12 and 13), where there are
259 notable differences in library sizes (Fig. 4a). This example illustrates the impact of library-size-
260 based normalization on single-cell data. Using a two-sample t-test, we compared gene expression
261 means between these groups with raw UMI counts and three other normalization methods (Fig.
262 4b) using absolute t-scores. While t-scores from CPM mirror those from UMI counts, albeit with
263 minor shrinkage, both VST and integration show substantial shrinkage. This normalization
264 process dampens the gene expression differences between the groups before deploying any DE
265 detection techniques.
266

267 Each method employs its unique filtering approach within the implemented function, resulting in
268 varying numbers of input genes. Specifically, Poisson-glm, Binomial-glm, and MAST utilized
269 nearly 4600 genes as input (Fig. 4c). In contrast, pseudo-bulk DESeq2 applied default quality
270 control criteria to both genes and cells, resulting in only 104 genes being retained. Pseudo-bulk
271 edgeR retained 9743 genes in the CPM data as inputs, while Muscat mixed models utilized 6732
272 genes. Notably, the Wilcox method from the Seurat package yielded no genes that passed the
273 default filtering procedure. However, when a more lenient filtering criterion was applied, the impact
274 on the differential expression results remained minimal.
275

276 In the volcano plots, both Poisson-glm and Binomial-glm display heavily imbalanced
277 expression patterns, aligning with the observations in the density plots (Fig. S3b). However, the
278 other methods do not reflect this observation, with fold change estimates appearing evenly spread.
279 The histograms of adjusted p-values for other methods are concentrated in large values (Fig.
280 S3c). Pseudo-bulk methods and mixed models from the Muscat package, in particular, exhibit p-
281 values that are clustered around one. Despite observing imbalanced expression patterns in
282 density plots and volcano plots in this comparison, only our GLMM methods identify a substantial
283 number of differentially expressed genes (DEGs) (Fig. 4c). The heatmaps of DEGs further
284 emphasize that raw counts can better capture the differences between groups compared to
285 integrated counts (Fig. 4d). Furthermore, 403 DEGs were excluded from the integrated data
286 before testing.
287

288 The DEGs prominently feature GO terms associated with actin cytoskeleton reorganization and
289 immune synapse formation (Fig. 4e). As T cells detect antigens on an antigen-presenting cell,
290 they establish an immunological synapse, necessitating substantial actin filament restructuring.
291 Actin polymerization within this synapse aids the transit of receptors and signaling molecules,
292 crucial for T cell activation. Our results hint that among these two CD8+ T cell groups, group 12
293 cells are actively recognizing antigens. Cell groups 12 and 13 had notable differences in library
294 sizes. While the DEGs we identified contributed to the disparity in measured RNA content
295 between the two groups, genes that were not differentially expressed had a much larger effect on
296 library size; consequently, normalization erased the contribution of the DEGs to differences in
297 expression patterns. Accordingly, in this example, only our GLMMs, which operate directly on
298 UMI counts, successfully identified DEGs.
299

300 *A Glimpse at CD4+ T Cells vs. NK Cells: No Striking Library Size Differences*

301

302 The second comparison is between CD4+ T cells and NK cells (clusters 2 and 19). In the density
303 plot, we observed similar library sizes based on UMI counts for the two clusters across donors
304 except for donor 7 (Fig. 5a). The zero-proportions of genes in these two clusters fit a Poisson
305 distribution well, indicating relative homogeneity within each cell cluster (Fig. 5a).

306
307 In this comparison, Poisson-glmm, Binomial-glmm, and MAST utilized nearly 4000 genes as input
308 (Fig. 5b). Methods implemented in the Muscat package, including pseudo-bulk methods DESeq2
309 and edgeR, as well as mixed models MMvst and MMpoisson, employed 1384, 9960, 5694, 5693
310 genes, respectively, in accordance with their filtering procedure. Notably, the Wilcox method from
311 the Seurat package includes only 47 genes as input due to the filtering based on the log2 fold
312 change between two groups of interest. The log2 fold change in the package is calculated using
313 the formula $\log_2(1 + \text{mean1}) / \log_2(1 + \text{mean2})$ on the input data, which can be
314 normalized/integrated data by Seurat or other packages. This transformation attenuates the ratio
315 of the two group means through the addition of 1 to each mean, resulting in the exclusion of a
316 substantial number of genes. Wilcox, MAST, pseudo-bulk methods, and MMvst each identified
317 fewer than 100 DEGs. In contrast, the methods that use UMI counts, Poisson-glmm, Binomial-
318 glmm, and MMpoisson, identified 273, 319, and 317 DEGs, respectively (Fig. 5b).

319
320 In the volcano plots, there are more positive estimates of log2 fold change by Poisson-glmm and
321 Binomial-glmm, signifying that genes are more expressed in cluster 19 (Fig. 5c). From the
322 pairwise comparisons of log2 fold change (Fig. S4b), MAST, Wilcox, and MMvst exhibit smaller
323 log2 fold change estimates, due to normalization processes that shrink the values. Pseudo-bulk
324 methods tend to yield more conservative p-values (Fig. 5c, S4c), as illustrated in the histograms
325 (Fig. S4d). While the log2 fold change estimates are consistent across our GLMMs, pseudo-bulk
326 methods, and MMpoisson, the presence of deviant p-values leads to significant disparities in the
327 identification of DEGs. Our GLMMs identified many more DEG candidates, surpassing the
328 thresholds of adjusted p-value and fold change.

329
330 In Figure 5d, we display gene expression from DEGs identified by Poisson-glmm alongside
331 heatmaps for VST, CPM, and integrated data. Notably, differences among these heatmaps are
332 subtler than those displayed in raw UMI counts. The integrated data displays elevated gene
333 expression across groups, obscuring distinctions. The heatmaps of DEGs from Poisson-glmm
334 and Binomial-glmm show the validity of DEGs (Fig. S4f), while most of the DEGs identified by
335 MMpoisson do not display any differential expression pattern in UMI counts (Fig. S4g). We
336 performed gene ontology (GO) enrichment analysis on DEGs from Poisson-glmm. The DEGs are
337 enriched for GO terms related to leukocyte activation, cell activation, and lymphocyte activation
338 (Fig. 5e), suggesting NK cells represented by cluster 19 are more active than the CD4+ T cells.

339
340 In summary, in the comparison of two cell clusters of similar library sizes, normalization continued
341 to obscure informative differences between the two clusters and hindered the identification of
342 potential DEGs.

343
344 Deciphering the Complexities of Heterogeneous Groups: Mature T Cells vs. CD4+ T Cells

345
346 Finally, by merging groups 8 and 17 and groups 2 and 19, we created two less homogenous
347 groups-of-interest: mature T cells and CD4+ T cells, respectively. The distribution of library sizes
348 between these clusters exhibits noticeable differences (Fig. 6a), and the zero proportions of these
349 groups deviate from a Poisson distribution (Fig. S5a).

350

351 In this comparison, Poisson-glmm, Binomial-glmm, and MAST used ~3480 genes as input.
352 Pseudo-bulk DEseq2, edgeR, and mixed models utilized 1937, 10483, 7099 genes, respectively.
353 For Wilcox, 123 genes passed the filtering procedure. The volcano plots revealed similar patterns
354 to previous comparisons across various methods (Fig. S5c). Our GLMM methods exhibited
355 predominantly positive estimates of fold change, suggesting higher expression of abundant genes
356 in CD4+ T cells (group 2&19). MAST, and MMvstn showed a somewhat similar tendency, but less
357 imbalanced. However, pseudo-bulk methods and MMpoisson provided evenly distributed
358 estimates in both directions.

359
360 The estimates of log2 fold change are not quite identical among different methods (Fig S5e). Both
361 pseudo-bulk methods exhibited a negative shift compared to Poisson-glmm, while MMpoisson
362 had a positive shift. MAST, Wilcox and MMvst showed shrinkage as before. Additionally, most
363 input genes for the Wilcox method displayed positive fold changes, albeit with small magnitudes.
364 This observation sheds light on how normalization and logarithmic transformation during pre-
365 processing influences the estimation of differences in gene expression.

366
367 When we examine the violin plots of gene expression frequency and log2 mean for the DEGs
368 identified by each method, it becomes apparent that MAST, Wilcox, and MMvst captured fewer
369 DEGs with lower gene expression frequency and smaller gene means than the remaining
370 methods (Fig. 6b, 6c). It is worth noting that MAST is a zero-inflated model, which incorporates
371 excessive zeros as an additional component. However, MAST might not effectively characterize
372 the zeros, as demonstrated in previous studies on UMI counts²⁶. Consequently, potential DEGs
373 that are lowly expressed may be masked by the model. The Wilcox method tends to filter out a
374 substantial number of genes, which poses challenges in identifying lowly expressed genes.
375 MMvst, despite having a considerable number of input genes (n=7099), only identified 35 DEGs.
376

377 The heatmap of DEGs in Poisson-glmm reveals distinct expression patterns between the two
378 groups (Fig. 6d (1)). However, in this example, the inherent heterogeneity within each group
379 impacts the fitness of Poisson model, potentially leading to false discoveries. To evaluate the
380 possibility of false discoveries by Poisson-glmm, we examined DEGs identified by other methods,
381 but not by Poisson-glmm (Fig. 6d (2)-(6)). The heatmaps make it evident the DEGs that
382 differentiate between the two groups are largely identified by Poisson-glmm only; the other
383 methods did not contribute additional valid DEGs that differentiate the two groups. Conversely,
384 most of the DEGs detected by Poisson-glmm exhibit differential expression despite the
385 heterogeneity within each group.

