

1 **Determination of essential phenotypic elements of clusters in high-**  
2 **dimensional entities - DEPECHE**

3

4 Short title: DEPECHE-a data-mining algorithm for mega-variate data

5

6 Authors

7 Axel Theorell<sup>1</sup>, Yenan Troi Bryceson<sup>2,3</sup>, Jakob Theorell<sup>2\*</sup>

8

9 Affiliations

10 <sup>1</sup> IBG-1: Biotechnology, Institute of Bio- and Geosciences, Forschungszentrum Jülich  
11 GmbH, Jülich, North Rhine-Westphalia, Germany

12 <sup>2</sup> Center for Hematology and Regenerative Medicine, Department of Medicine  
13 Huddinge, Karolinska Institutet, Stockholm, Sweden

14 <sup>3</sup> Broegelmann Research Laboratory, Department of Clinical Medicine, University of  
15 Bergen, Bergen, Norway

16

17 \* Corresponding author

18 Email: [Jakob.Theorell@ki.se](mailto:Jakob.Theorell@ki.se) (JT)

19

20 Author contributions

21 A.T. drafted mathematical models, coded the software implementation and co-wrote  
22 the manuscript, Y.T.B. co-wrote the manuscript, J.T. drafted mathematical models,  
23 coded software implementation and wrote the manuscript.

24

25

26

27  
28

## 29 **Abstract**

30 Technological advances have facilitated an exponential increase in the amount of  
31 information that can be derived from single cells, necessitating new computational  
32 tools that can make this highly complex data interpretable. Here, we introduce  
33 DEPECHE, a rapid, parameter free, sparse k-means-based algorithm for clustering of  
34 multi- and megavariable single-cell data. In a number of computational benchmarks  
35 aimed at evaluating the capacity to form biologically relevant clusters, including  
36 flow/mass-cytometry and single cell RNA sequencing data sets with manually curated  
37 gold standard solutions, DEPECHE clusters as well or better as the best performing  
38 state-of-the-art clustering algorithms. However, the main advantage of DEPECHE,  
39 compared to the state-of-the-art, is its unique ability to enhance interpretability of the  
40 formed clusters, in that it only retains variables relevant for cluster separation, thereby  
41 facilitating computational efficient analyses as well as understanding of complex  
42 datasets. An open source R implementation of DEPECHE is available at  
43 <https://github.com/theorell/DepecheR>.

44  
45

## 46 **Author summary**

47 DEPECHE-a data-mining algorithm for mega-variate data  
48 Modern experimental technologies facilitate an array of single cells measurements,  
49 *e.g.* at the RNA-level, generating enormous datasets with thousands of annotated  
50 biological markers for each of thousands of cells. To analyze such datasets,  
51 researchers routinely apply automated or semi-automated techniques to order the cells

52 into biologically relevant groups. However, even after such groups have been  
53 generated, it is often difficult to interpret the biological meaning of these groups since  
54 the definition of each group often depends on thousands of biological markers.  
55 Therefore, in this article, we introduce DEPECHE, an algorithm designed to  
56 simultaneously group cells and enhance interpretability of the formed groups.  
57 DEPECHE defines groups only with respect to biological markers that contribute  
58 significantly to differentiate the cells in the group from the rest of the cells, yielding  
59 more succinct group definitions. Using the open source R software DepecheR on  
60 RNA sequencing data and mass cytometry data, the number of defining markers were  
61 reduced up to 1000-fold, thereby increasing interpretability vastly, while maintaining  
62 or improving the biological relevance of the groups formed compared to state-of-the-  
63 art algorithms.  
64

## 65 **Introduction**

66 Since the introduction of the first single colour flow cytometers in the 1960s, there  
67 has been a remarkable increase in the complexity of data that can be generated with  
68 single-cell resolution. Currently, flow and mass cytometers able to simultaneously  
69 assess up to 40 cellular traits are becoming widely available [1]. In parallel, the  
70 development of high-throughput sequencing technology has facilitated deep single-  
71 cell transcriptomic analyses [2]. Furthermore, development of high-resolution single-  
72 cell proteomic analyses are underway [3].

73 These technological advances necessitate new computational approaches to  
74 analyses of multi- and megavariable single cell data [4–7]. Previous algorithms have  
75 contributed to automating analyses, thereby enhancing reproducibility and avoiding a  
76 need for *a priori* biological knowledge for design of manual gating analysis strategies.  
77 Automated analysis algorithms, not restricted to uni- or bivariate displays of the data,  
78 have also made it possible to display much more of the information embedded in  
79 multivariate data. To date, however, manual gating strategies are still dominantly  
80 used, which in part is likely due to that it is easy to interpret what population a certain  
81 gate refers to, as it is defined by few markers. In an attempt to combine the  
82 objectiveness and reproducibility of automated analysis pipelines with the high  
83 interpretability of manual gating strategies, we have developed an algorithm termed  
84 Determination of Essential Phenotypic Elements of Clusters in High-dimensional  
85 Entities (DEPECHE). DEPECHE simultaneously clusters and simplifies the data by  
86 identifying the variables that contribute to separate individual clusters from the rest of  
87 the data. We have implemented DEPECHE in R (in the open source package  
88 DepecheR), providing a complete software suite for statistical analysis and  
89 visualization of single cell omics data.

## 90 **Results and Discussion**

91 DEPECHE uses a penalized k-means clustering algorithm, related to the standard k-  
92 means algorithm [8]. In penalized k-means, a penalty term is introduced to the  
93 clustering algorithm. The value of the penalty,  $\lambda$ , determines the clustering resolution.  
94 Low clustering resolution implies that few clusters defined by few variables (high  
95 sparsity) are produced, and vice versa [9] (see online methods). Note that if  $k$  is high  
96 enough not to be limiting, the resolution of the emerging clusters depends entirely on  
97 the magnitude of  $\lambda$  and not on  $k$ , since DEPECHE annihilates all clusters that are  
98 pulled to the origin by the penalty. In DEPECHE, the penalty  $\lambda$  is tuned to identify  
99 the most *reproducible* clustering resolution, here termed the “optimal resolution”[10].  
100 To illustrate what we mean by reproducible, we constructed a show case, featuring a  
101 bi-variate dataset  $D$  (Fig 1a). Visually, the dataset  $D$  contains three clusters, where the  
102 centers of the two larger clusters are located close to either axis. For these two  
103 clusters, one variable is sufficient to define their position. Now, if multiple datasets  
104 were generated from the same data source as  $D$ , for example by repeated experiments,  
105 we assume that they would contain the same clusters. Hence, imposing the optimal  
106 penalty  $\lambda_i$  (that corresponds to the optimal resolution of  $D$ ) on all these datasets  
107 should ideally always result in the same clusters (high reproducibility). Contrarily,  
108 when clustering the same datasets with a penalty  $\lambda$  that differ significantly from the  
109 optimal penalty, the stochastic differences between the datasets are likely to induce  
110 solutions that deviate in cluster number, number of defining variables, and cluster  
111 center positions. In practice, DEPECHE tests a range of penalty values ( $\lambda_1 < \dots <$   
112  $\lambda_n$ ), each on a collection of dataset pairs which are generated by sampling  $N_R$  data  
113 points from  $D$  (Fig 1b) with resampling. The optimal resolution is defined as the  
114 penalty  $\lambda_i$ , which yields the lowest average variability within each dataset pair, as

115 measured by the Adjusted Rand Index (ARI) [11] . In our example, this corresponds  
116 to the penalty  $\lambda_i$  that yields 3 clusters, since 3 similar clusters are identified in each  
117 resampled dataset of  $D$  (Fig 1c). The penalties  $\lambda_1$  and  $\lambda_n$  are considered suboptimal,  
118 since with these penalties, the stochastic differences in the resampled datasets lead to  
119 less coherent clustering results compared to results obtained with the optimal penalty  
120  $\lambda_i$ . From here (Fig 1c) DEPECHE uses two alternative routes. If the number of data  
121 points in the dataset  $D$  is high (as default  $> 10^4$ ), the most generalizable cluster  
122 centers that were produced using the optimal penalty are chosen (see online methods)  
123 and the data points of  $D$  are allocated directly to their closest cluster center (Fig 1d). If  
124 the dataset  $D$  has few data points, the full dataset  $D$  is clustered using the optimal  
125 penalty  $\lambda_i$  (Fig 1e).

