

# *distinct*: a novel approach to differential distribution analyses

Simone Tiberi<sup>1\*</sup>, Helena L Crowell<sup>1</sup>, Lukas M Weber<sup>2</sup>, Pantelis Samartsidis<sup>3</sup> and Mark D Robinson<sup>1</sup>

<sup>1</sup>*Department of Molecular Life Sciences and Swiss Institute of Bioinformatics, University of Zurich, Zurich, Switzerland.*

<sup>2</sup>*Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA.*

<sup>3</sup>*MRC Biostatistics Unit, School of Clinical Medicine, University of Cambridge, Cambridge, UK.*

\* e-mail: [Simone.Tiberi@uzh.ch](mailto: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 high-  
6 dimensional 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  
13 only marginally affected. *distinct* is based on  
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  
18 well as more subtle variations that do not in-  
19 volve the mean. We performed extensive bench-  
20 marks across both simulated and experimen-  
21 tal datasets from single-cell RNA sequencing  
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 anal-  
29 yses; Differential state; High-throughput single-cell  
30 data; Single-cell RNA-seq; Single-cell flow and mass cy-  
31 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 high-  
36 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) [28], 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 [10]). Classical bulk differential expression  
45 methods have been shown to perform well when used  
46 on single-cell measurements [22, 23, 27] and on aggre-  
47 gated data (i.e., averages or sums across cells), also re-  
48 ferred to as pseudo-bulk (PB) [5, 28]. 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 [12, 27]. 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* [12],  
60 *SCDE* [11], *MAST* [8], *BASiCS* [26] and mixed mod-  
61 els [24]). Nevertheless, they all present significant limi-  
62 tations: *BASiCS* does not perform cell-type specific dif-  
63 ferential testing between conditions, *scDD* does not di-  
64 rectly handle covariates and biological replicates, while  
65 *PB*, *SCDE*, *MAST* and mixed models performed poorly  
66 in previous benchmarks when detecting differential pat-  
67 terns that do not involve the mean [5, 12].

## 68 Results

### 69 *distinct*'s full distribution approach

70 To overcome these challenges, we developed *distinct*, a  
71 flexible and general statistical methodology to perform  
72 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 dis-  
76 tribution function (ECDF) from the individual (e.g.,  
77 single-cell) measurements of each sample, and compares  
78 ECDFs to identify changes between full distributions,  
79 even when the mean is unchanged or marginally in-  
80 volved (Figure 1 and Supplementary Figure 1). First,  
81 we compute the ECDF of each individual sample; then,  
82 we build a fine grid and, at each cut-off, we average the  
83 ECDFs within each group, and compute the absolute  
84 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  
89 for the  $i$ -th sample in the  $j$ -th group, is denoted by  
90  $ecdf_i^{(j)}(\cdot)$ , for  $j \in \{A, B\}$  and  $i = 1, \dots, N_j$ . We  
91 then define  $K$  equally spaced cut-offs between the mini-  
92 mum,  $min$ , and maximum,  $max$ , values observed across  
93 all samples:  $b_1, \dots, b_K$ , where  $b_k = min + k \times l$ , for  
94  $k = 1, \dots, K$ , 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,  
100  $s^{obs}$ , is obtained by adding these differences across all  
101 cut-offs:

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

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

104 Intuitively,  $s^{obs}$ , which ranges in  $[0, \infty)$ , approximates  
105 the area between the average ECDFs, and represents  
106 a measure of distance between two groups of densities:  
107 the bigger  $s^{obs}$ , the greater the distance between groups.  
108 The number of cut-offs  $K$ , which can be defined by  
109 users, is set to 25 by default, because no detectable  
110 difference in performance was observed when further  
111 increasing it (data not shown). Note that, although at  
112 each cut-off we compute the average across each group's  
113 curves, ECDFs are computed separately for each indi-  
114 vidual sample, therefore our approach still accounts for  
115 the within-group variability; indeed, at a given thresh-  
116 old, the average of the sample-specific ECDFs differs  
117 from the group-level ECDF (i.e., the curve based on  
118 all individual measurements from the group). The null  
119 distribution of  $s^{obs}$  is then estimated via a hierarchical

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-  
125 sible in small samples, because there are many more  
126 cells than samples. In principle, this may lead to an  
127 inflation of false positives due to lack of exchangeabil-  
128 ity (see Methods); nonetheless, in our analyses, *distinct*  
129 provides good control of both false positive and false  
130 discovery rates.

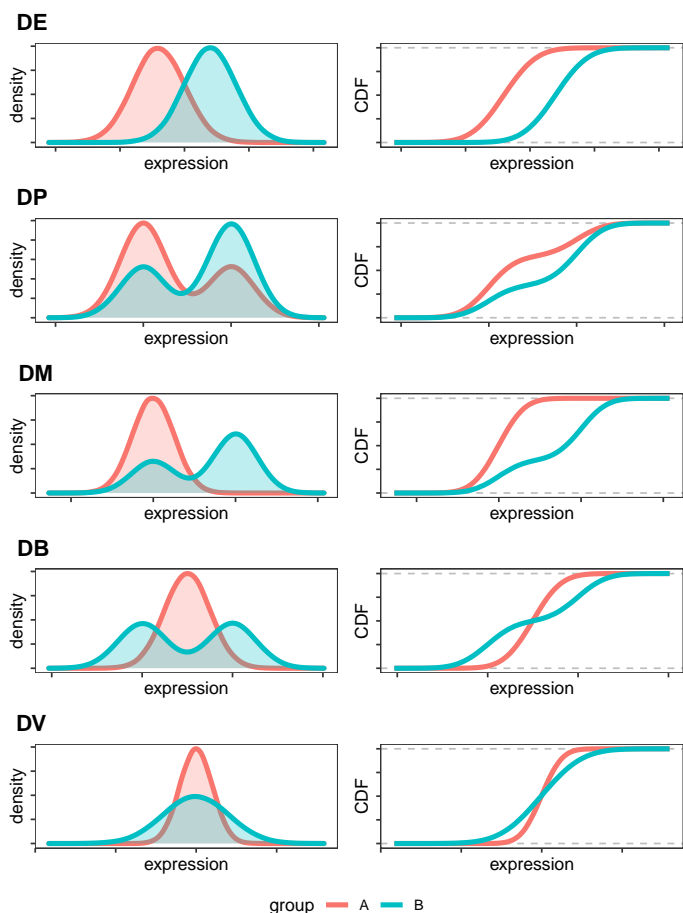
131 Importantly, *distinct* is general and flexible: it targets  
132 complex changes between groups, explicitly models bio-  
133 logical replicates within a hierarchical framework, does  
134 not rely on asymptotic theory, avoids parametric as-  
135 sumptions, and can be applied to arbitrary types of  
136 data. Additionally, *distinct* can also adjust for sample-  
137 level cell-cluster specific covariates (i.e., whose effect  
138 varies across cell clusters), such as batch effects; *dis-*  
139 *tinct* fits a linear model with the input data (e.g., CPMs  
140 or log2-CPMs) as response variable, and the covariates  
141 as predictors; the method then removes the estimated  
142 effect of covariates, and performs differential testing on  
143 these normalized values (see Methods).