386
387 Notably, MMpoisson mainly detected DEGs with small means (Fig. 6c), not showing clear
388 differences between different groups (Fig. 6d (6)). And the DEGs are mutually exclusive to those
389 identified by Poisson-glmm. Although Poisson-glmm and MMpoisson both use UMI counts,
390 MMpoisson includes group information as a random variable and involves library size as an offset;
391 our result underscores the significance of using an appropriate random effect in a mixed model
392 and suggests that the cell group information should be excluded from the random component.
393

394 The DEGs are enriched for GO terms related to peptide metabolic process and cytoplasmic
395 translation, indicating lower ribosomal RNA activities in mature T cells (Fig. S5g). Indeed, mature
396 T cells exhibit lower levels of ribosomal RNA activity compared to their immature counterparts,
397 mainly due to the state of activation and the metabolic requirements of the cells. On the other
398 hand, mature T cells, which are not rapidly proliferating, have less need for protein synthesis and
399 thus exhibit lower levels of rRNA activity. However, upon antigen recognition and activation,

400 mature T cells can rapidly upregulate rRNA activity and protein synthesis to support clonal
401 expansion and effector function. This differential regulation of rRNA activity is one of the many
402 ways in which cells regulate their metabolic activities to adapt to different physiological conditions.
403

404 In this example, Poisson-glmm detected more valid DEGs for heterogenous cell populations than
405 other methods. Normalization still diminished measurable differences between groups. We also
406 raise concerns about the masking of lowly expressed genes by the improper treatment of zeros,
407 as seen in MAST method and VST data.
408

409

410 **Case study 2 – DE analysis on different states of B cells**

411

412 In this case study, we applied our proposed DE framework to data collected by Kang et al⁴⁰; this
413 dataset consists of 29,065 cells and 7,661 genes from eight distinct cell types, collected from
414 peripheral blood mononuclear cells of eight lupus patients. Within each cell type, the cells are
415 evenly split into two groups for perturbation: unstimulated control and IFN- β stimulated (Fig. S6a).
416 UMAP plots (Fig. 7a) highlight that gene expression patterns are more differentiated between
417 stimulation states than between cell types. The zero-proportion plots fit better to Poisson
418 distribution when separated by stimulation states than only by cell types (Fig. S6b). This
419 observation motivated us to focus on DEGs between the cell states rather than between the cell
420 types.
421

422 Like the previous case study, we found that the distribution of library sizes underwent significant
423 changes after normalization (Fig. 7b). Raw UMI counts show that each cell type has a unique
424 library size distribution. However, these differences became less pronounced following
425 normalization, while library sizes remained relatively consistent between states within a single cell
426 type. Normalization seems to predominantly affect differences across cell types rather than
427 between states.
428

429 For the remainder of our case study, we focused on B cells. The cells from each donor were
430 divided approximately equally between the control and stimulated groups (Fig. 7c top), and the
431 library size distribution in these two groups is similar (Fig. 7c middle). The zero-proportion plot
432 suggests that the data does not perfectly fit the expected curve from the Poisson distribution,
433 indicating the presence of a mixture of subtypes within B cells (Fig. 7c bottom).
434

435 In our analysis of the subset comprising unstimulated and stimulated B cells, the majority of DE
436 methodologies used about 2,550 genes as inputs (Fig. S7a). However, the Wilcox approach within
437 Seurat selected only 144 genes. The estimates of fold change for the two states in B cells exhibit
438 an even spread across all methods, as depicted in the volcano plots (Fig. S7b). MAST and MMvst
439 struggled to identify differential patterns. Different from previous examples, our GLMM approach
440 flagged fewer DEGs than both pseudo-bulk techniques and MMpoisson. Notably, the DEGs that
441 were not shared between pseudo-bulk DESeq2 or MMpoisson and Poisson-glmm predominantly
442 belong to the extremely low expression category (Fig. 7d (2), (3)).
443

444 We hypothesized that this result could be explained by using fold change as a DEG criterion. In
445 bulk RNA-seq, a gene is typically labeled as a DEG if its adjusted p-value is below a certain
446 threshold, often 0.05, and the fold-change estimate exceeds a predetermined value, typically 1.5
447 or 2 (Fig. S8a). Most single-cell DE methods use the same criteria. However, in single-cell
448 datasets, the mean counts for many genes are exceedingly close to zero. Consequently, fold

449 change may not be a reliable metric to differentiate nuances in read counts. For instance, if gene
450 means are 2^{-3} for one group and 3^{-3} for another, the fold-change threshold of 1.5 is met, but the
451 actual difference is a mere 0.0625, which does not convey a significant disparity in expression,
452 especially when juxtaposed with genes having larger means. Moreover, near-zero values can
453 result in computational inaccuracies, causing ratio deviated from the underlying true value.
454

455 To overcome the limitation of using fold-change ratios on small counts, we established a new
456 criterion for DEGs based on absolute differences. Specifically, we mandated that the mean
457 difference between two groups exceeds a set threshold, such as -1. In the volcano plot, numerous
458 genes would be designated as DEGs when relying on ratio-defined fold change. Yet, as shown
459 in the mean vs. mean difference plot that many genes that meet the p-value criteria showcase
460 only modest changes in absolute means (Fig. S9a). This approach emphasizes genes with
461 significant absolute differences, yielding more biologically pertinent results.
462

463 We performed GO enrichment analysis on up-regulated and down-regulated genes separately
464 (Fig. 7e). We found IFN- β stimulated B cells have increased activities in interaction between
465 organisms, defense response, defense response to virus and defense response to symbiont,
466 while their activities in translation and other metabolic processes are decreased. Pseudo-bulk
467 technique detected similar GO terms while MMpoisson detected very different down-regulated
468 GOs (Fig. S9).

469
470 In this example, we demonstrated that conventional metrics to detect DEGs, especially fold
471 change based on ratios, are ill-suited for low-count data where the large fold changes reported
472 by current methods may be attributed to the ratio of two very small gene means. Careful post-
473 processing is needed to prioritize signals and manage false discoveries.
474

475 ***False discovery rates assessed under the null setting using permutation analysis***

476
477 To assess p-value calibration in empirical data, a permutation analysis was conducted within a
478 null dataset focusing on a group of interest. We specifically conducted the analysis on three
479 datasets: the control group of B cells, group 2, and group 13 in case study 1. Each underwent
480 random assignment to either the control or stimulus group. Subsequently, p-values for each gene
481 were computed employing various methods, with the gene set confined to those input into the
482 Poisson-glm model. To mitigate potential gene filtering, the threshold for the Wilcox method was
483 relaxed. This process was iterated 20 times, and on each iteration the proportion of p-values
484 below 0.05 was calculated along with the corresponding false discovery of differentially expressed
485 genes.
486

487 The analysis of the violin plot (Fig. 7f, S10) reveals that both our GLMM methods and the Wilcox
488 method exhibit consistently well-calibrated p-values among different choices of null datasets.
489 However, pseudo-bulk methods, and mixed models from Muscat appear excessively conservative,
490 with an overall proportion considerably below 0.05. The performance of MAST is conservative in
491 B cells but not in case study 1. The histograms of p-values across the 20 runs demonstrate a
492 consistently flat distribution for our glmm methods and the Wilcox method, indicative of adherence
493 to the null setting (Fig. S10). Conversely, other methods display overestimated p-values, yielding
494 conservative outcomes. Note that even though Wilcox performed well in the permutation analysis,
495 it is not powerful to detect real DEGs as shown in previous case studies. Under both the existing
496 criteria and our newly established criteria for determining DEGs, each method detected, at most,
497 one false discovery in each run.
498
499

500 **Discussion**

501
502 In this manuscript, we examined existing DE approaches to pre-processing, input values and test
503 statistics, and fold-change definitions in the context of single-cell DE analysis. We demonstrated
504 through extensive real-data examples the limitations and drawbacks of current practices. We
505 showed that current normalization and pre-processing techniques may obscure DEGs by an
506 overreliance on relative RNA abundance and ignoring or correcting for biological zeros. We also
507 illustrated how use of volcano plots in DE analysis, which also depends on relative RNA
508 abundance, leads to false discoveries in lowly expressed genes by prioritizing fold changes in
509 expression over absolute changes. We also argued that single-cell DE analysis suffers from false
510 discoveries due to the inappropriate handling of donor effects, as well as from biases that
511 accumulate as the consequence of sequential workflows.

512
513 We advocate a new paradigm, Poisson-glmm, which uses UMI counts as input and a generalized
514 Poisson mixed effect models to account for batch effects and within-sample variation. This
515 framework's use of UMI counts can significantly improve current practices by leveraging absolute
516 RNA expression. Poisson-glmm shows superior sensitivity and robustness toward model
517 misspecification when compared to current single-cell DE methods, which should ultimately lead
518 to new biological insights from single-cell data.

519
520 The use of UMI counts for DE analysis in scRNA-seq can significantly improve current practices,
521 potentially making some current practices (e.g., volcano plots as a diagnostic DE tool) obsolete.
522 However, relying on UMI counts as a representation of genuine RNA content predicates that
523 measurements are strictly single-cell based, underscoring the need for meticulous doublet and
524 triplet removal prior to DE analysis. Furthermore, seamlessly implementing this new paradigm
525 into existing popular tools remains a challenge. Given this significant shift from current practices,
526 a sustained effort will be required to educate and train researchers on these new alternatives and
527 to reshape existing practices accordingly.

528
529
530 **Methods and materials**

531
532 *Datasets and pre-processing*
533
534 In case study 1, a 10X scRNA-seq dataset of post-menopausal fallopian tubes, with 57,182 cells
535 sourced from five donors, covering 29,382 genes was analyzed. We obtained 20 clusters via
536 HIPPO algorithm. We did not apply a pre-filtering procedure on this dataset, except for built-in
537 filtering steps in each method. We used sctransform to get the VST data and the integration
538 workflow provided by Seurat to obtain the integrated data.

539
540 All integration or normalization processes were performed on the entire dataset, since cell types
541 are typically unknown during the pre-processing stage. In cross-batch integration, only the top
542 2,000 highly expressed genes were retained, which significantly reduced the number of genes for
543 downstream analysis. The dataset had already been fully analyzed and annotated with cell types.
544 We utilized the annotations to examine the effects of normalization/integration on distributions of
545 library sizes across cells.