126

127 **Fig 1: Illustration of the DEPECHE workflow:** a) The original dataset  $D$ . b)  $n$   
128 resampled datasets with  $N_R$  data points per dataset, are generated by sampling data  
129 points from  $D$  with resampling. Each resampled dataset has a corresponding penalty  
130  $\lambda_i$  ( $i = 1, \dots, n$ ). c) Each dataset in b is clustered with sparse k-means, using its  
131 corresponding penalty  $\lambda_i$ . The red frame highlights the clustering with the strongest  
132 attractors, *i.e.* the most generalizable solution (see online methods). d, e) Finally, the  
133 full dataset is clustered by allocating each data point to its closest cluster center, using  
134 the most generalizable cluster center solution produced in b.

135

136 To evaluate how biologically accurate DEPECHE clustering is on mass  
137 cytometry data, a 32-variate mass cytometry bone marrow dataset [12] was clustered,  
138 and the overlap to 14 manually pre-defined cell populations was quantified using the  
139 ARI. With this dataset, DEPECHE identified 7 clusters at the optimal resolution,

140 corresponding to all large pre-defined cell populations and to agglomerates of smaller  
141 cell populations, rendering an average ARI of 0.96, where an ARI of 1 corresponds to  
142 exact reproduction and an ARI of 0 means that the produced clusters are no more  
143 accurate than random allocation. (Fig 2a-b, S Fig 1a). Furthermore, using DEPECHE,  
144 the number of variables defining each cluster was reduced from 32 to a range from 8  
145 to 28, thereby enhancing interpretability (Fig 2d). When comparing to other state-of-  
146 the-art clustering algorithms [12–17], DEPECHE obtained similar ARI as the best  
147 algorithms for both the 32-variate dataset and another 14-variate, 24 population, mass  
148 cytometry dataset [18] (S Fig 2a, Table 1).

Table 1: background information on all datasets

Dataset	Data origin	n cells	n variables in analysis	n clusters in original
Levine	Mass cytometry	104184	32	14
Bendall	Mass cytometry	81747	14	24
Björklund	scRNAseq	648	35177	4
Biase	scRNAseq	56	19571	3
Deng	scRNAseq	268	13867	10
Goolam	scRNAseq	124	15487	5
Kolodziejczyk	scRNAseq	704	15117	3
Pollen	scRNAseq	301	13860	11
Yan	scRNAseq	90	13608	7

149

150

151 Single-cell transcriptomic datasets feature tens of thousands of variables.  
152 Thus, the need to exclude irrelevant variables is even more pressing, as compared to  
153 cytometry datasets. We therefore evaluated DEPECHE’s ability to cluster and extract  
154 the key transcripts defining clusters of a previously published single-cell  
155 transcriptomic dataset (Fig 2d-f) [19]. In this dataset, a total of 648 ILC1, ILC2, ILC3  
156 and NK cells from three donors’ tonsils were index-sorted prior to RNA sequencing.  
157 Hence, these cell types, manually defined by protein expression, can be compared to  
158 clusters unbiasedly determined by RNA expression profiles [19]. In the DEPECHE  
159 analysis, no pre-selection of transcripts was performed, and hence, 35177 unique

160 transcripts were included for each of the 648 cells. With the optimal penalty  $\lambda$ , four  
161 clusters were identified (Fig 2e). These corresponded well to the cell types as defined  
162 by protein expression; 84, 97, 91 and 97 percent of ILC1, ILC2, ILC3 and NK cells  
163 sorted into separate clusters, respectively (Fig 2d, e, S Fig 1b), leading to an average  
164 ARI of 0.78. Notably, cluster 1-4, corresponding to ILC1, ILC3, NK cells and ILC2,  
165 were defined by 27, 27, 108 and 10 transcripts, respectively (Fig 2f and Table 2),  
166 leading to a 99.9% average decrease in the number of variables. The transcripts  
167 identified to define the clusters in our analysis were among those most differentially  
168 expressed according to the original study [19] (Fig 2f). Thus, by identifying a finite  
169 number of variables, DEPECHE analysis increases interpretability and aides down-  
170 stream analyses. When DEPECHE clustering was compared to that of state-of-the-art  
171 algorithms [5–7] on the aforementioned dataset and six others (see Table 1), it  
172 performed consistently well as indicated by ARI (S Fig 2b). Thus, when applied to  
173 megavariate data, DEPECHE produces biologically relevant clusters and reduces the  
174 complexity of the result thousand-fold.

175

176 **Fig 2: DEPECHE performance with real datasets with 32 or 35177 variables.** a-  
177 b) bi-variate t-distributed stochastic neighbor embedding (tSNE) representation of the  
178 32-variate mass cytometry data. A: distribution of manually defined cell populations  
179 over the tSNE field. B: distribution of DEPECHE clusters over the tSNE field. c)  
180 Heatmap showing which variables that define each cluster. Red color indicates a  
181 higher expression in the cluster than the most common expression for all  
182 observations. Blue color conversely indicates lower expression than the geometric  
183 mean for all observations. Grey color indicates that the variable in question does not  
184 contribute to defining the cluster. For Fig a-c all, 104184 cells have been clustered. d-



185 e) tSNE representation of the 137-variate data subset that could efficiently distinguish  
186 the clusters in the 35177-variate single-cell transcriptome dataset. d) distribution of  
187 the cell types defined by index-sorting and manual gating on protein expression  
188 profiles shown over the tSNE field. e) distribution of DEPECHE clusters over the  
189 tSNE field. f) Violin plots illustrating the overlap between the original analysis by  
190 Björklund *et al* and the DEPECHE analysis. For each subplot, the left and right side  
191 illustrate the distribution of the transcripts defining the clusters, and all other  
192 transcripts, respectively. The y-axis shows the log<sub>10</sub> of the p-values in the original  
193 analysis adjusted for multiple comparisons. For Fig d-f, all 648 cells have been  
194 clustered.