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

149 Note that, although *distinct* and the Kolmogorov-  
150 Smirnov [15] (KS) test share similarities (they both  
151 compare distributions via non-parametric tests), the  
152 two approaches present several conceptual differences.  
153 Firstly, the KS considers the maximum distance be-  
154 tween two ECDFs, while our approach estimates the  
155 overall distance between ECDFs, which in our view is  
156 a more appropriate way to measure the difference be-  
157 tween distributions. Secondly, the KS test only com-  
158 pares two individual densities, while our framework  
159 compares groups of distributions. Thirdly, while the  
160 KS statistic relies on asymptotic theory, our framework  
161 uses a permutation test. Finally, a comparison between  
162 *distinct* and *scDD* [12] based on the KS test (labelled  
163 *scDD-KS*) shows that our method, compared to the KS  
164 test, has greater statistical power to detect differential  
165 effects and leads to fewer false discoveries (see Simula-  
166 tion studies).

## 167 Simulation studies

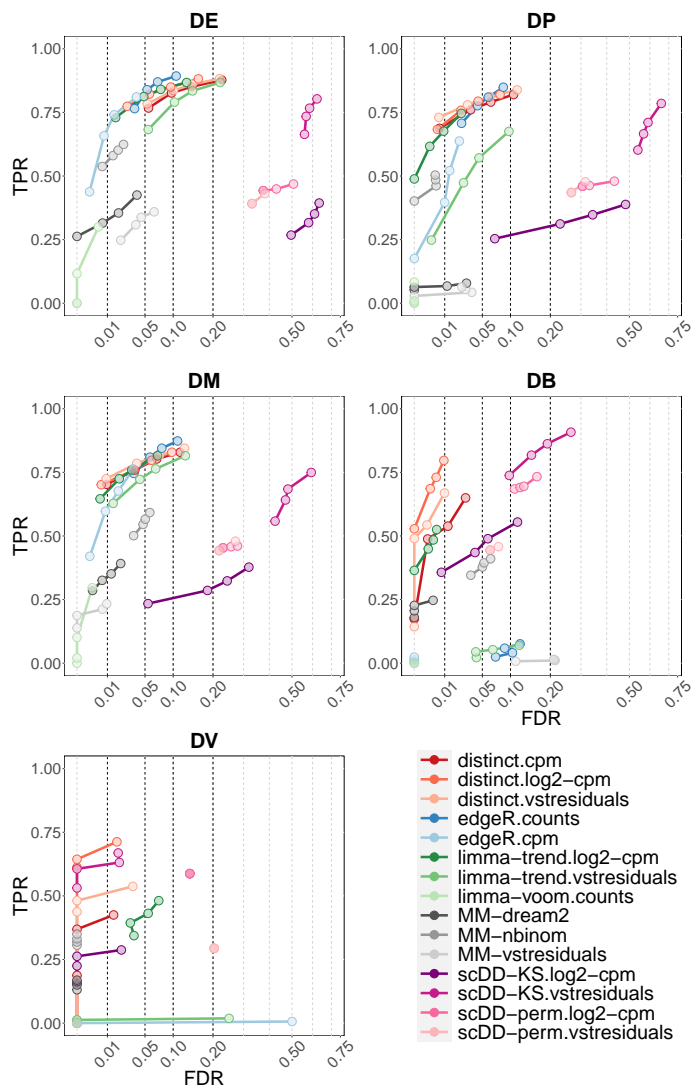
168 We conducted an extensive benchmark, based on  
169 scRNA-seq and mass cytometry simulated and experi-



**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. [12]: differential expression (DE), differential proportion (DP), differential modality (DM), and both differential modality and different component means (DB).

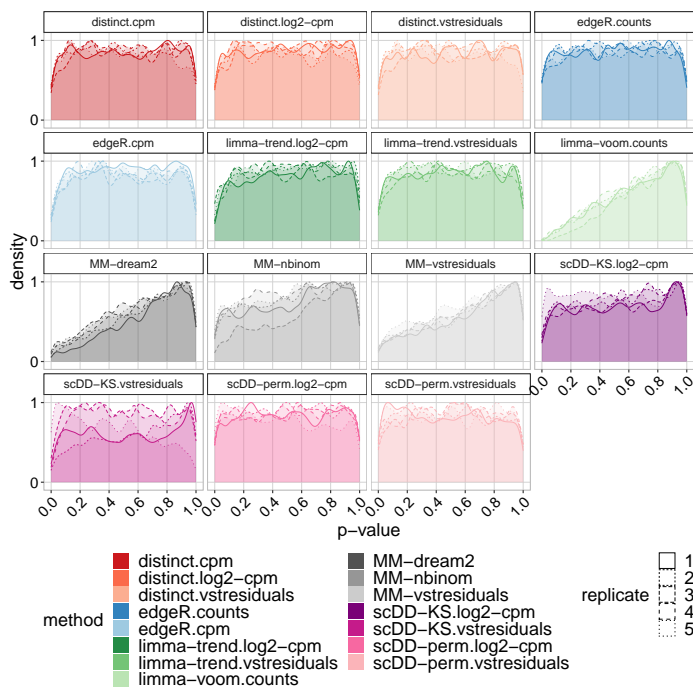
170 mental datasets to investigate *distinct*'s ability to iden-  
171 tify differential patterns in sub-populations of cells.

172 First, we simulated droplet scRNA-seq data via *mus-*  
173 *cat* [5] (see Methods). We ran five simulation repli-  
174 cates for each of the differential profiles in Figure 1,  
175 with 10% of the genes being differential in each cluster,  
176 where DE (differential expression) indicates a shift  
177 in the entire distribution, DP (differential proportion)  
178 implies two mixture distributions with different propor-  
179 tions of the two components, DM (differential modal-  
180 ity) assumes a unimodal and a bimodal distribution,  
181 DB (both differential modality and different component  
182 means) compares a unimodal and a bimodal distribu-  
183 tion with the same overall mean, and DV (differential  
184 variability) refers to two unimodal distributions with  
185 the same mean but different variance (Figure 1 and  
186 Supplementary Figure 1). Each individual simulation  
187 consists of 4,000 genes, 3,600 cells, separated into 3 clus-  
188 ters, and two groups of 3 samples each, corresponding  
189 to an average of 200 cells per sample in each cluster.



**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. Results are averages across the five simulation replicates. Circles indicate observed FDR for 0.01, 0.05, 0.1 and 0.2 significance thresholds. Two groups of 3 samples are compared and, on average, 200 cells are available for every sample in each of three clusters.