546
547 In case study 2, the dataset comprised 10X droplet-based scRNA-seq PBCM data from eight
548 Lupus patients obtained before and after 6h-treatment with IFN- β . After removing undetected and
549 lowly expressed genes (less than 10 cells expressing more than 1), the dataset consisted of

550 29065 cells and 7661 genes. The integrated data was replaced by log2-transformed normalized
551 expression values obtained via computeLibrarayFactors and logNormCounts functions in Muscat.
552

553 *Poisson-glmm and Binomial-glmm*

554
555 By default, we excluded any genes detected in fewer than 5% cells in the compared groups from
556 differential testing. The GLMMs were implemented with glmmPQL function of the MASS package.
557 We calculated adjusted p-values by using Benjamini-Hochberg correction. Each model fitting was
558 applied on one gene and the two compared groups.
559

560 We fit Poisson-glmm on UMI counts. Each count X_{cgk} sampled from cell c , donor k , and gene g ,
561 was modeled by

$$\begin{aligned} X_{cgk} | \lambda_{cgk} &\sim \text{Poisson}(\lambda_{cgk}) \\ \log(\lambda_{cgk}) &= \mu_g + X_c \beta_g + \epsilon_{gk}. \end{aligned}$$

562
563 We fit Binomial-glmm on the zero proportions. Each count X_{cgk} was modeled by
564

$$\begin{aligned} 1\{X_{cgk} = 0\} | p_{cgk} &\sim \text{Bernoulli}(p_{cgk}) \\ \log(p_{cgk}/(1 - p_{cgk})) &= \mu_g + X_c \beta_g + \epsilon_{gk} \end{aligned}$$

565 where X_c is the indicator for groups (e.g. cell types in case study 1, control/stimulus in case study
566 2), and $\epsilon_{gk} \sim N(0, \sigma_g^2)$ represents the random effects for donor k . Our goal was to test $H_0: \beta_g = 0$.
567

568 For both methods, we provided “log2 fold change” computed by $\log_2(\exp(\beta_g))$. In Poisson-glmm,
569 this estimate indicates the increment of $\log_2(\lambda_2)$ against $\log_2(\lambda_1)$, which is the conventional log2
570 fold change. However, this term in Binomial-glmm doesn’t represent the same meaning. It is the
571 difference between $\text{logit}(p_2)$ and $\text{logit}(p_1)$. The p-value and BH adjusted p-value are provided.
572

573 *Benchmarked methods*

574 Pseudo-bulk DESeq2 and pseudo-bulk edgeR are aggregation-based methods used in our
575 comparison. The input counts were summed up for a given gene over all cells in each group and
576 by donor. The pseudo-bulk data matrix has dimensions GxS, where S denotes the number of
577 interactions of donors and groups. For example, if there are two groups and 'a' and 'b' donors in
578 each group, then 'S' is equal to 2(a + b). We used raw counts as the input for DESeq2, while CPM
579 counts were used for edgeR. The log fold change was converted to log2 fold change in all the
580 comparisons. We implemented these two pseudo-bulk methods following the guided tutorial in
581 Muscat
582 package;
583 <https://www.bioconductor.org/packages/devel/bioc/vignettes/muscat/inst/doc/analysis.html>.

584 For MAST, we fitted a zero-inflated regression model (function zlm) for each gene and applied a
585 likelihood ratio test (function lrTest) to test for between-group differences in gene expression.
586 Besides the labels of groups and the cellular detection rate, we also included donor labels in the
587 covariates. This method was run on log (CPM+1) counts. We followed the tutorial
588 <https://github.com/RGLab/MAST>.

589 Wilcox, a rank sum test, is the default DE method in the FindMarkers function in the Seurat
590 package. We used integrated data and log counts as input. We computed the log fold change
591 given in the output as $\log(1+\text{mean1})/(1+\text{mean2})$. We applied the default filter in FindMarkers to
592 only test genes with a log fold change greater than 0.25. We calculated the adjusted p-value
593

598 provided from the function based on Bonferroni correction. We followed the guided tutorial found
599 here: https://satijalab.org/seurat/articles/de_vignette.
600

601 MMvst and MMpoisson are mixed models implemented in the Muscat package. MMvst fits linear
602 mixed models on variance-stabilizing transformation data. MMpoisson fits Poisson generalized
603 linear mixed models with an offset equal to the library size factors. In both models, we fit a $\sim 1 +$
604 group + (1|sample) model for each gene, where 'sample' denotes the experimental units (the
605 interaction of donors and groups). We followed the tutorial found at:
606 <https://www.bioconductor.org/packages-devel/bioc/vignettes/muscat/inst/doc/analysis.html>.
607

608 *The criteria to determine DEGs*
609

610 For the benchmarked methods, we adhered to conventional criteria for the identification of
611 Differentially Expressed Genes (DEGs). Specifically, a gene was classified as a DEG if its
612 absolute log2 fold change exceeded a predefined threshold, and the adjusted p-value was below
613 a specified cutoff. Typically, DEGs are visually represented in volcano plots. In the first dataset,
614 the log2 fold change threshold was set at log2(1.5), whereas in the second dataset, it was set at
615 1. The adjusted p-value threshold for both datasets was established at 0.05.
616

617 We proposed new criteria that are based on the convention plus the gene mean and the difference
618 in mean. If the log2 gene mean in two groups is lower than a certain value (-2.25 in case study 1)
619 and the log2 mean difference is smaller than a threshold (-1 in case study 1), the gene would not
620 be considered as a DEG. These can also be used as a filter before any DE analysis to speed up
621 the computation. Both criteria are adjustable, depending on the dataset's performance and
622 characteristics. An examination of heatmaps and mean difference against mean plot in advanced
623 can be helpful to determine the thresholds when analyzing a new dataset (Fig. S8b, c).
624

625 *Variation analysis*

626 To gain a deeper understanding of the donor effect and cell type effect concerning various types
627 of counts, we conducted a variation analysis across multiple group comparisons. To ensure the
628 consistency of our results, we restricted our analysis to genes presented in all datasets. For each
629 gene, we employed linear models ($\text{lm}(\text{count} \sim \text{donor} + \text{group})$) and computed the variances
630 attributed to three components: donor, group, and the residual. Logarithm transformation was
631 applied to UMI counts and CPM data to address skewness. The outcomes of this analysis were
632 then presented and compared based on the proportion of variation explained by the first two
633 components across different count types and various pairs. The results of the top 500 genes with
634 the lowest residual variations were exhibited.
635

636 *GO enrichment analysis*

637 GO over-representation analyses were performed using the enrichGO function in the R package
638 clusterProfiler with default parameters and the functional category for enrichment analysis set to
639 the GO 'Biological Processes' category.
640

641 *Data availability*
642

643 Both scRNA-seq datasets used in this study are publicly available. Processed and de-identified
644 human single-cell RNA sequencing data scRNA-seq dataset of post-menopausal fallopian tubes
645 has been deposited at Cellxgene under the following URL:
646 <https://cellxgene.cziscience.com/collections/d36ca85c-3e8b-444c-ba3e-a645040c6185>. The
647 droplet scRNA-seq data used in case study 2 is deposited under the Gene Expression Omnibus

648 under the accession number [GSE96583](#). The dataset is also available in R through the
649 Bioconductor ExperimentHub package.
650

651 **Code availability**

653 We provide an R package, LEMUR, implementing Poisson-glmm and Binomial-glmm methods
654 for DE analysis discussed in this study. The LEMUR package is available from GitHub
655 (<https://github.com/C-HW/LEMUR>). In addition, the R source code to reproduce all data analysis
656 in the study is available from GitHub at <https://c-hw.github.io/DEanalysis/index.html> .
657

658 **Acknowledgements**

660 The work was supported by National Institutes of Health grant R01 GM126553, R01 HG011883
661 and HG012927, and additional grant no. NSF 2016307 and Sloan Research Fellowship to M.C.
662

663 **Author information**

664 *Authors and Affiliations*

667 Department of Statistics, University of Chicago, Chicago, USA
668 Chih-Hsuan Wu
669

670 Department of Biostatistics, University of Michigan, Ann Arbor, USA
671 Xiang Zhou
672

673 Department of Human Genetics and Department of Medicine, University of Chicago, Chicago,
674 USA
675 Mengjie Chen
676

677 **Contributions**

678 M.C. conceived and led this work. C.W. and M.C. developed the methods and performed the
679 analyses. C.W. implemented the software. X.Z. participated in critically revising the draft. C.W.
680 and M.C. wrote the paper with feedback from X.Z. All authors read and approved the final
681 manuscript.
682

683 **Corresponding authors**

684 Correspondence to Mengjie Chen.
685

686 **Ethics declarations**

687 Ethics approval is not applicable to this study.
688

689

690 **Competing interests**

691 The authors declare no competing interests.
692

	Poisson-glmm	Binomial-glmm	Pb-DESeq2	Pb-edgeR	MAST	Wilcox	MMvst	MMpoisson
Package	LEMUR	LEMUR	Muscat	Muscat	MAST	Seurat	Muscat	Muscat
Input	UMI	Zero counts	UMI	CPM	CPM	Integrated/ Log normalized	VST	UMI
Model base	Poisson glmm	Binomial glmm	Negative binomial model	Negative binomial model	Zero-inflated model	Rank-sum test	LMM	Poisson glmm
Normalization	X	X	V	V	V	V	V	V
Normalization method			1. M median of ratio size factor and variance stabilizing transformation in the method	1. CPM normalization 2. Trimmed mean of M values (TMM) in the model	1. CPM normalization	1. Integration applied on log normalized data by Seurat in case study 1 package 2. Log2-transformed normalized data by Muscat in case study 2	1. VST normalization	1. Library size factor as offset in the model

693 Table 1. Comparison of DE methods used in this paper.
 694

695 **Figure legends:**

696 **Figure 1. Effects of normalization on library size, zero frequency, and gene count**

697 distributions.

698 a. Violin plots display library sizes based on raw UMI counts (top) and after data integration
699 (bottom), categorized by cell types and donors.

700 b. Violin plot illustrating the frequency of gene expression (non-zero counts) in raw UMI data.

701 c. Histograms representing the distribution of non-zero counts in raw UMI data across various
702 data transformations.

703 d. Histograms detailing the zero counts in raw UMI data, comparing VST with integrated data
704 where zeros are imputed or converted to non-zeros.

705 e. Histograms showing the distribution of gene RUNX3 across different data transformations.

706

707 **Figure 2: Cluster and Variation Analysis of Single-Cell Data from the Fallopian Tube in**

708 Case Study 1.

709 a. UMAP visualizing 20 clusters identified by HIPPO in case study 1.

710 b. Canonical markers delineate specific cell subtypes: clusters 9, 15, and 19 as NK cells; clusters
711 7, 10, 11, 14, 16, 18, and 20 as CD4+ T cells; clusters 4, 6, 12, and 13 as CD8+ T cells; clusters
712 8 and 17 as mature naive T cells.

713 c. Distribution of donors across the 20 identified clusters.

714 d. Comparative analysis of variation proportions attributable to donor and cell type effects across
715 different pairs and datasets.

716 e. Scatter plots comparing variation proportions due to donor and cell type effects across various
717 pairings and data sources.

718

719 **Figure 3. Comparison of established workflows and proposed paradigm for single-cell**

720 analysis.

721 Left: Under the current single-cell analysis pipeline, the raw UMI counts collected from multiple
722 donors are integrated to remove the batch effects and normalized for further analysis. It is
723 common to perform DE analysis on processed data.