Table 2: transcripts defining clusters in Björklund *et al* dataset

Cluster 1/NK cells			Cluster 4 (ILC2)			Cluster 3/ILC3			Cluster 3/ILC3		
Transcript	Log10 of adjusted original p-value	Cluster center	Transcript	Log10 of adjusted original p-value	Cluster center	Transcript	Log10 of adjusted original p-value	Cluster center	Transcript	Log10 of adjusted original p-value	Cluster center
CMC1	-8.95	0.23	A2M	-7.08	0.23	CD3G	-9.66	-0.06	PCDH9	-9.66	1.25
CST7	-8.95	0.17	AC092580.4	-9.17	-0.89	CD63	-6.03	0.21	PDCD4	-3.37	-0.02
GNLY	-8.95	1.82	CD2	-9.17	-2.69	CNN2	-9.66	-0.61	PECAM1	-9.66	0.06
GZMA	-8.95	0.23	CD300LF	-9.17	-0.06	COTL1	-8.04	-0.25	PRR5	-9.66	0.38
GZMK	-8.95	0.09	CD3E	-4.21	-0.40	CPNE7	-6.13	0.25	PTPN22	-4.79	0.05
KLRC1	-8.95	0.57	EMP3	-2.86	0.31	CTSA	-5.65	0.25	RBPJ	-5.14	0.07
KLRD1	-8.95	1.89	FCER1G	-7.08	-0.59	DCAF11	-6.48	0.04	RHOC	-9.66	0.76
KLRF1	-8.95	1.29	GATA3	-8.33	0.34	DHRS3	-3.72	0.02	RP11-264B17.3..1	-1.41	-0.18
NKG7	-8.95	1.93	GSN	-9.17	-0.40	DOCK5	-9.66	0.30	RP11-330A16.1	-9.47	0.01
PRF1	-8.95	0.14	HPGDS	-9.17	0.37	ELOVL6	-9.66	0.31	RP11-466H18.1	-0.35	-0.02
Cluster 2/ILC1			IL10RA	-8.84	0.03	EMP3	-7.26	-0.94	RP11-845M18.6	-9.66	0.44
Log10 of adjusted original p-value			IL17RB	-9.17	0.14	ENPP1	-9.66	0.17	RPS8	-1.54	-0.06
Cluster center			IL23R	-9.17	-0.47	FAIM3	-3.79	-0.23	S100A6	-6.94	-0.16
Transcript			IL2RB	-9.17	-0.64	FAM65B	-3.38	-0.47	S1PR1	-9.66	-0.08
Log10 of adjusted original p-value			IL32	-9.10	1.04	FCER1G	-9.66	1.02	SELL	-8.49	-0.97
Cluster center			KLRC1	-9.17	-0.51	FES	-9.66	0.11	SELPLG	-7.57	-0.37
AE000661.37	-8.67	-0.08	KLRG1	-9.17	0.50	GIMAP4	-7.49	-0.38	SERINC5	-9.66	-0.13
CCR7	-9.25	0.69	KRT1	-9.17	0.21	GIMAP7	-5.46	-0.39	SH2D1B	-9.66	0.90
CD27	-9.25	0.97	SH2D1B	-9.17	-0.64	GSN	-9.66	1.24	SLA	-5.90	0.16
CD3D	-9.25	3.50	TESPA1	-6.02	0.01	HCST	-0.34	-0.11	SLC38A1	-1.48	-0.06
CD3E	-4.13	0.48	TMIGD2	-9.03	-0.11	HDAC9	-9.66	0.30	SLC4A10	-9.66	0.32
CD3G	-9.25	2.56	TRAC	-9.17	-0.54	HLA-DRB1	-3.19	0.34	SORL1	-8.48	-0.23
CD4	-9.25	0.55	TXK	-8.53	-1.21	IL10RA	-9.66	-0.97	SPINK2	-9.66	0.67
CD6	-9.25	0.76	TYROBP	-9.17	-1.40	IL1R1	-9.66	0.42	SPRY1	-8.67	0.11
CNN2	-7.04	0.01	VWA5A	-9.17	-0.32	IL23R	-9.66	1.67	STARD3NL	-9.66	0.23
COTL1	-9.25	0.63	XCL1	-9.17	-0.62	IL32	-9.66	-0.85	TC2N	-6.72	-0.66
CTSW	-7.82	-0.02	XCL2	-9.17	-0.12	IL411	-9.66	1.50	TCIRG1	-7.11	0.10
FCER1G	-9.25	-0.81	Cluster 3/ILC3			ISG20	-2.10	-0.01	TLE3	-9.66	0.06
KLRB1	-9.25	-0.96	Transcript	Log10 of adjusted original p-value	Cluster center	ITGB2	-6.95	-0.76	TMIGD2	-9.66	1.09
LITAF	-9.25	0.08	Transcript	Log10 of adjusted original p-value	Cluster center	ITM2C	-9.66	0.41	TNFRSF18	-9.66	0.79
LST1	-9.25	-0.57	Transcript	Log10 of adjusted original p-value	Cluster center	KIAA1324	-9.66	1.07	TNFRSF25	-3.37	0.30
RNU2-6P	-7.67	-0.07	Transcript	Log10 of adjusted original p-value	Cluster center	KIT	-9.66	1.30	TNFRSF4	-7.13	0.04
RP11-466H18.1	0.00	0.20	A2M	-9.52	-0.04	KLRG1	-9.66	-0.26	TNFSF11	-9.66	0.31
SH2D1B	-9.25	-0.29	AC092580.4	-9.66	0.78	KRT81	-9.66	1.00	TNFSF13B	-9.66	1.02
SIT1	-9.25	1.06	ADAM10	-6.38	0.15	LAT2	-6.18	0.16	TOX	-9.66	0.35
TC2N	-9.08	1.02	ADAM28	-9.66	0.05	LDHB	-9.35	-0.18	TOX2	-9.66	0.69
TNFRSF18	-9.25	-0.21	AFF3	-9.66	0.80	LINC00299	-9.66	1.43	TRAC	-0.80	0.46
TOB1	-2.23	0.23	AHR	-8.27	0.24	LITAF	-9.66	-0.20	TRAJ45	-8.56	0.70
TRDC	-9.25	-1.71	AMICA1	-9.66	1.33	LST1	-9.66	1.68	TRAT1	-7.81	-0.12
TRDJ2	-9.25	-0.55	ARL4C	-4.37	-0.02	LTA4H	-9.66	0.89	TRDJ2	-6.30	0.07
TYROBP	-9.25	-0.71	ATP8B4	-9.66	0.10	LY6E	-3.30	-0.27	TRGJP1	-7.31	0.11
U2..37	-9.03	-0.30	BST2	-8.28	0.44	MPG	-8.65	0.49	TXK	-5.59	0.14
U2..55	-9.16	-0.39	C1orf162	-8.83	-0.69	NCR2	-9.66	1.00	TYROBP	-9.66	1.15
			CAT	-9.66	0.69	NKG7	-1.82	-0.04	VWA5A	-9.66	1.61
			CD2	-7.77	0.44	NRP1	-9.66	0.35	XCL1	-9.66	0.80
			CD300LF	-9.66	1.30	NSMCE1	-9.66	0.64	XCL2	-9.66	0.03
			CD3D	-9.65	-0.36	OTUD5	-9.66	0.73			

195

196

197

198

199

200

201

In conclusion, DEPECHE turns the penalized k-means methodology into a parameter free analysis technique guided by efficient calculation of the optimal clustering resolution. By doing so, it addresses the simultaneous problem of clustering and identification of biologically important variables that separate clusters. This is crucial in order to comprehend the noisy and often over-complicated data generated with current single cell technologies.

## 202 **Methods**

### 203 **Clustering with DEPECHE**

204 Clustering in DEPECHE is performed using a penalized version of the k-means  
205 algorithm, which is related to the k-means algorithm[8]. In this section, the k-means  
206 algorithm is outlined first, followed by an explanation of how it is extended to  
207 penalized k-means.