190 We considered three different normalizations: counts  
191 per million (CPMs), logarithm of CPMs to base 2 (log2-  
192 CPMs) and residuals from variance stabilizing normal-  
193 ization from *sctransform* (vstresiduals) [9]. We com-  
194 pared *distinct* to several PB approaches from *muscat*,  
195 based on *edgeR* [21], *limma-voom* and *limma-trend* [20],  
196 which emerged among the best performing methods for  
197 differential analyses from scRNA-seq data [5, 23]. We  
198 further considered three methods from *muscat* based  
199 on mixed models (MM), namely *MM-dream2*, *MM-*  
200 *vstresiduals* and *MM-nbinom* (see Methods). Finally,  
201 we included *scDD* [12], which is conceptually similar  
202 to our approach: *scDD* implements a non-parametric  
203 method to detect changes between individual distri-  
204 butions from scRNA-seq, based on the Kolmogorov-  
205 Smirnov test, *scDD-KS*, and on a permutation ap-

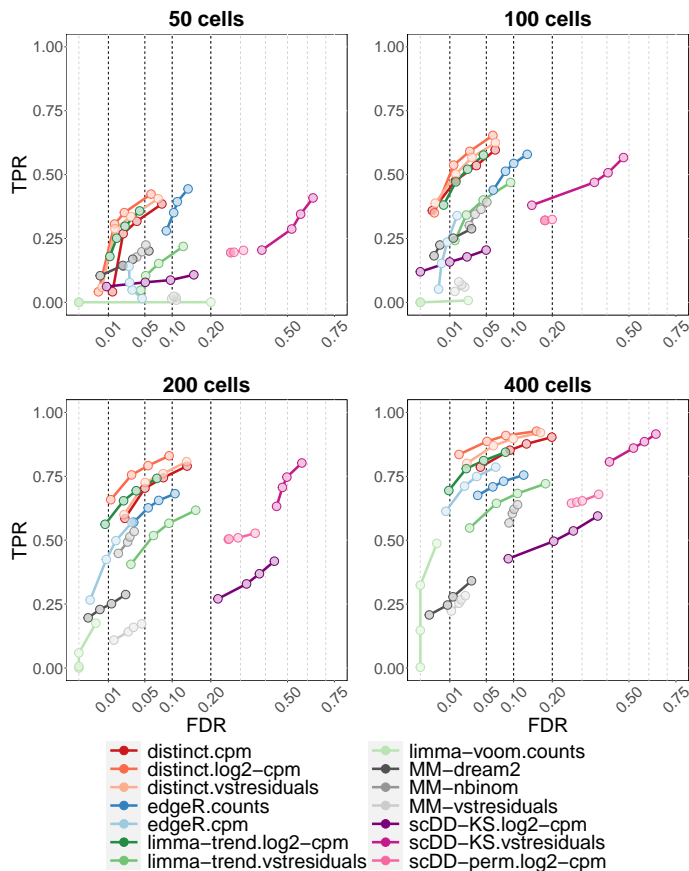


**Figure 3: *distinct* has uniform null p-values.** Density of raw p-values in *muscat* null simulated data; each replicate represents a different null simulation. Two groups of 3 samples are compared and, on average, 200 cells are available for every sample in each of three clusters.

206 proach, *scDD-perm*. For *scDD-perm* we used 100 per-  
207 mutations to reduce the computational burden.

208 In all scenarios and on all three input datasets, *dis-*  
209 *tinct* shows favourable performance: it has good sta-  
210 tistical power while controlling for the false discov-  
211 ery rate (FDR) (Figure 2). In particular, for DE,  
212 DP and DM, *distinct* has similar performance to the  
213 best performing competitors (*edgeR.counts* and *limma-*  
214 *trend.log2-CPMs*), while for DB and DV, it achieves  
215 significantly higher true positive rate (TPR), especially  
216 when using *log2-CPMs*. PB methods in general per-  
217 form well for differential patterns involving changes in  
218 the mean (DE, DP and DM), but struggle to identify  
219 DB and DV patterns. *scDD* provides good TPR across  
220 all patterns when using the KS test on vstresiduals  
221 (*scDD-KS.vstresiduals*), while the TPR is significantly  
222 reduced when using *log2-CPMs* and with the permu-  
223 tation approach(*scDD-perm*); however, *scDD* methods  
224 also show a significant inflation of the FDR. In contrast,  
225 MM methods provide good control of the FDR but have  
226 low statistical power in all differential scenarios.

227 We further simulated five null simulation replicates  
228 with no differential patterns; again with each simula-  
229 tion having 4,000 genes, 3,600 cells, 3 cell clusters and  
230 two groups of 3 samples each. In the null simulated  
231 data, no method presents an inflation of false positives,  
232 with *distinct*, *edgeR*, *limma-trend* and *scDD* showing



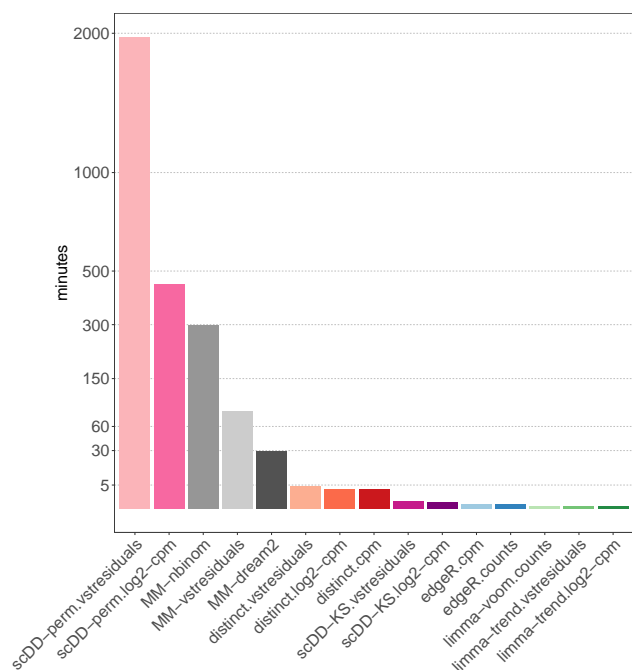
**Figure 4: *distinct* achieves good performance when varying the number of available cells.** TPR vs. FDR in *muscat* simulated data; with 50, 100, 200 and 400 cells per cluster-sample combination, corresponding to a total of 900, 1,800, 3,600 and 7,200 cells, respectively. Results are aggregated over 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.vstresiduals* was excluded from this analysis due to its computational cost.

233 approximately uniform p-values for all types of input  
234 data (Figure 3).