724 Right: Our new paradigm directly performs a generalized linear mixed model on raw UMI counts.
725 The random effect can account for the batch effect due to samples. The annotated cell types can
726 be obtained from existing pipeline or HIPPO algorithm which clusters cells based on the zero
727 proportions of UMI counts.

728

729 **Figure 4. DE analyses on CD8+ T cell subgroups**

730 a. Density plot of the library size for group 12 and 13.

731 b. Scatterplot comparisons of t-scores from mean difference tests between raw UMI counts and
732 other transformed data. Each gene's expression in two different groups is compared, showcasing
733 the pairwise absolute t-scores from various data sources.

734 c. Counts of input genes and DEGs in different DE methods.

735 d. Heatmaps visualize Poisson-glmm DEGs. Order: UMI counts (left), integrated data (middle),
736 and genes not included in the integrated data but shown in UMI counts (right). Heatmaps arrange
737 genes by descending Poisson-glmm fold change estimates and columns group cells by cell
738 clusters and donors.

739 e. GO analysis of the DEGs identified by Poisson-glmm.

740

741 **Figure 5. DE analyses on CD4+ T Cells vs. NK Cells**

742 a. Left: Density plot of the library size for group 2 and 19. Middle: Density plot of the library size
743 by different donors. Right: Zero proportion plots for each group and combined groups.

744 b. Counts of input genes and DEGs across different differential expression methods.

- 745 c. Volcano plots for each method, highlighting DEGs in blue. The signs of log2 fold change are
746 adjusted such that positive signs represent higher expression in group 19.
747 d. Heatmaps of Poisson-glmm DEGs shown in different data sources, with genes in integrated
748 data featured in the top block, and those absent in the lower block.
749 e. GO analysis of the DEGs identified by Poisson-glmm.

- 750
751 **Figure 6. DE analyses on heterogeneous groups: Mature T Cells vs. CD4+ T Cells**
752 a. Density plots comparing library sizes for combined groups 8 & 17 and 2 & 19.
753 b. Comparisons of the gene expression frequency of the DEGs from different methods.
754 c. Violin plot of log2 gene mean for DEGs from different methods.
755 d. Heatmaps of DEGs from different methods.

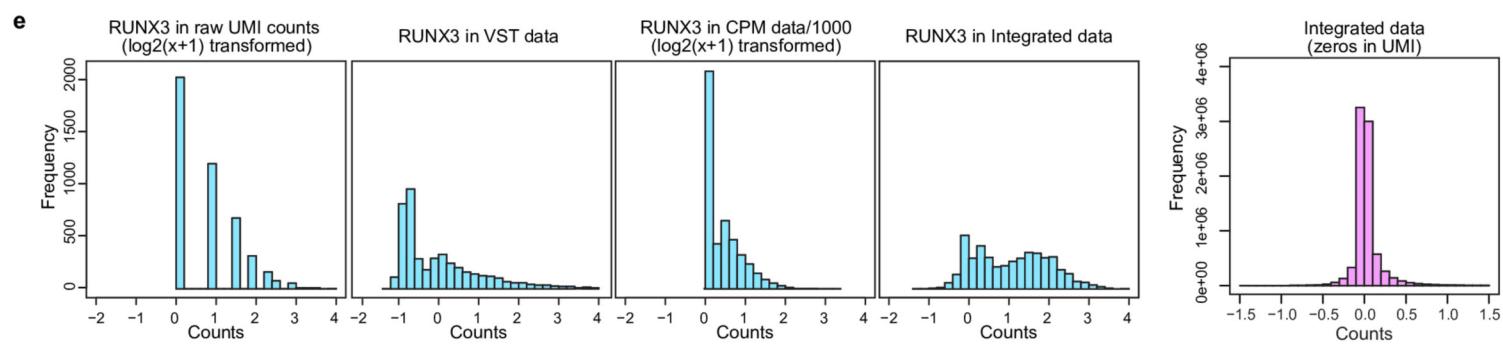
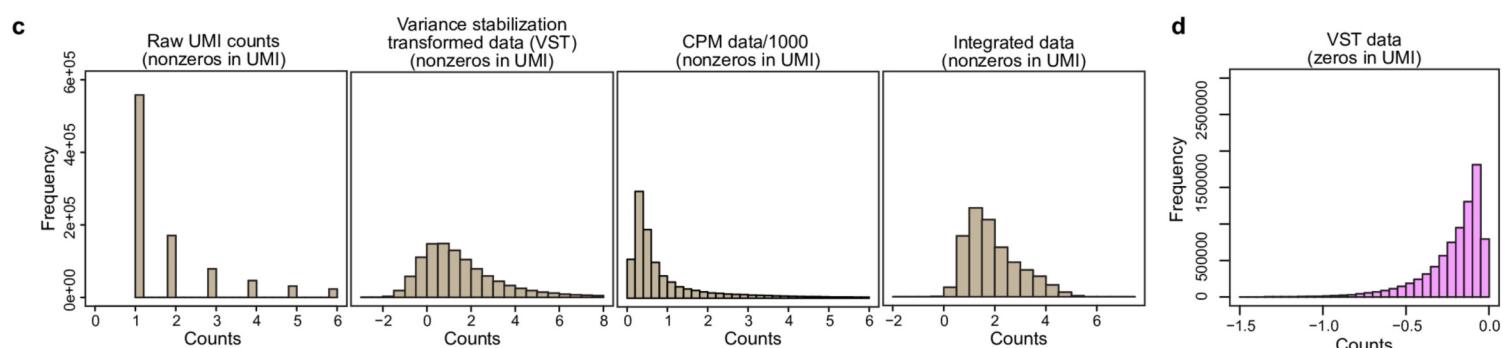
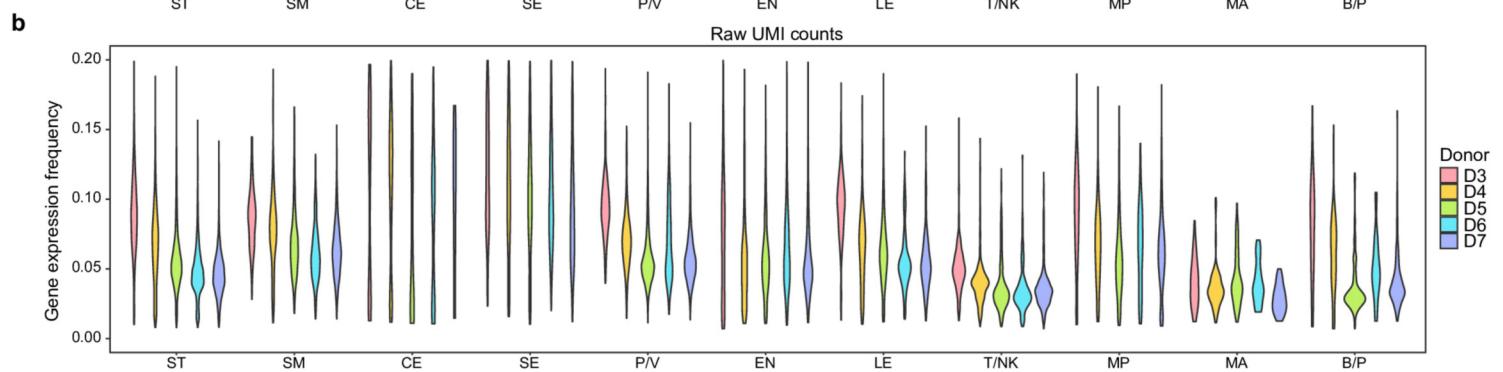
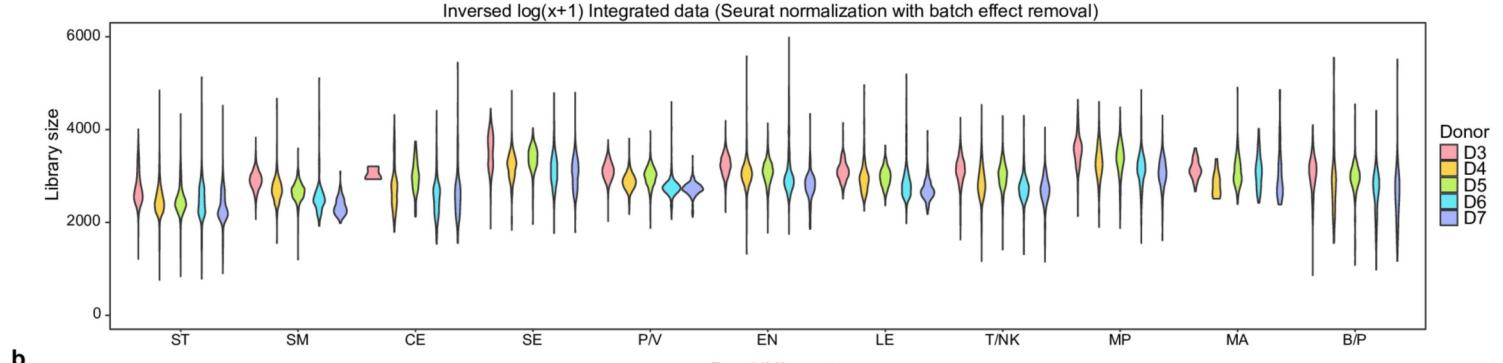
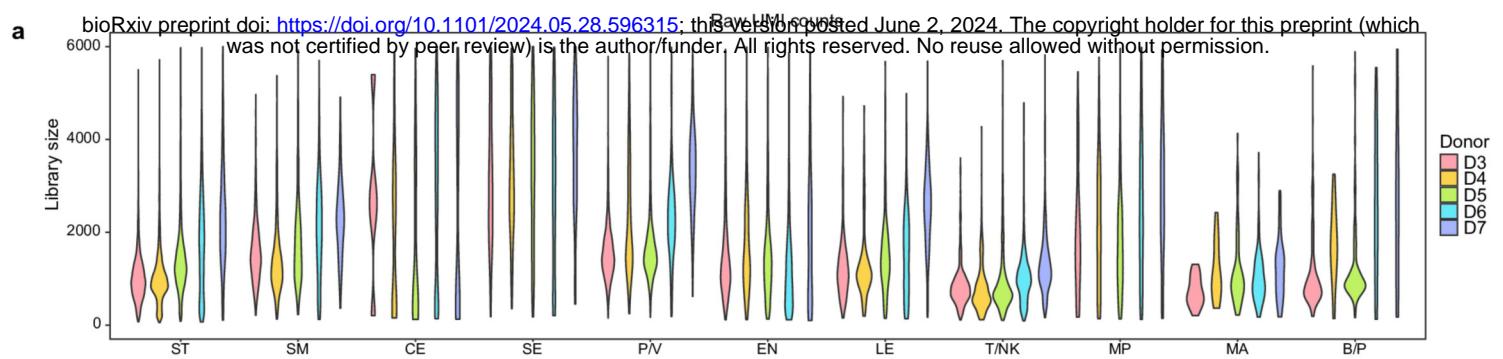
756
757 **Figure 7. Overview of case study 2 and DE analyses on different states in B cells**

- 758 a. UMAP showing groups and cell types for case study 2.
759 b. Library size comparisons before (raw UMI counts) and after normalization (log-normalized data)
760 by cell type.
761 c. Top: Donor distribution among B cells. Middle: Density plot of library size in different states.
762 Bottom: Zero proportion plots for different states and combined states.
763 d. Heatmaps of DEGs identified from different methods.
764 e. Violin plots depicting the proportion of p-values below 0.05 for each method.
765 f. GO analysis for up-regulated (left) and down-regulated genes (right).