208 The k-means algorithm clusters data by fitting a mixture of normal  
209 distributions to the data with k equal mixture components and unit variance.  
210 Formally,  $k$   $d$ -dimensional cluster centers, denoted  $\mu_{i,j}$  where  $i = 1 \dots k$  and  $j =$   
211  $1 \dots d$ , are fitted to the  $n d$ -dimensional datapoints  $x_{l,j}$ , where  $l = 1 \dots n$ , by  
212 maximizing the score function

213

$$Q = \sum_{i=1}^k \sum_{l=1}^n z_{i,l} \sum_{j=1}^d (x_{l,j} - \mu_{i,j})^2, \quad (1)$$

214

215 where  $z_{i,l}$  is 1 if the  $l$ th data point belongs to the  $i$ th cluster and zero otherwise. The  
216 score  $Q$  is optimized using an Expectation Maximization (EM) algorithm [20], *i.e.* so  
217 called E- and M-steps are iterated alternately until the score  $Q$  stops improving. In  
218 the E-step, the allocation variables  $z_{i,l}$  are updated so that each data point is allocated  
219 to its closest cluster. In the M-step, each cluster center  $\mu_{i,j}$  is moved to the center of  
220 the data points allocated to it. When no more reallocation occurs in the E-step, the  
221 algorithm has converged.

222 In order to reduce the influence of uninformative dimensions that only  
223 contribute with noise, penalized k-means introduces an L1-penalty for each element

224 of each cluster center  $\mu_{i,j}$  to the optimization objective. Formally, the score function

225  $Q$  in Eq. ( 1 ), is updated:

226

$$Q = \sum_{i=1}^k \sum_{l=1}^n z_{i,l} \sum_{j=1}^d (x_{i,j} - \mu_{i,j})^2 - \lambda \sum_{i=1}^k \sum_{j=1}^d |\mu_{i,j}|, \quad (2)$$

227 where  $\lambda$  is a positive penalty term. The additional term in the score function,

228 introduced in Eq. ( 2 ) results in a change in the M-step of the original EM-algorithm

229 of the k-means algorithm. Keeping  $z_{i,l}$  for all  $l$  fixed and optimizing  $Q$  with respect to

230  $\mu_{i,j}$ , the M-step is:

$$\mu_{i,j} = \text{sign} \left( \frac{\sum_{l=1}^n z_{i,l} x_{i,j}}{\sum_{l=1}^n z_{i,l}} \right) \cdot \max \left( \left| \frac{\sum_{l=1}^n z_{i,l} x_{i,j}}{\sum_{l=1}^n z_{i,l}} \right| - \frac{\lambda}{2 \sum_{l=1}^n z_{i,l}}, 0 \right). \quad (3)$$

231 Depending on the choice of the penalty parameter  $\lambda$ , some components of some

232 clusters centers will be set to 0 in the M-step. Note that penalized k-means with

233 penalty  $\lambda = 0$  reduces to the original k-means algorithm.

234 In DEPECHE, cluster centers that are moved to the origin in the M-step are

235 eliminated and not assigned any data points in the E-step. Due to the elimination of

236 clusters, the number of produced clusters is independent of  $k$  and dependent on the

237 penalty  $\lambda$  as long as at least one cluster is eliminated. In DEPECHE,  $k$  is always

238 chosen to be so large that at least one cluster is eliminated.

239 Eq. ( 2 ) is a special case of the penalized model based clustering algorithm by

240 Pan and Shen with unit variance and equal mixture components [9]. By imposing the

241 penalty for each dimension and each cluster, penalized k-means identifies the

242 dimensions that do not distinguish a particular cluster from the rest of the data, thus

243 leaving these dimensions out of the definition of that cluster. This differs from the

244 sparse k-means algorithm by Witten and Tibshirani [21] and the regularized k-means

245 algorithm by Sun et al [10], that only identify dimensions that do not contribute to  
246 distinguish any cluster from the rest of the data.

247 Penalized k-means, as well as k-means, relies on a procedure for initializing  
248 the positions of the cluster centers. Cluster initialization is particularly delicate in  
249 DEPECHE, due to the elimination of clusters at the origin in the E-step. Poor  
250 initialization of the clusters might lead to elimination of too many clusters in the early  
251 E-steps, yielding fewer clusters in the end result than necessary to optimize  $Q$ . To  
252 avoid early elimination of clusters, DEPECHE initializes the cluster positions using  
253 the seed generation algorithm of k-means++ by Arthur and Vassilvitskii[22] and  
254 always starts clustering with penalty  $\lambda = 0$ . The penalty is then increased linearly  
255 over a number of E-steps until it reaches the predetermined value.

256 The EM-algorithm guarantees convergence to an optimum of the score  $Q$ , but  
257 not necessarily to the global optimum. In order to diminish the influence of the  
258 starting state, the EM-algorithm is run several times with random initialization, and  
259 the solution with optimal  $Q$  is chosen. In addition,  $k$  is set considerably higher than  
260 the expected number of final clusters, which also diminishes the sensitivity to the  
261 starting state. In the extreme case where  $k$  is set equal to the number of data points  $n$ ,  
262 the outcome of penalized k-means is deterministic.

263

## 264 **Tuning the penalty**

265 In this section, we describe the optimization scheme which is used for tuning the  
266 linear penalty  $\lambda$ . The outline of the algorithm:

267

268 1. Choose a range of penalty terms  $\lambda_i, i = 1..N_\lambda$  that are considered for  
269 clustering the dataset  $D$

- 270 2. Create 2 datasets per penalty term  $\lambda_i$ , called  $D_{1,i}$  and  $D_{2,i}$ , by sampling  $N_r$  data  
271 points from  $D$  with replacement.
- 272 3. Run the penalized k-means algorithm on the datasets  $D_{1,i}$  and  $D_{2,i}$ , yielding  
273 sets of cluster centers, denoted  $M_{1,i}$  and  $M_{2,i}$ .
- 274 4. Create the partitions  $P_{1,i}$  and  $P_{2,i}$ , by allocating all data points of the dataset  $D$   
275 to their nearest cluster center of the sets  $M_{1,i}$  and  $M_{2,i}$ .
- 276 5. Determine the Adjusted Rand Index (ARI), denoted  $r(\lambda_i)$  from  $P_{1,i}$  to  $P_{2,i}$   
277 [11].
- 278 6. Repeat step 2-5 times and average the obtained ARIs  $r(\lambda_i)$  penalty wise until  
279 a stopping criteria regarding the statistical certainty of the obtained ARIs  $r(\lambda_i)$   
280 is met.
- 281 7. Choose the optimal penalty  $\lambda_i$ , which is the penalty with the largest  
282 ARI  $r(\lambda_i)$ .

283

284 Some remarks to the parameter tuning procedure: The repetition Step 6 is necessary,  
285 since the obtained ARI  $r(\lambda_i)$  is a random variable, due to the random procedure for  
286 creating the datasets  $D_{1,i}$  and  $D_{2,i}$  and the random procedure for initializing the  
287 penalized k-means algorithm. DEPECHE uses two stopping criteria: The first  
288 criterion creates an interval of width 2 standard errors around the obtained mean of  
289  $r(\lambda_i)$  and checks if the interval around the optimal ARI  $r(\lambda_i)$  has a zero overlap with  
290 the other intervals. The second criterion checks whether the standard error of the  
291 mean of  $r(\lambda_i)$  for the optimal penalty  $\lambda_i$  is below a threshold.

292 Step 2 requires a samples size  $N_r$ . A natural choice is to set  $N_r$  equal to the number of  
293 data points,  $n$ . However, in cases where  $n$  is very large, so that the computational load  
294 of the optimization scheme becomes limiting, it is preferable to choose a smaller  $N_r$ .

