

A Systematic Comparison of Differential Analysis Methods for CyTOF Data

Lis Arend^{1†*}, Judith Bernett^{1†}, Quirin Manz^{1†}, Melissa Klug^{1,2,3}, Olga Lazareva¹, Jan Baumbach^{4,5}, Dario Bongiovanni^{2,3,6}, Markus List^{1*}

***For correspondence:**

lis.arend@tum.de (LA);
markus.list@wzw.tum.de (ML)

[†]These authors contributed equally to this work

¹Chair of Experimental Bioinformatics, TUM School of Life Sciences, Technical University of Munich, Munich, Germany; ²Department of Internal Medicine I, School of Medicine, University hospital rechts der Isar, Technical University of Munich, Munich, Germany; ³German Center for Cardiovascular Research (DZHK), Partner Site Munich Heart Alliance, Munich, Germany; ⁴Chair of Computational Systems Biology, University of Hamburg, Hamburg, Germany; ⁵Institute of Mathematics and Computer Science, University of Southern Denmark, Odense, Denmark; ⁶Department of Cardiovascular Medicine, Humanitas Clinical and Research Center IRCCS and Humanitas University, Rozzano, Milan, Italy

Abstract Cytometry techniques are widely used to discover cellular characteristics at single-cell resolution. Many data analysis methods for cytometry data focus solely on identifying subpopulations via clustering and testing for differential cell abundance. For differential expression analysis of markers between conditions, only few tools exist. These tools either reduce the data distribution to medians, discarding valuable information, or have underlying assumptions that may not hold for all expression patterns.

Here, we systematically evaluated existing and novel approaches for differential expression analysis on real and simulated CyTOF data. We found that methods using median marker expressions compute fast and reliable results when the data is not strongly zero-inflated. Methods using all data detect changes in strongly zero-inflated markers, but partially suffer from overprediction or cannot handle big datasets. We present a new method, CyEMD, based on calculating the Earth Mover's Distance between expression distributions that can handle strong zero-inflation without being too sensitive.

Additionally, we developed CYANUS, a user-friendly R Shiny App allowing the user to analyze cytometry data with state-of-the-art tools, including well-performing methods from our comparison. A public web interface is available at <https://exbio.wzw.tum.de/cyanus/>.

Introduction

In conventional flow cytometry, single cells are passed through one or multiple lasers while being suspended in a liquid stream. Antibodies are labeled with fluorescent dyes and lasers produce both scattered and fluorescent light signals that are read by detectors (*McKinnon (2018)*). This enables the analysis, identification, and classification of cell populations which is required in multiple disciplines such as immunology, cancer biology, and virology. However, flow cytometry has some severe limitations restricting its utility. The number of parameters analyzed at one time is limited due to the overlap of the light emissions, the rupture of the stains, and the requirement of large

41 cell numbers. (*Gadalla et al. (2019)*)

42 Recently, high-dimensional time-of-flight mass cytometry (CyTOF) has emerged with the ability
43 to identify more than 40 parameters simultaneously. Its advantage over flow cytometry is that
44 antibodies are labeled with metal isotopes instead of fluorophores, allowing scientists to analyze
45 more antibodies in a single run while needing fewer cells per experiment. Traditional flow cytometry
46 would require multiple tubes with different antibody panels to cover the same number of markers
47 (*Gadalla et al. (2019)*). Consequently, CyTOF experiments are a powerful tool to unveil new cell
48 subtypes, functions, and biomarkers in many fields, e.g. the discovery of disease-associated
49 immunologic changes in cancer.

50 Cytometry experiments rely on a panel of antibodies that are associated with a specific ex-
51 perimental condition or phenotype of interest. Usually, the analysis of cytometry data starts by
52 clustering cells into cell subpopulations, followed by a differential expression analysis between
53 and within cell types (*Nowicka et al. (2019)*). Several methods have been developed for testing
54 clusters representing cell populations for differential abundance (DA) between conditions (*Bruggner*
55 *et al. (2014)*, *Arvaniti and Claassen (2017)*, *Weber et al. (2019)*). However, many experiments aim
56 to detect differential states (DS), i.e. differential expression of markers between conditions and
57 within cell populations (see Figure 1A).

58 Diffcyt (*Weber et al. (2019)*) presents two methods for differential expression detection, a linear
59 mixed effect model (LMM) and an adaptation of limma (*Ritchie et al. (2015)*). For both approaches,
60 the data is reduced to median marker expressions per sample and per cluster when comparing
61 conditions. Another recently developed method is CytoGLMM (*Seiler et al. (2021)*) which introduces
62 two multiple regression strategies for finding differential proteins in mass cytometry data: a
63 bootstrapped generalized linear model and a generalized linear mixed model allowing for random
64 and fixed effects.

65 Methods that rely solely on median marker expression and do not take other distribution
66 characteristics into account might be oblivious to certain marker expression patterns. At the
67 same time, the comparison of hundreds of thousands of cells per patient is computationally and
68 statistically tedious. In this study, we provide a clear overview of existing and novel approaches
69 and compare them in different scenarios, highlighting their strengths and weaknesses. As novel
70 approaches, we implemented three statistical tests relying on the medians, a logistic regression
71 using all expression data, two techniques modeling the expression distributions and a method
72 using the Earth Mover's distance (see Figure 1B). A similar approach to the latter, SigEMD, has
73 recently been introduced by *Wang and Nabavi (2018)* for single-cell RNA-seq data. All methods
74 are evaluated on one semi-simulated, one simulated, and two real datasets resembling several
75 experimental scenarios: globally visible differences in various magnitudes, patient-specific effects
76 on paired data, highly zero-inflated marker expressions, and an immune dataset composed of
77 multiple cell types (see Figure 1C).

78 In addition, we present CYANUS (CYtometry ANalysis Using Shiny), a user-friendly R Shiny App
79 available at <https://exbio.wzw.tum.de/cyanus/>. In contrast to existing cytometry analysis platforms
80 like Cytobank (*Kotecha et al. (2010)*) or OMIQ (*Belkina et al. (2019)*), we provide a free platform
81 allowing researchers to analyze normalized, gated cytometry data. To this end, we integrated
82 state-of-the-art methods from CATALYST (*Crowell et al. (2021)*) for preprocessing, visualization,
83 and clustering. Additionally, we integrated those methods for differential marker expression and
84 abundance which showed good performance in our benchmark.

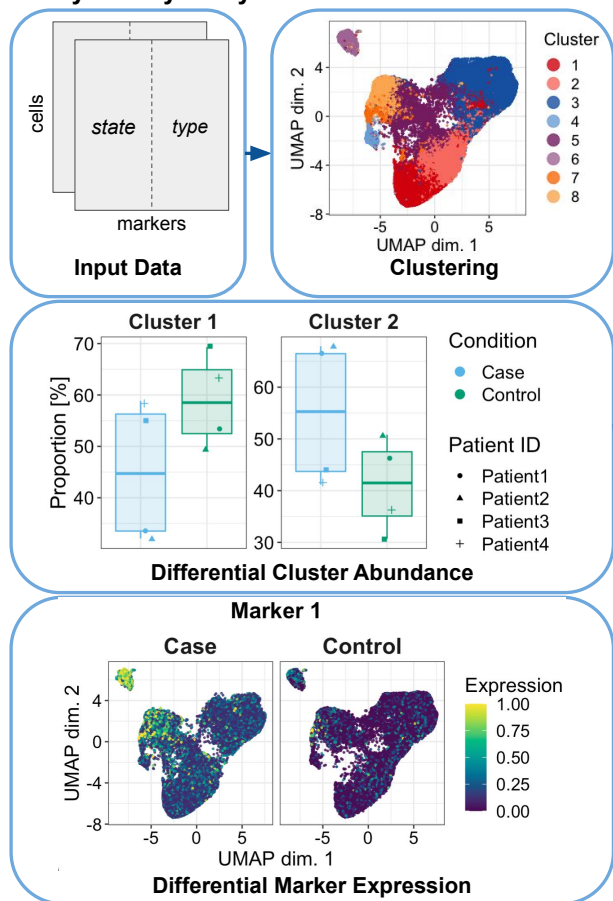
85 Results

86 In this study, we compared existing and novel approaches for detecting differentially expressed
87 markers in CyTOF datasets. The diffcyt package (*Weber et al. (2019)*) employs LMM and limma,
88 which both use median marker expressions per sample and cluster when comparing conditions.
89 In contrast, the methods from the CytoGLMM package (*Seiler et al. (2021)*) make use of the whole
90 data by modeling the condition with all marker expression values (per sample and cluster).

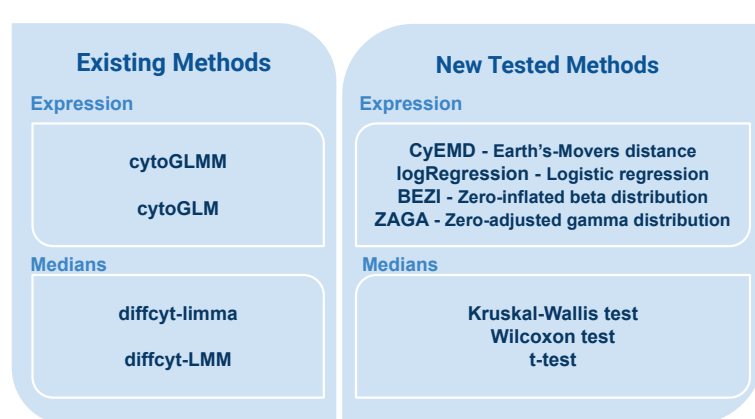
91 We hypothesize that when reducing the datasets to their medians as in diffcyt, simple statistical
 92 tests such as the Wilcoxon rank-sum/signed-rank test, Kruskal-Wallis or (paired) t-test could be
 93 effective. We further use a univariate logistic regression to examine whether the CytoGLMM
 94 approach could be simplified. To explore whether using the entire distribution of the dataset is
 95 beneficial, we modeled the expression data by fitting a zero-inflated beta distribution (BEZI) and a
 96 zero-adjusted gamma distribution (ZAGA), respectively. We further used the Earth Mover's Distance
 97 to compare normalized distributions for each marker (and cluster) between groups (CyEMD). For
 98 more details, please refer to the Methods section.