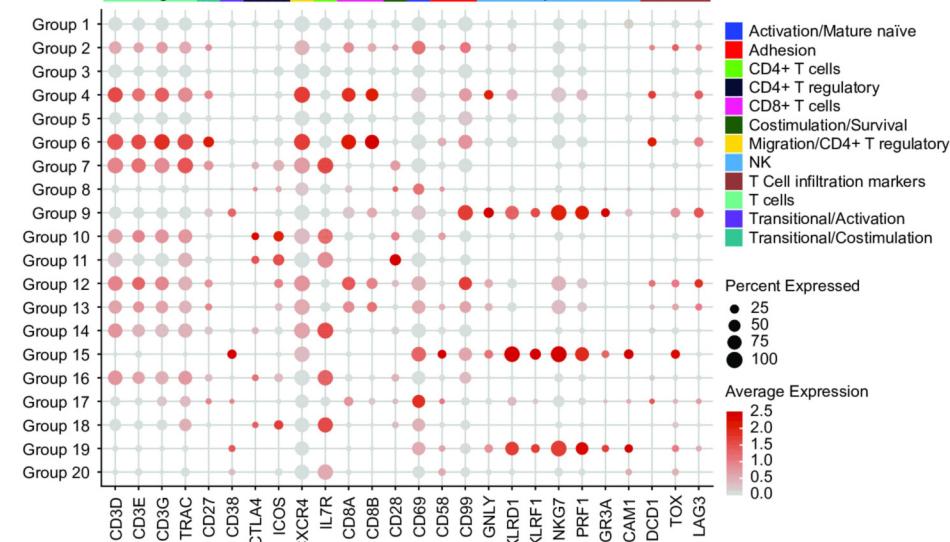
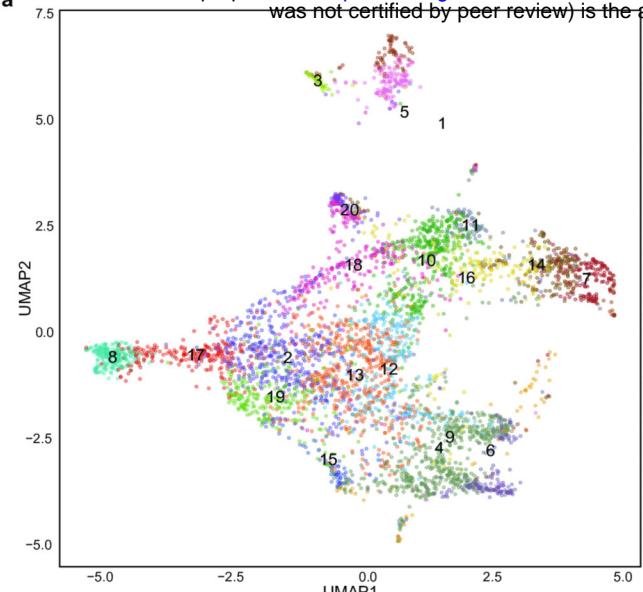
- 766
767
768
769
770
771 1. Saliba, A.-E., Westermann, A.J., Gorski, S.A. & Vogel, J. Single-cell RNA-seq: advances
772 and future challenges. *Nucleic acids research* **42**, 8845-8860 (2014).
773 2. Greenwald, W.W. et al. Pancreatic islet chromatin accessibility and conformation reveals
774 distal enhancer networks of type 2 diabetes risk. *Nature communications* **10**, 2078 (2019).
775 3. Grubman, A. et al. A single-cell atlas of entorhinal cortex from individuals with Alzheimer's
776 disease reveals cell-type-specific gene expression regulation. *Nature neuroscience* **22**,
777 2087-2097 (2019).
778 4. Lawlor, N. et al. Single-cell transcriptomes identify human islet cell signatures and reveal
779 cell-type-specific expression changes in type 2 diabetes. *Genome research* **27**, 208-222
780 (2017).
781 5. Squair, J.W. et al. Confronting false discoveries in single-cell differential expression.
782 *Nature communications* **12**, 1-15 (2021).
783 6. Das, S., Rai, A., Merchant, M.L., Cave, M.C. & Rai, S.N. A comprehensive survey of
784 statistical approaches for differential expression analysis in single-cell RNA sequencing
785 studies. *Genes* **12**, 1947 (2021).
786 7. Das, S., Rai, A. & Rai, S.N. Differential Expression Analysis of Single-Cell RNA-Seq Data:
787 Current Statistical Approaches and Outstanding Challenges. *Entropy* **24**, 995 (2022).
788 8. Lengyel, E. et al. A molecular atlas of the human postmenopausal fallopian tube and ovary
789 from single-cell RNA and ATAC sequencing. *Cell Reports* **41** (2022).
790 9. Li, P., Piao, Y., Shon, H.S. & Ryu, K.H. Comparing the normalization methods for the
791 differential analysis of Illumina high-throughput RNA-Seq data. *BMC bioinformatics* **16**, 1-
792 9 (2015).
793 10. Zyprych-Walczak, J. et al. The impact of normalization methods on RNA-Seq data
794 analysis. *BioMed research international* **2015** (2015).

- 795 11. Dillies, M.-A. et al. A comprehensive evaluation of normalization methods for Illumina high-
796 throughput RNA sequencing data analysis. *Briefings in bioinformatics* **14**, 671-683 (2013).
797 12. Robinson, M.D. & Oshlack, A. A scaling normalization method for differential expression
798 analysis of RNA-seq data. *Genome biology* **11**, 1-9 (2010).
799 13. Lytal, N., Ran, D. & An, L. Normalization methods on single-cell RNA-seq data: an
800 empirical survey. *Frontiers in genetics* **11**, 41 (2020).
801 14. Leek, J.T. et al. Tackling the widespread and critical impact of batch effects in high-
802 throughput data. *Nature Reviews Genetics* **11**, 733-739 (2010).
803 15. Korsunsky, I. et al. Fast, sensitive and accurate integration of single-cell data with
804 Harmony. *Nature methods* **16**, 1289-1296 (2019).
805 16. Chen, M. & Zhou, X. Controlling for Confounding Effects in Single Cell RNA Sequencing
806 Studies Using both Control and Target Genes. *Sci Rep* **7**, 13587 (2017).
807 17. Chen, M. et al. Alignment of single-cell RNA-seq samples without overcorrection using
808 kernel density matching. *Genome research* **31**, 698-712 (2021).
809 18. Hu, J., Chen, M. & Zhou, X. Effective and scalable single-cell data alignment with non-
810 linear canonical correlation analysis. *Nucleic Acids Research* (2021).
811 19. Schmid, R. et al. Comparison of normalization methods for Illumina BeadChip HumanHT-
812 12 v3. *BMC genomics* **11**, 1-17 (2010).
813 20. Tran, H.T.N. et al. A benchmark of batch-effect correction methods for single-cell RNA
814 sequencing data. *Genome biology* **21**, 1-32 (2020).
815 21. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-
816 seq data using regularized negative binomial regression. *Genome biology* **20**, 1-15 (2019).
817 22. Lause, J., Berens, P. & Kobak, D. Analytic Pearson residuals for normalization of single-
818 cell RNA-seq UMI data. *Genome biology* **22**, 1-20 (2021).
819 23. Argelaguet, R. et al. MOFA+: a statistical framework for comprehensive integration of
820 multi-modal single-cell data. *Genome biology* **21**, 1-17 (2020).
821 24. Wang, X., Park, J., Susztak, K., Zhang, N.R. & Li, M. Bulk tissue cell type deconvolution
822 with multi-subject single-cell expression reference. *Nature communications* **10**, 380 (2019).
823 25. Yang, Y. et al. Dimensionality reduction by UMAP reinforces sample heterogeneity
824 analysis in bulk transcriptomic data. *Cell reports* **36** (2021).
825 26. Kim, T.H., Zhou, X. & Chen, M. Demystifying "drop-outs" in single-cell UMI data. *Genome
826 biology* **21**, 196 (2020).
827 27. Qiu, P. Embracing the dropouts in single-cell RNA-seq analysis. *Nature communications*
828 **11**, 1169 (2020).
829 28. Svensson, V. Droplet scRNA-seq is not zero-inflated. *Nature Biotechnology* **38**, 147-150
830 (2020).
831 29. Gong, W., Kwak, I.-Y., Pota, P., Koyano-Nakagawa, N. & Garry, D.J. DrImpute: imputing
832 dropout events in single cell RNA sequencing data. *BMC bioinformatics* **19**, 1-10 (2018).
833 30. Li, W.V. & Li, J.J. An accurate and robust imputation method sclImpute for single-cell RNA-
834 seq data. *Nature communications* **9**, 997 (2018).
835 31. Tracy, S., Yuan, G.-C. & Dries, R. RESCUE: imputing dropout events in single-cell RNA-
836 sequencing data. *BMC bioinformatics* **20**, 1-11 (2019).
837 32. Chen, M. & Zhou, X. VIPER: variability-preserving imputation for accurate gene
838 expression recovery in single-cell RNA sequencing studies. *Genome Biol* **19**, 196 (2018).
839 33. Pierson, E. & Yau, C. ZIFA: Dimensionality reduction for zero-inflated single-cell gene
840 expression analysis. *Genome biology* **16**, 1-10 (2015).
841 34. Finak, G. et al. MAST: a flexible statistical framework for assessing transcriptional
842 changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome
843 biology* **16**, 1-13 (2015).
844 35. Kim, T.H., Zhou, X. & Chen, M. Demystifying "drop-outs" in single-cell UMI data. *Genome
845 Biol* **21**, 196 (2020).

