

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

Hypercluster: a python package and SnakeMake pipeline for flexible, parallelized  
unsupervised clustering optimization

Lili Blumenberg<sup>1,2</sup>, Kelly V. Ruggles<sup>1,2\*</sup>

<sup>1</sup>Institute of Systems Genetics, New York University School of Medicine, New York, NY  
10016, USA <sup>2</sup>Department of Medicine, New York University School of Medicine, New  
York, NY 10016, USA

\*Corresponding author

Email: [Kelly.Ruggles@nyulangone.org](mailto:Kelly.Ruggles@nyulangone.org)

24 Full title: Hypercluster: a python package and SnakeMake pipeline for flexible,  
25 parallelized unsupervised clustering optimization

26 Short title: Hypercluster: a tool for unsupervised clustering optimization

27

28 Abstract

29 Unsupervised clustering is a common and exceptionally useful tool for large  
30 biological datasets. However, clustering requires upfront algorithm and hyperparameter  
31 selection, which can introduce bias into the final clustering labels. It is therefore  
32 advisable to obtain a range of clustering results from multiple models and  
33 hyperparameters, which can be cumbersome and slow. To streamline this process, we  
34 present hypercluster, a python package and SnakeMake pipeline for flexible and  
35 parallelized clustering evaluation and selection. Hypercluster is available on bioconda;  
36 installation, documentation and example workflows can be found at:

37 <https://github.com/ruggleslab/hypercluster>.

38

39 Author summary

40 Unsupervised clustering is a technique for grouping similar samples within a  
41 dataset. It is extremely common when analyzing big data from patient samples, or high  
42 throughput techniques like single cell RNA-seq. When researchers use unsupervised  
43 clustering, they have to select parameters that affect the final result—for instance, how  
44 many groups they expect to find or what the smallest group is allowed to be. Some  
45 methods require setting even less intuitive parameters. For most applications, it is  
46 extremely challenging to guess what the values of these parameters should be;

47 therefore to prevent introducing bias into the final results, researchers should test many  
48 different parameters and methods to find the best groups. This process is cumbersome,  
49 slow and challenging to perform in a reproducible way. We developed hypercluster, a  
50 tool that automates this process, make it much faster, and presenting the results in a  
51 reproducible and helpful manner.

52

### 53 Introduction

54 Unsupervised clustering is commonly used for the interpretation of ‘omics  
55 datasets. It provides an objective and intuitive measure of similarity and difference  
56 between samples. Clustering can be used to determine biologically relevant subgroups  
57 of samples, find co-regulated molecular features, or provide objective support for the  
58 phenotypic similarity of biological perturbations. Moreover, clustering is a key step in the  
59 analysis of many emerging sequencing-based technologies. For example, a  
60 fundamental challenge in the analysis of single-cell measurement data, in particular  
61 single cell RNA-seq (scRNA-seq), is determining robust clusters of phenotypically  
62 similar cells (1–3). Clustering is also increasingly being used alongside traditional  
63 diagnostic techniques to establish new classifications of patient samples into disease-  
64 relevant subgroups (4–7) and for patient subgroup classification and risk stratification  
65 (6,8–12). The near-future of personalized medicine relies on researchers identifying  
66 robust unsupervised clustering-based disease subtypes. Therefore, it is essential that  
67 high-quality clustering results are easily and robustly obtainable, without user-selected  
68 hyperparameters introducing bias and impeding rapid analysis.

69           Currently, researchers robustly employing unsupervised clustering must choose  
70 specific algorithms and hyperparameters that are appropriate to their experiment type  
71 and data. Although some efforts have been made to advise researchers on optimal  
72 selection of both (13), biological datasets vary between batches, days, labs and  
73 researchers, underlining the importance of context- and experiment-dependent analysis  
74 tuning. Software packages for automatic hyperparameter tuning and model selection for  
75 regression and classification machine learning techniques exist, notably auto-sklearn  
76 from AutoML (14), but there are not yet packages for automated unsupervised  
77 clustering optimization.

