

ARCHIVIO ISTITUZIONALE DELLA RICERCA

Alma Mater Studiorum Università di Bologna Archivio istituzionale della ricerca

distinct: a novel approach to differential distribution analyses

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

Published Version: distinct: a novel approach to differential distribution analyses / Simone Tiberi; Helena L Crowell; Pantelis Samartsidis; Lukas M Weber; Mark D Robinson. - In: THE ANNALS OF APPLIED STATISTICS. - ISSN 1932-6157. - ELETTRONICO. - 17:2 (June)(2023), pp. 1681-1700. [10.1214/22-AOAS1689]

This version is available at: https://hdl.handle.net/11585/906917 since: 2022-11-28 *Published:* DOI: http://doi.org/10.1214/22-AOAS1689

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

(Article begins on next page)

This item was downloaded from IRIS Università di Bologna (https://cris.unibo.it/). When citing, please refer to the published version. This is the final peer-reviewed accepted manuscript of:

Simone Tiberi, Helena L. Crowell, Pantelis Samartsidis, Lukas M. Weber, Mark D. Robinson. (2023). *"distinct*: A novel approach to differential distribution analyses". *Annals of Applied Statistics*, Vol. 17, Issue 2, June 2023, pp. 1681-1700.

The final published version is available online at: https://doi.org/10.1214/22-AOAS1689

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<u>https://cris.unibo.it/</u>)

When citing, please refer to the published version.

distinct: a novel approach to differential distribution analyses

Simone Tiberi^{1*}, Helena L Crowell¹, Pantelis Samartsidis², Lukas M Weber³ and Mark D Robinson¹

¹Department of Molecular Life Sciences and Swiss Institute of Bioinformatics, University of Zurich, Zurich, Switzerland.

²MRC Biostatistics Unit, School of Clinical Medicine, University of Cambridge, Cambridge, UK.

³Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA.

* e-mail: Simone.Tiberi@uzh.ch

¹ Abstract

² We present *distinct*, a general method for dif-³ ferential analysis of full distributions that is 4 well suited to applications on single-cell data, 5 such as single-cell RNA sequencing and highdimensional flow or mass cytometry data. High-7 throughput single-cell data reveal an unprece-8 dented view of cell identity and allow com-⁹ plex variations between conditions to be discov-¹⁰ ered; nonetheless, most methods for differential ¹¹ expression target differences in the mean and ¹² struggle to identify changes where the mean is only marginally affected. *distinct* is based on 13 14 a hierarchical non-parametric permutation ap-¹⁵ proach and, by comparing empirical cumulative 16 distribution functions, identifies both differen-¹⁷ tial patterns involving changes in the mean, as well as more subtle variations that do not in-18 volve the mean. We performed extensive bench-19 marks across both simulated and experimen-20 tal datasets from single-cell RNA sequencing 21 22 and mass cytometry data, where *distinct* shows 23 favourable performance, identifies more differ-24 ential patterns than competitors, and displays ²⁵ good control of false positive and false discovery ²⁶ rates. *distinct* is available as a Bioconductor R 27 package.

28 keywords: Differential distribution; Differential anal29 yses; Differential state; High-throughput single-cell
30 data; Single-cell RNA-seq; Single-cell flow and mass cy31 tometry; Permutation tests.

32 Background

³³ Technology developments in the last decade have led to
³⁴ an explosion of high-throughput single-cell data, such
³⁵ as single-cell RNA sequencing (scRNA-seq) and high³⁶ dimensional flow or mass cytometry data, allowing re-

³⁷ searchers to investigate biological mechanisms at single-38 cell resolution. Single-cell data have also extended the 39 canonical definition of differential expression by dis-40 playing cell-type specific responses across conditions, ⁴¹ known as differential state (DS) [32], where genes or 42 proteins vary in specific sub-populations of cells (e.g., 43 a cytokine response in myeloid cells but not in other ⁴⁴ leukocytes [13]). Classical bulk differential expression 45 methods have been shown to perform well when used 46 on single-cell measurements [25, 26, 31] and on aggre-47 gated data (i.e., averages or sums across cells), also re-⁴⁸ ferred to as pseudo-bulk (PB) [7, 32]. However, most 49 bulk and PB tools focus on shifts in the means, and 50 may conceal information about cell-to-cell heterogene-51 ity. Indeed, single-cell data can show more complex ⁵² variations (Figure 1 and Supplementary Figure 1); such 53 patterns can arise due to increased stochasticity and 54 heterogeneity, for example owing to oscillatory and un-⁵⁵ synchronized gene expression between cells, or when 56 some cells respond differently to a treatment than oth-57 ers [15, 31]. In addition to bulk and PB tools, other 58 methods were specifically proposed to perform differ-⁵⁹ ential analyses on single-cell data (notably: *scDD* [15], 60 SCDE [14], MAST [11], BASiCS [10,29,30] and mixed ⁶¹ models [27]). Nevertheless, they all present significant 62 limitations: BASiCS does not perform cell-type spe-63 cific differential testing between conditions, scDD does 64 not directly handle covariates and biological replicates, ⁶⁵ while PB, SCDE, MAST and mixed models performed 66 poorly in previous benchmarks when detecting differ-⁶⁷ ential patterns that do not involve the mean [7, 15].

68 Results

⁶⁹ distinct's full distribution approach

To overcome these challenges, we developed *distinct*, a
flexible and general statistical methodology to perform
differential analyses between groups of distributions.

73 distinct is particularly suitable to compare groups of 74 samples (i.e., biological replicates) on single-cell data.

75 Our approach computes the empirical cumulative distribution function (ECDF) from the individual (e.g., 76 ⁷⁷ single-cell) measurements of each sample, and compares ECDFs to identify changes between full distributions, 78 even when the mean is unchanged or marginally involved (Figure 1 and Supplementary Figure 1). First, 80 we compute the ECDF of each individual sample; then, 81 we build a fine grid and, at each cut-off, we average the 82 ECDFs within each group, and compute the absolute 83 difference between such averages. A test statistic, s^{obs} , ⁸⁵ is obtained by adding these absolute differences.

86 More formally, assume we are interested in compar-⁸⁷ ing two groups, that we call A and B, for which N_A $_{88}$ and N_B samples are available, respectively. The ECDF so for the *i*-th sample in the *j*-th group, is denoted by 90 $ecdf_{i}^{(j)}(.)$, for $j \in \{A, B\}$ and $i = 1, ..., N_{j}$. We ⁹¹ then define K equally spaced cut-offs between the mini-⁹² mum, *min*, and maximum, *max*, values observed across 93 all samples: b_1, \ldots, b_K , where $b_k = min + k \times l$, for 94 k = 1, ..., K, with l = (max - min)/(K + 1) being 95 the distance between two consecutive cut-offs. We ex-96 clude *min* and *max* from the cut-offs because, trivially, 97 $ecdf_i^{(j)}(min) = 0$ and $ecdf_i^{(j)}(max) = 1, \forall j, i.$ At ev-98 ery cut-off, we compute the absolute difference between ⁹⁹ the mean ECDF in the two groups; our test statistic, s^{obs} , is obtained by adding these differences across all 101 cut-offs:

$$s^{obs} = \sum_{k=1}^{K} \left| \frac{\sum_{i=1}^{N_A} ecdf_i^{(A)}(b_k)}{N_A} - \frac{\sum_{i=1}^{N_B} ecdf_i^{(B)}(b_k)}{N_B} \right|.$$
(1)

¹⁰² Note that in differential state analyses, these operations are repeated for every gene-cluster combination. 103

105 107 108 users, is set to 25 by default, because no detectable 130 discovery rates. 109 difference in performance was observed when further 110 increasing it (data not shown). Note that, although at 131 Importantly, *distinct* is general and flexible: it targets 111 112 113 114 116 119 distribution of s^{obs} is then estimated via a hierarchical 139 ticular, distinct fits a linear mixed effects model with



Figure 1: Cumulative distribution functions (CDFs) unravel differences between distributions. Density (left panels) and CDF (right panels) of five differential patterns: differential variability (DV), and the four proposed by Korthauer et. al. [15]: differential expression (DE), differential proportion (DP), differential modality (DM), and both differential modality and different component means (DB).