99 We used four different datasets to evaluate method performance for different data distributions.
 100 The semi-simulated data contains a clear, globally visible artificial signal for four markers. In the
 101 simulated CytoGLMM dataset, five markers are differentially expressed but the differences are only
 102 present on patient-level, not overall. The dual platelet dataset contains strong zero-inflation for
 103 two platelet activation markers, leading to a median marker expression of zero. The PBMC dataset
 104 contains different cell types which is why a differential expression analysis should only be done cell
 105 type (cluster) -wise. The method performance will be discussed for each dataset to show strengths
 106 and weaknesses of the algorithms. Statistical test may report significant differences that are not

A Cytometry Analysis



B Methods Overview



C Datasets Overview

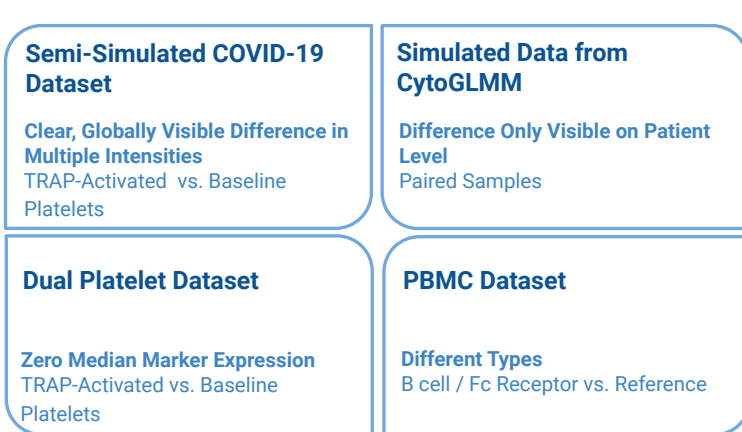


Figure 1. (A) Schematic overview of a differential analysis workflow for cytometry data. In a cytometry experiment, the abundance of state (condition) and type (lineage) markers are measured for each cell. Usually, cells are clustered using type markers to identify cell subpopulations. When differential cluster abundance is analysed, the proportion of cell types between conditions is compared (e.g. condition 2 stimulates the production of cell subpopulation 1). When differential marker expression is analysed, marker expression is compared between conditions within each cluster (e.g. Marker 1 is more highly expressed in condition 1 in clusters 6 and 8). **(B)** Overview of the methods compared in this study. **(C)** Overview of the datasets used in this study. One simulated, one semi-simulated, and two real CyTOF datasets were used to evaluate the methods.

107 meaningful due to their negligible effect size. To account for this, we computed, for all results, the
 108 overall (global) and grouped (accounting for patients or other groups) effect size. It should be noted
 109 that the grouped effect size must be treated with caution due to the small number of samples (see
 110 Methods).

111 Figure 1 shows a schematic representation of a differential analysis workflow for cytometry data
 112 as well as an overview of the methods and datasets investigated in this study.

113 Semi-Simulated COVID-19 Dataset With Clean, Globally Visible Difference Between 114 Conditions

115 For the semi-simulated COVID-19 platelet dataset, we expect to find the markers CD63, CD62P,
 116 CD107a, and CD154 differentially expressed between stimulated and non-stimulated samples
 117 because an artificial signal was only created for these markers specifically. In the original experiment,
 118 platelet expression from baseline samples was compared to expression measured for activated
 119 platelets after stimulation with thrombin receptor-activating peptide (TRAP). The four markers
 120 whose expression was used for creating the signal, hereinafter referred to as state markers, are
 121 known platelet activation markers (*Blair et al. (2018a)*).

122 All other markers (i.e. type markers) detected by any method can be classified as false positives,
 123 since the baseline expression values were not modified.

124 To examine the sensitivity of the methods, we reduced the differences in expression between
 125 the baseline and the spike condition step-wise via a parameter α (Equation 1) which indicates by
 126 what percentage the difference between the spiked-in expression and the baseline expression
 127 is reduced. In the datasets where α was set to 1, no marker should be classified as differentially
 128 expressed because the spiked-in expressions are equal to their corresponding baseline cell.

129 The differences are visible on a global level, as we can observe from the overall effect size which
 130 is large for CD63, CD62P, CD107a, and moderate for CD154. While all of the other 18 markers have
 131 a negligible overall effect size, six of them have a small grouped effect size and one has a moderate
 132 grouped effect size. Supplemental Figures 1 and 2 show the results containing all downsampled
 133 datasets for activation (state) markers and other (type) markers, respectively.

134 Table 1 gives an overview of the methods' performance across all COVID-19 datasets measured
 135 by the F1-score. Sensitivity, specificity, and precision on the same datasets can be found in Supple-
 136 mental Tables 1, 2, and 3, respectively. The methods relying on the median marker expression tend
 137 to perform better with an increasing number of cells. The opposite is the case for both methods
 138 from the CytoGLMM package, as well as BEZI, ZAGA, and the logistic regression.

139 The diffcyt methods can find all activation markers regardless of sample size and signal intensity.

Table 1. Methods' performance measured by F1 scores on the semi-simulated COVID-19 dataset. The means and standard deviations of the scores are reported across the multiple α values.

Number of Cells	1,000	2,000	5,000	10,000	15,000	20,000	4,052,622
diffcyt-DS-limma	0.89 +/- 0	0.89 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0
diffcyt-DS-LMM	0.62 +/- 0	0.67 +/- 0	0.73 +/- 0	0.73 +/- 0	0.73 +/- 0	0.73 +/- 0	0.73 +/- 0
t-test	0.89 +/- 0	0.89 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0
Wilcoxon test	0.96 +/- 0.07	0.85 +/- 0.07	0.96 +/- 0.07	0.96 +/- 0.07	0.96 +/- 0.07	0.96 +/- 0.07	0.96 +/- 0.07
Kruskal-Wallis test	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0
CytoGLM	0.79 +/- 0.05	0.69 +/- 0.14	0.5 +/- 0.07	0.47 +/- 0.09	0.48 +/- 0.13	0.48 +/- 0.12	0.4 +/- 0.05
CytoGLMM	0.74 +/- 0.06	0.47 +/- 0.03	0.38 +/- 0.02	0.36 +/- 0.03	0.34 +/- 0.04	0.33 +/- 0.03	0.37 +/- 0.04
logRegression	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	0.89 +/- 0
ZAGA	0.96 +/- 0.07	0.85 +/- 0.07	0.96 +/- 0.07	0.93 +/- 0.08	0.96 +/- 0.07	0.85 +/- 0.07	0.89 +/- 0
BEZI	0.93 +/- 0.08	0.96 +/- 0.07	0.96 +/- 0.07	0.88 +/- 0.16	0.96 +/- 0.07	0.58 +/- 0.06	0.28 +/- 0.09
CyEMD	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0

140 In the negative controls, both methods find markers in the small downsampled datasets (1000 and
141 2000 cells per patient).

142 Regarding the statistical tests on expression medians, the Kruskal-Wallis test correctly detects
143 all of the state markers and none of the type markers across all sample sizes and α values. The
144 Wilcoxon signed-rank test misses CD154 for $\alpha=0$ regardless of the sample size. In the negative
145 controls, the Wilcoxon test and the t-test find one type marker for $n=2000$ and the t-test finds one
146 type marker for $n=1000$. This observation can be made for all α values, except for $\alpha=1$.

147 The CytoGLMM methods find many false positive type markers across all α values (except
148 for $\alpha=1$). The number of false positives rises with increasing sample size. For the downsampled
149 datasets, more type markers are found for higher α values. Additionally, CD154 cannot be detected
150 by both CytoGLMM methods for $\alpha=0$, as well as by CytoGLM for $\alpha=0.25$.

151 BEZI fails to find different subsets of the state markers across all sample sizes and α values,
152 either due to convergence errors (for datasets bigger than 5000 cells/patient) or because they
153 did not pass the significance threshold of 0.05. Additionally, BEZI classifies PEAR as differentially
154 expressed for all α values in the datasets that were not subsampled. ZAGA and the univariate
155 logistic regression also find PEAR in this dataset but not for $\alpha=1$. Similar to the CytoGLMM methods,
156 ZAGA fails to find CD154 for smaller datasets when α is set to 0.25. Apart from that, ZAGA and the
157 univariate logistic regression do not make any further false predictions.

158 Finally, CyEMD was able to classify all markers correctly.

159 **Simulated Data from CytoGLMM Package With Differences Only Visible on Patient-** 160 **Level**

161 The design of the CytoGLMM simulation leads to patient-wise differences that are not visible globally.
162 The data was simulated in such a way that m01-m05 are differentially expressed and therefore
163 expected to be found. We observe that indeed, the grouped effect size is large for markers m01-m05
164 while the overall effect size is negligible (see Supplemental Figure 3).

165 Table 2 shows an overview of performance measurements on all subsets of this dataset. For
166 more detailed results, we refer to Supplemental Figure 3.

167 The two methods that cannot perform a paired analysis, CyEMD and the Kruskal-Wallis test on
168 marker expression medians, do not find any marker to be differentially expressed. Consequently,
169 they achieve a specificity of 1 and all other measurements are 0 or undefined. Overall, the diffcyt
170 methods have a high performance and gain power for greater numbers of cells. This effect can

Table 2. Methods' performance on the simulated CytoGLMM dataset. Sensitivity, specificity, precision, and F1 score are shown for each method. Means and standard deviations of the scores are reported across the multiple numbers of cells. If no positive classification was made, precision and F1 score cannot be computed and are marked as NaN in the table.