78           Typically, the effect of hyperparameter choice on the quality of clustering results  
79 cannot be described with a convex function, meaning that when searching the  
80 landscape of hyperparameter choices there are often local maxima that may appear to  
81 be the optimal results if broad choices of hyperparameters are not considered.  
82 Therefore it is unlikely that a sequential approach using for instance, gradient descent  
83 from a single initialized set of hyperparameters, would be able to select the optimal  
84 parameters for the majority of clustering challenges (15). Exhaustive (i.e. grid) search is  
85 the most likely to obtain optimal results from unsupervised clustering. However, grid  
86 search can be slow and cumbersome to perform for the multiple hyperparameters and  
87 clustering algorithms that are available from most clustering packages.

88           Here we present hypercluster, a python package and SnakeMake pipeline for  
89 parallelized clustering calculations and comparison. The hypercluster package allows  
90 users to calculate results from multiple hyperparameters using one or many algorithms,  
91 then easily calculate and visualize evaluation metrics for each result (16). The

92 accompanying SnakeMake pipeline allows parallelization on a single computer, across  
93 a high performance computing cluster, or on cloud based services (17,18), speeding up  
94 optimization, especially for large datasets. In addition, our pipeline has all the  
95 advantages of the SnakeMake framework, e.g. easily adding new datasets to analyze,  
96 keeping track of progress and simplified bug tracking. Currently, hypercluster can  
97 compare all clustering algorithms and evaluation metrics from scikit-learn (19), as well  
98 as non-negative matrix factorization (NMF) (20), Louvain and Leiden clustering (21,22).  
99 In addition, hypercluster can be extended to employ user-supplied clustering algorithm  
100 or evaluation metrics. Given a metric to maximize, hypercluster identifies “best” labels  
101 and optionally provides comparisons of labeling results. Even if no single metric can be  
102 used to select the best hyperparameters, hypercluster provides several visualizations  
103 that help users pick labels by balancing many metrics or picking the most reproducible  
104 clusters. Hypercluster provides researchers with a python package and pipeline for  
105 flexible, parallelized, distributed and user-friendly algorithm selection and hyper-  
106 parameter tuning for unsupervised clustering.