¹²⁰ non-parametric permutation approach (see Methods). 121 A major disadvantage of permutation tests, which of-122 ten restricts its usage on biological data, is that too ¹²³ few permutations are available from small samples. We 124 overcome this by permuting cells, which is still pos-104 Intuitively, s^{obs} , which ranges in $[0, \infty)$, approximates 125 sible in small samples, because there are many more the area between the average ECDFs, and represents 126 cells than samples. In principle, this may lead to an a measure of distance between two groups of densities: 127 inflation of false positives due to lack of exchangabilthe bigger s^{obs} , the greater the distance between groups. 128 ity (see Methods); nonetheless, in our analyses, distinct The number of cut-offs K, which can be defined by 129 provides good control of both false positive and false

each cut-off we compute the average across each group's 132 complex changes between groups, explicitly models biocurves, ECDFs are computed separately for each indi-133 logical replicates within a hierarchical framework, does vidual sample, therefore our approach still accounts for 134 not rely on asymptotic theory, avoids parametric asthe within-group variability; indeed, at a given thresh- 135 sumptions, and can be applied to arbitrary types of old, the average of the sample-specific ECDFs differs 136 data. Additionally, distinct can also adjust for samplefrom the group-level ECDF (i.e., the curve based on 137 level cell-cluster specific covariates (i.e., whose effect all individual measurements from the group). The null 138 varies across cell clusters), such as batch effects. In par140 the input data (e.g., normalized counts) as response 190 ters, and two groups of 3 samples each, corresponding variable, nuisance covariates as fixed effects, and sam- 191 to an average of 200 cells per sample in each cluster. 141 ples as random effects. The method then removes the estimated impact of fixed effect covariates, and per-143 forms differential testing on these normalized values 144 (see Methods). 145

146 Furthermore, to enhance the interpretability of differential results, *distinct* provides functionalities to compute 147 (log) fold changes between conditions, and to plot den-148 sities and ECDFs, both for individual samples and at 149 the group-level. 150

151 Note that, although *distinct* and the Kolmogorov-Smirnov [18] (KS) test share similarities (they both 152 compare distributions via non-parametric tests), the 153 wo approaches present several conceptual differences. 154 Firstly, the KS considers the maximum distance be-155 tween two ECDFs, while our approach estimates the 156 overall distance between ECDFs, which in our view is 157 more appropriate way to measure the difference be-158 tween distributions. Secondly, the KS test only com-159 pares two individual densities, while our framework 160 compares groups of distributions. Thirdly, while the 161 KS statistic relies on asymptotic theory, our framework 162 uses a permutation test. Finally, a comparison between 163 distinct and scDD [15] based on the KS test (labelled 164 scDD-KS) shows that our method, compared to the KS 165 test, has greater statistical power to detect differential 166 effects and leads to fewer false discoveries (see Simula-168 tion studies).

Simulation studies 169

170 We conducted an extensive benchmark, based on 171 scRNA-seq and mass cytometry simulated and experi-¹⁷² mental datasets to investigate *distinct*'s ability to identify differential patterns in sub-populations of cells. 173

First, we simulated droplet scRNA-seq data via mus-174 175 C 176 177 178 179 180 181 182 183 184 186 187 189 consists of 4,000 genes, 3,600 cells, separated into 3 clus- 241 10 simulations, while it failed to run within a week time

¹⁹² We considered six different normalization approaches: ¹⁹³ counts per million (CPMs), *scater*'s logcounts [19], 194 linnorm [34], BASiCS [10,29,30], SCnorm [3] and ¹⁹⁵ residuals from variance stabilizing normalization from 196 sctransform (vstresiduals) [12]. We compared dis-197 tinct to several PB approaches from muscat, based on 198 edgeR [24], limma-voom and limma-trend [23], which ¹⁹⁹ emerged among the best performing methods for differ-²⁰⁰ ential analyses from scRNA-seq data [7,26]. We further $_{201}$ considered three methods from *muscat* based on mixed 202 models (MM), namely MM-dream2, MM-vstresiduals ²⁰³ and *MM-nbinom* (see Methods). Finally, we included $204 \ scDD$ [15], which is conceptually similar to our ap-205 proach: *scDD* implements a non-parametric method to 206 detect changes between individual distributions from 207 scRNA-seq, based on the Kolmogorov-Smirnov test, ²⁰⁸ scDD-KS, and on a permutation approach, scDD-perm. ²⁰⁹ For *scDD-perm* we used 100 permutations to reduce the 210 computational burden.

211 In all scenarios and on all six input data, distinct shows 212 favourable performance: it has good statistical power ²¹³ while controlling for the false discovery rate (FDR) 214 (Figure 2). In particular, for DE, DP and DM, distinct ²¹⁵ has similar performance to the best performing com-216 petitors (edgeR.linnorm and limma-trend.logcounts). ²¹⁷ while for DB and DV, it achieves significantly higher ²¹⁸ true positive rate (TPR), especially when using log*counts.* PB methods in general perform well for differ-210 220 ential patterns involving changes in the mean (DE, DP ²²¹ and DM), but struggle to identify DB and DV patterns. 222 scDD provides good TPR across all patterns when us-²²³ ing the KS test on vstresiduals (*scDD-KS.vstresiduals*). ²²⁴ while the TPR is significantly reduced when using ²²⁵ other inputs and with the permutation approach(*scDD*-226 perm); however, scDD methods (in particular, scDDcat [7] (see Methods). We ran five simulation repli- 227 KS. vstresiduals) also show a significant inflation of the ates for each of the differential profiles in Figure 1, 228 FDR. In contrast, MM methods provide good control of with 10% of the genes being differential in each clus- 229 the FDR but have low statistical power in all differenter, where DE (differential expression) indicates a shift 230 tial scenarios. We also investigated how normalization in the entire distribution, DP (differential proportion) 231 influences each method's results (Supplementary Figimplies two mixture distributions with different propor- 232 ure 2): distinct appears to be the least affected method tions of the two components, DM (differential modal- 233 and displays the smallest variation across normalizaity) assumes a unimodal and a bimodal distribution, 234 tion inputs, possibly due to its non-parametric struc-DB (both differential modality and different component 235 ture, which can more flexibly accommodate various inmeans) compares a unimodal and a bimodal distribu- 236 puts. Given the computational cost of SCnorm, which tion with the same overall mean, and DV (differential 237 is significantly higher than the other normalizations, variability) refers to two unimodal distributions with 238 we only included this approach in the results from the the same mean but different variance (Figure 1 and 239 main simulations. Furthermore, among the 25 replicate Supplementary Figure 1). Each individual simulation 240 datasets in Figure 2. SCnorm ran in a few minutes on



Figure 2: distinct identifies various differential patterns and controls for the FDR. TPR vs. FDR in muscat simulated data; DE, DP, DM, DB and DV refer to the differential profiles illustrated in Figure 1. Circles indicate observed FDR for 0.01, 0.05, 0.1 and 0.2 significance thresholds. Results are averages across the five simulation replicates. Each individual replicate consists of 4,000 genes, 3,600 cells, separated into 3 clusters, and two groups of 3 samples each, corresponding to an average of 200 cells per sample in each cluster.

