



**University of
Zurich**^{UZH}

**Zurich Open Repository and
Archive**

University of Zurich
University Library
Strickhofstrasse 39
CH-8057 Zurich
www.zora.uzh.ch

Year: 2018

Bias, robustness and scalability in single-cell differential expression analysis

Soneson, Charlotte ; Robinson, Mark D

Abstract: Many methods have been used to determine differential gene expression from single-cell RNA (scRNA)-seq data. We evaluated 36 approaches using experimental and synthetic data and found considerable differences in the number and characteristics of the genes that are called differentially expressed. Prefiltering of lowly expressed genes has important effects, particularly for some of the methods developed for bulk RNA-seq data analysis. However, we found that bulk RNA-seq analysis methods do not generally perform worse than those developed specifically for scRNA-seq. We also present conquer, a repository of consistently processed, analysis-ready public scRNA-seq data sets that is aimed at simplifying method evaluation and reanalysis of published results. Each data set provides abundance estimates for both genes and transcripts, as well as quality control and exploratory analysis reports.

DOI: <https://doi.org/10.1038/nmeth.4612>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-186232>

Journal Article

Accepted Version

Originally published at:

Soneson, Charlotte; Robinson, Mark D (2018). Bias, robustness and scalability in single-cell differential expression analysis. *Nature Methods*, 15(4):255-261.

DOI: <https://doi.org/10.1038/nmeth.4612>

1 **Bias, robustness and scalability in differential expression analysis of**
2 **single-cell RNA-seq data**

3

4 Charlotte Sonesson^{1,2} and Mark D. Robinson^{1,2}

5

6 ¹ Institute of Molecular Life Sciences, University of Zurich, Zurich, Switzerland

7 ² SIB Swiss Institute of Bioinformatics, Zurich, Switzerland

8

9 Correspondence to C.S. (charlotte.soneson@uzh.ch) or M.D.R.

10 (mark.robinson@imls.uzh.ch)

11

12

13 **Abstract**

14 We perform an extensive evaluation of the performance and characteristics of 36
15 approaches for differential gene expression analysis in single-cell RNA-seq, using
16 both experimental and synthetic data. Considerable differences are found
17 between the methods in terms of the number and characteristics of the genes
18 that are called differentially expressed. Prefiltering of lowly expressed genes is
19 shown to have important effects on the results, particularly for some of the
20 methods originally developed for analysis of bulk RNA-seq data. Generally,
21 however, methods developed for bulk RNA-seq analysis do not perform notably
22 worse than those developed specifically for scRNA-seq. We also present *conquer*,
23 a repository of consistently processed, analysis-ready public single-cell RNA-seq
24 datasets, aimed at simplifying method evaluation and reanalysis of published
25 results. Each dataset provides abundance estimates for both genes and
26 transcripts, as well as quality control and exploratory analysis reports.

27

28

29 Keywords: single-cell RNA-seq, comparison, differential expression

30

31

32 Introduction

33 RNA-seq is used routinely to characterize transcriptomes, but until recently
34 sequencing libraries had to be prepared from pools of thousands or more cells,
35 and any measurement would represent an *average* across these cells. However,
36 recent advances enable library preparation from minute amounts of RNA and
37 thus profiling of the transcriptomes of individual cells¹⁻⁵. An increasing number
38 of such single-cell RNA-seq (scRNA-seq) datasets are being generated and
39 deposited in public repositories, which typically contain both raw read files and
40 processed data tables with, e.g., estimated gene abundances. Since the aims of
41 different studies vary widely, public datasets are often processed using very
42 different pipelines. Furthermore, the abundances may be represented in
43 different units and sometimes a fraction of the cells and/or genes are filtered
44 out. This can make reuse of the preprocessed public datasets, and especially
45 comparisons across datasets, challenging. To simplify this aspect, we have
46 developed *conquer*, a collection of consistently processed, analysis-ready public
47 scRNA-seq datasets. Each dataset has abundance estimates for all annotated
48 genes and transcripts, as well as quality assessment and exploratory analysis
49 reports to help users determine whether a particular dataset is suitable for their
50 purposes.

51
52 One of the most commonly performed computational tasks for RNA-seq data is
53 differential gene expression (DE) analysis. While well-established tools exist for
54 such analysis in bulk RNA-seq data⁶⁻⁸, methods for scRNA-seq data are just
55 emerging. Due to the special characteristics of scRNA-seq data, including
56 generally low library sizes, high noise levels and a large fraction of so-called
57 “dropout” events, it is unclear whether DE methods developed for bulk RNA-seq
58 are suitable also for scRNA-seq. A few recent studies have started to investigate
59 this question, suggesting that the optimal method choice may depend on the
60 number of cells and the strength of the signal⁹, and illustrating that also methods
61 that were not initially developed for RNA-seq analysis can perform well¹⁰. In this
62 study we use processed datasets, from *conquer* and other sources, to evaluate DE
63 methods in scRNA-seq data. Our study extends the previous comparisons to a
64 larger set of methods and a broader range of experimental datasets, and
65 additionally includes evaluations based on simulated data. We also investigate
66 the effect of filtering out lowly expressed genes and extend the set of employed
67 evaluation criteria. We focus on contrasting two predefined groups of cells since
68 this setup can be accommodated by all considered methods. However, it should
69 be noted that some scRNA-seq datasets contain cells from multiple subjects, or
70 from multiple plates, introducing a hierarchical variance structure that is not
71 accounted for by such a simple model¹¹. Moreover, single-cell measurements
72 allow additional questions that can not be addressed with bulk RNA-seq data,
73 such as testing whether different groups of cells show different levels of
74 variability or multimodality^{12,13}.

76 Results

77 Currently, *conquer* contains 36 datasets: 31 generated with full-length protocols
78 and 5 with 3'-end sequencing (UMI) protocols. With consistent processing and

79 representation of the datasets, we envision that *conquer* can be useful for a range
80 of applications. It can lower the barriers for evaluations and comparisons of
81 computational methods, for developers as well as end-users, and having easy
82 access to processed data is useful for teaching and tutorial construction. In
83 addition, *conquer* can be used for exploring the generality of biological
84 hypotheses across datasets from different species and cell types.

85
86 Seven datasets from *conquer* (six full-length and one UMI dataset) and two
87 additional UMI count datasets were used for the evaluation (Supplementary
88 Table 1, Supplementary Figures 1-2). We keep two predefined groups of cells
89 from each dataset, and generate multiple dataset instances with varying number
90 of cells. For eight datasets, we generate null datasets by subsampling from a
91 single group. Three datasets are used to simulate datasets with signal (10% of
92 the genes differentially expressed) as well as null datasets. For each instance, we
93 apply 36 DE approaches (Supplementary Table 2). Some methods failed to run
94 for certain datasets (Supplementary Figure 3), and these combinations are
95 excluded from the evaluations.

96

97 **Number of differentially expressed and non-tested genes**

98 Using all instances of the nine “signal” scRNA-seq datasets, we compare the
99 number of differentially expressed genes called by the different methods at an
100 adjusted p-value cutoff of 0.05 (Supplementary Figures 4-7). For full-length
101 datasets, SeuratBimod¹⁴ (without the default internal filtering) detects the
102 largest number of significant genes. edgeR/QLF^{7,15} detects large numbers of
103 genes if the dataset is not prefiltered to remove lowly expressed genes, but
104 shows the largest decrease in the number of significant genes after filtering
105 (Supplementary Figure 8). Conversely, SeuratBimod with non-zero expression
106 threshold, metagenomeSeq¹⁶ and scDD¹³ consistently detect few differentially
107 expressed genes. For UMI datasets, the performance of the methods based on the
108 voom transformation⁸ is highly variable without gene prefiltering.

109

110 Many DE methods implement internal filtering, which means that not all
111 quantified genes are actually being tested for DE. Such filtering is typically
112 performed to exclude lowly expressed genes and increase the power to detect
113 differences in the retained genes^{17,18}. For some methods, the model fitting
114 procedure can also fail to converge for some genes. While most evaluated
115 methods report valid results for all genes, some indeed exclude many genes if
116 run with default settings (Supplementary Figures 9-10). This is, however, not
117 specific to scRNA-seq data, and similar patterns can be seen if a subset of the
118 methods are applied to a large bulk RNA-seq dataset¹⁹ (Supplementary Figure
119 11). If the datasets are filtered before the DE analysis, the fraction of non-
120 reported results decreases, indicating that they mostly correspond to lowly
121 expressed genes.