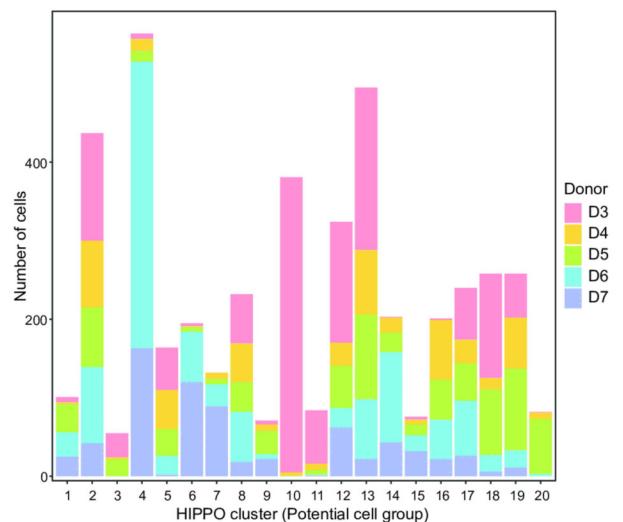
- 846 36. Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion
847 for RNA-seq data with DESeq2. *Genome biology* **15**, 1-21 (2014).
848 37. Robinson, M.D., McCarthy, D.J. & Smyth, G.K. edgeR: a Bioconductor package for
849 differential expression analysis of digital gene expression data. *bioinformatics* **26**, 139-140
850 (2010).
851 38. Clayton, D.G. Generalized linear mixed models. *Markov chain Monte Carlo in practice* **1**,
852 275-302 (1996).
853 39. Crowell, H.L. et al. Muscat detects subpopulation-specific state transitions from multi-
854 sample multi-condition single-cell transcriptomics data. *Nature communications* **11**, 6077
855 (2020).
856 40. Kang, H.M. et al. Multiplexed droplet single-cell RNA-sequencing using natural genetic
857 variation. *Nature biotechnology* **36**, 89-94 (2018).
858



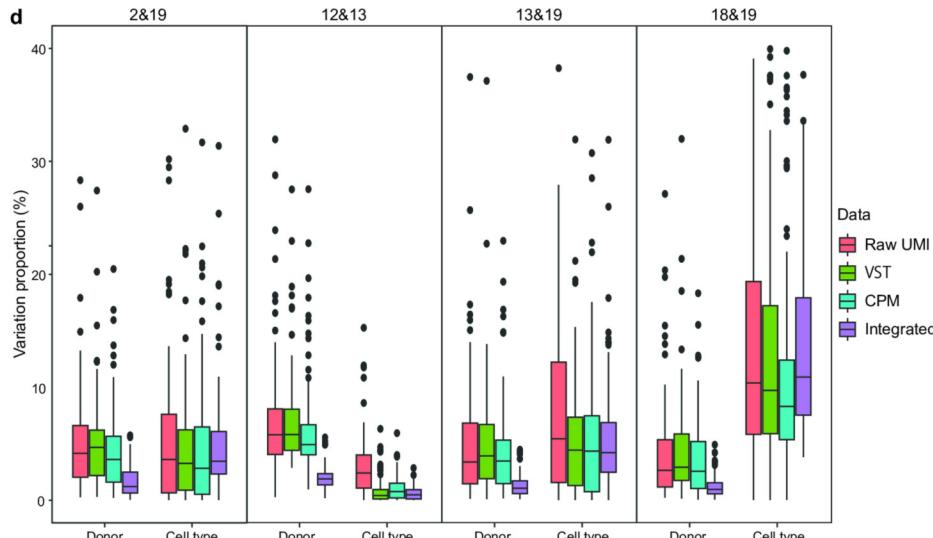
a



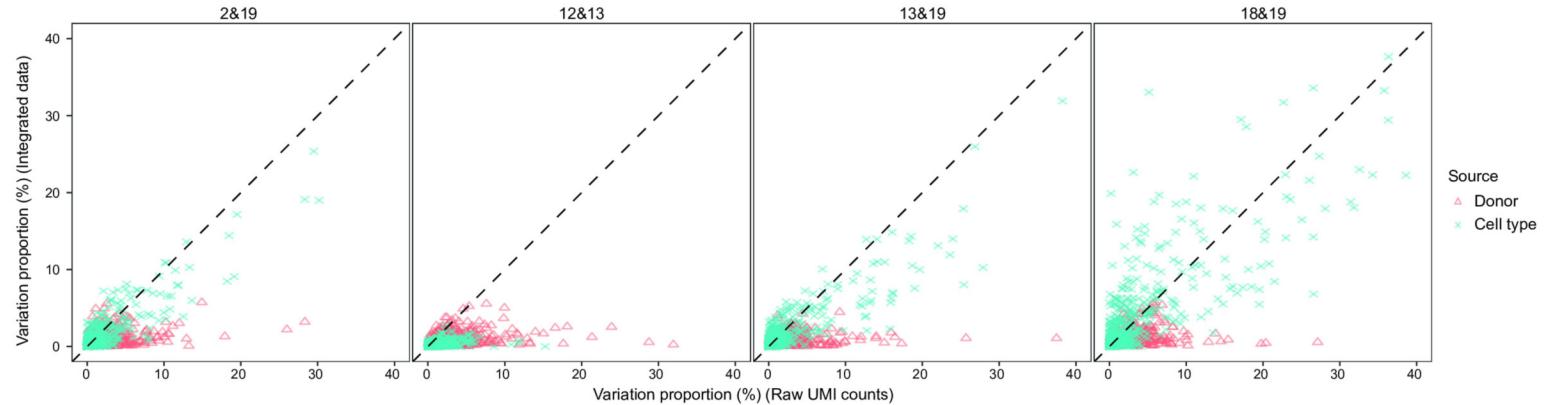
c



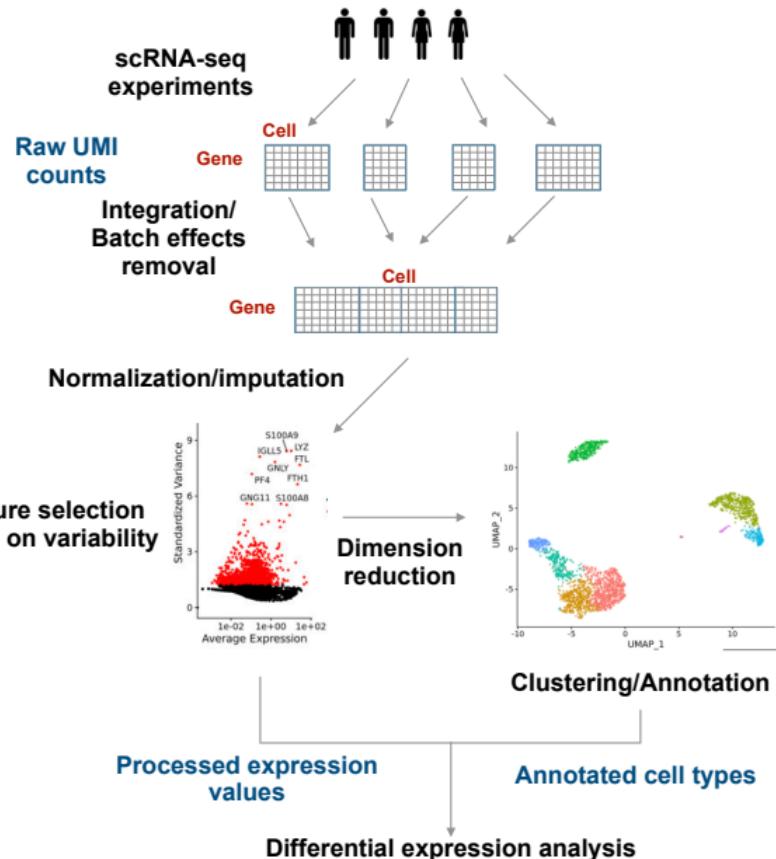
d



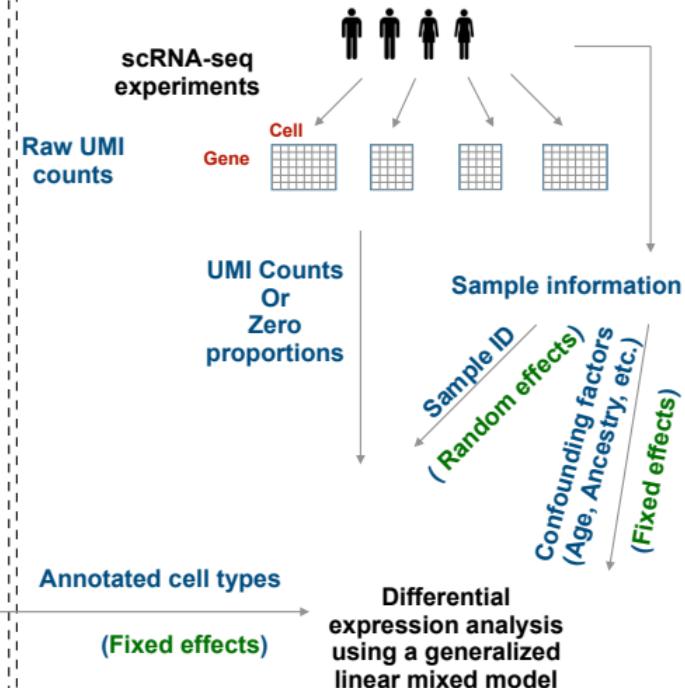
e



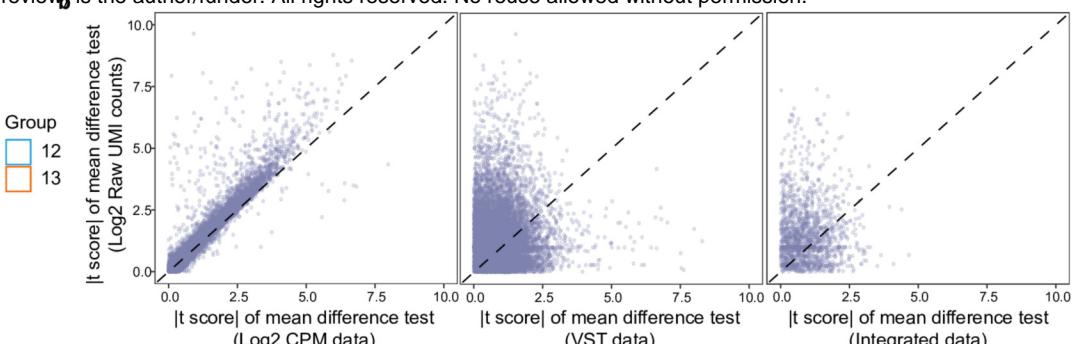
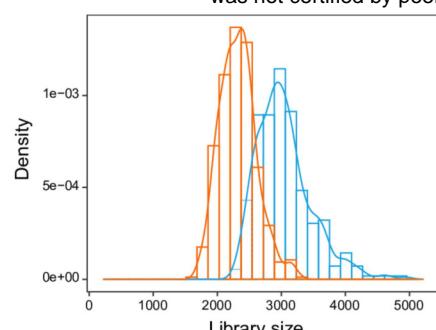
Current practice