243 245 246 247 ²⁴⁸ ticeable differences are detected between SCnorm and ²⁶⁶ ples from batch b_2 and one from b_1 . Differential results



Figure 3: distinct has uniform null p-values. Density of raw p-values in muscat null simulated data; each replicate represents a different null simulation. Each individual replicate consists of 4,000 genes, 3,600 cells, separated into 3 clusters, and two groups of 3 samples each, corresponding to an average of 200 cells per sample in each cluster.

the remaining normalization methods, while for *scDD*-250 KS SCnorm leads to a higher inflation of the FDR.

251 We further simulated five null simulation replicates ²⁵² with no differential patterns; again with each simulation ²⁵³ having 4,000 genes, 3,600 cells, 3 cell clusters and two ²⁵⁴ groups of 3 samples each. In the null simulated data, ²⁵⁵ only *limma-trend.basics* and *limma-trend.cpm* present a ²⁵⁶ mild inflation of false positives, while MM and, particu-²⁵⁷ larly, *edgeR.basics* lead to overly conservative p-values; 258 instead, distinct and scDD show approximately uni-²⁵⁹ form p-values for all types of input data (Figure 3).

242 (on 10 cores) on the remaining 15 datasets. Therefore, 260 We also extended previous simulations to add a cellwe excluded *SCnorm* from Figure 2 and, in Supple- 261 type specific batch effect (i.e., a batch effect that affects mentary Figures 3 and 4, we report a comparison of 262 differently each cell-type) [7,17]. In particular, we sim-SCnorm to the remaining normalization methods, on 263 ulated 2 batches, that we call b_1 and b_2 , with one group the subset of 10 simulations where all normalizations 264 of samples having two samples associated to b_1 and one successfully ran. For distinct, edgeR and limma, no no- $_{265}$ to b_2 , and the other group of samples having two sam-



Figure 4: distinct achieves good performance when varying the number of available cells. TPR vs. FDR in muscat simulated data; with 50, 100, 200 and 400 cells per cluster-sample combination, corresponding to a total of 900, 1,800, 3,600 and 7,200 cells, respectively. Results are aggregated over the five replicate simulations of each differential type (DE, DP, DM, DB and DV), contributing in equal fraction. Each individual simulation replicate consists of 4,000 genes, 3 cell clusters and two groups of 3 samples each. Circles indicate observed FDR for 0.01, 0.05, 0.1 and 0.2 significance thresholds. Note that *scDD-perm* and MM were excluded from this analysis due to their computational cost.

are substantially unchanged (Supplementary Figure 5), 267 which shows *distinct* can effectively remove nuisance 268 confounders. 269

270 Furthermore, we performed various sensitivity analyses and investigated how results are affected when varying: 271 i) the number of cells, ii) the library size, iii) the dis-272 persion parameter, iv) the fraction of significant genes, 273 and v) the sample sizes in each group. In particular, we 274 simulated 50, 100, 200 (as in the original simulation) 275 and 400 cells per sample in each cluster. We further 276 modified the library size and dispersion parameters of 277 the negative binomial model used by *muscat* to simu-278 late scRNA-seq data, influencing the mean expression 279 and cell-to-cell variability respectively, by considering 280 values 1/5, 1/2, 2 and 5 times as big as those used in 281 282 the original simulation. In addition, we varied the per- 303 which is higher than PB methods (0.1 to 0.2 minutes)



Figure 5: distinct requires more computational resources than PB and scDD-KS methods, but significantly less than MM and scDD-perm models. Average computing time, expressed in minutes, in *muscat* main simulations (Figures 2-3). For each method, times are averaged across simulation types (DE, DP, DM, DB, DV and null) and, for each type, across the five replicate simulations; in each replicate 3,600 cells are available (200, on average, per cluster-sample combination). distinct, MM and scDD models were run on 3 cores, while pseudo-bulk methods based on edgeR and limma used a single core because they do not allow for parellel computing. Note that scDD-perm requires much longer on vstresiduals than on the other normalized data, because scDD performs differential testing on non-zero values: vstresiduals, (unlike linnorm, cpm and basics normalized data) are not zero-inflated and, therefore, many more cells have to be used for differential testing.

283 centage of simulated differential genes as 1, 5, 10 (as in ²⁸⁴ the original simulation) and 20%, and considered var-²⁸⁵ ious unbalanced designs by comparing two groups of 286 different sample sizes: 3 vs. 2, 4 vs. 3, and 5 vs. 3. 287 Overall, increasing the number of cells or the library 288 size and decreasing the dispersion have a positive im-289 pact on the performance of all methods, by improving ²⁹⁰ their ability to detect differential effects (i.e., true pos-²⁹¹ itive rate); nonetheless, none of these factors seem to ²⁹² affect the relative ranking of methods, which remains ²⁹³ globally stable (Figure 4 and Supplementary Figures 294 6-7). In addition, changing the fraction of significant ²⁹⁵ genes and considering unbalanced designs does not ap-²⁹⁶ pear to introduce systematic changes in performance ²⁹⁷ (Supplementary Figures 8-9). Note that, in these sen-298 sitivity analyses, we excluded MM models due to the ²⁹⁹ high computational cost and low statistical power dis-300 played in the previous analyses.

³⁰¹ From a computational perspective, *distinct* required 302 an average time of 3.2 to 4.5 minutes per simulation,

and *scDD-KS* (0.5 to 0.7 minutes), but significantly lower than MM approaches (29.4 to 297.3 minutes) and *scDD-perm* (544.7 to 2085.6 minutes) (Figure 5 and Sor Supplementary Table 1). All methods were run on 3 because they do not allow for parellel computing.