235 We also extended previous simulations to add a cell-  
236 type specific batch effect (i.e., a batch effect that affects  
237 differently each cell-type) [5,14]. In particular, we simu-  
238 lated 2 batches, that we call  $b_1$  and  $b_2$ , with one group  
239 of samples having two samples associated to  $b_1$  and one  
240 to  $b_2$ , and the other group of samples having two sam-  
241 ples from batch  $b_2$  and one from  $b_1$ . Differential results  
242 are substantially unchanged (Supplementary Figure 2),  
243 which shows *distinct* can effectively remove nuisance  
244 confounders. Furthermore, by varying the number of  
245 cells in the simulated data, we show that, compared to  
246 PB, MM and *scDD* methods, *distinct* achieves higher  
247 overall TPR, while controlling for the FDR, regardless  
248 of the number of available cells (Figure 5 and Supple-  
249 mentary Figure 3).

250 From a computational perspective, *distinct* required

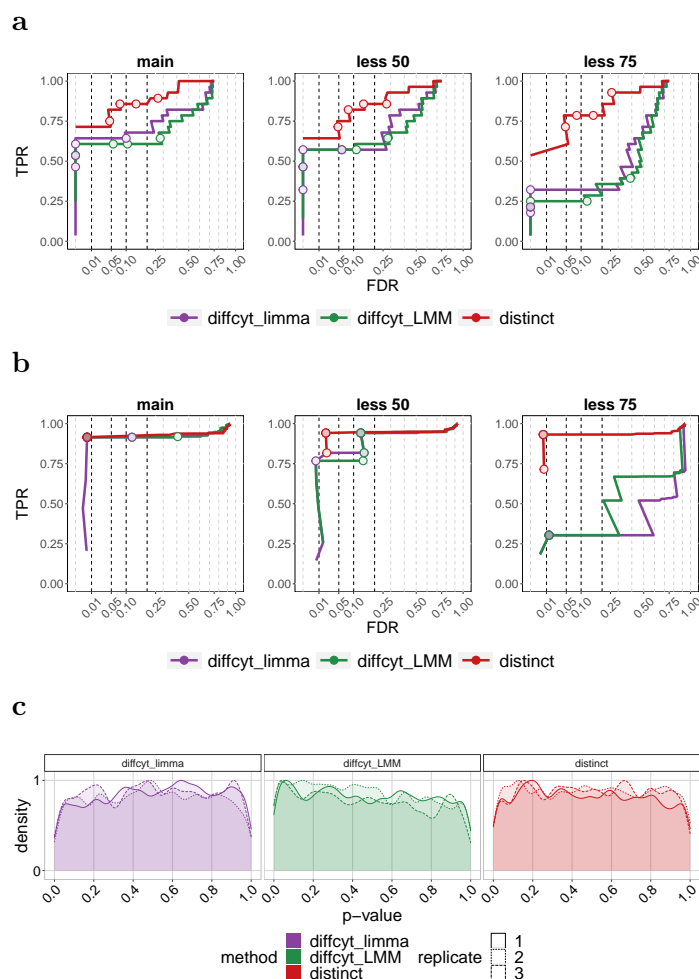




**Figure 5: *distinct* requires more computational resources than PB and *scDD-KS* methods, but significantly less than MM and *scDD-perm* models.** Average computing time, expressed in minutes, in *muscat* main simulations (Figures 2-3). For each method, times are averaged across simulation types (DE, 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 parallel computing.

251 an average time of 3.4 to 4.5 minutes per simulation,  
 252 which is higher than PB methods (0.1 to 0.2 minutes)  
 253 and *scDD-KS* (0.4 to 0.5 minutes), but significantly  
 254 lower than MM approaches (29.4 to 297.3 minutes) and  
 255 *scDD-perm* (447.5 to 1970.1 minutes) (Figure 4 and  
 256 Supplementary Table 1). All methods were run on 3  
 257 cores, except PB approaches, which used a single core,  
 258 because they do not allow for parallel computing.

259 We further considered the semi-simulated mass cytometry  
 260 data from Weber *et al.* [28] (labelled *diffcyt* sim-  
 261 ulation), where spike-in signals were computationally  
 262 introduced in experimental data [3], hence maintain-  
 263 ing the properties of real biological data while also  
 264 embedding a known ground truth signal. We evalu-  
 265 ated *distinct* and two methods from *diffcyt*, based on  
 266 *limma* [20] and linear mixed models (LMM), which out-  
 267 performed competitors on these same data [28]. In  
 268 particular, we considered three datasets from Weber  
 269 *et al.* [28]: the main DS dataset and two more where  
 270 differential effects were diluted by 50 and 75%. Each  
 271 dataset consists of 24 protein markers, 88,435 cells, and  
 272 two groups (with and without spike-in signal) of 8 sam-  
 273 ples each. Measurements were first transformed, and  
 274 then cells were grouped into sub-populations with two



**Figure 6: *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 on 8 populations based on manually annotated cell types [28]. (b) As in Weber *et al.* [28], 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.* [28], cells were clustered in an unsupervised manner.

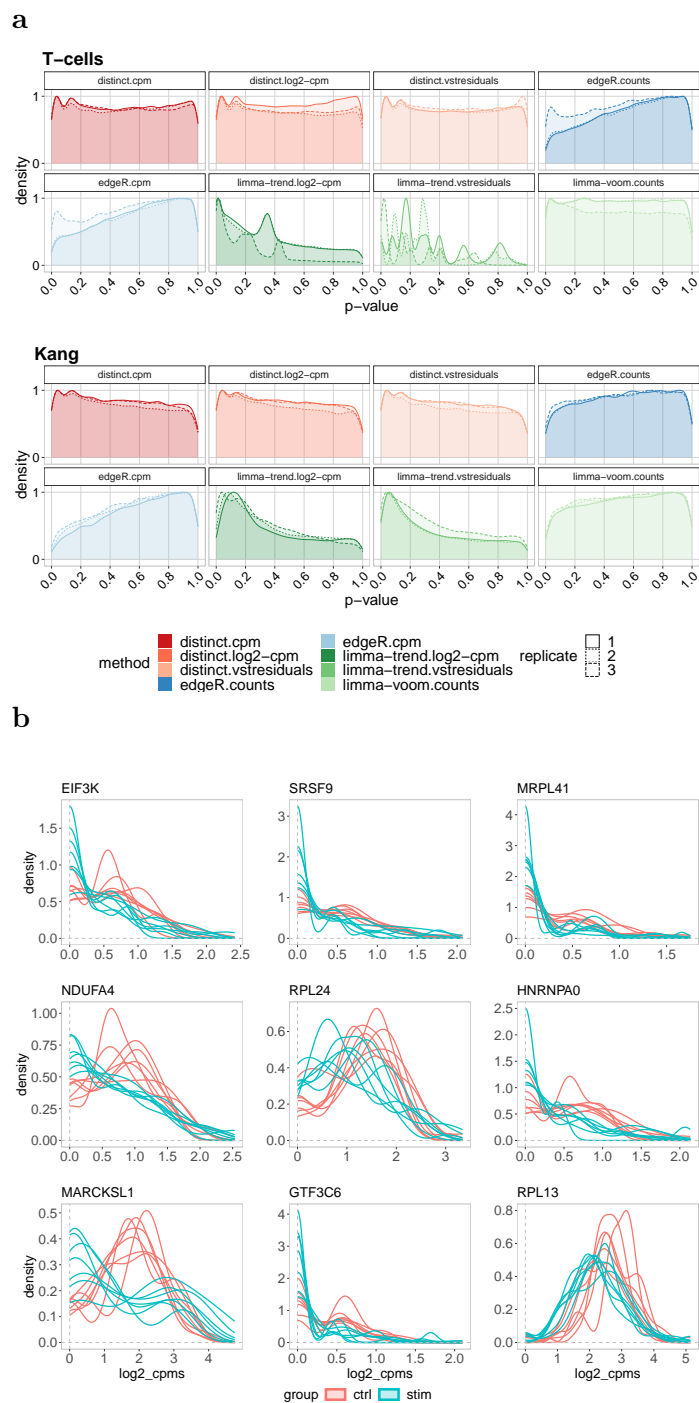
275 separate approaches (see Methods): i) similarly to the  
 276 *muscat* simulation study, cell labels were defined based  
 277 on 8 manually annotated cell types [28] (Figure 6a),  
 278 and ii) as in the original *diffcyt* study from Weber *et*  
 279 *al.* [28], cells were grouped into 100 high-resolution clus-  
 280 ters (based on 10 cell-type markers, see Methods) via  
 281 unsupervised clustering (Figure 6b). In the main simu-  
 282 lation, *distinct* achieves higher TPR when considering  
 283 cell-type labels (Figure 6a, ‘main’), while all methods  
 284 exhibit substantially overlapping performance when using  
 285 unsupervised clustering (Figure 6b, ‘main’). In both  
 286 clustering approaches, as the magnitude of the differ-  
 287 ential effect decreases, the distance between methods