	Sensitivity	Specificity	Precision	F1 Score
diffcyt-DS-limma	0.97 +/- 0.08	0.99 +/- 0.03	0.98 +/- 0.06	0.97 +/- 0.05
diffcyt-DS-LMM	1 +/- 0	0.96 +/- 0.04	0.9 +/- 0.09	0.95 +/- 0.05
t-test	0.97 +/- 0.08	1 +/- 0	1 +/- 0	0.98 +/- 0.04
Wilcoxon test	0.49 +/- 0.34	1 +/- 0	1 +/- 0	0.81 +/- 0.08
Kruskal-Wallis test	0 +/- 0	1 +/- 0	NaN	NaN
CytoGLM	0.8 +/- 0.38	1 +/- 0	1 +/- 0	0.96 +/- 0.1
CytoGLMM	0.97 +/- 0.08	0.97 +/- 0.05	0.94 +/- 0.12	0.95 +/- 0.07
logRegression	1 +/- 0	1 +/- 0	1 +/- 0	1 +/- 0
ZAGA	0.91 +/- 0.16	0.91 +/- 0.12	0.83 +/- 0.19	0.85 +/- 0.14
BEZI	0.86 +/- 0.15	0.93 +/- 0.07	0.84 +/- 0.16	0.83 +/- 0.1
CyEMD	0 +/- 0	1 +/- 0	NaN	NaN

171 also be observed for most other methods. CytoGLMM's and CytoGLM's scores are close to 1 except
172 for the sensitivity scores for CytoGLM which vary more strongly since some of the differentially
173 expressed markers cannot be detected for low cell counts. BEZI and ZAGA lose performance mostly
174 because the algorithms do not converge. Apart from that, they yield high scores. Only the univariate
175 logistic regression can correctly identify all differentially expressed markers without a false positive
176 discovery.

177 **Dual Platelet Dataset With Zero Median Marker Expression**

178 This dataset was generated by collecting two samples from each participant and stimulating one
179 of the two samples with TRAP to activate the platelets. Therefore, we expected to find platelet
180 activation (state) markers like CD63, CD62P, CD154 and CD107a to be differentially expressed
181 between the two conditions. Figure 4C shows that CD154 and CD107a have a median marker
182 expression of zero, posing a challenge for the methods using only marker medians.

183 We tested our methods twice on this dataset. For the first run, the patient ID was included as a
184 grouping variable while the second analysis was unpaired (see Figure 2). We used the Wilcoxon
185 rank-sum test and the Wilcoxon signed-rank test in the unpaired and paired design, respectively.

186 ZAGA, BEZI, and the univariate logistic regression classify all markers as significant or do not
187 converge. The issues of these three methods are examined thoroughly in the discussion.

188 The five algorithms (diffcyt-limma, diffcyt-LMM, t-test, Wilcoxon test, and Kruskal-Wallis test)
189 using only median expressions find the two state markers CD63 and CD62P but are not able to
190 find the zero-inflated markers CD154 and CD107a. In the unpaired run, no type markers are found
191 by these methods. In the analysis with patient ID as grouping variable, the Wilcoxon signed-rank
192 test, t-test, and both diffcyt methods find PAR1, PEAR, and CD69 to be significantly differentially
193 expressed between the non-stimulated and stimulated samples. CD42a was also found by the t-test
194 and both diffcyt methods. Each of the four markers has a large grouped effect size.

195 Two of the three methods using whole marker expression, CyEMD and CytoGLMM, classify all
196 four state markers as significant. CytoGLM only misses CD107a which has a small overall effect
197 size. Since CytoGLMM cannot be run without a random effect, its result for this run is not reliable.
198 Looking at the results for the type markers, CyEMD finds CD141 and CytoGLM finds PAR1 when no
199 paired analysis is performed. After including the patient ID as grouping variable, additional markers
200 are found by all methods able to incorporate this information. CytoGLMM, CytoGLM, and CyEMD all
201 detect CD141 (small overall effect size). The two methods of the CytoGLMM package find several
202 additional markers: CD41, CD61, PAR1, GPIIbIIIa, CD141, CD9, PEAR, CD47, CD31, and CD42a. While
203 PAR1, PEAR, and CD42a are also found by other methods (as mentioned above), some of these
204 markers (CD61, CD47, CD9, and CD31) have negligible effect sizes which is why we classify them as
205 false positives (see Supplemental Figure 4).

206 **PBMC Dataset With Different Cell Types**

207 Since our first real dataset, the dual platelets dataset, only contains one cell type, we also evaluated
208 the different approaches on the PBMC dataset by *Bodenmiller et al. (2012)* which contains eight
209 immune cell types annotated by *Nowicka et al. (2019)*. For each cluster of cell types, differential
210 expression analysis was performed, comparing the reference condition against the cells that were
211 cross-linked with B cell receptor/Fc receptor (BCR/FcR-XL). We expected to find pS6 differentially
212 expressed since these findings have been made in the original paper (*Bodenmiller et al. (2012)*).
213 Supplemental Figure 5 shows an overview of the results.

214 Many markers were significant across all clusters. Independent of the method, overall and
215 grouped effect sizes were large for numerous markers in all clusters (see Supplemental Figure 5).

216 Of all possible 192 marker-cluster combinations (24 markers in 8 cell types), the univariate
217 logistic regression, BEZI, and ZAGA find the most markers to be differentially expressed (168, 156,
218 and 155, respectively).

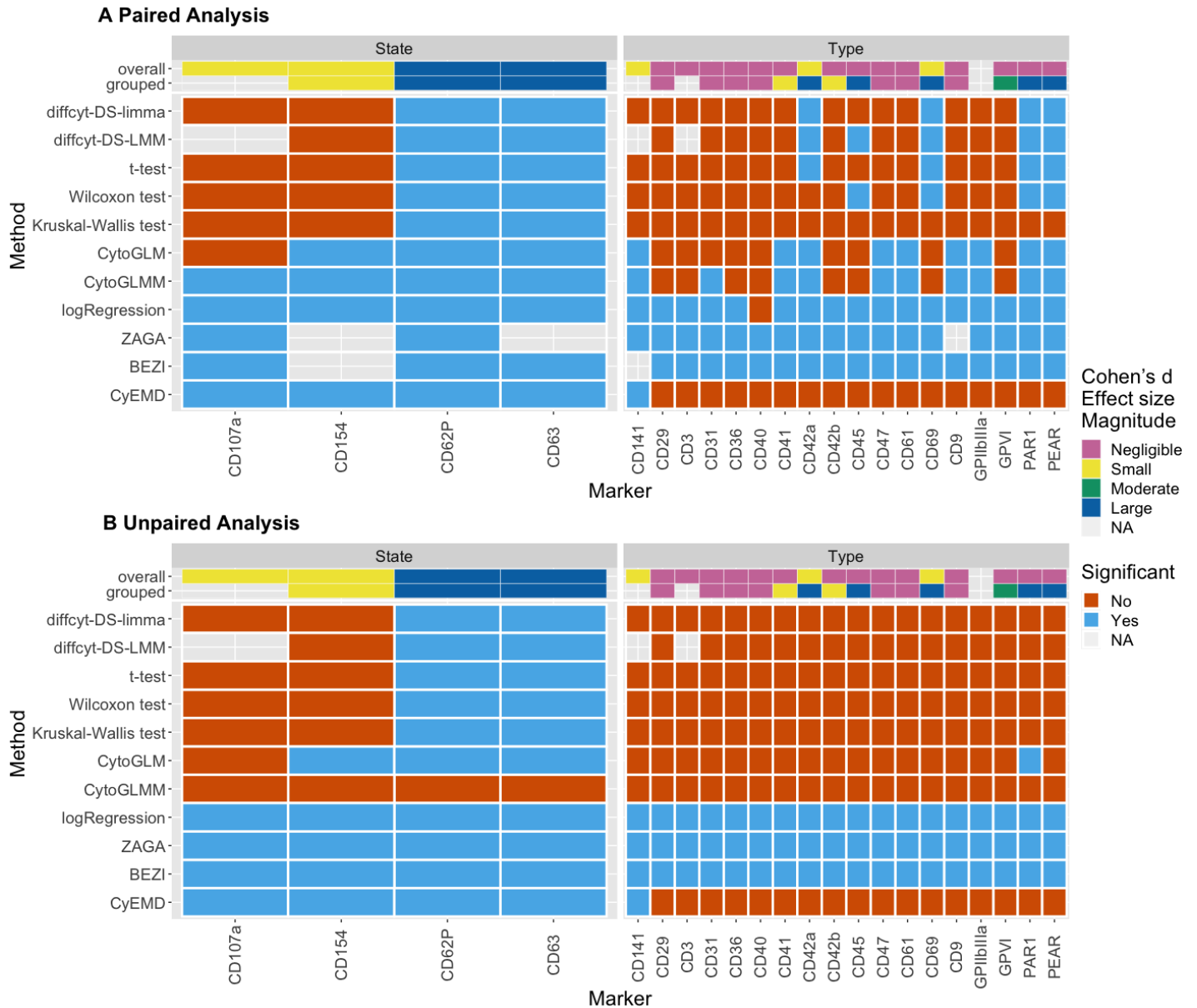


Figure 2. Method results for the dual dataset with patient id as grouping variable **(A)** and without any grouping variable **(B)**. Results colored in blue if the adjusted p-value < 0.05, else in red. Uncolored tiles mean convergence errors of the method for the specific marker. The overall and grouped effect size magnitudes per marker are shown at the top. Overall effect size refers to *Cohen's d* magnitudes using all expression data between two conditions. The magnitudes indicated by grouped effect size are computed in a paired fashion on the median marker expressions per sample. Wilcoxon test refers to the Wilcoxon signed-rank test and the Wilcoxon rank-sum test for the paired and the unpaired analysis, respectively. Markers are divided into their marker class (state and type).

219 The methods that are not able to include a grouping variable, CyEMD and the Kruskal-Wallis
 220 test, find the least markers to be differentially expressed (86 and 88, respectively). In contrast to the
 221 dual dataset results, CytoGLMM and CytoGLM do not produce more positive predictions than the
 222 statistical tests, CyEMD, or the diffcyt methods.

223 Runtime

224 The runtimes for the complete datasets are shown in Table 3. For the runtimes of the subsampled
 225 datasets, please refer to Supplemental Figure 6.

226 The diffcyt methods outperform all other methods in terms of runtime. All in all, methods that

227 use median marker expressions are fast independent of sample size. CytoGLMM and the unpaired
 228 logistic regression are quick as well, even though they take the whole distribution into account.
 229 The paired univariate logistic regression, CyEMD, and ZAGA have moderate runtimes while
 230 CytoGLM and BEZI often run more than six hours on the big datasets.