We also considered an alternative popular droplet 310 scRNA-seq data simulator, SplatPOP [2], which rep-311 ³¹² resents a generalization of *Splatter* [35], that allows nulti-sample multi-group synthetic data to be gener-313 ated. In particular, we simulated 20,345 genes from 314 a human genome with two groups of 4 samples each, 315 and 100 cells per sample, belonging to the same clus-316 er of cells, for a total of 800 cells across all samples. 317 We ran 8 differential simulations, with 10% of genes 318 truly differential between groups, by varying the lo-319 cation (*de.facLoc*) and scale (*de.facScale*) differential 320 parameters, mainly affecting the mean and variance, 321 respectively (see Methods). We considered the same 322 normalization and differential methods as in the mus-323 at simulation (except MM and *scDD-perm*, which were 324 not considered due to the high computational cost and 325 low statistical power displayed above). As expected, for 326 all methods, differential patterns are easier to detect as 327 the magnitude of the difference increases, with differen-328 tial location patterns having a higher true positive rate 329 than differential scale patterns. While all methods con-330 trol the FDR, in all simulations, distinct achieves sub-331 stantially higher TPR than competitors (Figure 6). We 332 also repeated the same simulations including a batch 333 effect, with two batches, with the same scale and lo-334 cation differential parameters for the batch and group 335 differences (i.e., increasing together from 0.2 to 1.5). 336 Again, we excluded *scDD* from these analyses because 337 it cannot handle covariates directly. Results agree with 338 those from the *muscat* batch effect simulation study: 339 FDR and TPRs are mostly unchanged when introduc-340 ing nuisance covariates, with only a minor decrease in 341 the TPR in stronger batch effects, i.e., when de.facLoc 342 and de.facScale are 1 and 1.5 (Supplementary Figure 343 10), which again indicates that *distinct* can effectively 344 control for nuisance covariates. C 345

We further considered the semi-simulated mass cytom-346 etry data from Weber et al. [32] (labelled diffcyt sim-347 ulation), where spike-in signals were computationally 348 introduced in experimental data [5], hence maintain-349 ing the properties of real biological data while also 350 embedding a known ground truth signal. We evalu-351 ated *distinct* and two methods from *diffcyt*, based on 352 *limma* [23] and linear mixed models (LMM), which out-353 performed competitors on these same data [32]. In 354 355 particular, we considered three datasets from Weber



Figure 6: distinct displays higher TPR than competitors. TPR vs. FDR in SplatPop simulated data, with various degrees of differential location (left) and scale (right) parameters, primarily affecting the mean and variance, respectively. Circles indicate observed FDR for 0.01, 0.05, 0.1 and 0.2 significance thresholds. Each simulation consists of 20,345 genes genes, 800 cells (belonging to the same cluster), and two groups of 4 samples each, corresponding to an average of 100 cells per sample.



Figure 7: distinct shows high power while controlling for false positive and false discovery rates. (a-b) TPR vs. FDR in diffcyt semi-simulated data. 'main', 'less 50' and 'less 75' indicate the main simulation, and those where differential effects are diluted by 50 and 75%, respectively. Each simulation consists of 88,435 cells and two groups of 8 samples each. Circles indicate observed FDR for 0.01, 0.05, 0.1 and 0.2 significance thresholds. (a) As in the muscat simulation study, cells were clustered into 8 populations based on manually annotated cell types [32]. (b) As in Weber et al. [32], cells were grouped in 100 high-resolution clusters via unsupervised clustering. (c) Density of raw p-values in *diffcyt* null semi-simulated data; each replicate represents a different null simulation. Each replicate consists of 88,438 cells and two groups of 8 samples each. As in Weber et al. [32], cells were clustered in an unsupervised manner.

357 358 359 360 361 362 363 364 365 366 367

lation, *distinct* achieves higher TPR when considering cell-type labels (Figure 7a, 'main'), while all methods 370 exhibit substantially overlapping performance when using unsupervised clustering (Figure 7b, 'main'). In both 372 clustering approaches, as the magnitude of the differen-373 tial effect decreases, the distance between methods in-374 creases: *diffcyt* tools show a significant drop in the true 375 positive rate whereas *distinct* maintains a higher TPR while effectively controlling for the false discovery rate (FDR) (Figures 7a-b and Supplementary Figure 11). 379 This indicates that *distinct* has good statistical power to detect even small changes between conditions. We 380 also considered the three replicate null datasets from 381 Weber et al. [32] (i.e., with no differential effect), con-382 383 taining 24 protein markers and 88,438 cells across 8 384 cell types, and found that all methods display approx-³⁸⁵ imately uniform p-values (Figure 7c).

Experimental data analyses 386

³⁸⁷ In order to investigate false positive rates (FPRs) in 388 real data, we considered two experimental scRNA-seq datasets where no differential signals were expected, by comparing samples from the same experimental con-390 dition. Given the high computational cost and low 391 power of MM, and the high FDR of *scDD* models, for 393 the real data analyses, we only included *distinct* and PB methods. We considered gene-cluster combinations 394 with at least 20 non-zero cells across all samples. The 395 ³⁹⁶ first dataset (labelled *T-cells*) consists of a Smart-seq2 scRNA-seq dataset of 19,875 genes and 11,138 T cells 397 ³⁹⁸ isolated from peripheral blood from 12 colorectal cancer patients [36]. We automatically separated cells in $_{400}$ 11 clusters (via *igraph* [1,8]), and generated replicate datasets, by randomly separating, three times, the 12 patients to two groups of size 6. The second dataset 402 (labelled Kang) contains 10x droplet-based scRNA-seq 403 peripheral blood mononuclear cell data from 8 Lupus 404 405 patients, before (controls) and after (stimulated) 6h-406 treatment with interferon- β (INF- β), a cytokine known 356 et al. [32]: the main DS dataset and two more where 407 to alter the transcriptional profile of immune cells [13]. differential effects were diluted by 50 and 75%. Each 408 The full dataset contains 35,635 genes and 29,065 cells, dataset consists of 24 protein markers, 88,435 cells, and 409 which are separated (via manual annotation [13]) into 8 two groups (with and without spike-in signal) of 8 sam- 410 cell types. One of the 8 patients was removed as it apples each. Measurements were first transformed, and 411 pears to be a potential outlier (Supplementary Figures then cells were grouped into sub-populations with two 412 12-14). Here we only included singlet cells and cells separate approaches (see Methods): i) similarly to the 413 assigned to a cell population, and considered control *muscat* simulation study, cell labels were defined based 414 samples only, resulting in 11,854 cells and 10,891 genes. on 8 manually annotated cell types [32] (Figure 7a), 415 Again, we artificially created three replicate datasets and ii) as in the original diffcyt study from Weber et 416 by randomly assigning the 7 retained control samples al. [32], cells were grouped into 100 high-resolution clus- 417 in two groups of size 3 and 4. In both null analyses, we ters (based on 10 cell-type markers, see Methods) via 418 found that *limma-trend*, particularly when using CPMs, ³⁶⁸ unsupervised clustering (Figure 7b). In the main simu- ⁴¹⁹ leads to an increase of FPRs, *distinct*'s p-values are only



Figure 8: On experimental scRNA-seq data, distinct has almost-uniform null p-values. Density of raw p-values in the null T-cells (top) and Kang (bottom) experimental data. Each replicate represents a random partition of samples in two groups. The *T*-cells data consists of 12 samples and 11,138 cells across 11 clusters. For the Kang dataset, we retained 7 samples and 11,854 cells across 8 clusters.

⁴²⁰ marginally inflated towards 0, while *edgeR* and *limmavoom* are the most conservative methods and provide 421 the best control of FPRs (Figure 8 and Supplemen-422 tary Tables 2-3). Regarding normalization, linnorm 423 and BASiCS lead to the most conservative p-values and 424 smallest false positive rates. 425

426 427 428 429 430 431 432 433 434 435 436 437 very coherent across different input data (Supplemen- 457 liver, and pancreatic cancer, KDELR2 for renal, head



Figure 9: distinct discovers non-canonical differential patterns. Density of logcounts for nine examples of differential patterns identified by *distinct* on all input data (adjusted p-values < 0.05), and not by any PB tool (adjusted p-values > 0.05), on the Kang dataset when comparing controls and stimulated samples. Gene RPL13 was identified in FCGR3A+ Monocytes (third row) and in NK cells (fourth row), while all other genes were detected in Dendritic cells. Each line represents a sample.