288 increases: *diffcyt* tools show a significant drop in the  
 289 true positive rate (TPR) whereas *distinct* maintains a  
 290 higher TPR while effectively controlling for the false  
 291 discovery rate (FDR) (Figures 6a-b and Supplemen-  
 292 tary Figure 4). This indicates that *distinct* has good  
 293 statistical power to detect even small changes between  
 294 conditions. We also considered the three replicate null  
 295 datasets from Weber *et al.* [28] (i.e., with no differential  
 296 effect), containing 24 protein markers and 88,438 cells  
 297 across 8 cell types, and found that all methods display  
 298 approximately uniform p-values (Figure 6c).

## 299 Experimental data analyses

300 In order to investigate false positive rates (FPRs) in  
 301 real data, we considered two experimental scRNA-seq  
 302 datasets where no differential signals were expected, by  
 303 comparing samples from the same experimental condi-  
 304 tion. Given the high computational cost and low  
 305 power of MM, and the high FDR of *scDD* models, for  
 306 the real data analyses, we only included *distinct* and  
 307 PB methods. We considered gene-cluster combinations  
 308 with at least 20 non-zero cells across all samples. The  
 309 first dataset (labelled *T-cells*) consists of a Smart-seq2  
 310 scRNA-seq dataset of 23,459 genes and 11,138 T cells  
 311 isolated from peripheral blood from 12 colorectal can-  
 312 cer patients [30]. We automatically separated cells in  
 313 11 clusters (via *igraph* [1, 6]), and generated replicate  
 314 datasets, by randomly separating, three times, the 12  
 315 patients to two groups of size 6. The second dataset  
 316 (labelled *Kang*) contains 10x droplet-based scRNA-seq  
 317 peripheral blood mononuclear cell data from 8 Lupus  
 318 patients, before (controls) and after (stimulated) 6h-  
 319 treatment with interferon- $\beta$  (INF- $\beta$ ), a cytokine known  
 320 to alter the transcriptional profile of immune cells [10].  
 321 The full dataset contains 35,635 genes and 29,065 cells,  
 322 which are separated (via manual annotation [10]) into  
 323 8 cell types. One of the 8 patients was removed as it  
 324 appears to be a potential outlier (Supplementary Fig-  
 325 ures 5-7). Here we only included singlet cells and cells  
 326 assigned to a cell population, and considered control  
 327 samples only, resulting in 11,854 cells. Again, we ar-  
 328 tificially created three replicate datasets by randomly  
 329 assigning the 7 retained control samples in two groups  
 330 of size 3 and 4. In both null analyses, we found that  
 331 *limma-trend* leads to a major increase of FPRs, *dis-*  
 332 *tinct*'s p-values are only marginally inflated towards 0,  
 333 while *edgeR* and *limma-voom* are the most conservative  
 334 methods and provide the best control of FPRs (Figure  
 335 7a and Supplementary Tables 2-3).

336 We then considered again the *Kang* dataset, and per-  
 337 formed a DS analysis between controls and stimulated  
 338 samples. Again, we removed one potential outlier pa-



**Figure 7: On experimental scRNA-seq data, *distinct* discovers non-canonical differential patterns, and has almost-uniform null p-values.** (a) 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. (b) Density of log2-CPMs 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 cells, while all other genes were detected in Dendritic cells. Each line represents a sample.

339 patient, and only considered singlet cells and cells as-  
 340 signed to a cell population, resulting in 35,635 genes,

341 23,571 cells across 8 cell types and 14 samples; we fur- 392  
342 ther filtered gene-cluster combinations with less than 20 393  
343 non-zero cells across all samples. We found that *distinct*  
344 identifies more differential patterns than PB methods,  
345 with *edgeR* and *limma-voom* being the most conser-  
346 vative methods, and that its results are very coherent  
347 across different input data (Supplementary Figure 8).  
348 When visually investigating the gene-cluster combina-  
349 tions detected by *distinct* (adjusted p-value < 0.05), on  
350 all input data (CPMs, log2-CPMs and vstresiduals),  
351 and not detected by any PB method (adjusted p-value  
352 > 0.05), we found several interesting non-canonical dif-  
353 ferential patterns (Figure 7b and Supplementary Fig-  
354 ures 9-17). In particular, gene MARCKSL1 displays  
355 a DB pattern, with stimulated samples having higher  
356 density on the tails and lower in the centre of the dis-  
357 tribution, gene RPL13 mirrors classical DE, while the  
358 other genes seem to emulate DP profiles. Interestingly,  
359 eight out of nine of these genes are known tumor prog-  
360 nostic markers: EIF3K for cervical and renal cancer,  
361 SRSF9 for liver cancer and melanoma, NDUFA4 for  
362 renal cancer, RPL24 for renal and thyroid cancer, HN-  
363 RNPA0 for renal and pancreatic cancer, MARCKSL1  
364 for liver and renal cancer, GTF3C6 for liver cancer and  
365 RPL13 for endometrial and renal cancer [25]. This is  
366 an interesting association, considering that INF- $\beta$  stim-  
367 ulation is known to inhibit and interfere with tumor  
368 progression [7, 19]. Finally, Supplementary Figures 9-  
369 17 show how *distinct* can identify differences between  
370 groups of distributions even when only a portion of the  
371 ECDF varies between conditions.