Table 3. Runtime of the methods on the different datasets. The methods that reduce the data to medians and CytoGLMM have very low runtime requirements while CytoGLM and BEZI are slow on big datasets. The univariate logistic regression, CyEMD, and ZAGA have moderate runtimes.

	Semi-Simulated COVID-19	Simulated CytoGLMM	Paired Dual Platelets	Unpaired Dual Platelets	PBMC
Number of Cells	4,052,622	4,400,000	4,491,504	4,491,504	906,815
diffcyt-DS-LMM	26 +/- 2 sec	29 sec	35 sec	33 sec	5 sec
diffcyt-DS-limma	29 +/- 5 sec	38 sec	47 sec	41 sec	5 sec
t-test	1.03 +/- 0.06 min	1.06 min	1.03 min	1.09 min	28 sec
Kruskal-Wallis test	1.04 +/- 0.09 min	1.01 min	1.08 min	1.04 min	31 sec
Wilcoxon test	1.04 +/- 0.06 min	1.07 min	1.11 min	1.07 min	29 sec
CytoGLMM	1.92 +/- 0.41 min	1.18 min	2.02 min	7.82 min	11 sec
CyEMD	1.9 +/- 0.2h	2.1h	2.0h	1.9h	6.66 min
logRegression	2.5 +/- 0.2h	2.2h	2.6h	3.62 min	49.73 min
ZAGA	2.7 +/- 0.4h	2.3h	2.1h	49.53 min	5.9 min
CytoGLM	6.5 +/- 1.6h	4.0h	6.5h	7.8h	15.92 min
BEZI	9.8 +/- 1.6h	9.7h	9.7h	4.7h	26.13 min

231 Discussion

232 Semi-Simulated COVID-19 Dataset With Clean, Globally Visible Difference Between 233 Conditions

234 Regarding the grouped effect sizes, we can observe a problem that occurs when the differences
 235 between the groups are very small, yielding a standard deviation close to zero (see Equation 3). By
 236 dividing by a value close to zero, bigger values for the effect sizes are obtained, hence seven of the
 237 markers that were not spiked in appear to have a small or even moderate grouped effect size. The
 238 difference in means between the two conditions is not significant, as shown by the paired t-test.
 239 Therefore, we recommend checking the effect size when a marker is classified as differentially
 240 expressed and additionally, checking the results of the paired t-test when the grouped effect size is
 241 not negligible because both methods compare paired means.

242 The diffcyt methods and the statistical tests perform well, especially for higher sample sizes. We
 243 hypothesize that the markers that are found in the small negative control datasets were detected
 244 because of noise in the measured data. Due to the law of large numbers, the median becomes more
 245 reliable for higher cell counts. Therefore, methods that reduce the expression data to medians
 246 become more stable with growing dataset size.

247 The CytoGLMM methods produce a high number of false positives for all α values except for
 248 $\alpha=1$, especially with rising sample size. A possible explanation could be that the multivariate
 249 generalized mixed effect models become too sensitive to small changes when there are only
 250 few bigger differences (here, CD62P and CD63) because all markers are included as explanatory
 251 variables. Therefore, the condition is modeled as a result of various small changes which are
 252 present because of the semi-simulated nature of the data. The increasing sample size seems to
 253 reduce the magnitude of the p-values.

254 The Wilcoxon signed-rank test and the CytoGLMM methods miss CD154 for $\alpha=0$ and $\alpha=0.25$ but
 255 find it for the other α values because of the way the artificial signal was created. Looking at the

256 medians per patient for the full number of cells (see Supplemental Figure 7), it is visible that the
257 medians of the spike condition are higher for $\alpha=0.25$, 0.5, and 0.75 than for $\alpha=0$ in two patients.
258 Additionally, the medians are extremely close to zero. Due to this strong zero-inflation, the spiked-in
259 expressions for CD154 were sometimes smaller than the baseline expressions. As Equation 1
260 leads to a convergence of the activated measurements towards the baseline measurements, the
261 spiked-in values become higher for the cells where the activated measurements were smaller than
262 the baseline measurements. Therefore, the expression is less zero-inflated in these two patients
263 and the marker can be found by the three methods for higher α values.

264 For BEZI, we see its high sensitivity for big sample sizes (see dual dataset), since PEAR is found
265 across all α settings in the big dataset, while it is not found for the downsampled datasets. ZAGA
266 and the univariate logistic regression yield reliable results for datasets with a clean, globally visible
267 difference, especially for smaller sample sizes.

268 **Simulated Data from CytoGLMM Package With Differences Only Visible on Patient-** 269 **Level**

270 Because this data is paired, we expect that only methods that can handle paired data can detect
271 the differentially expressed markers between the two conditions. This is confirmed as all methods
272 except for CyEMD and the unpaired Kruskal-Wallis test detect the differential expression and can
273 be sensitive to small changes in expression that are only detectable at the patient level.

274 CytoGLMM's false detection of one marker suggests an over-sensitivity further described in the
275 next section.

276 **Dual Platelet Dataset With Zero Median Marker Expression**

277 The results for this dataset clearly show the problem of reducing the data on median marker ex-
278 pressions to perform differential expression analysis. Methods taking the whole marker expression
279 into account find markers with zero-median marker expression, whereas methods working on the
280 medians are not able to find these.

281 Furthermore, the results differ depending on the applied method. The markers PAR1 and PEAR
282 are found by several methods. While PEAR has a higher expression in the stimulated condition, PAR1
283 is less expressed in this condition (see Supplemental Figure 8). In literature, the PEAR receptor has
284 been described to be increased on the platelet membrane after stimulation with several activators
285 (*Kauskot et al. (2012)*), while the effect on PAR1 expression after stimulation depends on the agonist.
286 Studies using a PAR1-AP are in line with our findings and show a decreased amount of the PAR1
287 receptor on the platelet surface after stimulation (*Ramström et al. (2008)*).

288 CD69 which is found by limma, LMM, the Wilcoxon test, and the t-test, shows a higher signal after
289 stimulation. Several studies have observed a similar trend for CD69 increase upon stimulation (*Testi*
290 *et al. (1990, 1992)*). CD42a is detected by the two diffcyt methods, the CytoGLM/M methods, and the
291 t-test, and shows a decreasing trend after TRAP stimulation. This also has been previously shown in
292 platelets using CyTOF (*Blair et al. (2018b)*). Several other studies examined a decrease of CD42a
293 expression after stimulation with activators ADP (*Braune et al. (2014)*) and collagen (*Hagberg et al.*
294 *(1997)*). The biological reason behind the differential expression of the two markers CD141 and
295 CD45 remains unclear. In general, CD141 is not found to be expressed on platelets (*Bongiovanni*
296 *et al. (2021)*) whereas CD45 has shown to be present on the surface of several platelets (*Gabbasov*
297 *et al. (2014)*).

298 The application of ZAGA, BEZI, and the univariate logistic regression is unfeasible for a real
299 dataset of this size. CytoGLMM and CytoGLM produce at least three false positives due to their high
300 sensitivity. Additionally, CytoGLM misses one of the two highly zero-inflated activation markers.
301 The diffcyt methods, the t-test, and the Wilcoxon signed-rank test perform fast and yield reliable
302 results but miss the two activation markers that have a median of zero. The Kruskal-Wallis test
303 performs worse than the Wilcoxon signed-rank test on this dataset because it is not able to handle
304 paired data and could therefore not detect markers like PAR1, PEAR, CD69, or CD42a. Lastly, CyEMD

305 detects the globally visible changes for the activation markers and CD141 but fails to detect any of
306 the changes that can only be seen on the patient level as seen in Figure 2.

307 **PBMC Dataset With Different Cell Types**

308 The evaluation of this dataset is limited by the number of cells per sample and cluster (see Sup-
309 plemental Figure 9). For cell types with less than 1000 cells per sample, noise is distorting the
310 analysis.

311 When *Weber et al. (2019)* evaluated their diffcyt methods on this dataset, they could confirm
312 that pS6 is differentially expressed in B-cells. From Figure 5, it becomes apparent that all methods
313 can find this marker. Moreover, *Nowicka et al. (2019)* showed that the diffcyt methods identify pS6
314 to be also differentially expressed in other cell types. All methods tested in this study confirm this
315 finding and find pS6 differentially expressed in all cell types except for dendritic cells.

316 In contrast to the dual platelet dataset, the univariate logistic regression, BEZI, and ZAGA were
317 not suffering from a clear over-identification of markers. Because the PBMC dataset is rather small
318 (172,791 cells in total vs. 4,491,504 in the dual platelet dataset), we hypothesize that the higher the
319 number of cells, the less suitable these three methods become. This is due to the influence of large
320 sample sizes on the magnitude of the p-values (*Lin et al. (2013)*).

321 Compared to the dual platelet dataset, the CytoGLM/M methods did not identify more markers
322 as significantly differentially expressed than the other methods, even though there are more
323 markers with a large effect size.

324 **Conclusion and Outlook**

325 Existing approaches for differential marker expression analysis were compared with simple and
326 advanced novel approaches that rely either on median or on full marker expression data using two
327 real, one semi-simulated, and one simulated dataset.

328 A limitation on the level of dataset evaluation is that we could not interpret the results obtained
329 on the PBMC dataset biologically. We could therefore not describe which markers were falsely
330 classified as differentially expressed and which markers were overlooked. Additionally, we did not
331 include a dataset with batch effects but assumed that the data had already been corrected for
332 it. Theoretically, it should be possible to include a batch effect as a random effect or additional
333 term in a model. This can be done for all the approaches we evaluated but the statistical tests and
334 CyEMD. Finally, the downsampling of the spiked and the CytoGLMM datasets was not repeated
335 multiple times. If that would have been done, the results would be more reliable and robust. In
336 this study, repeating the evaluations that many times was not feasible because of the high runtime
337 requirement of BEZI and ZAGA.