⁴³⁸ tary Figure 15). When visually investigating the gene-439 cluster combinations detected by *distinct* (adjusted pvalue < 0.1), on all five input data (CPMs, logcounts, ⁴⁴¹ linnorm, BASiCS and vstresiduals), and not detected $_{442}$ by any of the ten PB approaches (adjusted p-value > 443 0.1), we found several interesting non-canonical differ-444 ential patterns (Figure 9 and Supplementary Figures 445 16-27). In particular, gene MARCKSL1 displays a DB We then considered again the Kang dataset, and per- 446 pattern, with stimulated samples having higher density formed a DS analysis between controls and stimulated 447 on the tails and lower in the centre of the distribusamples. Again, we removed one potential outlier pa- 448 tion, gene RPL13 mirrors classical DE, while the other tient, and only considered singlet cells and cells as- 449 genes seem to emulate DP profiles. Interestingly, ten signed to a cell population; we further filtered gene- 450 out of eleven of these genes are known tumor progcluster combinations with less than 20 non-zero cells 451 nostic markers: H2AZ2 for cervical and renal cancer, across all samples, resulting in 12,045 genes and 23,571 452 SRSF9 for liver cancer and melanoma, RPL24 for recells across 8 cell types and 14 samples. We found 453 nal and thyroid cancer, HNRNPA0 for renal and panthat *distinct* identifies more differential patterns than 454 creatic cancer, MARCKSL1 for liver and renal cancer, PB methods, with edgeR and limma-voom being the 455 GTF3C6 for liver cancer, RPL13 for endometrial and most conservative methods, and that its results are 456 renal cancer, PGK1 for breast, head and neck, cervical,

method	% of unique results
distinct.logcounts	0.3
distinct.basics	0.8
limma-trend.logcounts	0.9
distinct.cpm	1.0
distinct.vstresiduals	1.1
edgeR.linnorm	1.2
limma-trend.vstresiduals	1.5
limma-trend.basics	1.5
edgeR.counts	1.7
edgeR.basics	3.0
distinct.linnorm	3.6
limma-trend.linnorm	3.7
limma-voom.counts	5.6
edgeR.cpm	10.4
limma-trend.cpm	26.8

Table 1: Percentage of unique gene/cell-type identifications that are unique to each method. Since methods return significantly different number of significant results, for each method, we selected the most significant 1,000 results. For every method, we then compute the fraction of such results that are unique, i.e., not in common with the top 1.000 results returned by any other method.

⁴⁵⁸ and neck and glioma cancer, and RPL11 for renal and breast cancer [28]. This is an interesting association, 459 considering that INF- β stimulation is known to inhibit 460 and interfere with tumor progression [9, 22]. Addition-461 ally, Supplementary Figures 16-27 show how distinct 462 can identify differences between groups of distributions 463 even when only a portion of the ECDF varies between 464 conditions. Finally, we computed the fraction of detected genes that are unique by each method. Given 466 that a ground truth is absent, we speculate that gene-467 cluster combinations detected by multiple methods are 468 more likely to be truly differential, while those detected 469 by a single method are more likely to be false posi-470 tive detections. Since methods return widely different 471 472 number of significant genes, for each method, we considered the top (i.e., smallest p-value) 1,000 genes per 473 cell-type. We then computed the percentage of results 474 that are unique to each method (Table 1), i.e., not in 475 common with the top 1,000 results returned by any 476 other method. Overall, *distinct* displays a lower frac-477 478 479 480 (i.e., distinct.logcounts and limma-trend.logcounts).

483 Discussion

485 differential patterns; nonetheless, most methods for dif- 537 ferent scenarios. For instance, by suitably modifying

ferential expression fail to identify changes where the 486 mean is not affected. To overcome the limitations of 487 present differential tools, we have developed *distinct*, a novel method to identify differential patterns between groups of distributions, which is particularly well suited 490 to perform differential analyses on high-throughput 491 ⁴⁹² single-cell data. *distinct* is based on a flexible hierarchical multi-sample full-distribution non-parametric approach. In order to compare it to state-of-the-art 494 differential methods, we ran extensive benchmarks on both simulated and experimental datasets from scRNAseq and mass cytometry data, where our approach ex-497 hibits favourable performance, provides good control of the FNR and FDR, and is able to identify more patterns 400 of differential expression compared to canonical tools, even when the overall mean is unchanged. In particular, 501 our approach displays a higher statistical power (i.e., TPR) not only than PB methods, but also compared 503 to other non-parametric frameworks from scDD, based 505 on the Kolmogorov-Smirnov test statistic (scDD-KS) 506 and on permutation tests (*scDD-perm*). distinct also ⁵⁰⁷ allows for biological replicates, does not rely on asymp-⁵⁰⁸ totic theory, which could be inaccurate in small sample ⁵⁰⁹ sizes (typical of biological data), and avoids parametric ⁵¹⁰ assumptions, that may be challenging to meet in single-511 cell data. Additionally, *distinct* can also effectively ad-⁵¹² just for sample-level cell-cluster specific covariates (i.e., ⁵¹³ whose effect varies across cell clusters), such as batch ⁵¹⁴ effects (Supplementary Figure 5). Importantly, distinct 515 is a very general test that, due to its non-parametric 516 nature, can be applied to various types of data, even ⁵¹⁷ beyond the single-cell applications shown here. Fur-518 thermore, thanks to its flexible form, we have shown in 519 our simulations that *distinct* has the most consistent ⁵²⁰ performance across normalization approaches (Supple-⁵²¹ mentary Figure 2 and 4).

522 However, these advantages come at the expense of a 523 higher computational burden, particularly when com-⁵²⁴ pared to PB methods or KS approaches (Figure 5). 525 Nonetheless, by employing clever computational tech-⁵²⁶ niques (i.e., parallel computing and C++ coding within ⁵²⁷ R), the method runs within minutes on a laptop, even tion of unique results (1.4% on average across all input 528 for large datasets. Overall, we believe that distinct data) compared to edgeR (4%) and limma (6.7%). It is 529 represents a valid alternative for differential detections also interesting to note that *scater*'s logcounts normal- 530 from single-cell data, particularly when interest lies beization lead to the 2 smallest fractions of unique values 531 yound canonical differences in means, as it allows to en-532 hance statistical power at the cost of a reasonable in-533 crease in the computational time.

⁵³⁴ Finally, although we have focused here on comparing 535 two groups of samples, several future extensions are 484 High-throughput single-cell data can display complex 536 possible to allow our framework to be applied to dif538 the test statistics in (1), one may ideally extend our ap- 584 Competing interests proach to perform a joint differential test between three 539 540 of more groups of samples. Although, it is worth noting that, in the presence of three or more experimental 541 conditions, at present, it is still possible to run pairwise 542 comparisons between pairs of conditions. While a joint 543 test across all groups may certainly be of interest in 544 some cases, from our experience, comparisons between pairs of groups are usually more used among scientists. 546 In addition, as we were suggested by a user, *distinct* 547 could be employed to compare cell clusters instead of 548 experimental conditions, hence discovering differential 549 genes between cell clusters (e.g., cell types), even from 550 551 individual samples.