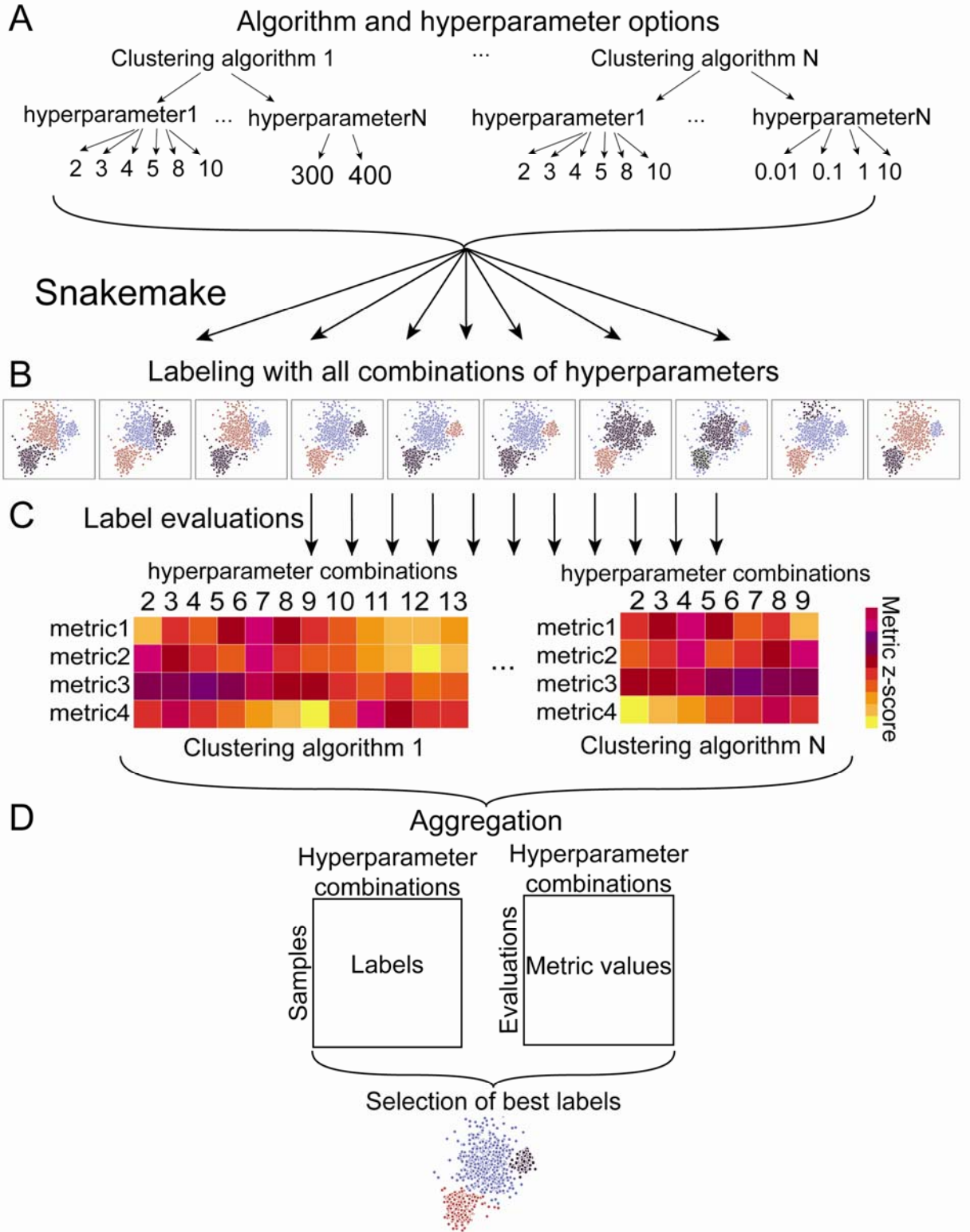
107

## 108 Design and Implementation

### 109 *Requirements and structure*

110 The hypercluster package uses scikit-learn (19), python-igraph (23), leidenalg  
111 (24) and louvain-igraph (25) to assign cluster labels and uses scikit-learn and custom  
112 metrics to compare clustering algorithms and hyperparameters to find optimal clusters  
113 for any given input data (Fig. 1). Hypercluster requires python3, pandas (26), numpy  
114 (27), scipy (28), matplotlib (29), seaborn (30), scikit-learn (19), python-igraph (23),

115 leidenalg (24), louvain-igraph (25) and SnakeMake (17). Hypercluster can be run  
116 independently of SnakeMake, as a standalone python package. Inputs, outputs and an  
117 example workflow are described below, but additional example workflows are provided  
118 at <https://github.com/ruggleslab/hypercluster/tree/master/examples>.



119

120 *Fig. 1 Hypercluster workflow schematic*

- 121 a) Clustering algorithms and their respective hyperparameters are user-specified.  
122 Hypercluster then uses those combinations to create exhaustive configurations, and if  
123 selected a random subset is chosen.  
124 b) Snakemake is then used to distribute each clustering calculation into different jobs.  
125 c) Each set of clustering labels is then evaluated in a separate job by a user-specified  
126 list of metrics.  
127 d) All clustering results and evaluation results are aggregated into tables. Best labels  
128 can also be chosen by a user-specified metric.

129

### 130 *Modes*

131 Hypercluster takes pandas DataFrames as input. For local running,  
132 AutoClusterer and MultiAutoClusterer objects can be instantiated with default or user-  
133 defined values. To run through hyperparameters for a dataset, users simply provide a  
134 pandas DataFrame to the “fit” method on either object. Users evaluate the labeling  
135 results by running the “evaluate” method.

### 136 *config.yml*

137 SnakeMake allows users to parallelize clustering calculations. To configure the  
138 SnakeMake pipeline, users edit a config.yml file (Table 1). In that file, users can specify  
139 input and output directories and files (Table 1, lines 1-3, 5-7) and the hyperparameter  
140 search space (Fig 1A, Table 1, line 18). Users can specify whether to use exhaustive  
141 grid search or random search; if random search is selection, they can specify probability  
142 weights for each hyperparameter (Table 1, line 9). Snakemake then schedules  
143 performing each clustering algorithm and evaluating the results as a separate job (Fig.