122

123 **Type I error control**

124 Using the eight real null datasets, where no truly differential genes are expected,
125 we evaluate the type I error control by recording the fraction of tested genes that

126 are assigned a nominal p-value below 0.05 (Figure 1A). For unfiltered datasets,
127 many methods struggle to correctly control the type I error, and the best
128 performance is obtained by ROTS^{20,21} and SeuratTobit. Several of the other
129 methods are too liberal, with SeuratBimod and edgeR/QLF standing out with a
130 large number of false positive findings. Setting a non-zero expression threshold
131 in Seurat (SeuratBimodIsExpr2) improves the error control, but at the price of
132 detecting much fewer significant genes (Supplementary Figures 4-7). Conversely,
133 metagenomeSeq, scDD, SCDE²² and DESeq2⁶ on Census counts²³ instead control
134 the false positive rate well below the imposed level. Methods based on voom
135 mostly perform well, but sometimes the number of false positives is very high
136 (Supplementary Figure 12). For UMI datasets, monocle²⁴ performs best when
137 applied to transcript counts (monoclecount), whereas converting these values to
138 TPMs and applying a tobit model (monocle) deteriorates performance. For full-
139 length datasets, however, the TPM values lead to a slightly better performance
140 than the read counts. After filtering out lowly expressed genes (Figure 1B) the
141 performance of voom-limma, ROTSvoom and edgeR/QLF stabilizes and
142 improves, along with most other methods, while SeuratBimod still assigns low p-
143 values to a large fraction of the tested genes. P-value histograms further
144 illustrate that without filtering, few methods return uniformly distributed p-
145 values while after the applied filtering, results are considerably improved
146 (Supplementary Figures 13-14). The results are largely similar for the three
147 simulated datasets (Supplementary Figure 15).
148

149 **Characteristics of false positive genes**

150 To investigate the presence of biases in the DE calling, we use the eight
151 unfiltered real null datasets to characterize the set of genes that are (falsely)
152 called significant by the different methods. For each gene in each dataset
153 instance, we estimate the average, variance and coefficient of variation of the
154 CPM values across all cells as well as the fraction of cells in which the gene is
155 undetected. For each instance, and for each method calling at least five genes DE,
156 we calculate a signal-to-noise statistic comparing the values of each of the four
157 gene characteristics between the significant and non-significant (including non-
158 tested) genes (Figure 2, Supplementary Figure 16). The results show striking
159 differences between the types of genes detected by the different methods. False
160 positives of NODES²⁵, ROTS, SAMseq²⁶ and SeuratBimod have few zeros, high
161 expression and mostly a relatively low coefficient of variation. Conversely, false
162 positives of edgeR/QLF, SeuratTobit, MAST²⁷ and metagenomeSeq have
163 relatively many zeros. The same evaluation performed on the simulated datasets
164 shows largely similar results (Supplementary Figure 17).
165

166 **Between-method similarity**

167 Using the nine real scRNA-seq “signal” datasets we quantify the concordance
168 between gene rankings returned by different methods (within-method
169 consistency is investigated in Supplementary Figure 18). For each dataset we
170 calculate the area under the concordance curve (AUCC) for the top-ranked 100
171 genes for each pair of methods (Online Methods). Averaging the AUCCs across all
172 datasets and clustering based on the resulting similarities (Figure 3) shows, for

173 example, that while the four MAST modes give overall similar rankings, the
174 inclusion of the detection rate as a covariate has a larger effect on the rankings
175 than changing the type of expression values from CPMs to TPMs. Moreover, the
176 count-based bulk RNA-seq methods cluster together, as do some of the general
177 non-parametric methods (the Wilcoxon test and D3E²⁸), which are also similar to
178 the robust count-based methods and several approaches based on log-like
179 transformations of the data. The methods using Census transcript counts as
180 input give similar rankings. The degree of similarity between any given pair of
181 methods can vary widely across the dataset instances (Supplementary Figure
182 19), but for most method pairs, it is somewhat positively associated with the
183 number of cells per group (Supplementary Figure 20).
184

185 **FDR control and power**

186 Using the simulated datasets, we evaluate the false discovery rate control and
187 statistical power of the methods. Several methods, such as voom/limma,
188 ROTStpm, MAST, the methods applied to Census counts, SeuratTobit,
189 SeuratBimod with non-zero expression cutoff and SAMseq, robustly control the
190 FDR close to the imposed level (Figure 4A). SCDE, scDD, the t-test, D3E, limma-
191 trend^{8,29}, the Wilcoxon test, and the other variants of ROTS control the FDR at a
192 lower level than imposed. The worst FDR control for the unfiltered data is
193 obtained by monocle, SeuratBimod and edgeR/QLF. After filtering, edgeR/QLF
194 improves dramatically (Figure 4B), whereas MAST and SCDE yield even lower
195 false discovery proportions (FDPs). Most methods perform closer to the optimal
196 level for large sample sizes (Supplementary Figure 21). Adjusting the nominal p-
197 values for multiple testing using independent hypothesis weighting¹⁸ with the
198 average expression as covariate rather than using the values returned by the
199 respective methods has only minor impact (Supplementary Figure 22).
200

201 Practically all methods show increased power with increased sample size
202 (Figure 4C-D, Supplementary Figure 23). Among the methods with good, robust
203 FDR control after filtering, edgeR/QLF, SAMseq, DEsingle³⁰ and voom-limma
204 achieve high power, whereas for methods like metagenomeSeq, SeuratTobit,
205 SeuratBimodIsExpr2 and the methods applied to Census counts, the FDR control
206 comes at the price of reduced power. The power to detect true differences is
207 weakly related to the fraction of genes that are excluded by internal filtering
208 procedures (Supplementary Figure 24). However, DESeq2 and NODES achieve
209 high power despite strong filtering. The area under the ROC curve (AUROC),
210 indicating whether the methods are able to rank truly differentially expressed
211 genes ahead of truly non-differential ones, shows favourable performance of
212 edgeR, followed by MAST, limma (voom and trend), SCDE, DEsingle, DESeq2 and
213 SeuratBimod without filtering and the non-parametric methods (Figure 4E).
214 After prefiltering the rankings of most methods are improved (Figure 4F), and
215 the AUROC is typically higher for datasets with more cells (Supplementary
216 Figure 25).
217

218 Other aspects

219 As the number of cells that are studied in a dataset increases, computational
220 efficiency becomes important for method selection. For comparative purposes,
221 we ran all methods on a single core in this study. However, DESeq2, BPSC³¹,
222 MAST, SCDE, scDD and monocle all feature explicit arguments to take advantage
223 of parallelization, and methods that perform gene-wise tests without
224 information sharing between genes, such as the Wilcoxon test, the t-test and
225 D3E, can be run in parallel after splitting the data into chunks. Four dedicated
226 single-cell methods, namely BPSC, DEsingle, D3E and SCDE, are the slowest for
227 most datasets, while the bulk methods (edgeR, DESeq2 and especially the limma
228 variants) are generally faster (Supplementary Figure 26A). Most single-cell
229 methods (with the exception of SCDE) scale well with increasing number of cells,
230 while the computational time required for the bulk RNA-seq methods is more
231 sample size dependent (Supplementary Figures 26B, 27-31).

232
233 While the evaluations in this study are centered on the simplest experimental
234 situation, comparing two groups of cells, many real studies require a more
235 complex experimental design, which not all evaluated methods can
236 accommodate. Specifically, the Wilcoxon test, the t-test, scDD, NODES, SCDE,
237 Seurat, ROTS, DEsingle and D3E are limited to two-group comparisons, while
238 SAMseq can perform a limited number of analysis types. The remaining methods
239 implement statistical frameworks that can accommodate more complex (fixed
240 effect) designs, including comparisons across multiple groups and adjustments
241 for batch effects and other covariates.

242
243 Other important aspects are the availability and documentation of the software
244 packages. Most methods are available either via Bioconductor³² or CRAN, or via a
245 public GitHub repository (Supplementary Table 2). NODES was obtained via a
246 Dropbox link provided by the authors. The Bioconductor packages have
247 extensive documentation, including help pages for individual functions and a
248 vignette to guide the user through a typical workflow, all tested to work with the
249 current version of the package. Some packages, such as Seurat, D3E, monocle and
250 SCDE, have dedicated webpages with instructions for users, examples and
251 tutorials.

252

253 Discussion

254 We have presented an extensive evaluation and comparison of methods for DE
255 analysis of scRNA-seq data, using mainly real datasets from *conquer*, a repository
256 of consistently processed public single-cell RNA-seq datasets. The fact that
257 *conquer* provides gene expression estimates in multiple units allowed us to
258 compare methods requiring different types of input values, and also to
259 investigate the effect of using different input values for the same method. We
260 have shown that prefiltering of genes is essential to obtain good, robust
261 performance for several of the evaluated methods, most notably edgeR/QLF,
262 which tends to call lowly expressed genes with many zeros significant if these
263 are present in the data but otherwise performs well, and voom-limma, which
264 also performs more robustly after filtering out lowly expressed genes.

265
266 We noted a large variability among the number of genes called differential with
267 the different methods, as well as in the ability to control the type I error rate and
268 the false discovery rate. After appropriate filtering, a subset of the methods
269 managed to control the FDR and FPR close to the imposed level while achieving a
270 high power while for many other methods, appropriate error control was
271 associated with a lack of power.

272
273 We also showed that the DE methods are biased in different ways in terms of the
274 types of genes they preferentially detect as differential, which can have
275 important implications in practical applications. In agreement with previous
276 evaluations, methods originally developed for bulk RNA-seq analysis did not
277 perform worse than methods specifically developed for scRNA-seq data, but
278 sometimes showed a stronger dependence on the data being appropriately
279 prefiltered.

280
281 Figure 5 summarizes the performance of the different methods across the main
282 evaluation criteria in our study. For each evaluation aspect, each method was
283 classified as “good”, “intermediate” or “poor” (Online Methods). While it is
284 difficult to capture the full complexity of the evaluation in a crude categorization,
285 the table provides a convenient summary of our results and can be used to select
286 an appropriate method based on the criteria that are most important for a
287 specific application.

288
289 The number of cells per group ranged between 6 and 400 in our datasets. While
290 these are relatively small numbers compared to the thousands of cells that can
291 be sequenced in an actual experiment, DE analysis is typically used to compare
292 sets of homogeneous cells (e.g., from given, well-defined cell types), and these
293 collections are likely to be much smaller. Thus, we believe that the range of
294 sample sizes considered in our comparisons are relevant for real applications
295 and that it is important to know how the methods perform under these
296 circumstances.
297

298 **Acknowledgements**

299 The authors would like to acknowledge M Love and V Svensson for helpful online
300 instructions regarding automated download of raw data from ENA. This study
301 was supported by the Forschungskredit of the University of Zurich, grant no. FK-
302 16-107 to C.S.

303 **Author Contributions**

304 C.S. and M.D.R. designed analyses and wrote the manuscript. C.S. performed
305 analyses. Both authors have read and approved the final manuscript.