## 372 Discussion

373 High-throughput single-cell data can display complex  
374 differential patterns; nonetheless, most methods for dif-  
375 ferential expression fail to identify changes where the  
376 mean is not affected. To overcome present limitations,  
377 we have introduced *distinct*, a general method to iden-  
378 tify differential patterns between groups of distribu-  
379 tions, which is particularly well suited to perform differ-  
380 ential analyses on high-throughput single-cell data. We  
381 ran extensive benchmarks on both simulated and ex-  
382 perimental datasets from scRNA-seq and mass cytom-  
383 etry data, where our method exhibits favourable per-  
384 formance, provides good control of the FPR and FDR,  
385 and is able to identify more patterns of differential ex-  
386 pression compared to canonical tools, even when the  
387 overall mean is unchanged. Furthermore, *distinct* al-  
388 lows for biological replicates, can adjust for covariates  
389 (e.g., batch effects), and does not rely on asymptotic  
390 theory. Finally, note that *distinct* is a very general test  
391 that, due to its non-parametric nature, can be applied

to various types of data, beyond the single-cell applica-  
tions shown here.

## 394 Availability

395 *distinct* is freely available as a Bioconductor R pack-  
396 age at: <https://bioconductor.org/packages/distinct>.  
397 The scripts used to run all analyses are avail-  
398 able on GitHub ([https://github.com/SimoneTiberi/  
399 distinct\\_manuscript](https://github.com/SimoneTiberi/distinct_manuscript), version v2) and Zenodo (DOI:  
400 10.5281/zenodo.4739098). The *diffcyt* simulated data  
401 is available via FlowRepository (accession ID FR-FCM-  
402 ZYL8 [28]) and *HDCytoData* R Bioconductor pack-  
403 age [29]; the *Kang* dataset can be accessed via *musc-*  
404 *Data* R Bioconductor package [4]; the *T-cells* dataset  
405 is deposited on the European Genome-phenome (acces-  
406 sion id EGAD00001003910 [30]).

## 407 Acknowledgements

408 We acknowledge Almut Luetge and the entire Robinson  
409 lab for precious comments and suggestions. This work  
410 was supported by Forschungskredit to ST (grant num-  
411 ber FK-19-113) as well as by the Swiss National Sci-  
412 ence Foundation to MDR (grants 310030\_175841, CR-  
413 SII5\_177208). MDR acknowledges support from the  
414 University Research Priority Program Evolution in Ac-  
415 tion at the University of Zurich.

## 416 Author contributions

417 ST conceived the method, implemented it, performed  
418 all analyses and wrote the manuscript. ST and MDR  
419 designed the study. HLC and LMW contributed to  
420 *muscat* and *diffcyt* simulation studies, respectively. PS  
421 contributed to the computational development of *dis-*  
422 *tinct*. All authors read, contributed to, and approved  
423 the final article.

## 424 Competing interests

425 The authors declare no competing interests.

## 426 Methods

### 427 Permutation test

428 In order to test for differences between groups, we em-  
429 ploy a hierarchical permutation approach: to estimate  
430 the null distribution of  $s^{obs}$ , we permute the individual  
431 observations (e.g., single-cell measurements) instead of  
432 the samples. Note that this violates the exchangeability  
433 assumption of permutation tests and, hence, p-values  
434 are not guaranteed to be uniformly distributed under



435 the null hypothesis; nonetheless, in our simulated and  
436 experimental analyses, we empirically show that *dis-*  
437 *tinct* provides good control of both false positive and  
438 false discovery rates. We randomly permute individual  
439 observations  $P$  times across all samples and groups, by  
440 retaining the original sample sizes. We denote by  $s_p$   
441 the test statistic computed at the  $p$ -th permutation,  
442  $p = 1, \dots, P$ . A p-value,  $\tilde{p}$ , is obtained as [18]:

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

443 where  $\mathbf{1}(cond)$  is 1 if *cond* is true, and 0 otherwise. In  
444 order to accurately infer small p-values, when  $\tilde{p}$  is below  
445 some pre-defined thresholds, the number of permuta-  
446 tions are automatically increased and  $\tilde{p}$  is re-computed.  
447 By default, *distinct* initially computes 100 permuta-  
448 tions; when  $\tilde{p} \leq 0.1$  these are increased to 500; when  
449 the new  $\tilde{p} \leq 0.01$  we use 2,000 permutations, which  
450 are further increased to 10,000 if  $\tilde{p} \leq 0.001$ . Note that  
451 the number of permutations (i.e., 100, 500, 2,000 and  
452 10,000) can be specified by the user.

## 453 Covariates

Assume we observe  $Z$  nuisance covariates, and that  $N$   
samples are available across all groups, where for the  
 $i$ -th sample we observe  $C_i$  values (e.g., single-cell mea-  
surements). We fit the following linear model:

$$y_c^{(i)} = \beta_0 + \sum_{z=1}^Z \beta_z X_z^{(i)} + \epsilon_c^{(i)}, \text{ for } i = 1, \dots, N, \\ \text{and } c = 1, \dots, C_i, \quad (3)$$

454 where  $y_c^{(i)}$  represents the  $c$ -th observation for the  $i$ -th  
455 sample,  $\beta_0$  is the intercept of the model,  $X_z^{(i)}$  indi-  
456 cates the  $z$ -th covariate in the  $i$ -th sample,  $\beta_z$  repre-  
457 sents the coefficient for the  $z$ -th covariate, and  $\epsilon_c^{(i)}$  is  
458 the residual for the  $c$ -th observation in the  $i$ -th sample.  
459 We infer parameters  $\beta_0, \dots, \beta_Z$  via least squares regres-  
460 sion, with the estimated values denoted by  $\hat{\beta}_0, \dots, \hat{\beta}_Z$ .  
461 We then remove the estimated effect of covariates as  
462  $y_c^{(i)} - \sum_{z=1}^Z \hat{\beta}_z X_z^{(i)}$ ; differential testing is performed, as  
463 described above, on these normalized values. For DS  
464 analyses, model (3) is fit, separately, for every gene-  
465 cluster combination, hence accommodating for cell-type  
466 specific effects of covariates.

## 467 Normalization

468 In scRNA-seq datasets, CPMs and log2-CPMs were  
469 computed via *scater* Bioconductor R package [16],  
470 while vstresiduals were calculated via *sctransform* R

471 package [9] (except for the *T-cells* data, where, due to  
472 a failure of *sctransform*'s variance stabilizing normal-  
473 ization, we used *DESeq2*'s vst transformation [13]).