144 1B). Users can specify which evaluation metrics to apply (Fig. 1C, Table 1, line 10) and  
 145 add keyword arguments to tune several steps in the process (Table 1, lines 4, 8-9, 11-  
 146 16). Clustering and evaluation results are then aggregated into final tables (Fig. 1D).  
 147 Other than the location and names of the input files, everything has a predefined default  
 148 that allows the pipeline to be used “out of the box.” Users can reference the  
 149 documentation and examples for more information.

**Table 1 Parameters in SnakeMake configuration file**

config.yml parameter	Explanation	Example
1 input_data_folder	Path to folder in which input data can be found.	/input_data
2 input_data_files	List of prefixes of data files.	['input_data1', 'input_data2']
3 gold_standard_file	File name of gold_standard_file, must be in input_data_folder	{'input_data': 'gold_standard_file.txt'}
4 read_csv_kwargs	pandas.read_csv keyword arguments for input data.	{'test_input': {'index_col': [0]}}
5 output_folder	Path to folder into which results should be written.	/results
6 intermediates_folder	Name of subfolder to put intermediate results.	clustering_intermediates
7 clustering_results	Name of subfolder to put aggregated results.	clustering
8 clusterer_kwargs	Additional arguments to pass to clusterers.	KMeans: {'random_state': 8}}
9 generate_parameters_addtl_kwargs	Additional keyword arguments for the hypercluster.AutoClusterer class.	{'KMeans': {'random_search': true}}
10 evaluations	Names of evaluation metrics to use.	['silhouette_score', 'number_clustered']
11 eval_kwargs	Additional kwargs per evaluation metric function.	{'silhouette_score': {'random_state': 8}}
12 metric_to_choose_best	Which metric to maximize to choose the labels.	silhouette_score
13 metric_to_compare_labels	Which metric to use to	adjusted_rand_score

	compare label results to each other.	
14	compare_samples	Whether to made a table and figure with counts of how often each two samples are in the same cluster. "true"
15	output_kwargs	pandas.to_csv and pandas.read_csv keyword arguments for output tables. {'evaluations': {'index_col':[0]}, 'labels': {'index_col':[0]}}
16	heatmap_kwargs	Arguments for seaborn.heatmap for pairwise visualizations. {'vmin':-2, 'vmax':2}
17	optimization_parameters	Which algorithms and corresponding hyperparameters to try. {'KMeans': {'n_clusters': [5, 6, 7]}}

150

151 *Table 1 Line-by-line explanation of the config.yml for SnakeMake*

152

153 *Input data and execution*

154 After specifying the config.yml file, users provide a data table with samples to be  
 155 clustered as the rows and features as the columns, with the location specified in the  
 156 config.yml file (Table 1, line2). Users can then simply run “snakemake -s  
 157 hypercluster.smk --configfile config.yml” in the command line, with any additional  
 158 SnakeMake flags appropriate for their system. Applying the same configuration to new  
 159 files or adjusting algorithms and hyperparameter options simply requires editing the  
 160 config.yml file and rerunning SnakeMake.