306 **Competing interests**

307 The authors declare that they have no competing interests.

309 **References**

- 310 1. Tang, F. *et al.* mRNA-Seq whole-transcriptome analysis of a single cell. *Nat*
311 *Methods* **6**, 377–382 (2009).
- 312 2. Picelli, S. *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in
313 single cells. *Nat Methods* **10**, 1096–1098 (2013).
- 314 3. Klein, A. M. *et al.* Droplet barcoding for single-cell transcriptomics applied to
315 embryonic stem cells. *Cell* **161**, 1187–1201 (2015).
- 316 4. Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of
317 Individual Cells Using Nanoliter Droplets. *Cell* **161**, 1202–1214 (2015).
- 318 5. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single
319 cells. *Nat Commun* **8**, 14049 (2017).
- 320 6. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and
321 dispersion for RNA-seq data with DESeq2. *Genome Biol* **15**, 550 (2014).
- 322 7. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package
323 for differential expression analysis of digital gene expression data.
324 *Bioinformatics* **26**, 139–140 (2010).
- 325 8. Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: precision weights unlock
326 linear model analysis tools for RNA-seq read counts. *Genome Biol* **15**, R29
327 (2014).
- 328 9. Miao, Z. & Zhang, X. Differential expression analyses for single-cell RNA-Seq:
329 old questions on new data. *Quant Biol* **4**, 243–260 (2016).
- 330 10. Jaakkola, M. K., Seyednasrollah, F., Mehmood, A. & Elo, L. L. Comparison of
331 methods to detect differentially expressed genes between single-cell
332 populations. *Brief Bioinform* bbw057 (2016).

- 333 11. Lun, A. T. L. & Marioni, J. C. Overcoming confounding plate effects in
334 differential expression analyses of single-cell RNA-seq data. *Biostatistics*
335 (2017).
- 336 12. Vallejos, C. A., Richardson, S. & Marioni, J. C. Beyond comparisons of
337 means: understanding changes in gene expression at the single-cell level.
338 *Genome Biol* **17**, 70 (2016).
- 339 13. Korthauer, K. D. *et al.* A statistical approach for identifying differential
340 distributions in single-cell RNA-seq experiments. *Genome Biol* **17**, 222 (2016).
- 341 14. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial
342 reconstruction of single-cell gene expression data. *Nat Biotechnol* **33**, 495–502
343 (2015).
- 344 15. Lun, A. T. L., Chen, Y. & Smyth, G. K. It's DE-licious: A Recipe for
345 Differential Expression Analyses of RNA-seq Experiments Using Quasi-
346 Likelihood Methods in edgeR. in *Statistical Genomics* (eds. Mathé, E. & Davis,
347 S.) 391–416 (Springer New York, 2016).
- 348 16. Paulson, J. N., Stine, O. C., Bravo, H. C. & Pop, M. Differential abundance
349 analysis for microbial marker-gene surveys. *Nat Methods* **10**, 1200–1202
350 (2013).
- 351 17. Bourgon, R., Gentleman, R. & Huber, W. Independent filtering increases
352 detection power for high-throughput experiments. *Proc Natl Acad Sci U A* **107**,
353 9546–9551 (2010).
- 354 18. Ignatiadis, N., Klaus, B., Zaugg, J. B. & Huber, W. Data-driven hypothesis
355 weighting increases detection power in genome-scale multiple testing. *Nat*
356 *Methods* **13**, 577–580 (2016).

- 357 19. Lappalainen, T. *et al.* Transcriptome and genome sequencing uncovers
358 functional variation in humans. *Nature* **501**, 506–511 (2013).
- 359 20. Elo, L. L., Filén, S., Lahesmaa, R. & Aittokallio, T. Reproducibility-optimized
360 test statistic for ranking genes in microarray studies. *IEEEACM Trans Comput
361 Biol Bioinform* **5**, 423–431 (2008).
- 362 21. Seyednasrollah, F., Rantanen, K., Jaakkola, P. & Elo, L. L. ROTS:
363 reproducible RNA-seq biomarker detector-prognostic markers for clear cell
364 renal cell cancer. *Nucleic Acids Res* **44**, e1 (2016).
- 365 22. Kharchenko, P. V., Silberstein, L. & Scadden, D. T. Bayesian approach to
366 single-cell differential expression analysis. *Nat Methods* **11**, 740–742 (2014).
- 367 23. Qiu, X. *et al.* Single-cell mRNA quantification and differential analysis with
368 Census. *Nat Methods* **14**, 309–315 (2017).
- 369 24. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are
370 revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381–
371 386 (2014).
- 372 25. Sengupta, D., Rayan, N. A., Lim, M. & Lim, B. Fast, scalable and accurate
373 differential expression analysis for single cells. *bioRxiv* (2016).
374 doi:10.1101/049734
- 375 26. Li, J. & Tibshirani, R. Finding consistent patterns: a nonparametric
376 approach for identifying differential expression in RNA-Seq data. *Stat Methods
377 Med Res* **22**, 519–536 (2013).
- 378 27. Finak, G. *et al.* MAST: A flexible statistical framework for assessing
379 transcriptional changes and characterizing heterogeneity in single-cell RNA-
380 seq data. *Genome Biol* **16**, 278 (2015).

- 381 28. Delmans, M. & Hemberg, M. Discrete distributional differential expression
382 (D3E)–a tool for gene expression analysis of single-cell RNA-seq data. *BMC*
383 *Bioinformatics* **17**, 110 (2016).
- 384 29. Smyth, G. K. Linear Models and Empirical Bayes Methods for Assessing
385 Differential Expression in Microarray Experiments. *Stat Appl Genet Mol Biol* **3**,
386 Article 3 (2004).
- 387 30. Miao, Z. & Zhang, X. DEsingle: A new method for single-cell differentially
388 expressed genes detection and classification. *bioRxiv* (2017).
389 doi:10.1101/173997
- 390 31. Vu, T. N. *et al.* Beta-Poisson model for single-cell RNA-seq data analyses.
391 *Bioinformatics* **32**, 2128–2135 (2016).
- 392 32. Huber, W. *et al.* Orchestrating high-throughput genomic analysis with
393 Bioconductor. *Nat Methods* **12**, 115–121 (2015).

394

395 **Figure legends**

396

397 **Figure 1**

398 Type I error control across several instances from eight single-cell null datasets,
399 with a range of sample sizes. Values are split between full-length and UMI
400 datasets, and the methods are ordered by the median FPR across all datasets
401 (separately for unfiltered and prefiltered datasets). A. Without any prefiltering of
402 genes (only excluding genes with zero counts across all cells). B. After filtering,
403 retaining only genes with an estimated expression above 1 TPM in more than
404 25% of the cells. Only methods returning nominal p-values are included. The
405 black line indicates the target FPR=0.05, and the y-axis is square-root
406 transformed for increased visibility. Center line, median; hinges, first and third
407 quartiles; whiskers, most extreme values within 1.5 IQR from the box; *n*, number
408 of data set instances.
409

410 **Figure 2**

411 Characteristics of genes falsely called significant by the evaluated methods. For
412 each instance of the eight real scRNA-seq null datasets, we record characteristics
413 of each gene (average CPM, variance and coefficient of variation of CPM, fraction
414 zeros across all cells) and use a signal-to-noise statistic to compare each of these
415 characteristics between genes called significant and the rest of the genes. A
416 positive statistic indicates that the corresponding characteristic is more
417 pronounced in the set of genes called significant than in the remaining genes.
418 Note that ROTSvroom, D3E, limma-trend, the t-test and the Wilcoxon test did not
419 return enough false positive findings to be included in the evaluation. Center line,
420 median; hinges, first and third quartiles; whiskers, most extreme values within
421 1.5 IQR from the box; n , number of data set instances.
422

423 **Figure 3**

424 Dendrogram illustrating the average similarities between the gene rankings
425 obtained by the evaluated methods. The dendrogram is obtained by complete-
426 linkage hierarchical clustering based on the matrix of average AUCC values
427 across all datasets. The labels of the internal nodes represent their stability
428 across datasets, in terms of the fraction of instances where they are observed.
429 Only nodes with stability scores of at least 0.1 are labeled. The colored boxes
430 below the methods represent characteristics of the methods.
431

432 **Figure 4**

433 Differential expression detection performance, summarized across all instances
434 of the three simulated datasets. The methods are stratified by their ability to
435 control the FDR at the 0.05 level across the datasets. A method where more than
436 75% of the observed FDPs are above 0.05 or where the median FDP is above
437 0.15 is considered to have “high FDP”, whereas a method where more than 75%
438 of the observed FDPs are below 0.05 or where the median FDP is below 0.0167 is
439 considered to have “low FDP”. A-B. Observed FDP at an adjusted p-value cutoff at
440 0.05. The horizontal line represents the target FDR of 0.05, and the y-axis is
441 square-root transformed for increased visibility. C-D. Observed TPR at an
442 adjusted p-value cutoff at 0.05. E-F. Observed area under the ROC curve. Center
443 line, median; hinges, first and third quartiles; whiskers, most extreme values
444 within 1.5 IQR from the box; n , number of data set instances.
445

446 **Figure 5**

447 Summary of the performance of the evaluated methods across all major
448 evaluation criteria in the current study. A description of the criteria and the
449 cutoff values for assigning a method to a performance category is available in the
450 Online Methods. The methods are ranked by their average performance across
451 the criteria, with the numerical encoding good=2, intermediate=1, poor=0.
452 NODES and SAMseq do not return nominal p-values and are therefore not
453 evaluated in terms of the FPR.
454

455 Online Methods

456 *conquer*

457 The *conquer* pipeline processes (sc)RNA-seq datasets using the steps outlined in
458 Supplementary Table 3, including quality control, abundance estimation,
459 exploratory analysis and summarization.