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

## 477 *diffcyt* simulation

478 The *diffcyt* semi-simulated data originates from a real  
479 mass cytometry dataset of healthy peripheral blood  
480 mononuclear cells from two paired groups of 8 samples  
481 each [3]; one group contains unstimulated cells, while  
482 the other was stimulated with B cell receptor/Fc recep-  
483 tor cross-linker. The original dataset contains a total  
484 of 172,791 cells and 24 protein markers: 10 of these  
485 are cell-type markers used for cell clustering, while 14  
486 are cell state markers used for differential state anal-  
487 yses; the distinction between cell state and cell-type  
488 markers is based on prior biological knowledge [28].  
489 In Weber *et al.* [28], semi-simulated data were gener-  
490 ated by separating the cells of each unstimulated sam-  
491 ple in two artificial samples; a differential signal was  
492 then computationally introduced by replacing, in one  
493 group, unstimulated B cells with B cells from stimu-  
494 lated samples. Measurements were transformed and  
495 cells clustered via *diffcyt*'s *transformData* (which ap-  
496 plies an *arcsinh* transformation) and *generateClusters*  
497 functions, respectively. For the DS simulation in Fig-  
498 ure 6b, as in Weber *et al.* [28], we evaluated methods'  
499 performance in terms of detecting DS for phosphory-  
500 lated ribosomal protein S6 (pS6) in B cells, which is  
501 the strongest differential signal across the cell types in  
502 this dataset [17, 28]. For the DS simulation in Figure  
503 6a, we considered previously manually annotated cell  
504 types [28] and included all 14 cell state markers. *dif-*  
505 *fcyt*'s *limma* and LMM methods were applied via *dif-*  
506 *fcyt*'s *testDS\_limma* and *testDS\_LMM* functions, re-  
507 spectively [28]. We accounted for the paired design by  
508 modelling the patient id as a covariate.

## 509 *muscat* simulation and *Kang* data

510 In all *muscat* simulations, we used the control samples  
511 of the *Kang* dataset as a anchor data; as in the real  
512 data analyses, we excluded one sample as it emerged  
513 as a potential outlier (Supplementary Figures 5-7), and  
514 only considered singlet cells and cells assigned to a cell  
515 population. In *muscat*'s simulation studies, we con-  
516 sidered gene-cluster combinations with simulated ex-  
517 pression mean greater than 0.2; for DB patterns, we  
518 increased this threshold to 1 because with low expres-  
519 sion values differences are not visible by eye. For every



520 simulations, five replicates were simulated, and results  
521 were averaged across replicates. In the main simulation  
522 (Figure 2) and the batch effect simulation (Supplemen-  
523 tary Figure 3), we simulated from a paired design 2  
524 groups of 3 samples each, with 4,000 genes, and 3,600  
525 cells distributed in 3 clusters (corresponding to an av-  
526 erage of 200 cells per sample in each cluster). For the  
527 simulation study when varying the number of cells (Fig-  
528 ure 5 and Supplementary Figure 3), the total numbers  
529 of available cells were 900, 1,800, 3,600 and 7,200, cor-  
530 responding to an average of 50, 100, 200 and 400 cells  
531 per sample in every cluster. For the differential sim-  
532 ulations, we used log<sub>2</sub>-FC values of 1 for DE, 1.5 for  
533 DP and DM, and 3 for DB and DV. For the batch  
534 effect simulation study we used a modified version of  
535 *muscat*, developed by Almut Luetge at the Robinson  
536 lab (available at: [https://github.com/SimoneTiberi/  
537 distinct\\_manuscript](https://github.com/SimoneTiberi/distinct_manuscript)), which allows simulating cluster-  
538 specific batch effects [5,14]. All *muscat* simulation stud-  
539 ies, as well as the *Kang* non-null data analysis, were  
540 performed by editing the original snakemake workflow  
541 from Crowell *et al.* [5]. PB methods were applied on  
542 aggregated data by summing cell-level measurements;  
543 for differential testing, we used *muscat*'s *pbDS* function  
544 [5]. Mixed model methods were implemented, via *mus-*  
545 *cat*'s *mmDS* function, using the same approaches as in  
546 Crowell *et al.* [5]: in *MM-dream2* and *MM-vstresiduals*  
547 linear mixed models were applied to log-normalized  
548 data with observational weights and variance-stabilized  
549 data, respectively, while in *MM-nbinom* generalized lin-  
550 ear mixed models were fitted directly to raw counts. In  
551 the *muscat* simulations and in the *Kang* non-null data  
552 analysis, we accounted for the paired design by mod-  
553 elling the patient id as a covariate in all methods that  
554 allow for covariates (i.e., *distinct*, PB and MM).

## 555 P-values adjustment

556 All p-values were adjusted via Benjamini-Hochberg cor-  
557 rection [2]. In *diffcyt* simulations we used globally ad-  
558 justed p-values for all methods, i.e., p-values from all  
559 clusters are jointly adjusted once. However, since PB  
560 methods were found to be over-conservative when glob-  
561 ally adjusting p-values [5], in *muscat* simulations and  
562 *Kang* discovery analyses, we used locally adjusted p-  
563 values for all methods.