295 In DepecheR,  $N_r = 10^4$  by default, in case  $n \geq 10^4$ . Notice that when an optimal  
296 penalty  $\lambda_i$  is discovered using sample size  $N_r \neq n$ , the corresponding optimal penalty  
297 when sampling the full dataset  $D$  with magnitude  $n$  is (approximately)  $\lambda_i \cdot \frac{n}{N_r}$ , since  
298 the attraction force of a cluster is proportional to the number of data points in it.  
299 Exact calculation of the ARI in step 5 is computationally intractable for large datasets.  
300 Therefore, DEPECHE relies on an approximate ARI computation, based on  $10^4$   
301 random pairs of data points.

302

### 303 **Simultaneous Clustering and Parameter Tuning**

304 For very large datasets ( $n > 10^8$ ), not only the penalty optimization, but also the  
305 final clustering once the optimal penalty has been found may be computationally  
306 intractable. However, increasing the size of the dataset, does not necessarily lead to an  
307 increase in number of clusters at the optimal resolution. In this case, it is feasible to  
308 cluster a subset of the full dataset  $D$  to obtain cluster centers  $M$  and then allocate the  
309 remaining data points of  $D$  to their closest clusters in  $M$ . This boosts computational  
310 efficiency since allocation imposes a much smaller computational load than  
311 clustering. Since several subsets of  $D$  are produced and clustered during the tuning of  
312 the penalty parameter  $\lambda$ , it seems natural to retrieve cluster centers  $M$  that were  
313 produced during the parameter tuning and use them to cluster  $D$ .

314 When picking a set of cluster centers  $M$  from the penalty tuning, the question arises  
315 which set of centers  $M$  to take, since several sets of centers, denoted  $M_{i,l}$  ( $l = 1, \dots, p$ ),  
316 are produced for the optimal penalty  $\lambda_i$ . In DEPECHE, the centers  $M_{i,j}$  that have the  
317 strongest similarity (on average) to the remaining  $p - 1$  centers is chosen and is  
318 referred to as the most generalizable cluster set. The level of similarity between the

319 centers  $M_{i,j}$  and  $M_{i,l}$  is quantified using the ARI between the partitions  $P_{i,j}$  and  $P_{i,l}$ ,  
320 induced by allocating each data point of  $D$  to its closest cluster center in  $M_{i,j}$  and  $M_{i,l}$   
321 respectively.

322

323 **Empiric Performance of the Penalty Tuning Scheme.** Roughly speaking,  
324 DEPECHE combines a flavored penalized k-means algorithm with a parameter tuning  
325 scheme, which identifies an optimal resolution. A naturally arising question is then  
326 whether the parameter tuning scheme is able to determine a biologically relevant  
327 resolution or if other penalized k-means clustering resolutions outperform the  
328 resolution chosen by DEPECHE. Using a range of datasets (Table 3), the biological  
329 relevance (measured in ARI to the manually curated solution) of the optimized  
330 DEPECHE partitions were compared to the biologically optimal partition among all  
331 partitions generated with 20 repetitions on each of a range of 11 penalties per dataset.  
332 Overall, the DEPECHE resolution-selection showed close to optimal performance, as  
333 the selected solutions only had a median of 0.02 lower ARI to the gold standard  
334 (range 0-0.065) than the best possible solution with all penalties (Table 3).

335

Table 3: ARI between DEPECHE partitions and golden standard partitions

Dataset	Median ARI in supplementary figure 2	Maximal ARI with any penalty	Difference
Levine	0.961	0.975	0.015
Bendall	0.841	0.873	0.032
Biase	1	1	0
Björklund	0.782	0.842	0.06
Deng	0.827	0.848	0.021
Goolam	0.629	0.639	0.009
Kolodziejczyk	0.992	1	0.008
Pollen	0.863	0.928	0.065
Yan	0.626	0.691	0.064

336

337

338



### 339 **Scaling and Centering the Data**

340 The clusters produced by DEPECHE, as well as their interpretation, depends on the  
341 scaling and centering of the data. The scaling determines the relative importance of  
342 the measured variables, where variables with a larger spread have stronger influence  
343 on the clustering. The centering defines where zero occurs in each variable, thereby  
344 influencing the clustering results due to the linear penalty.

345 DEPECHE is applicable to a large range of datasets where the numbers of  
346 dimensions,  $d$ , and the number of data points,  $n$ , can vary with many orders of  
347 magnitude. The differing characteristics of these datasets require different treatments  
348 with respect to scaling and centering.

349

350 **Scaling.** Empirically, a majority of single-cell transcriptome datasets tend to have a  
351 few variables where the variance is many orders of magnitude greater than in the  
352 other variables. In this case, the high-variance variables will *de-facto* determine the  
353 clustering, implying that the clustering will fail to take the majority of the measured  
354 information into account. To even out the influence of these high-variance variables  
355 on the clustering outcome, the data is log transformed when such variables are  
356 present. In DepecheR, this data behavior is detected automatically by concatenating  
357 all variables into a one dimensional vector, for which the kurtosis is calculated. A  
358 high kurtosis, indicates that the variables differ greatly in their internal variance. For  
359 datasets with low kurtosis, refraining from the log transform is preferable, since  
360 transformation distorts the information.

361

362

363 **Centering.** Centering the origin to be close to the bulk of the data is preferable, in  
364 order to have all biological clusters at approximately the same distance from the

365 origin. Having some biological clusters close to the origin and some far off is often  
366 unwanted, since the linear penalty then imposes a preference for creating clusters  
367 close to the origin. Apart from influencing the clustering, the centering also  
368 determines the interpretation of the obtained sparsity. Just as for scaling, which  
369 centering scheme to apply depends on the dataset.

370 For low dimensional datasets ( $n > 100$ ), DEPECHE applies maximal density  
371 centering, which sets the zero in each dimension to coincide with the highest data  
372 density. The density is computed by collecting the data in equally spaced bins (default  
373 number of bins in DepecheR is the number of data points  $n$  divided by 50), where the  
374 bin with the highest number of data points has the highest density. Using this scheme,  
375 sparsity (*i.e.* that a variable is non-contributing to the definition of a cluster) is  
376 interpreted as that the data points in the cluster do not deviate from the most common  
377 outcome. It also ensures that the origin is relatively close to the bulk of the data, since  
378 it is located at the most common outcome for each variable respectively. The benefit  
379 of this scheme is that it boosts sparsity, by declaring the most common outcome non-  
380 defining. However, for high dimensional datasets ( $n \geq 100$ ), maximal density  
381 centering can push the origin so far away from the center of mass of the dataset, that  
382 the penalty starts to impose an unwanted, artificial influence on the clustering,  
383 hampering the biological relevance of the clusters. To avoid this, DEPECHE imposes  
384 a mean centering scheme for such datasets, which locates the origin at the center of  
385 mass of the dataset.

386 A potential complication, related to centering, occurs when a biologically relevant  
387 cluster is located very close to the origin, since DEPECHE creates no clusters in the  
388 origin and will then force the cluster to merge with other clusters. However, this  
389 scenario was never detected in real data.

390

## 391 **Experimental procedures**

392 **Preprocessing of mass cytometry data.** The benchmark datasets from Levine *et al*  
393 [12] and Bendall *et al*[18] were transformed using the flowTrans package [23] before  
394 used in any clustering algorithm.

395 **Preprocessing of single-cell transcriptomic data.** The dataset from Björklund *et al*  
396 [19] was normalized using the sva package [24] as in the original manuscript. For this  
397 dataset, doublet variables were removed, lowering the number of variables from  
398 64443 to 35177.