460
461 Many of the processed datasets contain not only scRNA-seq samples (single
462 cells), but also bulk RNA-seq samples for comparison, or technical control
463 samples. Whenever these could be identified, they are excluded from the
464 processed data. A list of the excluded samples is provided in the online
465 repository. Cells belonging to the same SRA/GEO dataset but sequenced on
466 different platforms are separated into different repository entries. No filtering
467 based on poor quality or low abundance is performed, since that may introduce
468 unwanted biases for certain downstream analyses and since no universally
469 adopted filtering approach or threshold currently exists. However, the provided
470 quality control and exploratory analysis reports can be used to determine
471 whether some cells need to be excluded for specific applications. The Ensembl
472 catalog (v38)³³ was used as reference when processing the currently available
473 datasets. Information about the underlying reference is also included as
474 metadata in the processed datasets and displayed in the exploratory report.
475 Since TPMs and read counts are estimated using the same reference annotation,
476 with the same software and using the same data, the *conquer* datasets can be
477 used to compare computational methods that require different types of input,
478 with minimal bias. The processed datasets and the resulting reports can be
479 browsed and downloaded from <http://imlspenticton.uzh.ch:3838/conquer/>, and
480 the underlying code used to process all datasets is available from
481 <https://github.com/markrobinsonuzh/conquer>.

482
483

484 Evaluation of differential expression methods

485 Experimental and simulated data

486 Seven of the real datasets from *conquer*, with a large number of cells, are selected
487 as the basis for the evaluation of DE analysis methods. For each of the datasets,
488 we retain only cells from two of the annotated cell groups (Supplementary Table
489 1), attempting to select large and relatively homogeneous populations among the
490 ones annotated by the data generators. The selected datasets span a wide
491 spectrum of signal strengths and population homogeneities (Supplementary
492 Figures 1 and 2). For each dataset, we then generate one instance of “maximal”
493 size (with the number of cells per group equal to the size of the smallest of the
494 two selected cell populations) and several subsets with fewer cells per group by
495 random subsampling from the maximal size subset (see Supplementary Table 1
496 for exact group sizes). For each non-maximal sample size, we generate five
497 replicate dataset instances, and thus each original dataset contribute 11-21
498 separate instances, depending on the number of different sample sizes
499 (Supplementary Table 1). Moreover, for each dataset with enough cells we
500 generate *null* datasets with different sample sizes (again, five instances per

501 sample size except for the maximal size) by sampling randomly from one of the
502 two selected cell populations. Finally, three of the datasets (GSE45719,
503 GSE74596 and GSE60749-GPL13112) are used as the basis for simulation of data
504 using a slightly modified version of the *powsim* R package³⁴. Individual reports
505 generated by *countsimQC*³⁵ and verifying the similarity between the simulated
506 and real datasets across a range of aspects are provided as Supplementary Data.
507 As for the original, experimental datasets, we subsample dataset instances with
508 varying number of cells per group, and further generate null datasets by random
509 sampling from one of the simulated groups. In each simulated dataset, 10% of
510 the genes are selected to be differentially expressed between the two groups,
511 with fold changes sampled from a Gamma distribution with shape 4 and rate 2.
512 The direction of the DE is randomly determined for each gene, with equal
513 probability of up- and downregulation. Mean and dispersion parameters used as
514 basis for the simulations are estimated from the respective real datasets using
515 edgeR⁷. For each of the three datasets, the rounded length-scaled TPMs for all
516 genes with at least two non-zero counts are used as input to the simulator, and a
517 dataset with the same number of genes is generated. The counts for each
518 simulated gene are based on one of the original genes (however, the same
519 original gene can be the basis for more than one simulated gene), and by
520 retaining this information we can link average transcript lengths (calculated by
521 *tximport*³⁶ for the original data) to each simulated gene, and thus estimate
522 approximate TPMs also for the simulated data.

523
524 In addition to the seven datasets from *conquer*, we downloaded and processed
525 two additional UMI datasets. First, the UMI counts corresponding to the GEO
526 entry GSE59739³⁷ were downloaded from <http://linnarssonlab.org/drg/>
527 (accessed December 18, 2016). The provided UMI RPMs were used in the place
528 of TPMs, and were combined with the provided information about the total
529 number of reads per cell to generate gene counts. Empty wells were filtered out.
530 Second, we downloaded UMI count matrices for C14+ monocytes and cytotoxic
531 T-cells processed with the 10X Genomics GemCode protocol⁵
532 (<https://support.10xgenomics.com/single-cell-gene-expression/datasets>,
533 accessed September 17, 2017). For this dataset, as well as for the UMI dataset
534 obtained from *conquer* (GSE62270-GPL17021), the UMI counts were used as
535 “raw counts” in the DE analysis, and since these counts are supposed to be
536 proportional to the concentration of transcript molecules, we estimated the TPM
537 by scaling the UMI counts to sum to 1 million. Although this may be suboptimal
538 due to the low capture efficiency of single-cell protocols, it allows us to apply
539 methods consistently across full-length and UMI datasets.

540
541 For comparison, we also downloaded a bulk RNA-seq dataset from the Geuvadis
542 project¹⁹ from <http://www.ebi.ac.uk/arrayexpress/experiments/E-GEUV-1/> and
543 estimated gene expression levels using the same pipeline as for the *conquer*
544 datasets. For this dataset, we perform DE analysis using a subset of the methods
545 applied to the single-cell RNAseq datasets, comparing samples from the CEU and
546 YRI populations generated at the University of Geneva.

547
548 For each real and simulated dataset, we perform the DE analysis evaluation both
549 on the full, “unfiltered”, dataset (excluding only genes with 0 counts in all

550 considered cells) and on a filtered dataset, where we retain only genes with an
551 estimated TPM above 1 in more than 25% of the considered cells. Depending on
552 the dataset and the number of considered cells, between 4 and 50% of the genes
553 are retained after this filtering (Supplementary Figure 32).
554

555 Differential expression analysis methods

556
557 For each of the real and simulated scRNA-seq datasets, we apply 36 statistical
558 approaches for DE analysis to compare the expression levels in the two groups of
559 cells (Supplementary Table 2). As representatives for methods developed for
560 differential analysis of bulk RNA-seq data, we include edgeR⁷, DESeq2⁶, voom-
561 limma⁸ and limma-trend⁸. For edgeR, we apply both the likelihood ratio test
562 (LRT)³⁸ and the more recent quasi-likelihood approach (QLF)¹⁵. For the LRT, in
563 addition, we use both the default dispersion estimates³⁹ and the robust
564 dispersion estimates developed to address outlier counts⁴⁰, and we apply edgeR
565 both with the default TMM normalization⁴¹ and with the recently developed
566 deconvolution normalization approach for scRNA-seq⁴². In addition, we run
567 edgeR/QLF including the cellular detection rate (the fraction of detected genes
568 per cell) as a covariate. DESeq2 is run in three modes, after rounding the length-
569 scaled TPM values to integers: with default settings, without the log-fold change
570 shrinkage (beta prior), and after disabling the internal independent filtering and
571 outlier detection and replacement. Additionally, both edgeR/LRT and DESeq2
572 are applied to both the read counts (length-scaled TPMs as described above) and
573 Census transcript counts²³, aimed at converting relative abundances such as
574 TPMs into transcript counts, based on the assumption that the most common
575 signal among the genes detectable with current single-cell library preparation
576 protocols corresponds to a single molecule. The Census counts are calculated
577 from the estimated TPMs using monocle²⁴ with default settings. We note that it is
578 possible that modifications of these settings, optimized for the library
579 preparation parameters for each individual dataset, would lead to different
580 absolute count values, and thus potentially altered performance, in some of the
581 datasets.

582
583 Three non-parametric methods are included in the comparison: SAMseq²⁶, the
584 Wilcoxon test⁴³ and NODES²⁵. SAMseq is applied to the length-scaled TPMs,
585 while the Wilcoxon test is applied to TPM estimates after applying TMM
586 normalization to address the compositionality of the TPMs. NODES was initially
587 run in two modes: with default settings, and after disabling the internal filtering
588 steps. However, disabling the internal filtering caused the method to fail in
589 subsequent steps, and thus we retain only the runs with default settings.

590
591 We include a broad range of methods developed specifically for scRNA-seq DE
592 analysis. BPSC³¹ is applied to CPMs (calculated using edgeR) as suggested by the
593 package authors. D3E²⁸ is run with the method-of-moments approach to
594 parameter estimation, the non-parametric Cramer-von Mises test to compare
595 distributions and without removing zeros before the analysis. MAST²⁷ is applied
596 to both $\log_2(\text{CPM}+1)$ and $\log_2(\text{TPM}+1)$ values, both with and without including
597 the cellular detection rate (the fraction of genes that are detected with non-zero

598 counts) as a covariate in the model. For monocle²⁴, the input is either TPM
599 estimates (with a tobit model), raw counts (read counts or UMI counts,
600 depending on the dataset, with a Negative Binomial model) or Census counts
601 (with a Negative Binomial model), calculated from the TPMs as for edgeR and
602 DESeq2 above. SCDE²² is applied to rounded length-scaled TPMs, following the
603 instructions provided in the package documentation, and p-values are calculated
604 from the provided z-scores. Seurat¹⁴ is applied using either the default “bimod”
605 likelihood ratio test⁴⁴ (applied to the length-scaled TPMs, which are log-
606 normalized internally), both with default settings and disabling the internal
607 filtering steps, as well as after setting the internal expression threshold to 2
608 instead of the default of 0, or the “tobit” test²⁴ (applied to the TPMs). scDD¹³ was
609 applied to counts normalized with the median normalization, and using the
610 default “fast” procedure based on the Kolmogorov-Smirnov test, without
611 permutations. We applied DEsingle³⁰ to rounded counts.

612

613 Given the similarities between single-cell RNA-seq data and operational
614 taxonomic unit (OTU) count data from 16S marker studies in metagenomics
615 applications, we also apply metagenomeSeq¹⁶ to the count values, fitting the
616 zero-inflated log-normal model using the *fitFeatureModel* function from the
617 metagenomeSeq package and testing for differences in abundance.