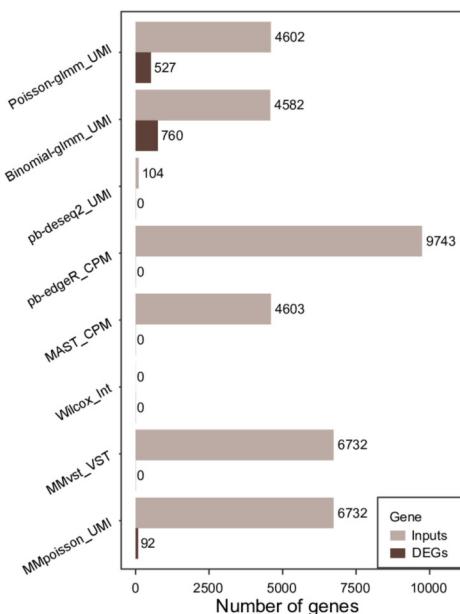
New paradigm



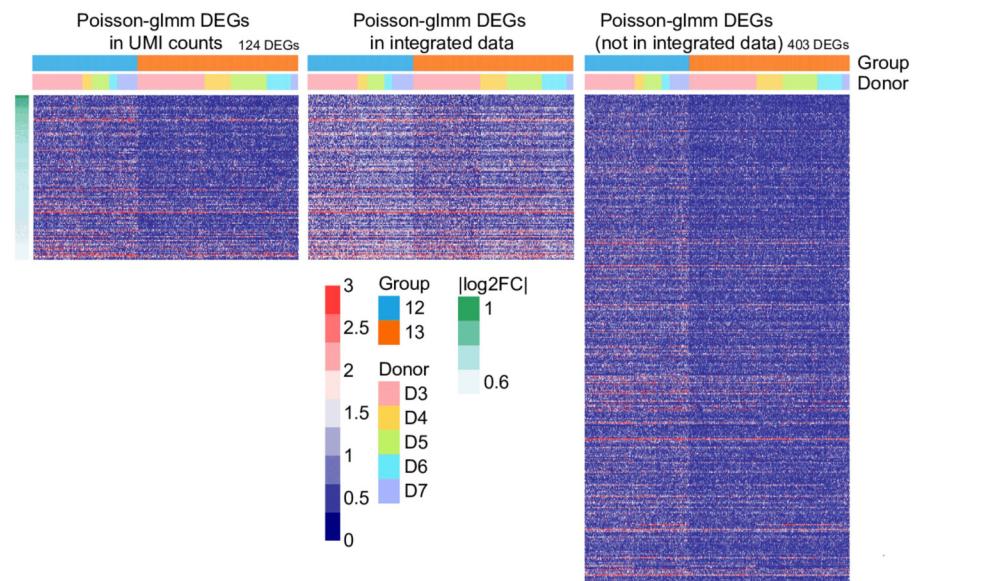
a



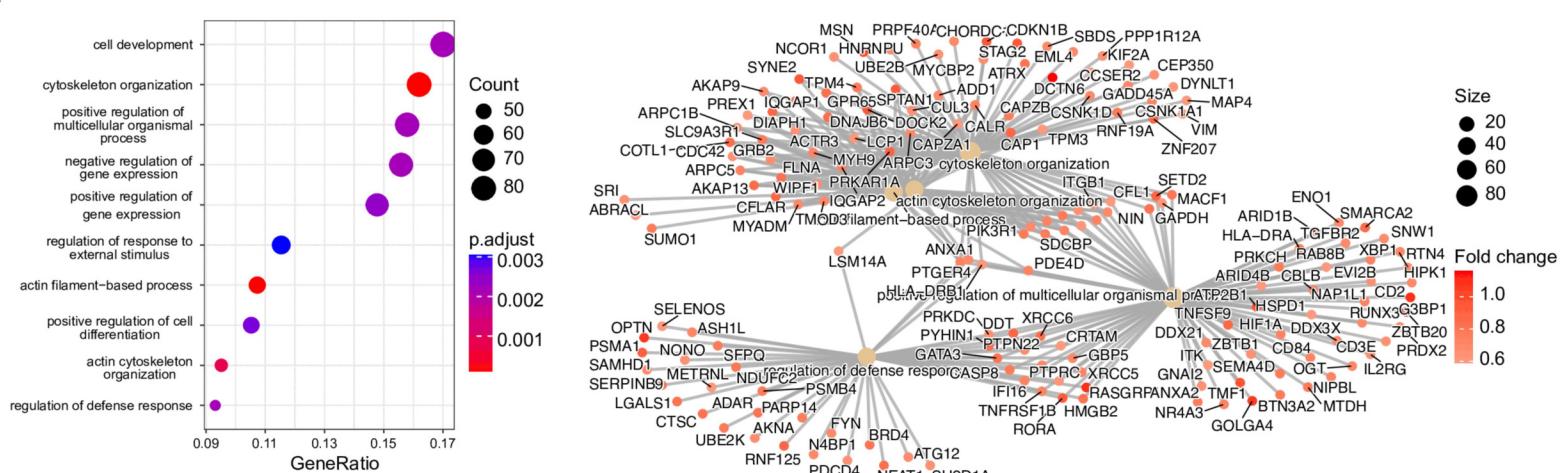
c

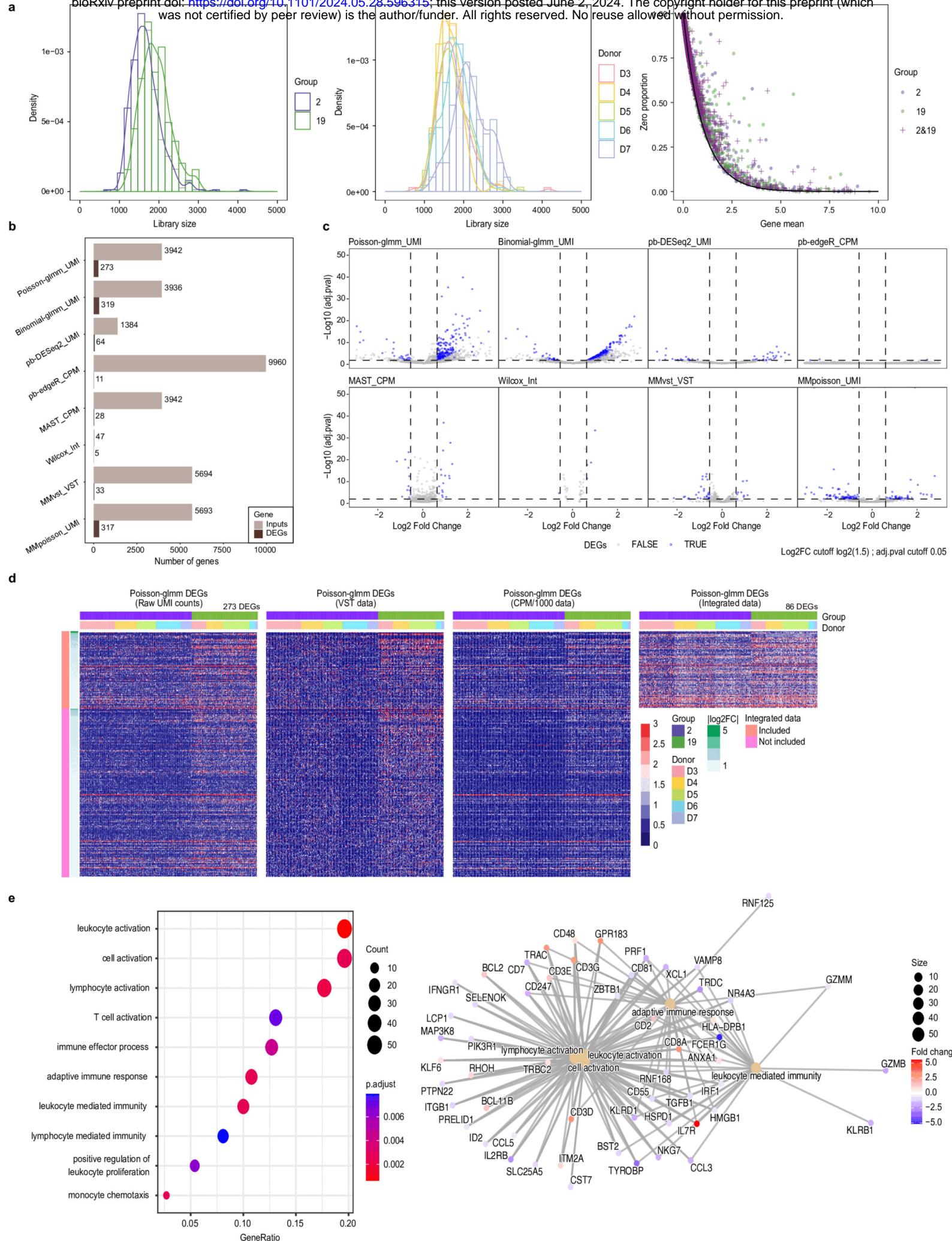


d

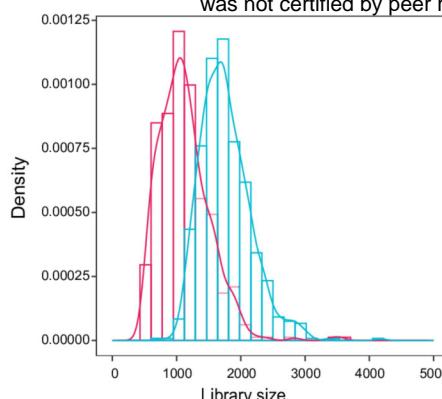


e

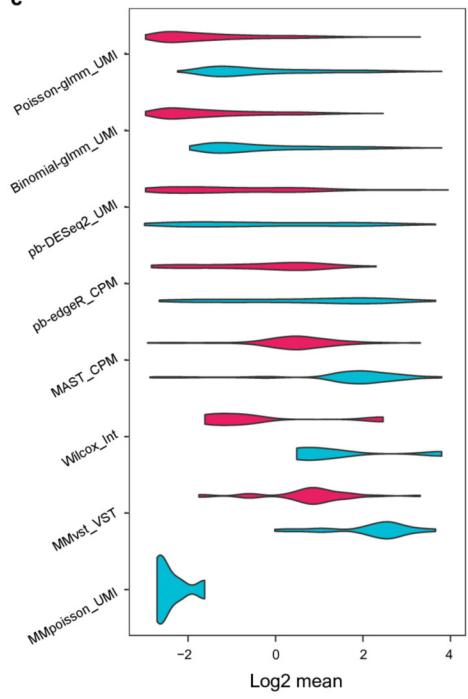




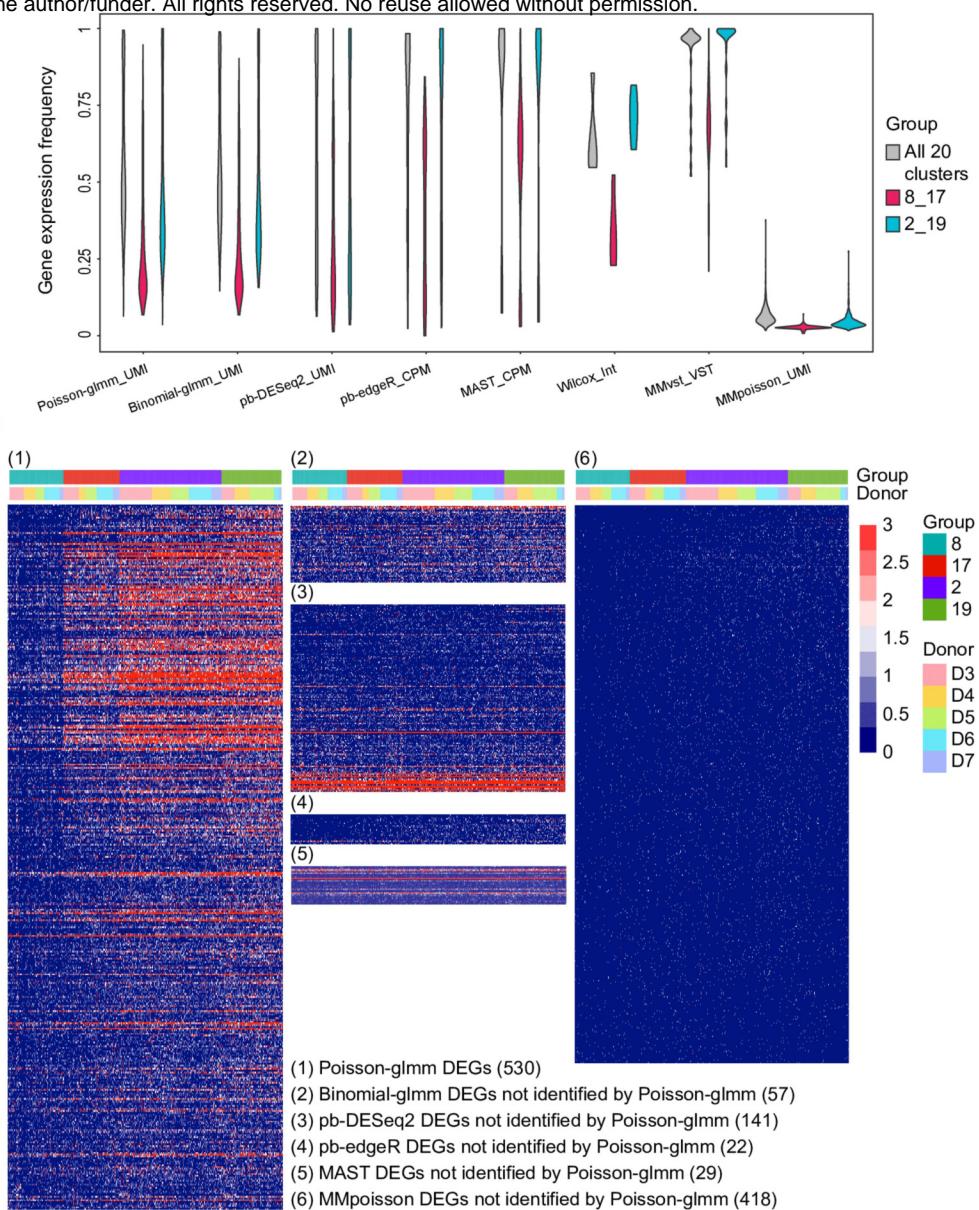
a



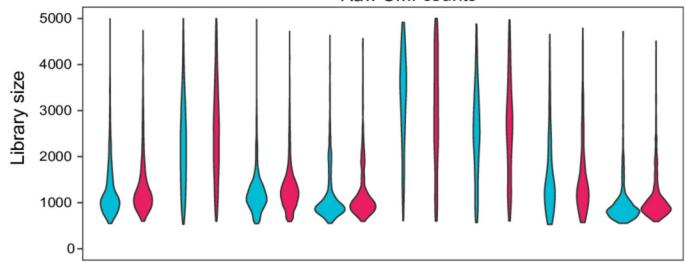
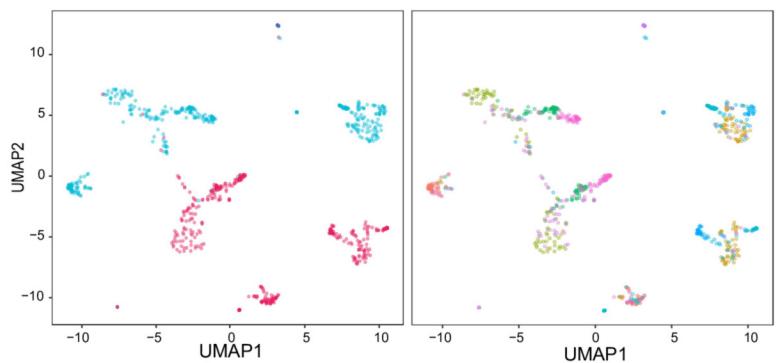
c