399 The gold-standard datasets used for benchmarking in the publication by  
400 Kiselev *et al* [5] were obtained in a pre-processed state. Before clustering with any  
401 algorithm, the gene filter function used in the sc3 package was used [5], with settings  
402 removing the genes that were expressed in more than 90% of the cells. This resulted  
403 in the number of transcripts presented in Table 1 (range 13608-19571 transcripts).

404

## 405 **Code availability**

406 All code necessary to generate the figures and tables in the manuscript are included in  
407 supporting code 1. The software package DepecheR is available for download at  
408 (<https://github.com/theorell/depecher>).

409

## 410 **Acknowledgments**

411 The authors are grateful for all important input that has come from the initial users of  
412 the DepecheR software, especially Dr Niklas Björkström, Dr Jakob Michaëlsson and  
413 Sigrun Stultz. Other colleagues that have contributed seminally to the framework of

414 ideas that have led to the creation of DEPECHE are Dr Bruce Bagwell and Dr  
415 Geoffrey Hart.

416

417

## 418 **Author contributions**

419 A.T. drafted mathematical models, co-wrote the software implementation and the  
420 manuscript, Y.T.B. co-wrote the manuscript, J.T. co-wrote mathematical models and  
421 drafted software implementation and the manuscript.

422

## 423 **References**

- 424 1. Spitzer MH, Nolan GP. Mass Cytometry: Single Cells, Many Features. *Cell*.  
425 2016;165: 780–791. doi:10.1016/j.cell.2016.04.019
- 426 2. Tanay A, Regev A. Scaling single-cell genomics from phenomenology to  
427 mechanism. *Nature*. 2017;541: 331–338. doi:10.1038/nature21350
- 428 3. Budnik B, Levy E, Slavov N. Mass-spectrometry of single mammalian cells  
429 quantifies proteome heterogeneity during cell differentiation. *bioRxiv*. 2017;  
430 102681. doi:10.1101/102681
- 431 4. Saeys Y, Gassen SV, Lambrecht BN. Computational flow cytometry: helping to  
432 make sense of high-dimensional immunology data. *Nat Rev Immunol*. 2016;16:  
433 449–462. doi:10.1038/nri.2016.56
- 434 5. Kiselev VY, Kirschner K, Schaub MT, Andrews T, Yiu A, Chandra T, et al. SC3:  
435 consensus clustering of single-cell RNA-seq data. *Nat Methods*. 2017;14: 483–  
436 486. doi:10.1038/nmeth.4236
- 437 6. Guo M, Wang H, Potter SS, Whitsett JA, Xu Y. SINCERA: A Pipeline for  
438 Single-Cell RNA-Seq Profiling Analysis. *PLoS Comput Biol*. 2015;11:  
439 e1004575. doi:10.1371/journal.pcbi.1004575
- 440 7. Žurauskienė J, Yau C. pcaReduce: hierarchical clustering of single cell  
441 transcriptional profiles. *BMC Bioinformatics*. 2016;17: 140.  
442 doi:10.1186/s12859-016-0984-y
- 443 8. MacQueen J. Some methods for classification and analysis of multivariate  
444 observations. *Proceedings of the Fifth Berkeley Symposium on Mathematical  
445 Statistics and Probability*. Berkeley, California: University of California Press;

- 446 1967. pp. 281–297. Available:  
447 <https://projecteuclid.org/euclid.bsmmsp/1200512992>
- 448 9. Pan W. Penalized model-based clustering with application to variable selection. *J*  
449 *Mach Learn Res.* 2007;8: 1145–1164.
- 450 10. Sun W, Wang J, Fang Y. Regularized k-means clustering of high-dimensional  
451 data and its asymptotic consistency. *Electron J Stat.* 2012;6: 148–167.  
452 doi:10.1214/12-EJS668
- 453 11. Hubert L, Arabie P. Comparing partitions. *J Classif.* 1985;2: 193–218.  
454 doi:10.1007/BF01908075
- 455 12. Levine JH, Simonds EF, Bendall SC, Davis KL, Amir ED, Tadmor MD, et al.  
456 Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that  
457 Correlate with Prognosis. *Cell.* 2015;162: 184–197.  
458 doi:10.1016/j.cell.2015.05.047
- 459 13. Lo K, Hahne F, Brinkman RR, Gottardo R. flowClust: a Bioconductor package  
460 for automated gating of flow cytometry data. *BMC Bioinformatics.* 2009;10:  
461 145. doi:10.1186/1471-2105-10-145
- 462 14. Aghaeepour N. flowMeans: Non-parametric Flow Cytometry Data Gating. 2010.
- 463 15. Sørensen T, Baumgart S, Durek P, Grützkau A, Häupl T. immunoClust--An  
464 automated analysis pipeline for the identification of immunophenotypic  
465 signatures in high-dimensional cytometric datasets. *Cytom Part J Int Soc Anal*  
466 *Cytol.* 2015;87: 603–615. doi:10.1002/cyto.a.22626
- 467 16. Ge Y, Sealfon SC. flowPeaks: a fast unsupervised clustering for flow cytometry  
468 data via K-means and density peak finding. *Bioinforma Oxf Engl.* 2012;28:  
469 2052–2058. doi:10.1093/bioinformatics/bts300
- 470 17. Zare H, Shooshtari P, Gupta A, Brinkman RR. Data reduction for spectral  
471 clustering to analyze high throughput flow cytometry data. *BMC Bioinformatics.*  
472 2010;11: 403. doi:10.1186/1471-2105-11-403
- 473 18. Bendall SC, Simonds EF, Qiu P, Amir ED, Krutzik PO, Finck R, et al. Single-  
474 Cell Mass Cytometry of Differential Immune and Drug Responses Across a  
475 Human Hematopoietic Continuum. *Science.* 2011;332: 687–696.  
476 doi:10.1126/science.1198704
- 477 19. Björklund ÅK, Forkel M, Picelli S, Konya V, Theorell J, Friberg D, et al. The  
478 heterogeneity of human CD127(+) innate lymphoid cells revealed by single-cell  
479 RNA sequencing. *Nat Immunol.* 2016;17: 451–460. doi:10.1038/ni.3368
- 480 20. Dempster AP, Laird NM, Rubin DB. Maximum Likelihood from Incomplete Data  
481 via the EM Algorithm. *J R Stat Soc Ser B Methodol.* 1977;39: 1–38.
- 482 21. Witten DM, Tibshirani R. A framework for feature selection in clustering. *J Am*  
483 *Stat Assoc.* 2010;105: 713–726. doi:10.1198/jasa.2010.tm09415

- 484 22. Arthur D, Vassilvitskii S. K-means++: The Advantages of Careful Seeding.  
485 Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete  
486 Algorithms. Philadelphia, PA, USA: Society for Industrial and Applied  
487 Mathematics; 2007. pp. 1027–1035. Available:  
488 <http://dl.acm.org/citation.cfm?id=1283383.1283494>
- 489 23. Finak G, Manuel-Perez J, Gottardo R. flowTrans: Parameter Optimization for  
490 Flow Cytometry Data Transformation. 2010.
- 491 24. Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD. The sva package for  
492 removing batch effects and other unwanted variation in high-throughput  
493 experiments. *Bioinforma Oxf Engl*. 2012;28: 882–883.  
494 doi:10.1093/bioinformatics/bts034

495

496

## 497 **Supporting information**

498 **S1 Fig Heatmaps comparing the golden standard partitions to the DEPECHE**  
499 **partitions for a) the 32-variate Levine dataset and b) the 35166-variate**  
500 **Björklund dataset.** Red color indicates large overlap, blue color indicates low  
501 overlap between a gold standard-vs-depeche cluster pair.