338 All in all, the diffcyt methods perform fast and yield good, trustworthy results when the median
339 of the differentially expressed marker is unequal to zero. Nevertheless, they did not outperform
340 a simple, Wilcoxon signed-rank test or t-test on the medians, meaning that a more complicated
341 model is not certainly necessary to detect significant differences in CyTOF marker medians. The
342 comparison with the Kruskal-Wallis test on marker medians shows that the clear advantage of the
343 Wilcoxon/t- test is the ability to compute a paired test statistic.

344 Regarding the cytoGLMM methods, we observe that small, individual changes can be detected as
345 well as globally visible changes on very clean data, even when it is strongly zero-inflated. Additionally,
346 cytoGLMM is fast even though it takes the whole distribution into account. On the other hand, the
347 methods are extremely sensitive to changes even without a small, grouped effect size and classify
348 many markers to be differentially expressed, especially with growing dataset size. Therefore, we
349 recommend checking for overlaps between cytoGLMM and other methods, making diagnostic plots
350 and looking at the effect size magnitude when running cytoGLMM on larger, real datasets.

351 BEZI, ZAGA, and the univariate logistic regression proved to be infeasible for larger, real datasets.
352 While the performance on the completely artificial CytoGLMM dataset was acceptable, the method

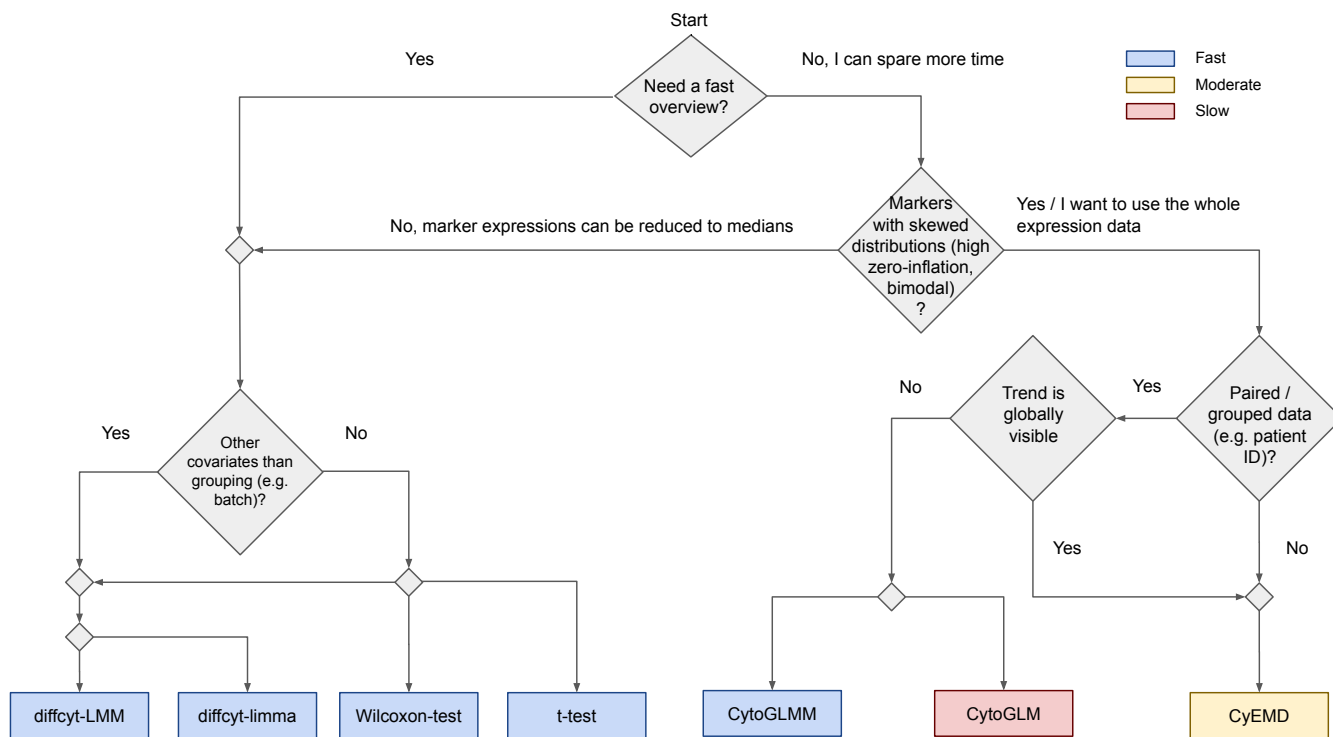


Figure 3. Overview of the methods suitable for CyTOF data. Several scenarios can occur while analyzing CyTOF data. This graph helps to identify the most suitable method and includes the runtime of the different methods.

353 performance dropped for the semi-simulated spike dataset and eventually only produced positive
 354 predictions on the real, dual platelets dataset. Additionally, BEZI is unacceptably slow.

355 Finally, our novel method CyEMD exploits the advantages of taking the whole marker expressions
 356 into account and still performs well on big datasets because it partitions the distribution into
 357 bins and computes p-values via permutation tests. We showed that the EMD approach can
 358 detect differentially expressed markers that are strongly zero-inflated in an acceptable amount of
 359 time. Additionally, the approach should be able to find differences in bimodal or skewed marker
 360 expressions, even when the medians are similar. A disadvantage to the EMD approach is that it
 361 cannot detect differentially expressed markers when the changes are only visible by comparing
 362 expressions group- or patient-wise.

363 Our results across datasets with different properties show that each of the tested methods
 364 comes with its own strengths and weaknesses. Taking factors like runtime, zero-inflation and
 365 skewness and sample groups into account, we offer a guideline for users to choose optimal
 366 methods for their analysis (Figure 3). However, often several methods are suitable for a given
 367 scenario and should be compared to obtain robust and interpretable results.

368 To make such a comparative analysis easily accessible, we integrated the diffcyt methods, the
 369 Wilcoxon rank-sum and signed-rank test, the t-test, the cytoGLMM methods, and CyEMD into a user-
 370 friendly R Shiny App CYANUS available at <https://exbio.wzw.tum.de/cyanus/>. CYANUS (CYtometry
 371 ANalysis Using Shiny) allows the user to analyze gated and normalized cytometry data (i.e. flow
 372 cytometry as well as CyTOF) with state-of-the-art methods from CATALYST (Crowell et al. (2021)).
 373 For differential abundance analysis, we integrated the methods included in the diffcyt package. All
 374 differential analysis methods can be easily compared to each other, enabling thorough analysis of
 375 cytometry data exploiting the advantages of the various approaches.

376 **Methods**

377 **Data Description**

378 For the evaluation of the differential expression methods, we worked with four different datasets.
 379 The methods were tested on one semi-simulated, one simulated, and two real CyTOF datasets (see
 380 Figure 1).

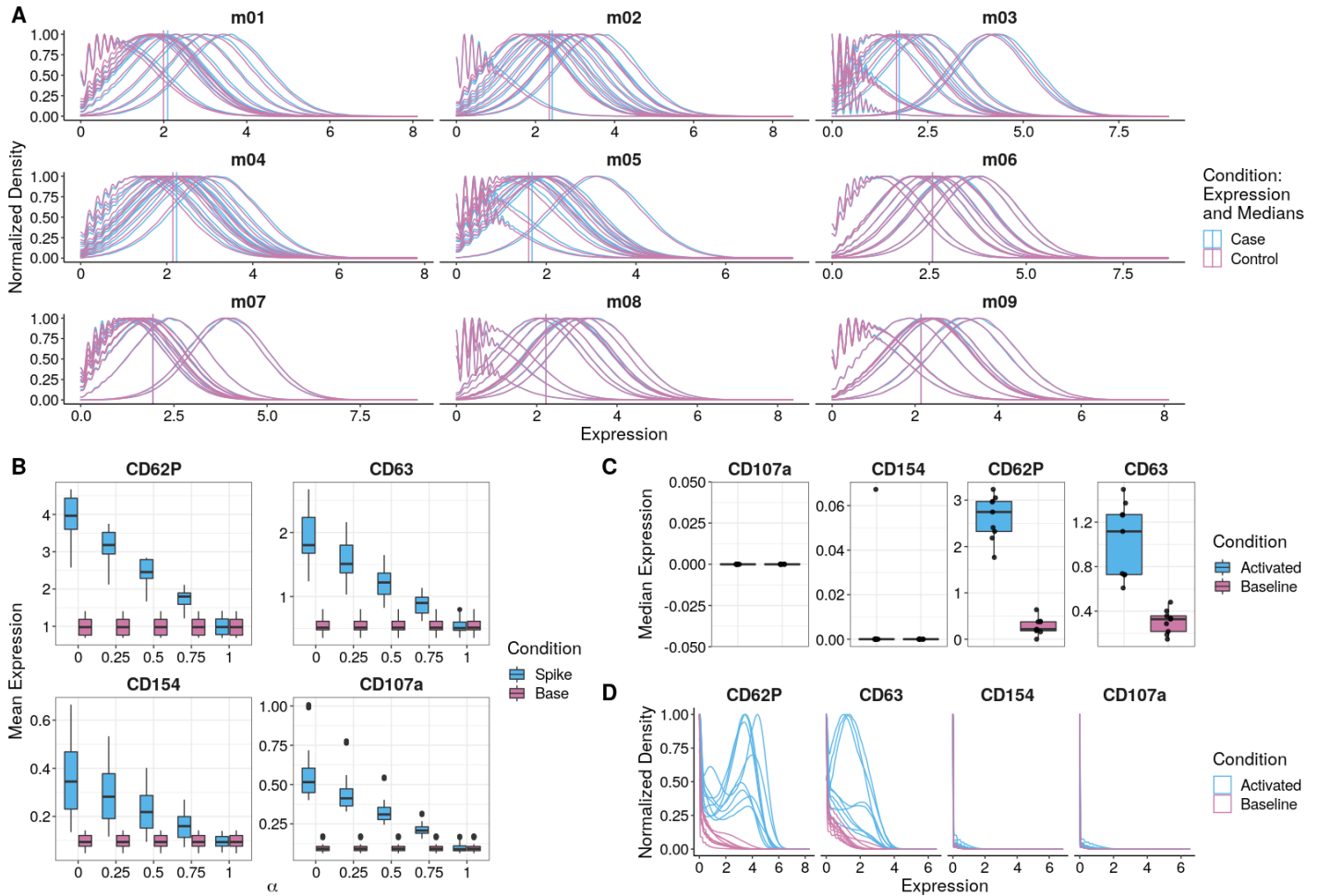


Figure 4. Marker Expressions of the simulated CytoGLMM (A), semi-simulated COVID-19 (B), and the dual platelet datasets (C,D). **(A)** Normalized density of the markers m01-m09 of the dataset simulated using the CytoGLMM data generation process by *Seiler et al. (2021)*. The markers m01-m05 are simulated to be differentially expressed in such a way that the expression differs slightly but consistently for each patient. Meanwhile, the median marker expressions of the whole dataset, marked by the vertical lines, do not differ significantly. **(B)** Mean expressions for the four spiked-in activation markers at different intensities. For $\alpha=0$ (full intensity), the originally measured expressions of the corresponding activated sample were used. Subsequently, α was repeatedly increased by 0.25 in order to reduce the difference between the spiked and the base condition so that the differences would become harder to detect. $\alpha=1$ was used as control dataset. **(C)** Median expression of state markers of the dual platelet dataset. Markers CD62P and CD63 are higher expressed in the activated condition. The median marker expression of CD107a and CD154 is zero, except for one sample. **(D)** Normalized density of state markers of the dual platelet dataset. CD107a and CD154 show a small difference in the expression.