618

619 Finally, we include ROTS (reproducibility-optimized test statistic)^{20,21}, which is a
620 general test, originally developed for microarray data, in which a t-like test
621 statistic is optimized for reproducibility across bootstrap resamplings. We apply
622 ROTS to CPM and TPM values, as well as to the log-transformed CPM values
623 calculated by the voom function in the limma package⁸. For comparison, we also
624 apply a Welch t-test⁴⁵ to TMM-normalized TPM values, after adding 1 and
625 applying a log-transformation.

626

627 All code used for the DE analysis and evaluation is accessible via
628 https://github.com/csoneson/conquer_comparison.

629

630 Evaluation strategies

631 Most of the evaluations in this study are performed using real, experimental
632 data, where no independently validated truth is available. The advantage of this
633 approach is that no assumptions or restrictions are made regarding data
634 distributions or specific structures of the data. However, the set of evaluation
635 measures is more limited than in situations where the ground truth is accessible.
636 Our first battery of evaluation approaches aim to catalog the number of genes
637 found to be significantly differentially expressed, as well as the number and
638 characteristics of the false positive detections from each method. For the latter
639 evaluations we use the null datasets, where no truly differential genes are
640 expected and thus all significant genes are false positives. First, we investigate
641 the fraction of genes for which no interpretable test results are returned by the
642 applied methods (e.g., due to internal filtering or convergence failure of fitting
643 procedures). Then, for all methods returning nominal p-values, we calculate the
644 fraction of performed tests that give a nominal p-value below 0.05. For a well-
645 calibrated test, this fraction should be around 5%. Next, we calculate

646 characteristics such as the expression level (CPM), the fraction of zero counts
647 and the expression variability (variance and coefficient of variation for CPM
648 estimates) for all genes, and compare these characteristics between genes called
649 differentially expressed (with an adjusted p-value/FDR threshold of 0.05) and
650 genes not considered DE, for each of the methods. More precisely, for each
651 characteristic and for each method detecting at least five differentially expressed
652 genes at this threshold, we calculate a signal-to-noise statistic:
653

$$\frac{\mu_S - \mu_{NS}}{\sigma_S + \sigma_{NS}}$$

654
655 where μ_S (μ_{NS}) and σ_S (σ_{NS}) represent the mean and standard deviation of the
656 gene characteristic among the significant (nonsignificant) genes. Genes with non-
657 interpretable test results (e.g., NA adjusted p-values) are considered non-
658 significant in this evaluation. This approach gives insights into the inherent
659 biases of the different methods, in the sense of the type of genes that are
660 preferentially called significantly differential. Note that since the evaluation is
661 done on the null datasets, the results are not confounded by the characteristics
662 of *truly* differentially expressed genes.

663
664 The second type of evaluations focus on *robustness* of methods when applied to
665 different subsets of the same dataset. In a dataset where there is a true
666 underlying signal (i.e., truly differential genes between cell populations), ideally,
667 this signal will be detected regardless of the set of cells that are sampled for the
668 analysis. Thus, a high concordance between results obtained from different
669 subsets of the cells is positive, and indicative of robust performance. For a
670 dataset without truly differential genes, however, any detections should be
671 random, and a high similarity between results obtained from different subsets
672 can rather indicate a bias in the DE calling. Thus, we first calculate a measure of
673 concordance between the gene rankings from each pair of instances of a dataset
674 with the same number of cells per group (five such instances were generated for
675 each group size, giving 10 pairwise comparisons). Then, we match “signal” and
676 null instances from the same original dataset and with the same number of cells
677 per group, and compare the robustness values between signal and null instances.
678 A large difference indicates a significant difference between the cross-instance
679 concordance in a dataset with a true underlying signal and a dataset without a
680 true signal, suggesting that the method is able to robustly detect underlying
681 effects, and that this robustness is not due to a strong bias in the significance
682 testing. As a measure of concordance, we use the area under the *concordance*
683 *curve* for the top- K genes ranked by significance, with $K=100$ (cf. Irizarry *et al.*⁴⁶).
684 More precisely, for each dataset instance and each DE method, we rank the genes
685 by statistical significance (nominal p-value or adjusted p-value). Then, for each
686 pair of dataset instances with the same sample size, for $k=1, \dots, K$, we count the
687 number of genes that are ranked among the top k in both the corresponding
688 rankings. Plotting the number of shared genes against k gives a curve, and the
689 area under this curve is used as a measure of the concordance. To obtain more
690 interpretable values, we divide the calculated area with the maximal possible
691 value ($K^2/2$). Thus, a normalized value of 1 indicates that the two compared
692 rankings are identical, whereas a value of 0 indicates that the sets of top- K genes

693 from the two rankings don't share any genes. The rationale for using this type of
694 concordance index to evaluate robustness is that it is independent of the number
695 of genes that are actually called significant (which can vary widely across
696 methods), and it is applicable to situations where not all compared rankings
697 have interpretable results for the same sets of genes (e.g., due to different
698 internal filtering criteria), which would cause a problem for e.g. overall
699 correlation estimation. Furthermore, as opposed to a simple intersection of the
700 top- K genes in the two rankings, the concordance score incorporates the actual
701 ranking of these top- K genes.

702
703 A similar approach is used to evaluate similarities between methods. Briefly, for
704 each dataset instance, we rank the genes by significance using each of the DE
705 methods. Then, for each pair of methods, we construct a concordance curve and
706 calculate the area under this curve as a measure of similarity between the results
707 from the two methods. This evaluation is only performed on the “signal”
708 datasets.

709
710 Finally, we use the simulated data to evaluate false discovery rate (FDR) control
711 and true positive rate (TPR, power), as well as the area under the receiver
712 operating characteristic (ROC) curve, indicating the ability of a method to rank
713 truly differential genes ahead of truly non-differential ones. For the prefiltered
714 datasets, we limit the evaluation to the genes retained after the filtering.

715
716 An interesting aspect, although not strictly related to performance, is the
717 computational time requirement for the different methods. We investigate two
718 aspects of this: first, the actual time required to run each method using a single
719 core. Since this depends on the size of the dataset, we normalize all times for a
720 given dataset instance so that the maximal value across all methods is 1. Thus, a
721 “relative” computational time of 1 for a given method and a given dataset
722 instance means that this method was the slowest one for that particular instance,
723 and a value of, e.g., 0.1 means that the time requirement was 10% of that for the
724 slowest method. Second, we investigate how the computational time
725 requirement scales with the number of cells. This is particularly important for
726 scRNA-seq data, since the number of cells sequenced per study is now increasing
727 rapidly⁴⁷. For this, we consider all instances of all datasets (“signal” and null, as
728 well as simulated data), and divide them into 10 equally sized bins depending on
729 the total number of tested genes. Within each such bin, we model the required
730 time T as a function of the number of cells per group (N) as

$$731 \quad \quad \quad 732 \quad \quad \quad T = aN^p,$$

733
734 and record the estimated value of p .

735

736 **Performance summary criteria**

737 Figure 5 summarizes the performance of the evaluated methods across the range
738 of evaluation metrics. For each metric, the performance of each method is
739 considered either “good”, “intermediate” or “poor”. Metrics that are mainly

740 descriptive rather than quantitative are excluded from the summary. Here, we
741 list the criteria used to categorize the methods for each evaluation metric:

742

743 **MedianFDP.** Evaluated after filtering, across all simulated signal datasets

- 744 - Good: no more than 75% of FDPs on one side (above or below) of 0.05
745 and $0.0167 < \text{median FDP} < 0.15$
- 746 - Intermediate: $0.15 \leq \text{median FDP} < 0.25$ or $0.01 < \text{median FDP} \leq 0.0167$,
747 or $0.0167 < \text{median FDP} < 0.15$ but more than 75% of FDPs on one side of
748 0.05
- 749 - Poor: $\text{median FDP} \geq 0.25$ or $\text{median FDP} \leq 0.01$

750

751 **MaxFDP.** Evaluated after filtering, across all simulated signal datasets

- 752 - Good: maximal FDP < 0.15
- 753 - Intermediate: $0.15 \leq \text{maximal FDP} < 0.35$
- 754 - Poor: maximal FDP ≥ 0.35

755

756 **TPR.** Evaluated after filtering, across all simulated signal dataset instances with
757 more than 20 cells

- 758 - Good: median TPR > 0.8
- 759 - Intermediate: $0.6 < \text{median TPR} \leq 0.8$
- 760 - Poor: median TPR ≤ 0.6

761

762 **AUROC.** Evaluated after filtering, across all simulated signal datasets

- 763 - Good: median AUC > 0.8
- 764 - Intermediate: $0.65 < \text{median AUC} \leq 0.8$
- 765 - Poor: median AUC ≤ 0.65

766

767 **MedianFPR.** Evaluated after filtering, across all real null datasets, separately for
768 full-length and UMI datasets

- 769 - Good: $|\log_2(\text{median FPR}/0.05)| < \log_2(1.5)$
- 770 - Intermediate: $\log_2(1.5) \leq |\log_2(\text{median FPR}/0.05)| < 2$
- 771 - Poor: $|\log_2(\text{median FPR}/0.05)| \geq 2$

772

773 **MaxFPR.** Evaluated after filtering, across all real null datasets, separately for full-
774 length and UMI datasets

- 775 - Good: maximal FPR < 0.1
- 776 - Intermediate: $0.1 \leq \text{maximal FPR} < 0.25$
- 777 - Poor: maximal FPR ≥ 0.25

778

779 **Scalability.** Evaluated based on all datasets

- 780 - Good: median exponent in power model of timing vs number of cells < 0.5
- 781 - Intermediate: $0.5 \leq \text{median exponent in power model of timing vs number}$
782 of cells < 1
- 783 - Poor: median exponent in power model of timing vs number of cells ≥ 1