502

503 **S2 Fig Algorithm comparisons.** For all graphs, the x-axis shows the algorithms and  
504 the y-axis shows the Adjusted Rand Index comparing the clustering result with the  
505 golden standard clustering. a) Subsamples with 20000 unique cells from two mass  
506 cytometry datasets published by Levine *et al* and Bendall *et al* were clustered with  
507 DEPECHE and six previously published algorithms. For each dataset and algorithm,  
508 clustering was performed on 20 unique subsamples. For flowClust, flowPeaks and  
509 SamSPECTRAL, that do not perform internal parameter tuning, a range of parameter  
510 values were evaluated and the parameter value sets generating the highest ARI values  
511 were selected for display. b) The full Björklund dataset, as well as six other datasets  
512 previously used for benchmarking by Kiselev *et al* were clustered 20 times with

513 DEPECHE and three other algorithms. The Björklund dataset was normalized to  
514 reduce batch effects, with the procedure described in the original publication. These  
515 six datasets were also automatically log<sub>2</sub>-transformed within DEPECHE, and thus,  
516 log<sub>2</sub>-transformation was applied also for Sincera and pcaReduce, whereas sc3 was fed  
517 both log<sub>2</sub>- and untransformed data. The lower and upper hinges of all boxplots  
518 extend to the 25:th and 75:th percentile, whereas the line in the middle describes the  
519 median. The whiskers extend to the lowest and highest value no further than 1.5 times  
520 the distance between the 25:th and 75:th percentile. Outside of this range, the  
521 observations are considered outliers and are shown as dots.