552 Availability

553 distinct is freely available as a Bioconductor R packhttps://bioconductor.org/packages/distinct. 554 age at: The scripts used to run all analyses are avail-555 able on GitHub (https://github.com/SimoneTiberi/ 556 distinct manuscript, version v3) and Zenodo (DOI: 10.5281/zenodo.6397114). The *diffcyt* simulated data 558 ⁵⁵⁹ is available via FlowRepository (accession ID FR-FCM-ZYL8 [32]) and HDCytoData R Bioconductor pack-560 age [33]; the Kang dataset can be accessed via musc-561 $Data \ R \ Bioconductor \ package \ [6]; the \ T-cells \ dataset$ 562 ⁵⁶³ is deposited on the European Genome-phenome (acces-564 sion id EGAD00001003910 [36]).

Acknowledgements 565

566 We acknowledge Almut Luetge, Brian D M Tom, 567 Christina Azodi, Davis McCarthy, Reinhard Furrer, and the entire Robinson lab for precious com-568 ⁵⁶⁹ ments and suggestions. This work was supported by 570 Forschungskredit to ST (grant number FK-19-113) as ⁵⁷¹ well as by the Swiss National Science Foundation to ⁵⁷² MDR (grants 310030 175841, CRSII5 177208). MDR 573 acknowledges support from the University Research 574 Priority Program Evolution in Action at the Univer-575 sity of Zurich.

576 Author contributions

577 ST conceived the method, implemented it, performed 578 all analyses and wrote the manuscript. ST and MDR designed the study. HLC and LMW contributed to *muscat* and *diffcyt* simulation studies, respectively. PS 580 581 contributed to the computational development of dis-582 *tinct* and to the revision process. All authors read, 583 contributed to, and approved the final article.

585 The authors declare no competing interests.

586 Methods

587 Permutation test

588 In order to test for differences between groups, we em-589 ploy a hierarchical permutation approach: to estimate ⁵⁹⁰ the null distribution of s^{obs} , we permute the individual ⁵⁹¹ observations (e.g., single-cell measurements) instead of ⁵⁹² the samples. Note that this violates the exchangeability ⁵⁹³ assumption of permutation tests and, hence, p-values are not guaranteed to be uniformly distributed under 594 the null hypothesis; nonetheless, in our simulated and 595 596 experimental analyses, we empirically show that dis-597 tinct provides good control of both false positive and ⁵⁹⁸ false discovery rates. We randomly permute individual $_{599}$ observations P times across all samples and groups, by ⁶⁰⁰ retaining the original sample sizes. We denote by s_p 601 the test statistic computed at the *p*-th permutation, 602 $p = 1, \ldots, P$. A p-value, \tilde{p} , is obtained as [21]:

$$\tilde{p} = \frac{\sum_{p=1}^{P} \mathbf{1} \left(s_p \ge s^{obs} \right) + 1}{P+1},$$
(2)

603 where 1(cond) is 1 if cond is true, and 0 otherwise. In order to accurately infer small p-values, when \tilde{p} is below ⁶⁰⁵ some pre-defined thresholds, the number of permutations are automatically increased and \tilde{p} is re-computed. 607 By default, distinct initially computes 100 permutations; when $\tilde{p} < 0.1$ these are increased to 500; when 609 the new $\tilde{p} \leq 0.01$ we use 2,000 permutations, which 610 are further increased to 10,000 if $\tilde{p} < 0.001$. Note that 611 the number of permutations (i.e., 100, 500, 2,000 and 612 10,000) can be specified by the user.

613 Covariates

Assume we observe Z nuisance covariates, and that Nsamples are available across all groups, where for the *i*-th sample we observe C_i values (e.g., single-cell measurements). We fit the following linear mixed effects model:

$$y_{c}^{(i)} = \beta_{0} + \sum_{z=1}^{Z} \beta_{z} X_{z}^{(i)} + \alpha_{i} + \epsilon_{c}^{(i)}, \text{ for } i = 1, \dots, N,$$

and $c = 1, \dots, C_{i}, \quad (3)$

614 where $y_c^{(i)}$ represents the *c*-th observation for the *i*-th ⁶¹⁵ sample, β_0 is the intercept of the model, $X_z^{(i)}$ indicates 616 the z-th covariate in the *i*-th sample, β_z denotes the α_i fixed effect coefficient for the z-th covariate, α_i rep- α_i aged across replicates. In the main simulation (Figure 618 resents the random effect term for the *i*-th sample, 667 2) and the batch effect simulation (Supplementary Figand $\epsilon_c^{(i)}$ is the (zero-mean) residual for the c-th obser- 668 ure 5), we simulated from a paired design 2 groups of 621 623 624 terms, observations from the same sample are posi- 673 ure 4), the total numbers of available cells were 900. 625 tively correlated while, observations between different 674 1,800, 3,600 and 7,200, corresponding to an average of 626 samples are independent. We infer model parameters 675 50, 100, 200 and 400 cells per sample in every clus-627 via maximum likelihood, with the estimated values for 676 ter. For the differential simulations, we used log2-FC 628 the fixed effect terms denoted by $\hat{\beta}_0, \ldots, \hat{\beta}_Z$. We then 677 values of 1 for DE, 1.5 for DP and DM, and 3 for DB 630 631 633 634 specific effects of covariates.

635 Normalization

637 638 639 ure of *sctransform*'s variance stabilizing normalization, 641 642 643 644 646 ter. 647

648 In mass cytometry datasets, measurements were trans-649 formed via diffcyt's transformData function, which ap-650 plies an *arcsinh* transformation.

651 muscat simulation and Kang data

652 In all muscat simulations, we used the control samples 703 https://www.gencodegenes.org/human/release 19.html. 653 of the Kang dataset as a anchor data; as in the real 704 We ran a total of 16 simulations: 8 with and 8 without 654 data analyses, we excluded one sample as it emerged as 705 batch effects as nuisance covariate. In each case, we 655 656 657 658 659 660 662 663 tary Figure 6), we filtered gene-clusters combinations 714 and the differential location and scale parameters 665 five replicates were simulated, and results were aver- 716 respectively) matched those between groups of samples