784

785 **Speed.** Evaluated based on all datasets

- 786 - Good: median relative computation time requirement (relative to slowest
787 method) < 0.1

- 788 - Intermediate: $0.1 \leq$ median relative computation time requirement
789 (relative to slowest method) < 0.7
790 - Poor: median relative computation time requirement (relative to slowest
791 method) ≥ 0.7
792

793 **BiasDEG.** Evaluated based on all unfiltered real null datasets

- 794 - Good: No false positive genes detected, or $|\text{median SNR}| < 0.5$ for all four
795 SNR statistics (for fraction of zeros, CV(CPM), $\log_2(\text{average CPM})$ and
796 $\log_2(\text{variance(CPM)})$)
797 - Intermediate: $|\text{median SNR}| \geq 0.5$ for at least one statistic, but $|\text{median}$
798 SNR < 1 for all four statistics
799 - Poor: $|\text{median SNR}| \geq 1$ for at least one statistic
800

801 **Consistency.** Evaluated after filtering

- 802 - Good: The t-statistic of robustness values between signal and null
803 datasets is > 2 for GSE60749-GPL13112 and 10XMonoCytoT, and all t-
804 statistics are ≥ 0
805 - Intermediate: Any of the t-statistics for GSE60749-GPL13112 or
806 10XMonoCytoT is ≤ 2 , but all t-statistics (across all real datasets for which
807 both signal and datasets are available) are ≥ 0
808 - Poor: The t-statistic for any dataset is < 0
809

810 **ComplexDesign**

- 811 - Good: The method allows arbitrary complex (fixed) designs
812 - Intermediate: The method can accommodate a limited set of designs
813 - Poor: The method only performs two-group comparisons
814

815 **FailureRate.** Evaluated across all datasets

- 816 - Good: Average failure rate < 0.01
817 - Intermediate: $0.01 \leq$ Average failure rate < 0.25
818 - Poor: Average failure rate ≥ 0.25
819

820 **Software specifications and code availability**

821 The datasets currently available in the *conquer* repository were processed with
822 Salmon v0.6.0-v0.8.2⁴⁸, FastQC v0.11.6.devel and MultiQC v0.8⁴⁹. All analyses for
823 the method evaluation were run in R v3.3⁵⁰, with Bioconductor v3.4³², except for
824 scDD and DEsingle, which required R 3.4 and Bioconductor v3.5. Performance
825 indices were calculated with iCOBRA v1.2.0⁵¹ when applicable, and results were
826 visualized using ggplot2 v2.2.1⁵². All code used to process the datasets for
827 *conquer* can be accessed via GitHub:

828 <https://github.com/markrobinsonuzh/conquer>. The code used to perform the
829 evaluation of the DE analysis methods is also available from GitHub:

830 https://github.com/csoneson/conquer_comparison. The results of the
831 evaluation can be browsed in a shiny application available at

832 http://imlspenticton.uzh.ch:3838/scrnaseq_de_evaluation/.
833

834 Data availability

835 All public datasets included in *conquer* can be downloaded from
836 <http://imlspenticton.uzh.ch:3838/conquer/>. The processed abundances for the
837 UsoskinGSE59739 dataset were downloaded from <http://linnarssonlab.org/drg/>
838 on December 18, 2016. The UMI count matrices for the 10X MonoCytoT dataset
839 were downloaded from [https://support.10xgenomics.com/single-cell-gene-](https://support.10xgenomics.com/single-cell-gene-expression/datasets)
840 [expression/datasets](https://support.10xgenomics.com/single-cell-gene-expression/datasets) on September 17, 2017. All processed datasets used for the
841 evaluation (listed in Supplementary Table 1) can be downloaded as a
842 compressed archive from the accompanying website:
843 http://imlspenticton.uzh.ch/robinson_lab/conquer_de_comparison/. Figures 1,
844 2, 4 and 5 have associated source data.

846 References

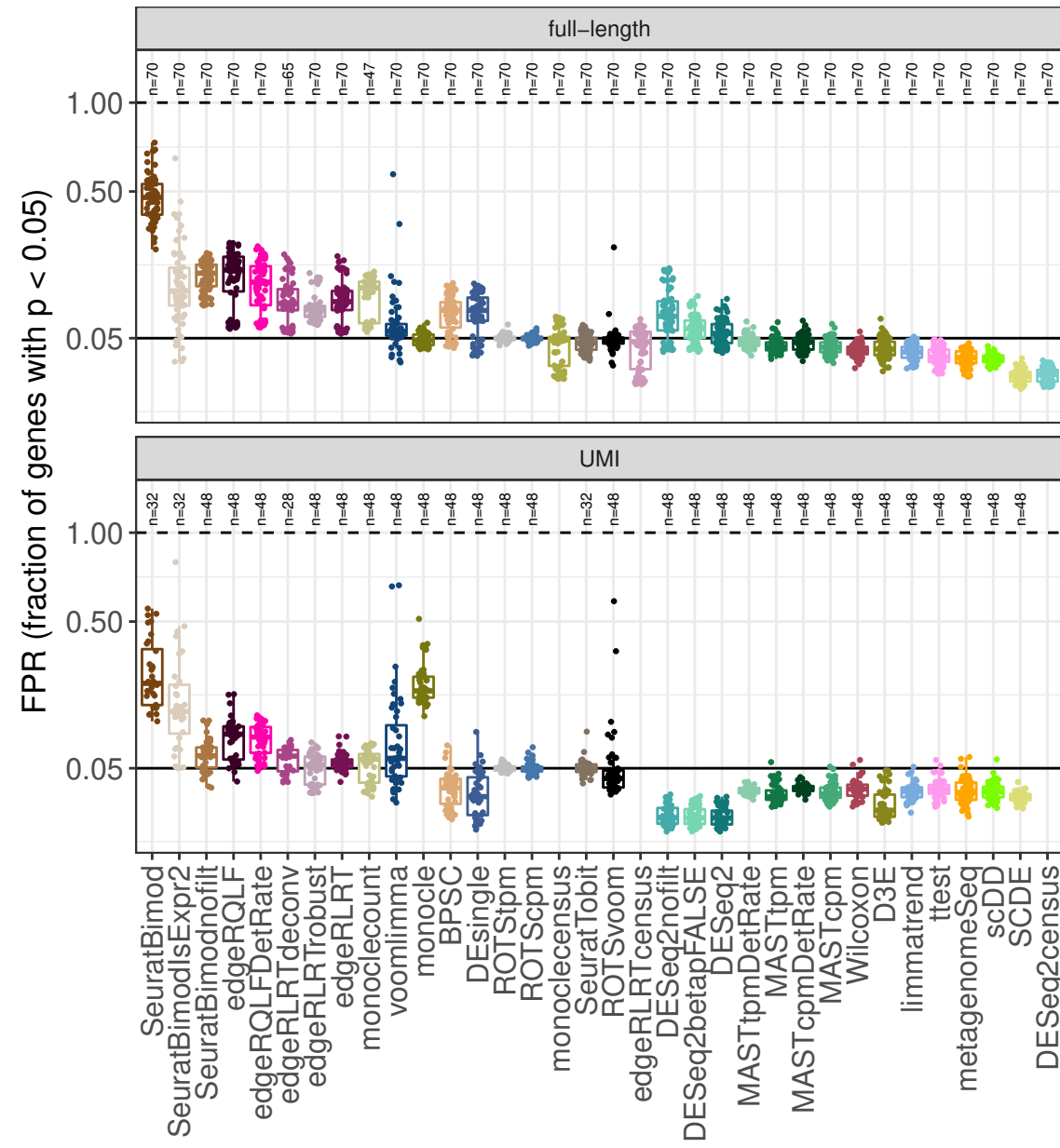
- 847 33. Aken, B. L. *et al.* The Ensembl Gene Annotation System. *Database* baw093
848 (2016).
- 849 34. Vieth, B., Ziegenhain, C., Parekh, S., Enard, W. & Hellmann, I. powsim:
850 Power analysis for bulk and single cell RNA-seq experiments. *bioRxiv* (2017).
851 doi:10.1101/117150
- 852 35. Soneson, C. & Robinson, M. D. Towards unified quality verification of
853 synthetic count data with countsimQC. *Bioinformatics* btx631-btx631 (2017).
854 doi:10.1093/bioinformatics/btx631
- 855 36. Soneson, C., Love, M. I. & Robinson, M. D. Differential analyses for RNA-
856 seq: transcript-level estimates improve gene-level inferences. *F1000Res* **4**,
857 1521 (doi: 10.12688/f1000research.7563.1) (2015).
- 858 37. Usoskin, D. *et al.* Unbiased classification of sensory neuron types by large-
859 scale single-cell RNA sequencing. *Nat Neurosci* **18**, 145–153 (2015).
- 860 38. McCarthy, D. J., Chen, Y. & Smyth, G. K. Differential expression analysis of
861 multifactor RNA-Seq experiments with respect to biological variation. *Nucleic*
862 *Acids Res* **40**, 4288–4297 (2012).