381 **Semi-Simulated COVID-19 Data**

382 The semi-simulated COVID-19 dataset originates from the University Hospital rechts der Isar, Munich,
 383 Germany (*Bongiovanni et al. (2021)*). The original dataset comprises CyTOF data of 8 symptomatic
 384 SARS-CoV-2-infected patients, hospitalized between March and May 2020. Additionally, 11 healthy
 385 donors were included in the study. A baseline sample (non-stimulated platelets) and one sample
 386 stimulated with TRAP was prepared for each donor.

387 In order to study the sensitivity of the methods to changes in the expression patterns, we
388 performed the following data generation procedure. Firstly, the baseline healthy samples were
389 randomly split in half. Half of a sample was used for randomly spiking in the expression values for
390 the four known activation markers (CD62P, CD63, CD107a, CD154) from the activated sample of the
391 corresponding patient. Because this leads to very clear, well distinguishable results, we reduced
392 the differences in expression between baseline and spike expressions for the four markers using
393 the following formula:

$$c_{m,x_i} := 5 \sinh \left[\operatorname{asinh} \left(\frac{c_{m,y_i}}{5} \right) - \alpha \left(\operatorname{asinh} \left(\frac{c_{m,y_i}}{5} \right) - \operatorname{asinh} \left(\frac{c_{m,x_i}}{5} \right) \right) \right] \quad (1)$$

394 where m is the marker, c_{m,x_i} is the raw value measured for the baseline sample for cell x_i , c_{m,y_j} is
395 the raw value measured for the activated sample for cell y_j , $X = x_1, \dots, x_{N/2}$ are the indices of the
396 baseline cells whose expression was randomly replaced and $Y = y_1, \dots, y_{N/2}$ are the indices of the
397 activated cells whose values were used for spiking. Since we wanted to observe the differences in
398 the asinh transformed expression values, the reduction was made on the level of the transformed
399 values. Using the formula, five datasets were produced by setting α to 0 (full intensity), 0.25, 0.5,
400 0.75, and 1.0 (control) (see Figure 4B). Each dataset contains eleven paired samples with 4,052,622
401 cells in total (see Supplemental Table 4 for the number of cells per sample). This approach was
402 inspired by the diffcyt benchmarking strategy ([Weber et al. \(2019\)](#)). In contrast to their approach, we
403 did not use differences in means and standard deviations between the two conditions for reducing
404 the signal but the actual differences between c_{m,x_i} and c_{m,y_j} .

405 Simulated CytoGLMM Data

406 To investigate the methods' handling of data without global but paired differences in expression,
407 we used a customized version of the data simulation process described by [Seiler et al. \(2021\)](#). The
408 algorithm samples from a Poisson GLM with an underlying hierarchical model combining effects on
409 cell and donor-level for two conditions. Figure 4A shows that this leads to expression differences
410 on a patient level, but not overall. The resulting simulated dataset used in this study consists of 20
411 markers, of which 5 are differentially expressed, in 22 paired samples from 11 patients with 200,000
412 cells per sample.

413 Dual Platelet Data

414 We used a CyTOF platelet dataset originating from the University Hospital rechts der Isar, Munich,
415 Germany, consisting of platelet heterogeneity measurements of patients with chronic coronary
416 syndrome receiving dual anti-thrombotic therapy. The dataset contains 4,491,504 cells and includes
417 18 paired samples from 9 donors in two conditions: non-stimulated and stimulated (TRAP). For the
418 exact number of cells per sample, refer to Supplemental Table 5. The panel containing 22 protein
419 markers (see Supplemental Table 6) includes four well-known platelet activation markers ([Blair
420 et al. \(2018a\)](#)). Two of the platelet activation markers, CD63 and CD62P, are known to be highly
421 upregulated after TRAP stimulation, whereas CD107a and CD154 are upregulated less strongly (see
422 Figures 4C, D).

423 PBMC Data

424 The peripheral blood mononuclear cells (PBMCs) dataset originating from [Bodenmiller et al. \(2012\)](#)
425 consists of samples from 8 healthy donors in 12 conditions. [Nowicka et al. \(2019\)](#) performed a
426 complete CyTOF analysis on a subset of this data containing the reference and one stimulated
427 condition. In the stimulated condition, the cells were cross-linked with B cell receptor/Fc receptor
428 for 30 minutes. This subset consists of 172,791 cells in 16 paired samples from 8 patients (see
429 Supplemental Table 7). [Nowicka et al. \(2019\)](#) manually merged 20 clusters obtained via meta
430 clustering into 8 cell populations which were made publicly available. In this study, this annotated
431 and well-described subset was used.

432 Downsampling of Artificial Datasets

433 To review changes in the methods' power and runtime with respect to sample size, we downsampled
434 the two simulated datasets. Both the spiked COVID-19 and the simulated CytoGLMM data was
435 subsampled to contain 1000, 2000, 5000, 10000, 15000, and 20000 cells per patient. For both
436 datasets, we sampled in such a way that the smaller sets are always subsets of the bigger ones. The
437 same cells were used in the COVID-19 dataset for different α values, to ensure a fair comparison.

438 Effect Size

439 To quantify the difference between marker expressions, we computed *Cohen's d* (*Cohen (1977)*) for
440 each marker in every dataset using the `rstatix` R package (*Kassambara (2021)*). The thresholds for
441 the absolute value of d to consider the magnitude of the effect size at least *small*, *moderate*, and
442 *large* are 0.2, 0.5, and 0.8, respectively. Values smaller than 0.2 are referred to as *negligible*. The
443 effect size was calculated overall (on the whole expression) and grouped (based on the median
444 marker expression of the paired samples). The overall effect size compares marker intensities
445 between two conditions by using their mean and (shared) standard deviation:

$$d = \frac{\mu_1 - \mu_2}{\sigma} \quad (2)$$

446 Differences on the patient-level can be captured with a paired effect-size estimation, defined
447 as grouped effect size. To obtain paired data points, the expression median was computed for
448 each sample. Additionally, the paired effect-size allows us to check whether significant results from
449 the paired t-test have considerable effect-sizes and whether effect-sizes of higher magnitudes are
450 statistically significant because both methods investigate differences in population means. Since
451 the sample size for the paired calculation is limited to the number of patients n , which is smaller
452 than 20 for all datasets, and there is a known upwards bias for small sample sizes, we used Hedges'
453 correction to adjust for that. The grouped effect size with Hedges correction is computed as follows
454 (*Hedges and Olkin (1985)*):

$$d_z = \frac{\mu_z}{\sigma_z} \times \frac{n-2}{n-1.25}, \quad (3)$$

455 where x and y are the median marker expressions of two groups with paired samples and z is
456 their difference $z = x - y$.

457 Differential Analysis

458 In this work, we compared the existing approaches for differential marker expression analysis from
459 *Weber et al. (2019)* (`diffcyt-limma`, `diffcyt-LMM`) and *Seiler et al. (2021)* (`CytoGLM`, `CytoGLMM`) with
460 simple and advanced novel approaches that rely either on median or on full marker expression data.
461 The simple approaches consisted of a t-test, a Wilcoxon test (both paired and unpaired) on median
462 marker expressions, a Kruskal-Wallis test on median marker expressions, and a univariate logistic
463 regression predicting the condition from the whole marker expression profiles. More advanced
464 approaches comprised modeling the marker expression distributions with a zero-inflated beta
465 distribution (BEZI) and a zero-adjusted gamma distribution (ZAGA). Furthermore, we developed
466 CyEMD, a method which compares the normalized distributions using the Earth Mover's Distance
467 (EMD). All p-values mentioned in this study have been adjusted per method and dataset to control
468 the false-discovery rate using the Benjamini-Hochberg procedure at a significance level of 0.05.

469 Diffcyt Methods

470 The `diffcyt limma` method fits a linear model for each marker-cluster combination, predicting the
471 sample medians from the corresponding conditions. The LMM method builds a linear mixed-effects
472 model and can therefore handle random effects in contrast to the `limma` method where a grouping
473 variable can be included only as an additional fixed effect (*Weber et al. (2019)*). The `diffcyt` methods
474 can easily incorporate other covariates such as batch effects in their model as additional terms.

475 CytoGLMM Methods

476 Instead of predicting the expression from the conditions, the CytoGLMM methods fit a generalized
477 mixed model predicting the conditions from the whole expression vectors. The package contains
478 two methods, CytoGLMM and CytoGLM. The former can only handle grouped data since it relies on
479 a random effect like patient ID whereas the latter can also handle unpaired data. CytoGLM builds a
480 bootstrapped generalized linear model while CytoGLMM builds a generalized linear mixed model
481 (*Seiler et al. (2021)*). In this study, 500 bootstrap replications were used. These methods can also
482 include additional terms in their model.