## 564 Software versions

565 All analyses were performed via R software version  
566 4.0.0, with Bioconductor packages from release 3.11.

## 567 References

- 568 [1] R. A. Amezcua, A. T. Lun, E. Becht, V. J. Carey, L. N. Carpp,  
569 L. Geistlinger, F. Marini, K. Rue-Albrecht, D. Risso, C. Soneson, et al.  
570 Orchestrating single-cell analysis with bioconductor. *Nature methods*,  
571 17(2):137–145, 2020.
- 572 [2] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a prac-  
573 tical and powerful approach to multiple testing. *Journal of the Royal sta-  
574 tistical society: series B (Methodological)*, 57(1):289–300, 1995.
- 575 [3] B. Bodenmiller, E. R. Zunder, R. Finck, T. J. Chen, E. S. Savig, R. V.  
576 Bruggner, E. F. Simonds, S. C. Bendall, K. Sachs, P. O. Krutzik, et al.  
577 Multiplexed mass cytometry profiling of cellular states perturbed by small-  
578 molecule regulators. *Nature biotechnology*, 30(9):858–867, 2012.
- 579 [4] H. L. Crowell. *muscData: Multi-sample multi-group scRNA-seq data*,  
580 2020. R package version 1.1.2.
- 581 [5] H. L. Crowell, C. Soneson, P.-L. Germain, D. Calini, L. Collin, C. Ra-  
582 poso, D. Malhotra, and M. D. Robinson. *muscat* detects subpopulation-  
583 specific state transitions from multi-sample multi-condition single-cell tran-  
584 scriptomics data. *Nature Communications*, 11(1):1–12, 2020.
- 585 [6] G. Csardi and T. Nepusz. The igraph software package for complex network  
586 research. *InterJournal, Complex Systems*:1695, 2006.
- 587 [7] M. R. Doherty, H. Cheon, D. J. Junk, S. Vinayak, V. Varadan, M. L. Telli,  
588 J. M. Ford, G. R. Stark, and M. W. Jackson. Interferon-beta represses  
589 cancer stem cell properties in triple-negative breast cancer. *Proceedings of  
590 the National Academy of Sciences*, 114(52):13792–13797, 2017.
- 591 [8] G. Finak, A. McDavid, M. Yajima, J. Deng, V. Gersuk, A. K. Shalek, C. K.  
592 Slichter, H. W. Miller, M. J. McElrath, M. Prlic, et al. MAST: a flexible  
593 statistical framework for assessing transcriptional changes and characterizing  
594 heterogeneity in single-cell RNA sequencing data. *Genome biology*, 16(1):1–  
595 13, 2015.
- 596 [9] C. Hafemeister and R. Satija. Normalization and variance stabilization  
597 of single-cell RNA-seq data using regularized negative binomial regression.  
598 *Genome biology*, 20(1):1–15, 2019.
- 599 [10] H. M. Kang, M. Subramaniam, S. Targ, M. Nguyen, L. Maliskova, E. Mc-  
600 Carthy, E. Wan, S. Wong, L. Byrnes, C. M. Lanata, et al. Multiplexed  
601 droplet single-cell RNA-sequencing using natural genetic variation. *Nature  
602 biotechnology*, 36(1):89, 2018.
- 603 [11] P. V. Kharchenko, L. Silberstein, and D. T. Scadden. Bayesian approach to  
604 single-cell differential expression analysis. *Nature methods*, 11(7):740–742,  
605 2014.
- 606 [12] K. D. Korthauer, L.-F. Chu, M. A. Newton, Y. Li, J. Thomson, R. Stewart,  
607 and C. Kendziorski. A statistical approach for identifying differential dis-  
608 tributions in single-cell RNA-seq experiments. *Genome biology*, 17(1):222,  
609 2016.
- 610 [13] M. I. Love, W. Huber, and S. Anders. Moderated estimation of fold change  
611 and dispersion for RNA-seq data with DESeq2. *Genome biology*, 15(12):550,  
612 2014.
- 613 [14] A. Lütge, J. Zypych-Walczak, U. B. Kunzmann, H. L. Crowell, D. Calini,  
614 D. Malhotra, C. Soneson, and M. D. Robinson. Cellmixs: quantifying and  
615 visualizing batch effects in single-cell rna-seq data. *Life science alliance*,  
616 4(6), 2021.
- 617 [15] F. J. Massey Jr. The kolmogorov-smirnov test for goodness of fit. *Journal  
618 of the American statistical Association*, 46(253):68–78, 1951.
- 619 [16] D. J. McCarthy, K. R. Campbell, A. T. Lun, and Q. F. Wills. Scater:  
620 pre-processing, quality control, normalization and visualization of single-cell  
621 RNA-seq data in R. *Bioinformatics*, 33(8):1179–1186, 2017.
- 622 [17] M. Nowicka, C. Krieg, L. M. Weber, F. J. Hartmann, S. Guglietta,  
623 B. Becher, M. P. Levesque, and M. D. Robinson. CyTOF workflow: differ-  
624 ential discovery in high-throughput high-dimensional cytometry datasets.  
625 *F1000Research*, 6, 2017.
- 626 [18] B. Phipson and G. K. Smyth. Permutation p-values should never be zero:  
627 Calculating exact p-values when permutations are randomly drawn. *Statistical  
628 applications in genetics and molecular biology*, 9:Article39, 2010.
- 629 [19] X.-Q. Qin, N. Tao, A. Dergay, P. Moy, S. Fawell, A. Davis, J. M. Wilson, and  
630 J. Barsoum. Interferon- $\beta$  gene therapy inhibits tumor formation and causes  
631 regression of established tumors in immune-deficient mice. *Proceedings of  
632 the National Academy of Sciences*, 95(24):14411–14416, 1998.
- 633 [20] M. E. Ritchie, B. Phipson, D. Wu, Y. Hu, C. W. Law, W. Shi, and G. K.  
634 Smyth. limma powers differential expression analyses for RNA-sequencing  
635 and microarray studies. *Nucleic acids research*, 43(7):e47–e47, 2015.
- 636 [21] M. D. Robinson, D. J. McCarthy, and G. K. Smyth. edgeR: a Bioconductor  
637 package for differential expression analysis of digital gene expression data.  
638 *Bioinformatics*, 26(1):139–140, 2010.
- 639 [22] C. Soneson and M. D. Robinson. Bias, robustness and scalability in single-  
640 cell differential expression analysis. *Nature methods*, 15(4):255, 2018.
- 641 [23] J. W. Squair, M. Gautier, C. Kathe, M. A. Anderson, N. D. James, T. H.  
642 Hutson, R. Hudelle, T. Qaiser, K. J. Matson, Q. Barraud, et al. Confronting  
643 false discoveries in single-cell differential expression. *bioRxiv*, 2021.

- 644 [24] P.-Y. Tung, J. D. Blischak, C. J. Hsiao, D. A. Knowles, J. E. Burnett, J. K.  
645 Pritchard, and Y. Gilad. Batch effects and the effective design of single-cell  
646 gene expression studies. *Scientific reports*, 7:39921, 2017.
- 647 [25] M. Uhlén, L. Fagerberg, B. M. Hallström, C. Lindskog, P. Oksvold,  
648 A. Mardinoglu, Å. Sivertsson, C. Kampf, E. Sjöstedt, A. Asplund, et al.  
649 Tissue-based map of the human proteome. *Science*, 347(6220), 2015.
- 650 [26] C. A. Vallejos, J. C. Marioni, and S. Richardson. BASiCS: Bayesian analysis  
651 of single-cell sequencing data. *PLoS Comput Biol*, 11(6):e1004333, 2015.
- 652 [27] T. Wang, B. Li, C. E. Nelson, and S. Nabavi. Comparative analysis of dif-  
653 ferential gene expression analysis tools for single-cell RNA sequencing data.  
654 *BMC bioinformatics*, 20(1):40, 2019.
- 655 [28] L. M. Weber, M. Nowicka, C. Sonesson, and M. D. Robinson. diffcyt: Differ-  
656 ential discovery in high-dimensional cytometry via high-resolution cluster-  
657 ing. *Communications biology*, 2(1):1–11, 2019.
- 658 [29] L. M. Weber and C. Sonesson. Hdcytodata: Collection of high-  
659 dimensional cytometry benchmark datasets in bioconductor object formats.  
660 *F1000Research*, 8, 2019.
- 661 [30] Y. Zhang, L. Zheng, L. Zhang, X. Hu, X. Ren, and Z. Zhang. Deep single-cell  
662 RNA sequencing data of individual T cells from treatment-naive colorectal  
663 cancer patients. *Scientific data*, 6(1):1–15, 2019.