- 863 39. Chen, Y., Lun, A. T. L. & Smyth, G. K. Differential Expression Analysis of
864 Complex RNA-seq Experiments Using edgeR. in *Statistical Analysis of Next*
865 *Generation Sequencing Data* (eds. Datta, S. & Nettleton, D.) 51–74 (Springer
866 International Publishing, 2014).
- 867 40. Zhou, X., Lindsay, H. & Robinson, M. D. Robustly detecting differential
868 expression in RNA sequencing data using observation weights. *Nucleic Acids*
869 *Res* **42**, e91 (2014).
- 870 41. Robinson, M. D. & Oshlack, A. A scaling normalization method for
871 differential expression analysis of RNA-seq data. *Genome Biol* **11**, R25 (2010).
- 872 42. Lun, A. T. L., Bach, K. & Marioni, J. C. Pooling across cells to normalize
873 single-cell RNA sequencing data with many zero counts. *Genome Biol* **17**, 75
874 (2016).
- 875 43. Wilcoxon, F. Individual comparisons by ranking methods. *Biom. Bull.* **1**,
876 80–83 (1945).
- 877 44. McDavid, A. *et al.* Data exploration, quality control and testing in single-
878 cell qPCR-based gene expression experiments. *Bioinformatics* **29**, 461–467
879 (2013).
- 880 45. Welch, B. A. The generalization of Student's problem when several
881 different population variances are involved. *Biometrika* **34**, 28–35 (1947).
- 882 46. Irizarry, R. A. *et al.* Multiple-laboratory comparison of microarray
883 platforms. *Nat Methods* **2**, 345–350 (2005).
- 884 47. Svensson, V., Vento-Tormo, R. & Teichmann, S. A. Moore's Law in Single
885 Cell Transcriptomics. *arXiv:1704.01379v1* (2017).

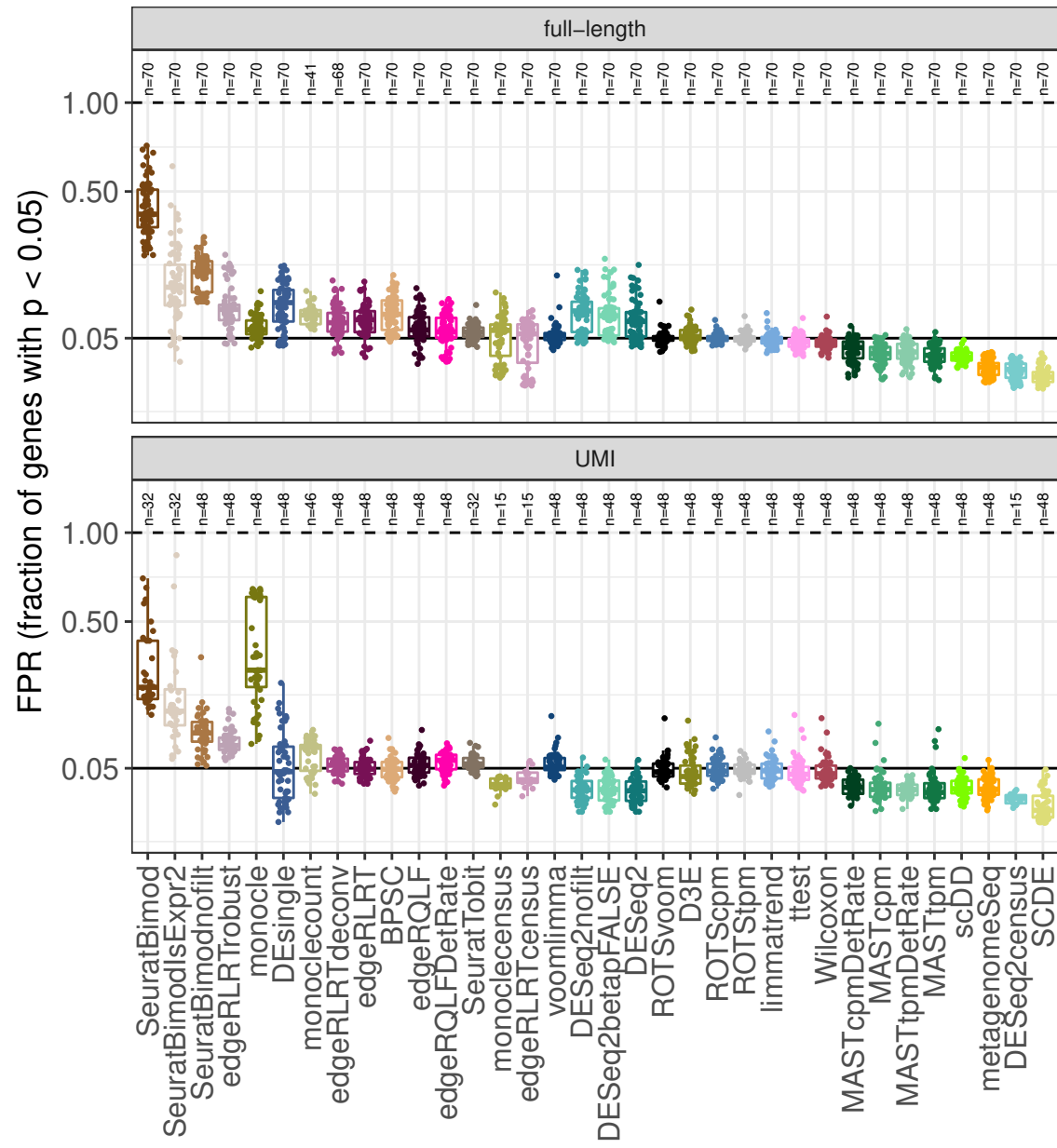
- 886 48. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon
887 provides fast and bias-aware quantification of transcript expression. *Nat*
888 *Methods* **14**, 417–419 (2017).
- 889 49. Ewels, P., Magnusson, M., Lundin, S. & Källér, M. MultiQC: Summarize
890 analysis results for multiple tools and samples in a single report.
891 *Bioinformatics* btw354 (2016).
- 892 50. R Core Team. *R: A Language and Environment for Statistical Computing*. (R
893 Foundation for Statistical Computing, 2016).
- 894 51. Sonesson, C. & Robinson, M. D. iCOBRA: open, reproducible, standardized
895 and live method benchmarking. *Nat Methods* **13**, 283 (2016).
- 896 52. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag
897 New York, 2009).
- 898