483 Logistic Regression

484 In order to find out whether the CytoGLMM approach based on the whole marker expression could
485 be simplified, we fitted univariate logistic regression models per marker and cluster and extracted
486 the p-value from the regression model. A multivariate approach was omitted since the markers
487 are not statistically independent by design. CytoGLMM partially evades this problem by fitting a
488 hierarchical model containing random slopes and intercepts for the grouping variable (patient ID)
489 which assumes dependent errors.

490 Approaches Modeling the Expression: BEZI, ZAGA

491 As CyTOF data can be strongly zero-inflated (*Papoutsoglou et al. (2019)*), we fit a zero-inflated
492 beta distribution (BEZI) as well as a zero-adjusted gamma distribution (ZAGA) to our expression
493 data. As a basis, we chose the gamma distribution for modeling non-zero expressions because
494 it was demonstrated that the gamma distribution fits expression data more often than other
495 non-Normal distributions (*de Torrenté et al. (2020)*). A common choice for single cell RNA-seq data
496 is the negative binomial distribution (*He et al. (2021)*) which is not suitable for CyTOF data as it
497 requires discrete values. Therefore, we selected the beta distribution as a conjugate to negative
498 binomial distribution, i.e. it belongs to the same probability distribution family. To model the marker
499 expression distributions of CyTOF data, the condition was used as an explanatory variable and its
500 model coefficient was tested for equality to zero. For this, we used the `gamlss` and the `gamlss.dist`
501 packages (*Rigby and Stasinopoulos (2005)*; *Stasinopoulos and Rigby (2021)*). To model changes on
502 a patient level for paired data, random intercepts can be included.

503 The zero-adjusted gamma distribution, $ZAGA(\mu, \sigma, \nu)$, is a continuous distribution on $(0, \infty)$.
504 The response variable $Y \sim ZAGA(\mu, \sigma, \nu) \in [0, \infty)$ is modeled using the mixed probability function
505 $f_Y(y|\mu, \sigma, \nu)$:

$$f_Y(y|\mu, \sigma, \nu) = \begin{cases} \nu & \text{if } y = 0 \\ (1 - \nu)f_{Y_1}(y|\mu, \sigma) & \text{if } y > 0 \end{cases} \quad (4)$$

506 where $y \geq 0$, $\mu > 0$, $\sigma > 0$, and $0 < \nu < 1$, and where $Y_1 \sim GA(\mu, \sigma)$. The parameter ν is the
507 non-zero probability for $Y = 0$. For $Y \in (0, \infty)$, Y is gamma-distributed.

508 The zero-inflated beta distribution, $BEZI(\mu, \sigma, \nu)$, is defined on $[0, 1)$. The response variable
509 $Y \sim BEZI(\mu, \sigma, \nu) \in [0, 1)$ is modeled using the mixed probability function $f_Y(y|\mu, \sigma, \nu)$:

$$f_Y(y|\mu, \sigma, \nu) = \begin{cases} \nu & \text{if } y = 0 \\ (1 - \nu)f_W(y|\mu, \sigma) & \text{if } 0 < y < 1 \end{cases} \quad (5)$$

510 where $0 < \sigma < 1$, $\mu > 0$, $0 < \nu < 1$. The beta distribution $f_W(y|\mu, \sigma)$ is based on the work of *Ospina*
511 *and Ferrari (2012)*. To fit a zero-inflated beta distribution on CyTOF data, the marker expressions
512 were first scaled to the range $[0, 1)$. For further details regarding the implementation of `gamlss`,
513 please refer to *Rigby et al. (2020)*.

514 CyEMD

515 Our novel approach, CyEMD, uses the Earth Mover's Distance to compare normalized distributions
516 for each marker (and cluster) between groups.

517 For two normalized histograms P and Q , the EMD is calculated by minimizing the cost of
518 transforming one into the other. The histograms are represented as $P = \{(p_1, w_{p1}), \dots, (p_n, w_{pn})\}$ and
519 $Q = (q_1, w_{q1}), \dots, (q_n, w_{qn})$, where p_i/q_j is the center of the i th/ j th histogram bin and w_{p_i}/w_{q_j} describes
520 the height of the corresponding bin for P/Q .

521 To transform histogram P into histogram Q , certain proportions of the bins p_i, q_j differing
522 between P and Q need to be moved to other bins. The optimization problem for this task is how
523 much has to be transferred from one bin to another bin (defined as flow $F = [f_{ij}]$) in order to
524 minimize the cost. The flow is weighted according to the distances d_{ij} between the bins such that
525 transporting a high amount of a bin over a long distance is penalized (**Rubner et al. (1998)**):

$$COST(P, Q, F) = \sum_{i=1}^n \sum_{j=1}^n f_{ij} d_{ij} = \sum_{i=1}^n \sum_{j=1}^n |w_{p_i} - w_{q_j}| \cdot |p_i - q_j| \quad (6)$$

526 After normalizing the minimal cost by the overall flow we get

$$EMD(P, Q) = \frac{\min COST}{\sum_{i=1}^n \sum_{j=1}^n f_{ij}}. \quad (7)$$

527 Since the expression densities in CyTOF data can have different ranges for distinct values, we
528 use a flexible bin width estimated by the Freedman–Diaconis rule evaluated on all nonzero values:

$$\text{Bin width} = 2 \frac{IQR(x)}{\sqrt[3]{n}} \quad (8)$$

529 where $IQR(x)$ is the interquartile range of nonzero marker expressions and n is the number of
530 observed expressions (**Freedman and Diaconis (1981)**).

531 To determine the significance of the EMD between two marker expressions, a permutation
532 test (500 permutations) that permutes the condition labels sample-wise is performed to obtain a
533 p-value for each marker.

534 As opposed to **Wang and Nabavi (2018)**, we compute the EMD on normalized histograms, which
535 can be done in linear time. In order to speed up the computationally intensive EMD computation,
536 we implemented this part in C++. Furthermore, SigEMD permutes the labels cell-wise instead of
537 sample-wise which proved to be infeasible for big datasets since the empirical p-values become
538 smaller with growing dataset size.

539 Data Availability

540 The scripts for the analysis and the code for the Shiny App are available at <https://github.com/biomedbigdata/cyanus>
541 under the GPL-3 license.

542 The original COVID-19 dataset is publicly available at flowrepository.org, accessible at repository
543 ID FR-FCM-Z4AE. The script for producing the semi-simulated COVID-19 data is provided in the
544 Github repository. The simulated CytoGLMM data can be reproduced using a script of the Github
545 repository. Access to the dual dataset (9 patients) is granted upon request. The original PBMC
546 dataset is published at www.cytobank.org/nolanlab. We followed the CyTOF workflow by **Nowicka**
547 **et al. (2019)** and downloaded the data using `HDCytoData` (**Weber and Soneson (2019)**). The manual
548 cluster annotation of the CyTOF workflow can be downloaded from http://imlspenticton.uzh.ch/robinson_lab/cytofWorkflow/.
549

550 Acknowledgments

551 The authors thank Simona Ursu and Sarah Warth at the Core Facility Cytometry of the Ulm University
552 Medical Facility for their support acquiring both platelet datasets. We thank Marc Rosenbaum,
553 Dries van Hemelen, Gloria Martrus and Mayur Bakshi for excellent technical assistance, and Kilian
554 Kirmes for testing and evaluating our app.