vation in the *i*-th sample. We assume that random 669 3 samples each, with 4,000 genes, and 3,600 cells disterms are normally distributed as $\alpha_i \sim \mathcal{N}(0, \sigma_i^2)$, where 670 tributed in 3 clusters (corresponding to an average of $\mathcal{N}(a,b)$ denotes the normal distribution with mean a 671 200 cells per sample in each cluster). For the simuand variance b. Note that, due to the random effect 672 lation study when varying the number of cells (Figremove the estimated effect of nuisance covariates as 678 and DV. For the batch effect simulation study we used a $y_c^{(i)} - \sum_{z=1}^Z \hat{\beta}_z X_z^{(i)}$; differential testing is performed, as 679 modified version of *muscat*, developed by Almut Luetge described above, on these normalized values. In DS 680 at the Robinson lab (available at: https://github.com/ analyses, model (3) is fit, separately, for every gene- 681 SimoneTiberi/distinct manuscript), which allows simcluster combination, hence accommodating for cell-type 682 ulating cluster-specific batch effects [7, 17]. All mus-683 cat simulation studies, as well as the Kang non-null 684 data analysis, were performed by editing the original 685 snakemake workflow from Crowell et al. [7]. PB meth-686 ods were applied on aggregated data by summing cell-636 In scRNA-seq datasets, CPMs and logcounts were com- 687 level measurements; for differential testing, we used puted via scater Bioconductor R package [19], vstresid- 688 muscat's pbDS function [7]. Mixed model methods uals were calculated via sctransform R package [12] 689 were implemented, via muscat's mmDS function, us-(except for the *T*-cells data, where, due to a fail- 690 ing the same approaches as in Crowell et al. [7]: in 691 MM-dream2 and MM-vstresiduals linear mixed models we used *DESeq2*'s vst transformation [16]), while *lin-* 692 were applied to log-normalized data with observational norm, BASiCS and SCnorm, normalized data were 693 weights and variance-stabilized data, respectively, while calculated with the respective Bioconductor R pack- 694 in MM-nbinom generalized linear mixed models were ages [3, 10, 29, 30, 34]. For SCnorm, following the au- 695 fitted directly to raw counts. In the muscat simulations thor's suggestions, we normalized each cell cluster (3 in 696 and in the Kang non-null data analysis, we accounted total) separately, using samples as *Conditions* parame- 697 for the paired design by modelling the patient id as a ⁶⁹⁸ covariate in all methods that allow for covariates (i.e., 699 *distinct*, PB and MM).

700 *splatPop* simulation

SplatPOP simulated data, 701 In we used a huversion 702 man genome, 19. downloaded from a potential outlier (Supplementary Figures 12-14), and 706 ran 4 differential location ("de.facLoc" parameter) only considered singlet cells and cells assigned to a cell 707 and 4 differential scale ("de.facScale" parameter) population. In *muscat*'s simulation studies, we con- 708 simulations, with differential parameters equals to sidered gene-cluster combinations with simulated ex- 709 0.2, 0.5, 1 and 1.5. In every simulation, 10% of pression mean greater than 0.2; for DB patterns, we 710 genes were differential between groups, and a total increased this threshold to 1 because with low expres- 711 of 20,345 genes and 800 cells were simulated (100 sion values differences are not visible by eye. In the 712 per sample). In the simulation with batch effects, simulation when varying the library size (Supplemen- 713 the 8 samples were randomly assigned to 2 batches, with at least 50 non-zero cells. For every simulations, 715 between batches ("batch.facLoc" and "batch.facScale",

717 ("de.facLoc" and "de.facScale"). For more details on 764 References ⁷¹⁸ how *SplatPOP*'s data is simulated, please refer to the 765 766 ⁷¹⁹ original manuscript [2] and vignettes. 767

720 diffcyt simulation

772 773 721 The *diffcyt* semi-simulated data originates from a real 774 722 mass cytometry dataset of healthy peripheral blood 723 mononuclear cells from two paired groups of 8 samples 775 776 ⁷²⁴ each [5]; one group contains unstimulated cells, while 777 725 the other was stimulated with B cell receptor/Fc recep-778 tor cross-linker. The original dataset contains a total 780 726 of 172,791 cells and 24 protein markers: 10 of these 727 are cell-type markers used for cell clustering, while 14 728 783 729 are cell state markers used for differential state anal-784 730 yses; the distinction between cell state and cell-type 787 ⁷³¹ markers is based on prior biological knowledge [32]. 788 789 732 In Weber et al. [32], semi-simulated data were gener-733 ated by separating the cells of each unstimulated sam-790 791 734 ple in two artificial samples; a differential signal was then computationally introduced by replacing, in one 793 735 736 group, unstimulated B cells with B cells from stimu-794 737 lated samples. Measurements were transformed and 796 738 cells clustered via *diffcyt*'s *transformData* (which ap-739 plies an arcsinh transformation) and generateClusters 799 ⁷⁴⁰ functions, respectively. For the DS simulation in Fig-⁷⁴¹ ure 7b, as in Weber *et al.* [32], we evaluated methods' 802 performance in terms of detecting DS for phosphory-742 804 ⁷⁴³ lated ribosomal protein S6 (pS6) in B cells, which is 805 ⁷⁴⁴ the strongest differential signal across the cell types in 807 808 ⁷⁴⁵ this dataset [20, 32]. For the DS simulation in Figure 746 7a, we considered previously manually annotated cell 809 747 types [32] and included all 14 cell state markers. dif-811 748 fcyt's limma and LMM methods were applied via dif-749 fcyt's testDS limma and testDS_LMM functions, re-814 ⁷⁵⁰ spectively [32]. We accounted for the paired design by ⁷⁵¹ modelling the patient id as a covariate.

752 P-values adjustment

753 All p-values were adjusted via Benjamini-Hochberg cor-⁷⁵⁴ rection [4]. In *diffcyt* simulations we used globally ad-755 justed p-values for all methods, i.e., p-values from all ⁷⁵⁶ clusters are jointly adjusted once. However, since PB 757 methods were found to be over-conservative when glob-⁷⁵⁸ ally adjusting p-values [7], in *muscat* simulations and Kang discovery analyses, we used locally adjusted p-759 values for all methods. 760