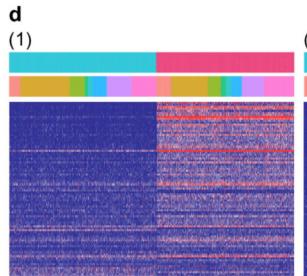
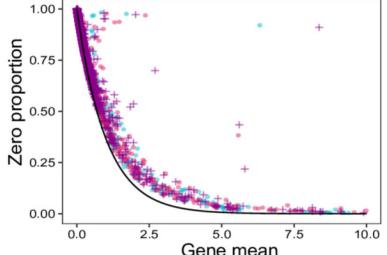
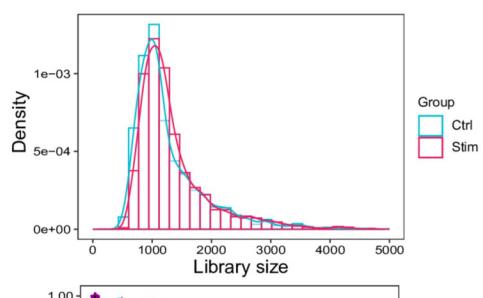
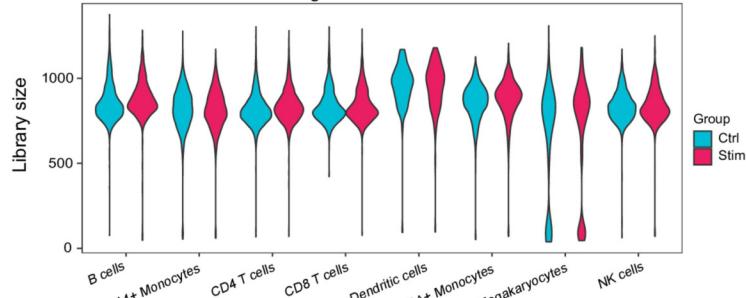
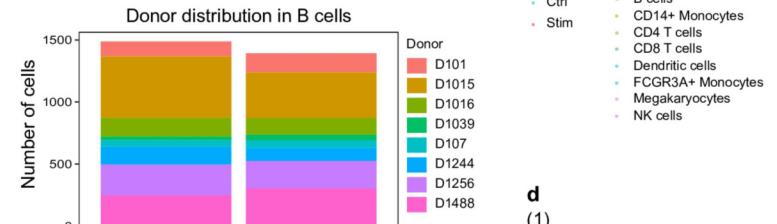
d



a

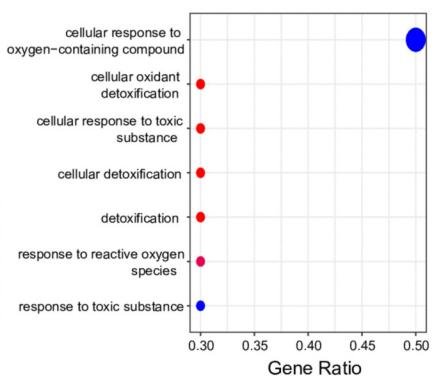
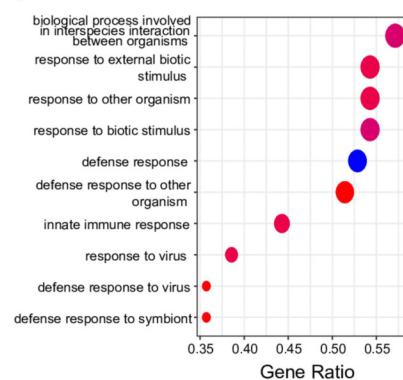


c



- (1) Poisson-glm DEGs (88)
- (2) pb-DESeq2 DEGs not identified by Poisson-glm (283)
- (3) MMpoisson DEGs not identified by Poisson-glm (220)

e



f