References

- 555
556 **Arvaniti E**, Claassen M. Sensitive detection of rare disease-associated cell subsets via representation learn-
557 ing. *Nature Communications*. 2017 Apr; 8(1):14825. <https://www.nature.com/articles/ncomms14825>, doi:
558 <https://doi.org/10.1038/ncomms14825>.
- 559 **Belkina AC**, Ciccolella CO, Anno R, Halpert R, Spidlen J, Snyder-Cappione JE. Automated optimized parameters
560 for T-distributed stochastic neighbor embedding improve visualization and analysis of large datasets. *Nature*
561 *communications*. 2019; 10(1):1–12. doi: <https://doi.org/10.1038/s41467-019-13055-y>.
- 562 **Blair TA**, Michelson AD, Frelinger AL. Mass Cytometry Reveals Distinct Platelet Subtypes in Healthy Subjects
563 and Novel Alterations in Surface Glycoproteins in Glanzmann Thrombasthenia. *Scientific Reports*. 2018 Jul; 8.
564 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6037710/>, doi: 10.1038/s41598-018-28211-5.
- 565 **Blair TA**, Michelson AD, Frelinger AL. Mass cytometry reveals distinct platelet subtypes in healthy subjects and
566 novel alterations in surface glycoproteins in glanzmann thrombasthenia. *Scientific reports*. 2018; 8(1):1–13.
567 doi: <https://doi.org/10.1038/s41598-018-28211-5>.
- 568 **Bodenmiller B**, Zunder ER, Finck R, Chen TJ, Savig ES, Bruggner RV, Simonds EF, Bendall SC, Sachs K, Krutzik PO,
569 Nolan GP. Multiplexed mass cytometry profiling of cellular states perturbed by small-molecule regulators.
570 *Nature Biotechnology*. 2012 Sep; 30(9):858–867. doi: <https://doi.org/10.1038/nbt.2317>.
- 571 **Bongiovanni D**, Klug M, Lazareva O, Weidlich S, Biasi M, Ursu S, Warth S, Buske C, Lukas M, Spinner CD, Scheidt
572 Mv, Condorelli G, Baumbach J, Laugwitz KL, List M, Bernlochner I. SARS-CoV-2 infection is associated with a
573 pro-thrombotic platelet phenotype. *Cell Death & Disease*. 2021 Jan; 12(1):1–10. [https://www.nature.com/](https://www.nature.com/articles/s41419-020-03333-9)
574 [articles/s41419-020-03333-9](https://www.nature.com/articles/s41419-020-03333-9), doi: <https://doi.org/10.1038/s41419-020-03333-9>.
- 575 **Braune S**, Walter M, Schulze F, Lendlein A, Jung F. Changes in platelet morphology and function during 24 hours
576 of storage. *Clinical hemorheology and microcirculation*. 2014; 58(1):159–170. doi: [https://doi.org/10.3233/ch-](https://doi.org/10.3233/ch-141876)
577 [141876](https://doi.org/10.3233/ch-141876).
- 578 **Bruggner RV**, Bodenmiller B, Dill DL, Tibshirani RJ, Nolan GP. Automated identification of stratifying signatures
579 in cellular subpopulations. *Proceedings of the National Academy of Sciences*. 2014; 111(26):E2770–E2777.
580 <https://www.pnas.org/content/111/26/E2770>, doi: <https://doi.org/10.1073/pnas.1408792111>.
- 581 **Cohen J**. *Statistical power analysis for the behavioral sciences*. Academic press; 1977. doi:
582 <https://doi.org/10.1016/C2013-0-10517-X>.
- 583 **Crowell HL**, Zanotelli VRT, Chevrier S, Robinson MD. CATALYST: Cytometry dATa anALYSIS Tools; 2021, <https://github.com/HelenaLC/CATALYST>, doi: <https://doi.org/10.18129/B9.bioc.CATALYST>, r package version 1.14.1.
584
- 585 **Freedman D**, Diaconis P. On the histogram as a density estimator: L 2 theory. *Zeitschrift für Wahrscheinlichkeit-*
586 *s-theorie und verwandte Gebiete*. 1981; 57(4):453–476. doi: <https://doi.org/10.1007/BF01025868>.
- 587 **Gabbasov Z**, Ivanova O, Kogan-Yasny V, Ryzhkova E, Saburova O, Vorobyeva I, Vasilieva E. Activated platelet
588 chemiluminescence and presence of CD45+ platelets in patients with acute myocardial infarction. *Platelets*.
589 2014; 25(6):405–408. doi: <https://doi.org/10.3109/09537104.2013.829211>.
- 590 **Gadalla R**, Noamani B, MacLeod BL, Dickson RJ, Guo M, Xu W, Lukhele S, Elsaesser HJ, Razak ARA, Hirano N,
591 et al. Validation of CyTOF against flow cytometry for immunological studies and monitoring of human cancer
592 clinical trials. *Frontiers in oncology*. 2019; 9:415. doi: [10.3389/fonc.2019.00415](https://doi.org/10.3389/fonc.2019.00415).
- 593 **Hagberg IA**, Roald HE, Lyberg T. Platelet activation in flowing blood passing growing arterial thrombi. *Arterioscle-*
594 *rosis, thrombosis, and vascular biology*. 1997; 17(7):1331–1336. doi: <https://doi.org/10.1161/01.ATV.17.7.1331>.
- 595 **He L**, Davila-Velderrain J, Sumida TS, Hafler DA, Kellis M, Kulminski AM. NEBULA is a fast negative binomial mixed
596 model for differential or co-expression analysis of large-scale multi-subject single-cell data. *Communications*
597 *biology*. 2021; 4(1):1–17. doi: NEBULA is a fast negative binomial mixed model for differential or co-expression
598 analysis of large-scale multi-subject single-cell data.
- 599 **Hedges LV**, Olkin I. *Statistical methods for meta-analysis*. Academic press; 1985. doi:
600 <https://doi.org/10.1016/C2009-0-03396-0>.
- 601 **Kassambara A**. rstatix: Pipe-Friendly Framework for Basic Statistical Tests; 2021, [https://CRAN.R-project.org/](https://CRAN.R-project.org/package=rstatix)
602 [package=rstatix](https://CRAN.R-project.org/package=rstatix), r package version 0.7.0.

- 603 **Kauskot A**, Di Michele M, Loyen S, Freson K, Verhamme P, Hoylaerts MF. A novel mechanism of sustained
604 platelet α IIb β 3 activation via PEAR1. *Blood, The Journal of the American Society of Hematology*. 2012;
605 119(17):4056–4065. doi: <https://doi.org/10.1182/blood-2011-11-392787>.
- 606 **Kotecha N**, Krutzik PO, Irish JM. Web-based analysis and publication of flow cytometry experiments. *Current*
607 *Protocols in Cytometry*. 2010 Jul; Chapter 10:Unit10.17. doi: <https://doi.org/10.1002/0471142956.cy1017s53>.
- 608 **Lin M**, Lucas Jr HC, Shmueli G. Research commentary—too big to fail: large samples and the p-value problem.
609 *Information Systems Research*. 2013; 24(4):906–917. doi: <https://psycnet.apa.org/doi/10.1287/isre.2013.0480>.
- 610 **McKinnon KM**. Flow Cytometry: An Overview. *Current protocols in immunology*. 2018 Feb; 120:5.1.1–5.1.11.
611 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5939936/>, doi: <https://doi.org/10.1002/cpim.40>.
- 612 **Nowicka M**, Krieg C, Crowell HL, Weber LM, Hartmann FJ, Guglietta S, Becher B, Levesque MP, Robin-
613 son MD. CyTOF workflow: differential discovery in high-throughput high-dimensional cytometry
614 datasets. *F1000Research*. 2019 May; 6. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5473464/>, doi:
615 <https://doi.org/10.12688/f1000research.11622.3>.
- 616 **Ospina R**, Ferrari SLP. A general class of zero-or-one inflated beta regression models. *Computational Statistics &*
617 *Data Analysis*. 2012; 56(6):1609–1623. <https://www.sciencedirect.com/science/article/pii/S0167947311003628>,
618 doi: <https://doi.org/10.1016/j.csda.2011.10.005>.
- 619 **Papoutsoglou G**, Lagani V, Schmidt A, Tsirlis K, Cabrero DG, Tegnér J, Tsamardinos I. Challenges in the Multivari-
620 ate Analysis of Mass Cytometry Data: The Effect of Randomization. *Cytometry Part A*. 2019; 95(11):1178–1190.
621 <https://onlinelibrary.wiley.com/doi/abs/10.1002/cyto.a.23908>, doi: <https://doi.org/10.1002/cyto.a.23908>.
- 622 **Ramström S**, Öberg KV, Åkerström F, Enström C, Lindahl TL. Platelet PAR1 receptor density—correlation to
623 platelet activation response and changes in exposure after platelet activation. *Thrombosis research*. 2008;
624 121(5):681–688. doi: <https://doi.org/10.1016/j.thromres.2007.06.010>.
- 625 **Rigby RA**, Stasinopoulos DM. Generalized additive models for location, scale and shape,(with discussion).
626 *Applied Statistics*. 2005; 54:507–554. doi: <https://doi.org/10.1111/j.1467-9876.2005.00510.x>.
- 627 **Rigby RA**, Stasinopoulos MD, Heller GZ, De Bastiani F. Distribution for modelling location, scale, and shape:
628 using GAMLSS in R. Chapman & Hall/CRC : the R series, Boca Raton, Florida: CRC Press; 2020. doi:
629 <https://doi.org/10.1201/9780429298547>.
- 630 **Ritchie ME**, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK. limma powers differential expression
631 analyses for RNA-sequencing and microarray studies. *Nucleic acids research*. 2015; 43(7):e47–e47. doi:
632 <https://doi.org/10.1093/nar/gkv007>.
- 633 **Rubner Y**, Tomasi C, Guibas LJ. A metric for distributions with applications to image databases. In:
634 *Sixth International Conference on Computer Vision (IEEE Cat. No. 98CH36271)* IEEE; 1998. p. 59–66. doi:
635 <https://doi.org/10.1109/ICCV.1998.710701>.
- 636 **Seiler C**, Ferreira AM, Kronstad LM, Simpson LJ, Le Gars M, Vendrame E, Blish CA, Holmes S. CytoGLMM:
637 conditional differential analysis for flow and mass cytometry experiments. *BMC bioinformatics*. 2021; 22(1):1–
638 14. doi: <https://doi.org/10.1186/s12859-021-04067-x>.
- 639 **Stasinopoulos M**, Rigby R. *gamlss.dist: Distributions for Generalized Additive Models for Location Scale and*
640 *Shape*; 2021, <https://CRAN.R-project.org/package=gamlss.dist>, r package version 5.3-2.
- 641 **Testi R**, Pulcinelli F, Frati L, Gazzaniga PP, Santoni A. CD69 is expressed on platelets and mediates
642 platelet activation and aggregation. *The Journal of experimental medicine*. 1990; 172(3):701–707. doi:
643 <https://doi.org/10.1084/jem.172.3.701>.
- 644 **Testi R**, Pulcinelli FM, Cifone MG, Botti D, Del Grosso E, Rioldino S, Frati L, Gazzaniga PP, Santoni A. Preferential
645 involvement of a phospholipase A2-dependent pathway in CD69-mediated platelet activation. *The Journal of*
646 *Immunology*. 1992; 148(9):2867–2871. <https://www.jimmunol.org/content/148/9/2867.long>.
- 647 **de Torrenté L**, Zimmerman S, Suzuki M, Christopheit M, Greally JM, Mar JC. The shape of gene expression distri-
648 butions matter: how incorporating distribution shape improves the interpretation of cancer transcriptomic
649 data. *BMC bioinformatics*. 2020; 21(21):1–18. doi: <https://doi.org/10.1186/s12859-020-03892-w>.
- 650 **Wang T**, Nabavi S. SigEMD: A powerful method for differential gene expression analysis in single-cell RNA
651 sequencing data. *Methods*. 2018; 145:25–32. doi: <https://doi.org/10.1016/j.ymeth.2018.04.017>.

- 652 **Weber LM**, Nowicka M, Sonesson C, Robinson MD. diffcyt: Differential discovery in high-dimensional cytometry
653 via high-resolution clustering. *Communications biology*. 2019; 2(1):1–11. doi: [https://doi.org/10.1038/s42003-](https://doi.org/10.1038/s42003-019-0415-5)
654 [019-0415-5](https://doi.org/10.1038/s42003-019-0415-5).
- 655 **Weber LM**, Sonesson C. HDCytoData: collection of high-dimensional cytometry benchmark datasets in Biocon-
656 ductor object formats. *F1000Research*. 2019; 8. doi: <https://doi.org/10.18129/B9.bioc.HDCytoData>.