761 Software versions

762 All analyses were performed via R software version 763 4.0.0, with Bioconductor packages from release 3.11.

768

769

770 771

779

782

785

792

795

798

- R. A. Amezquita, A. T. Lun, E. Becht, V. J. Carey, L. N. Carpp, L. Geistlinger, F. Marini, K. Rue-Albrecht, D. Risso, C. Soneson, et al. Orchestrating single-cell analysis with bioconductor. Nature methods, 17(2):137-145, 2020.
- C. B. Azodi, L. Zappia, A. Oshlack, and D. J. McCarthy. splatPop: simu-[2] lating population scale single-cell RNA sequencing data. Genome biology, 22(1):1-16, 2021.
- R. Bacher, L.-F. Chu, N. Leng, A. P. Gasch, J. A. Thomson, R. M. Stewart, M. Newton, and C. Kendziorski. Scnorm: robust normalization of single-cell [3] rna-seq data. Nature methods, 14(6):584-586, 2017
- [4] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a prac-tical and powerful approach to multiple testing. *Journal of the Royal sta-tistical society: series B (Methodological)*, 57(1):289–300, 1995.
- B. Bodenmiller, E. R. Zunder, R. Finck, T. J. Chen, E. S. Savig, R. V. [5] Bruggner, E. F. Simonds, S. C. Bendall, K. Sachs, P. O. Krutzik, et al. Multiplexed mass cytometry profiling of cellular states perturbed by small-molecule regulators. *Nature biotechnology*, 30(9):858-867, 2012.
- [6] H. L. Crowell. muscData: Multi-sample multi-group scRNA-seq data, 2020. R package version 1.1.2.
- H. L. Crowell, C. Soneson, P.-L. Germain, D. Calini, L. Collin, C. Raposo, D. Malhotra, and M. D. Robinson. muscat detects subpopulationspecific state transitions from multi-sample multi-condition single-cell transcriptomics data. Nature Communications, 11(1):1-12, 2020
- [8] G. Csardi and T. Nepusz. The igraph software package for complex network research. InterJournal, Complex Systems:1695, 2006
- M. R. Doherty, H. Cheon, D. J. Junk, S. Vinayak, V. Varadan, M. L. Telli, [9] J. M. Ford, G. R. Stark, and M. W. Jackson. Interferon-beta represses cancer stem cell properties in triple-negative breast cancer. Proceedings of the National Academy of Sciences, 114(52):13792-13797, 2017.
- N. Eling, A. C. Richard, S. Richardson, J. C. Marioni, and C. A. Vallejos. [10] Correcting the mean-variance dependency for differential variability testing using single-cell RNA sequencing data. Cell systems, 7(3):284-294, 2018.
- 797 [11] G. Finak, A. McDavid, M. Yajima, J. Deng, V. Gersuk, A. K. Shalek, C. K. Slichter, H. W. Miller, M. J. McElrath, M. Prlic, et al. MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. Genome biology, 16(1):1-800 13, 2015
 - [12] C. Hafemeister and R. Satija. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. Genome biology, 20(1):1-15, 2019.
 - [13]H. M. Kang, M. Subramaniam, S. Targ, M. Nguyen, L. Maliskova, E. Mc-Carthy, E. Wan, S. Wong, L. Byrnes, C. M. Lanata, et al. Multiplexed droplet single-cell RNA-sequencing using natural genetic variation. *Nature* biotechnology, 36(1):89, 2018.
- [14]P. V. Kharchenko, L. Silberstein, and D. T. Scadden. Bayesian approach to 810 single-cell differential expression analysis. Nature methods, 11(7):740-742. 2014.
- 812 [15] K. D. Korthauer, L.-F. Chu, M. A. Newton, Y. Li, J. Thomson, R. Stewart, and C. Kendziorski. A statistical approach for identifying differential distributions in single-cell RNA-seq experiments. Genome biology, 17(1):222, 815 2016.
- 816 [16]M. I. Love, W. Huber, and S. Anders. Moderated estimation of fold change 817 and dispersion for RNA-seq data with DESeq2. Genome biology, 15(12):550, 2014. 818
- 819 A. Lütge, J. Zyprych-Walczak, U. B. Kunzmann, H. L. Crowell, D. Calini, D. Malhotra, C. Soneson, and M. D. Robinson. Cellmixs: quantifying and [17]820 821 visualizing batch effects in single-cell rna-seq data. Life science alliance, 822 4(6), 2021.
- 823 [18] F J. Massey Jr. The kolmogorov-smirnov test for goodness of fit. Journal 824 of the American statistical Association, 46(253):68-78, 1951.
- D. J. McCarthy, K. R. Campbell, A. T. Lun, and Q. F. Wills. 825 [19] Scater: pre-processing, quality control, normalization and visualization of single-cell 826 RNA-seg data in R. Bioinformatics, 33(8):1179-1186, 2017. 827
- 828 [20] M. Nowicka, C. Krieg, L. M. Weber, F. J. Hartmann, S. Guglietta, B. Becher, M. P. Levesque, and M. D. Robinson. CyTOF workflow: dif-829 830 ferential discovery in high-throughput high-dimensional cytometry datasets. 831 F1000Research, 6, 2017
- 832 [21]B. Phipson and G. K. Smyth. Permutation p-values should never be zero: Calculating exact p-values when permutations are randomly drawn. Statis-833 834 tical applications in genetics and molecular biology, 9:Article39, 2010.
- 835 [22] X.-Q. Qin, N. Tao, A. Dergay, P. Moy, S. Fawell, A. Davis, J. M. Wilson, and 836 J. Barsoum. Interferon- β gene therapy inhibits tumor formation and causes 837 regression of established tumors in immune-deficient mice. Proceedings of the National Academy of Sciences, 95(24):14411-14416, 1998. 838
- 839 [23] M. E. Ritchie, B. Phipson, D. Wu, Y. Hu, C. W. Law, W. Shi, and G. K. Smyth. limma powers differential expression analyses for RNA-sequencing 840 841 and microarray studies. Nucleic acids research, 43(7):e47-e47, 2015.

- 842 [24] M. D. Robinson, D. J. McCarthy, and G. K. Smyth. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. Bioinformatics, 26(1):139-140, 2010.
- 845 [25] C. Soneson and M. D. Robinson. Bias, robustness and scalability in singlecell differential expression analysis. *Nature methods*, 15(4):255, 2018.
- 847 [26] J. W. Squair, M. Gautier, C. Kathe, M. A. Anderson, N. D. James, T. H.
 848 Hutson, R. Hudelle, T. Qaiser, K. J. Matson, Q. Barraud, et al. Confronting
 849 false discoveries in single-cell differential expression. *bioRxiv*, 2021.
- 850 [27] P.-Y. Tung, J. D. Blischak, C. J. Hsiao, D. A. Knowles, J. E. Burnett, J. K.
 851 Pritchard, and Y. Gilad. Batch effects and the effective design of single-cell
 852 gene expression studies. Scientific reports, 7:39921, 2017.
- 853 [28] M. Uhlén, L. Fagerberg, B. M. Hallström, C. Lindskog, P. Oksvold,
 854 A. Mardinoglu, Å. Sivertsson, C. Kampf, E. Sjöstedt, A. Asplund, et al.
 855 Tissue-based map of the human proteome. *Science*, 347(6220), 2015.
- 856 [29] C. A. Vallejos, J. C. Marioni, and S. Richardson. BASiCS: Bayesian
 857 analysis of single-cell sequencing data. PLoS computational biology, 11(6):e1004333, 2015.
- 859 [30] C. A. Vallejos, S. Richardson, and J. C. Marioni. Beyond comparisons of means: understanding changes in gene expression at the single-cell level.
 861 Genome biology, 17(1):1-14, 2016.
- 862 [31] T. Wang, B. Li, C. E. Nelson, and S. Nabavi. Comparative analysis of differential gene expression analysis tools for single-cell RNA sequencing data.
 864 BMC bioinformatics, 20(1):40, 2019.
- 865 [32] L. M. Weber, M. Nowicka, C. Soneson, and M. D. Robinson. diffcyt: Differential discovery in high-dimensional cytometry via high-resolution clustering. Communications biology, 2(1):1-11, 2019.
- 868 [33] L. M. Weber and C. Soneson. Hdcytodata: Collection of high-dimensional cytometry benchmark datasets in bioconductor object formats.
 870 F1000Research, 8, 2019.
- 871 [34] S. H. Yip, P. Wang, J.-P. A. Kocher, P. C. Sham, and J. Wang. Linnorm: improved statistical analysis for single cell RNA-seq expression data. Nucleic acids research, 45(22):e179-e179, 2017.
- 874 [35] L. Zappia, B. Phipson, and A. Oshlack. Splatter: simulation of single-cell
 875 RNA sequencing data. *Genome biology*, 18(1):1–15, 2017.
- 876 [36] Y. Zhang, L. Zheng, L. Zhang, X. Hu, X. Ren, and Z. Zhang. Deep single-cell
 877 RNA sequencing data of individual T cells from treatment-naive colorectal
 878 cancer patients. Scientific data, 6(1):1-15, 2019.