A

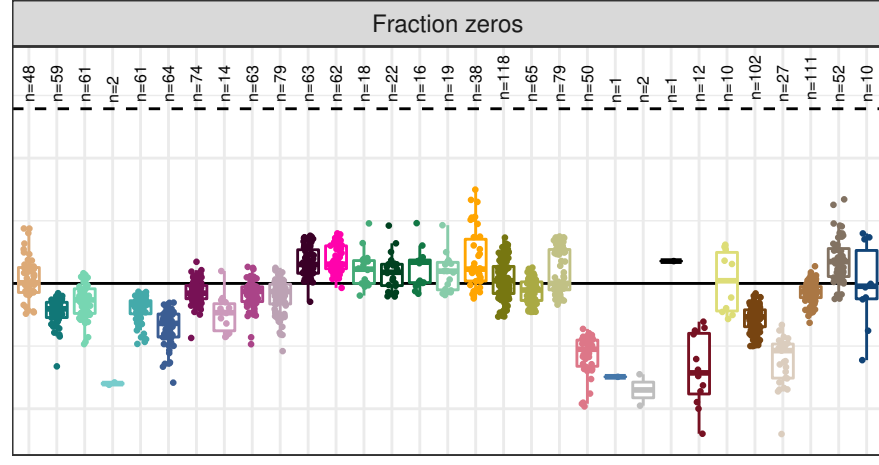
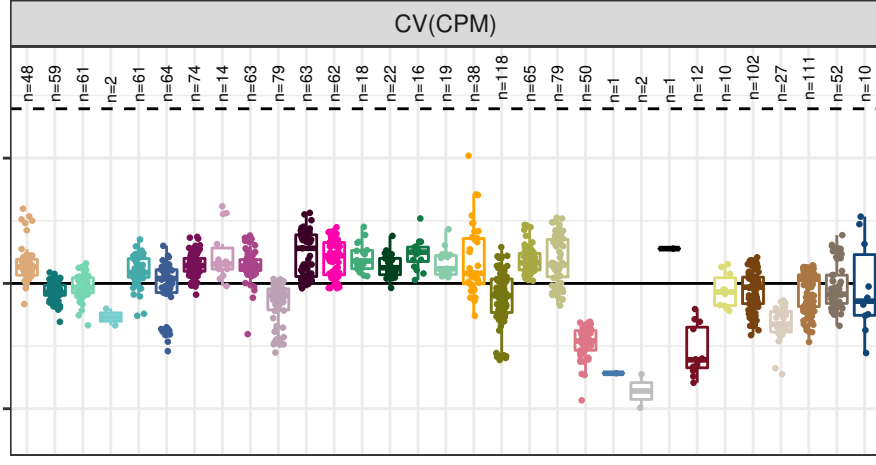
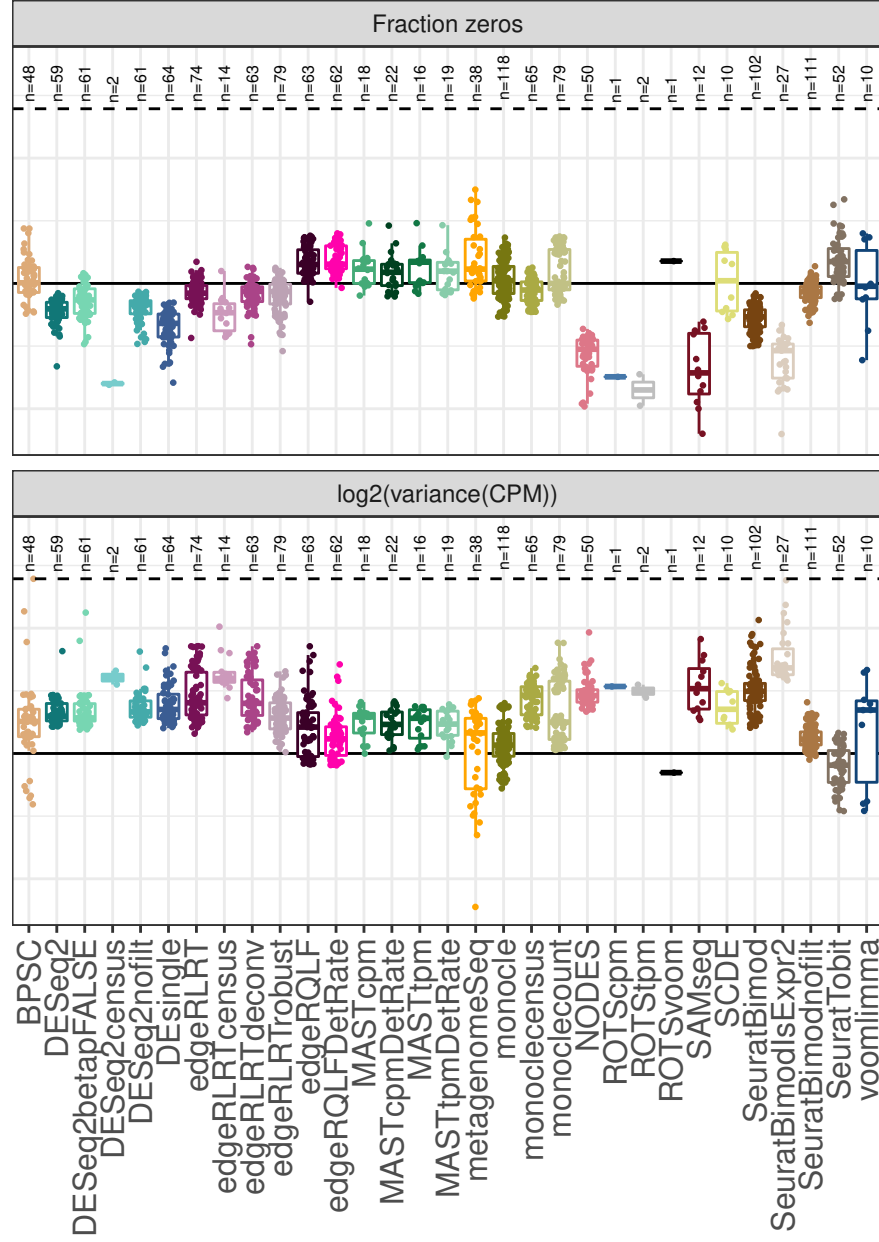
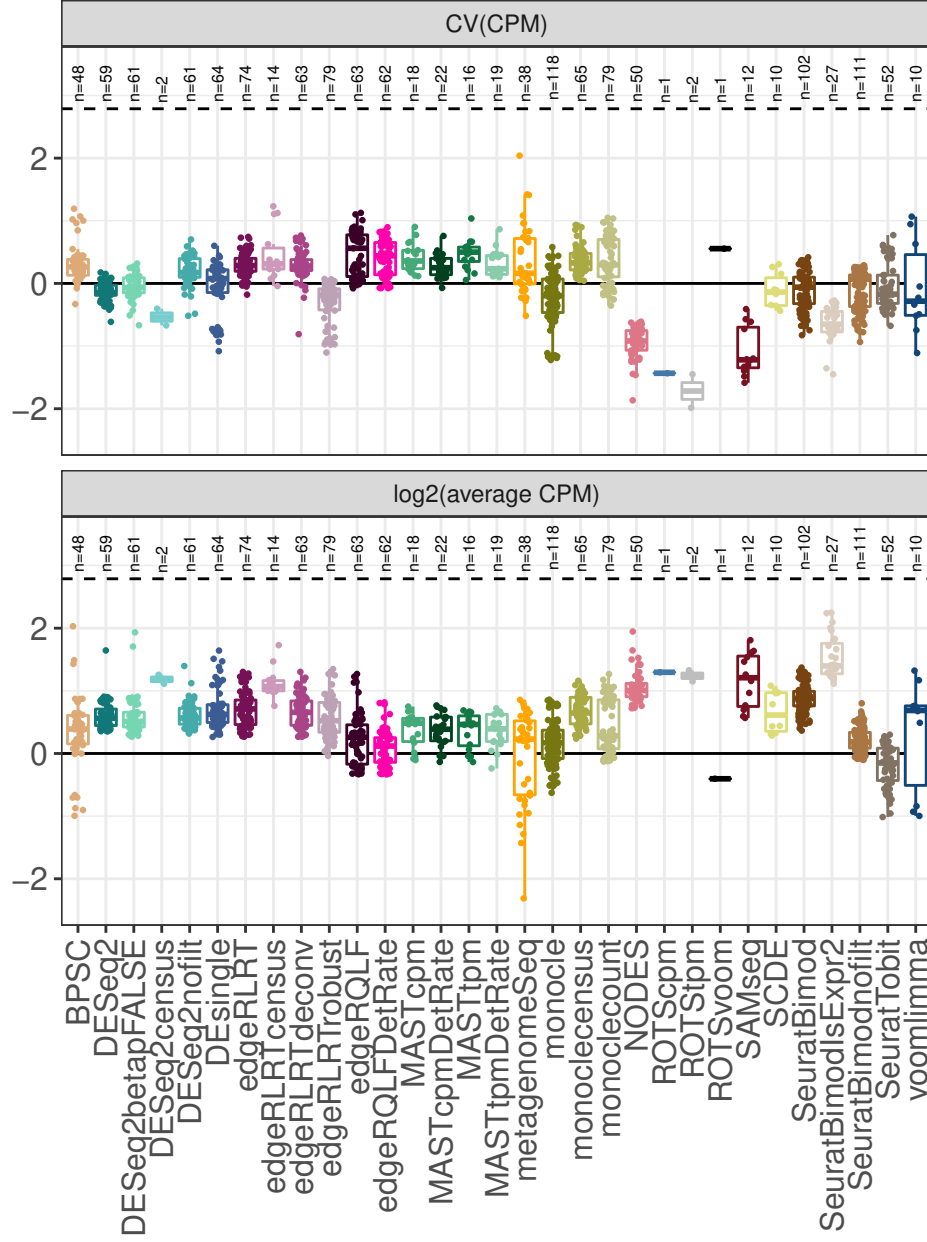
Without filtering

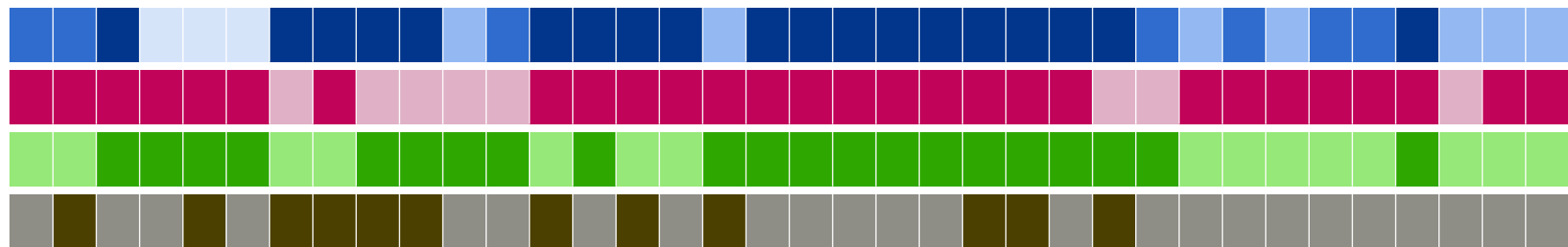
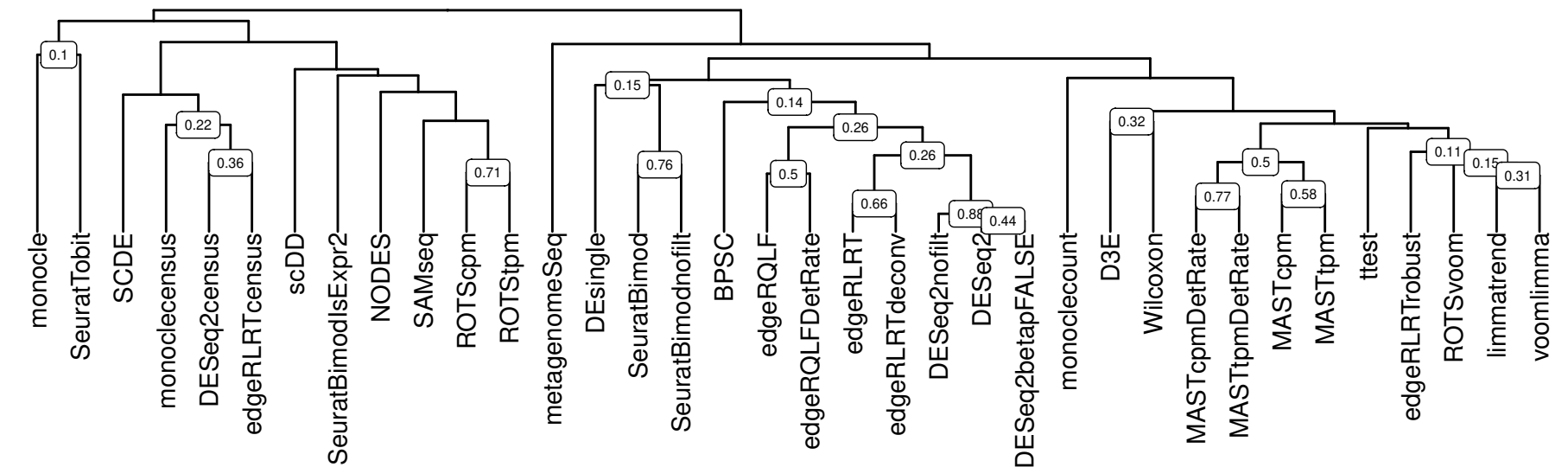
**B**

After filtering

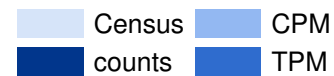


signal-to-noise statistic comparing significant and non-significant genes

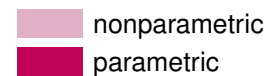




Input



Modeling



Transformation



NA values