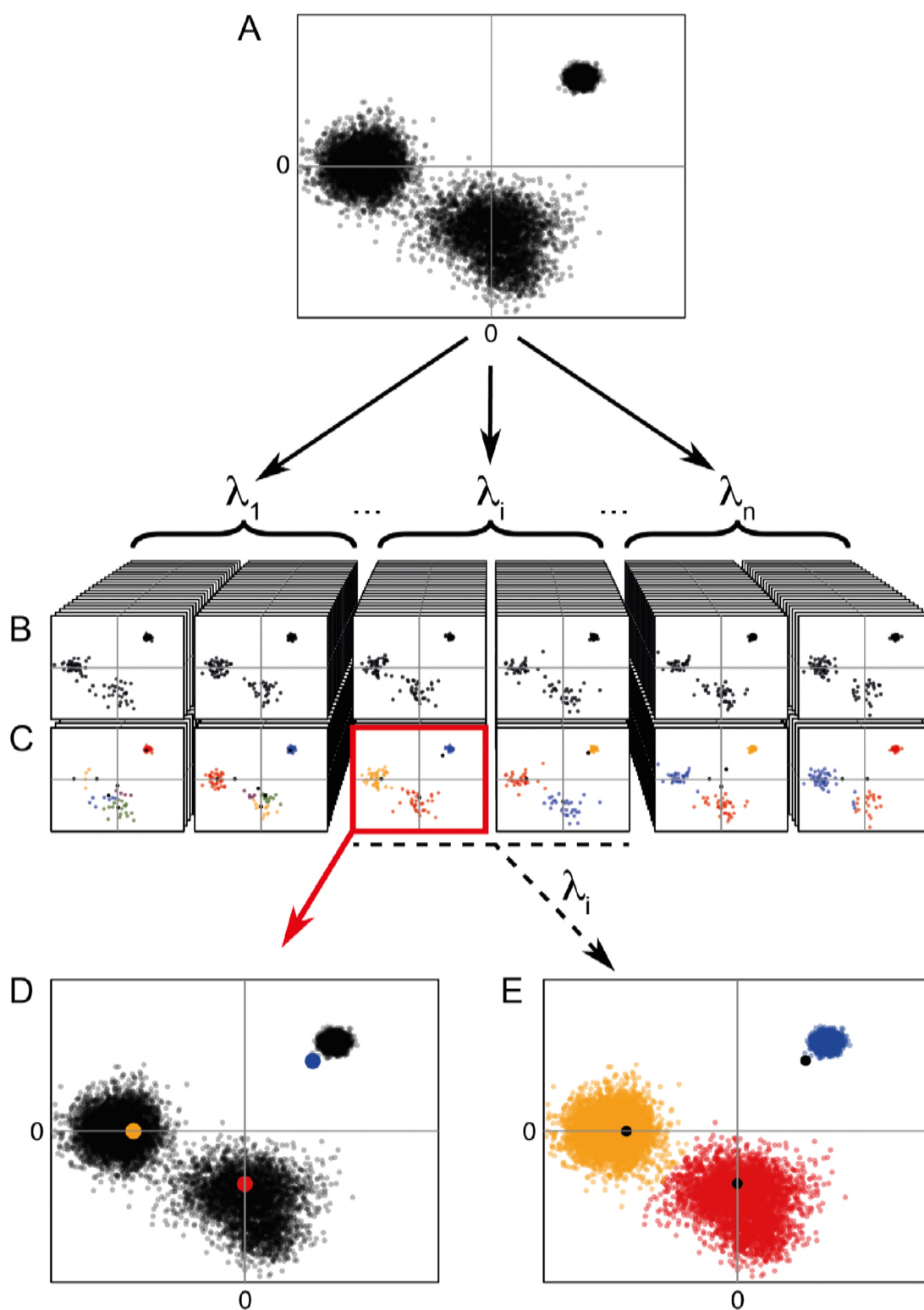
522

523 **S1 File. The DepecheR software, for the review phase.**

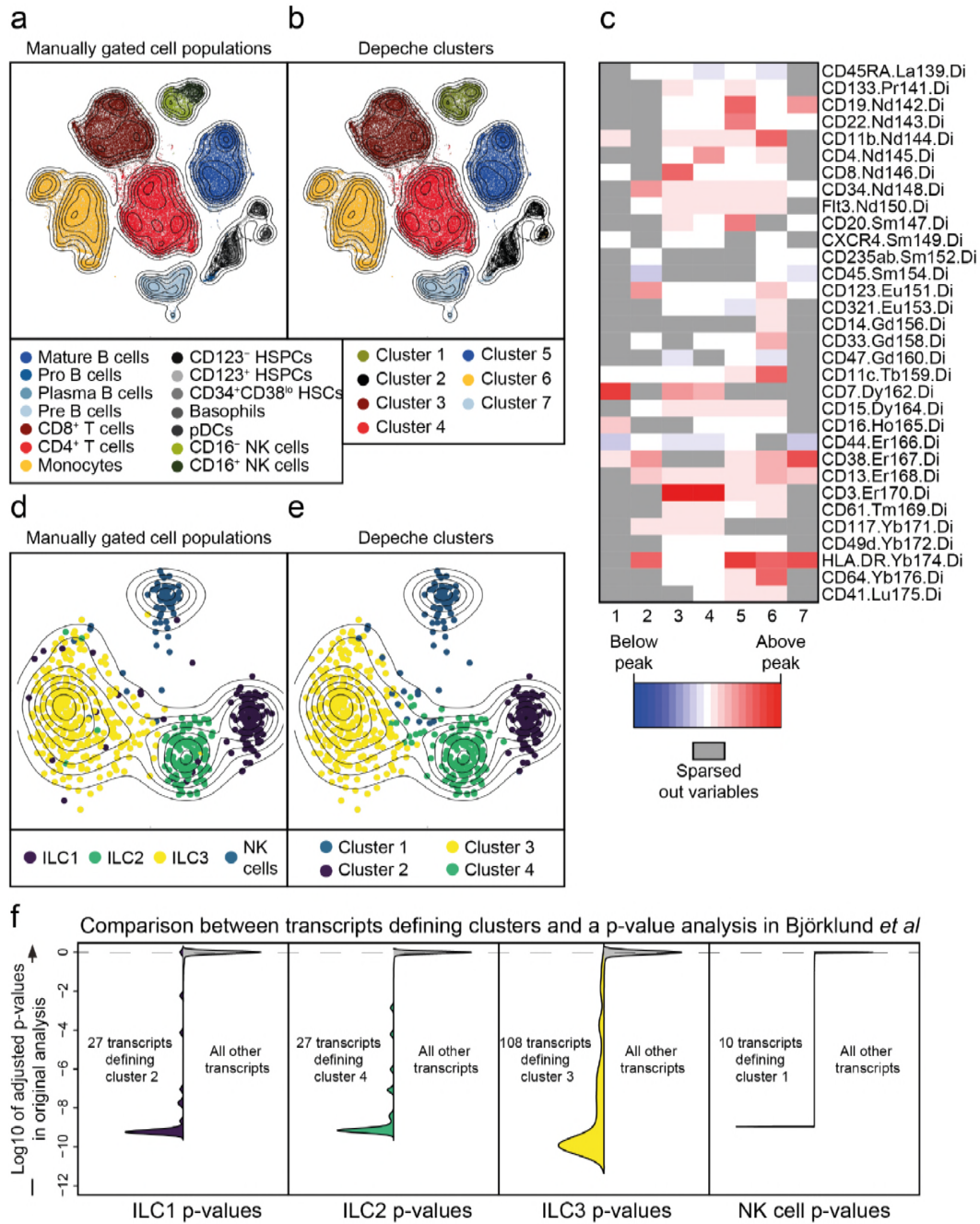
524

525 **S2 File. The code needed to generate all figures, for the review phase.**

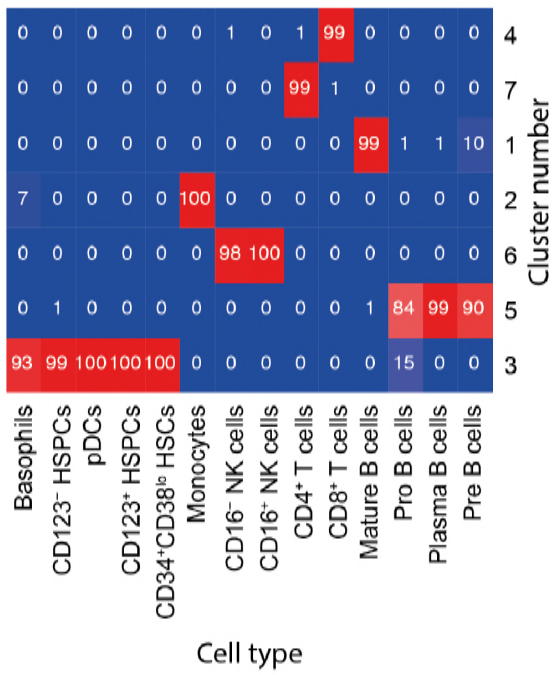
526



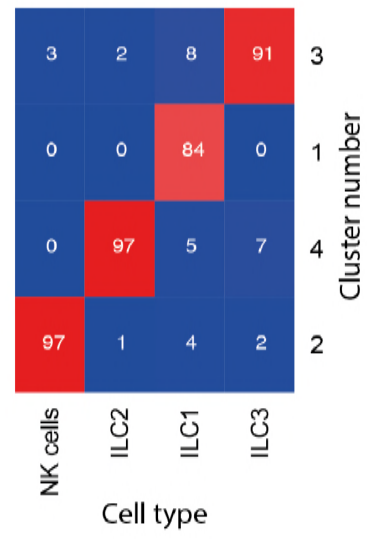


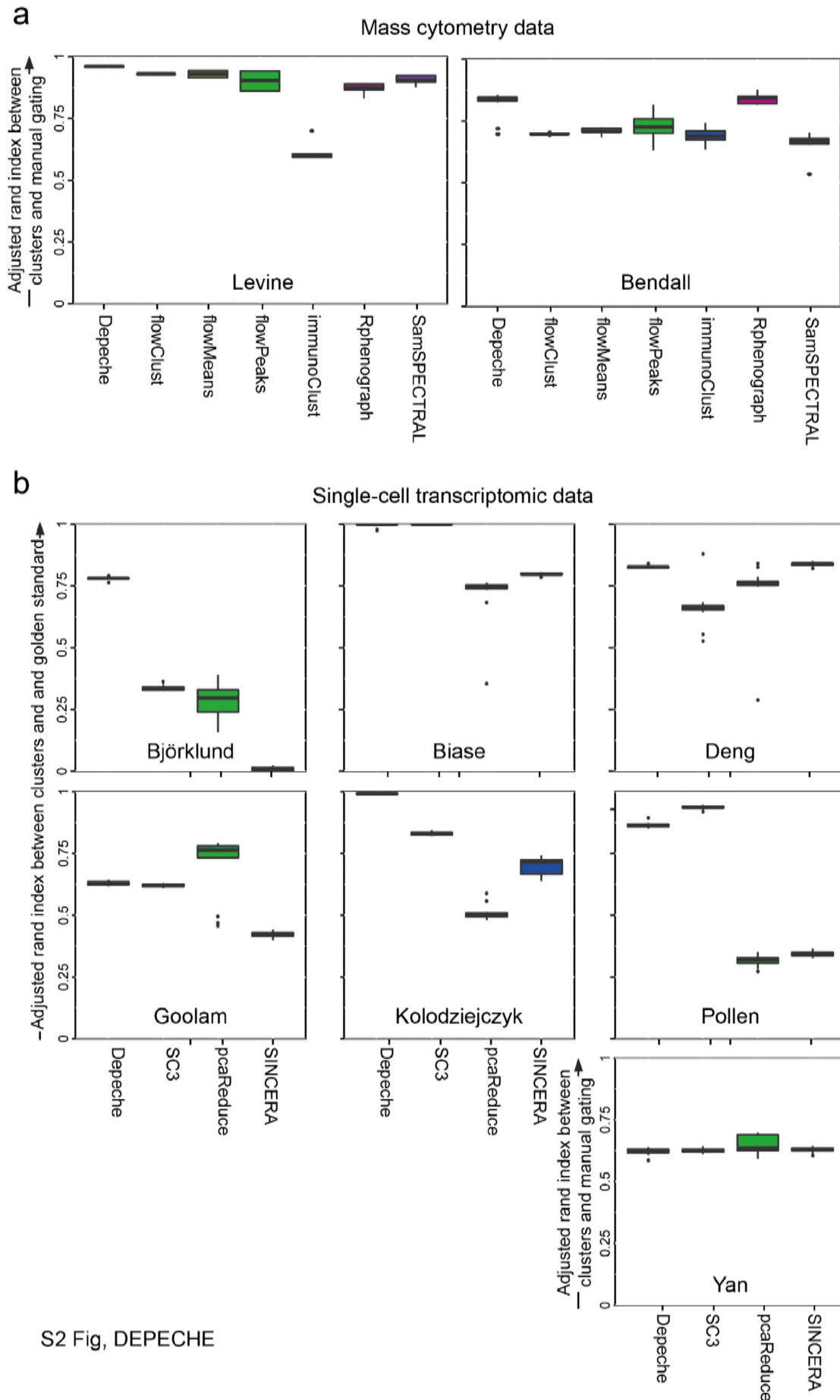


a



b





S2 Fig, DEPECHE