161 *Extending hypercluster*

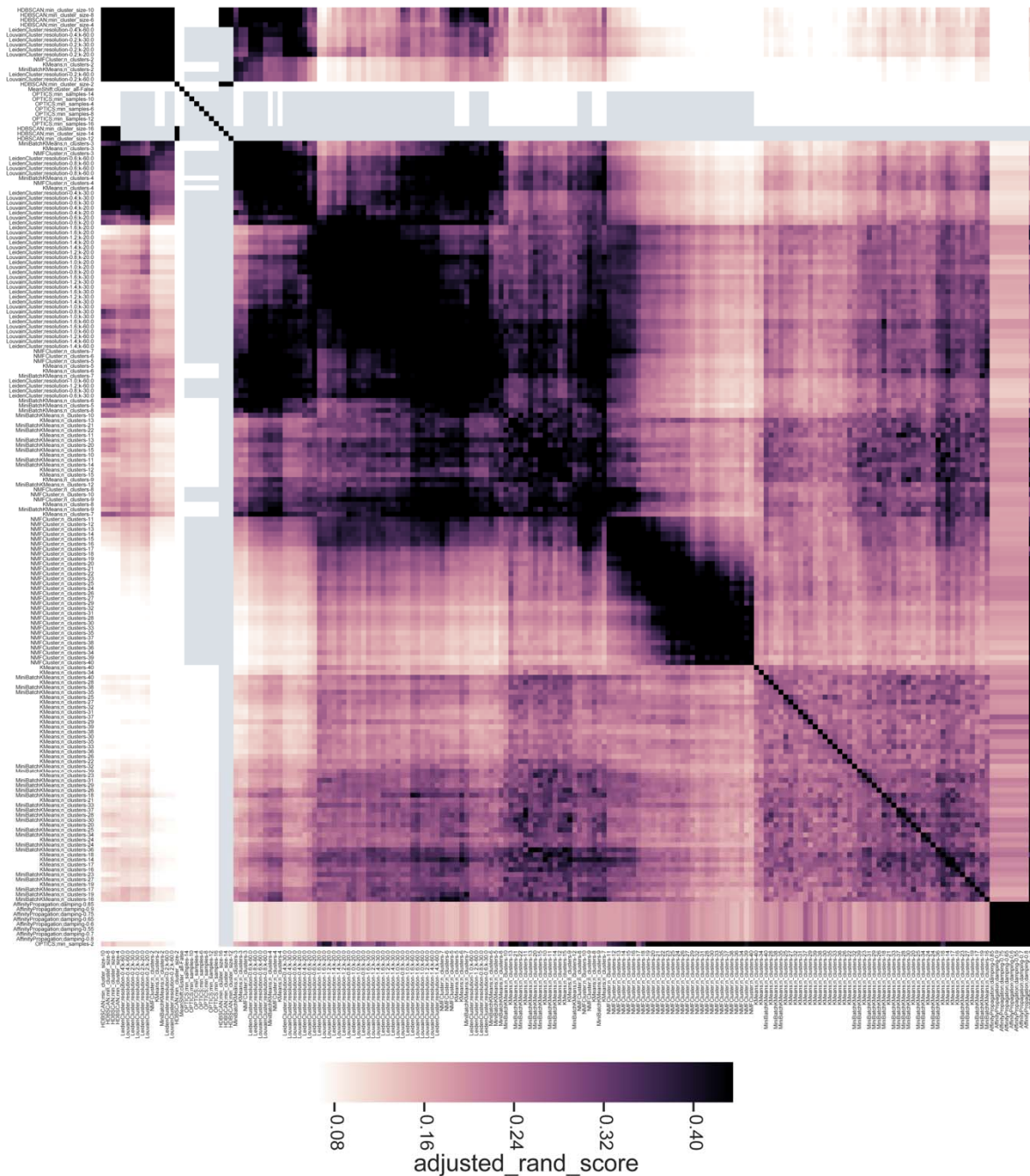
162 Currently, hypercluster can optimize any clustering algorithm and calculate any  
 163 evaluation available in scikit-learn (19,31), as well as NMF, Louvain and Leiden

164 clustering. Additional clustering classes and evaluation metric functions can be added  
165 by users in the `additional_clusterer.py` and `additional_metrics.py` files, respectively, if  
166 written to accommodate the same input, outputs and methods (see  
167 `additional_clusterers.py` and `additional_metrics.py` for examples).

### 168 *Outputs*

169 By default, hypercluster outputs a yaml file containing all configurations of the  
170 clustering algorithms and hyperparameters that are being searched. For each set of  
171 labels, it generates a file containing labels and a file containing evaluations. It also  
172 outputs aggregated tables of all labels and evaluations. Finally, given a metric to  
173 maximize, hypercluster writes files containing the optimal labels. Optionally,  
174 hypercluster will also output a table and heatmap of pairwise comparisons of labeling  
175 similarities with a user-specified metric (Figure S1). This figure is particularly useful for  
176 finding labels that are robust to differences in hyperparameters. It can also optionally  
177 output a table and heatmap showing how often each pair of samples were assigned the  
178 same cluster (Figure S2).

## Pairwise similarity between labels for breast cancer RNA-seq



179

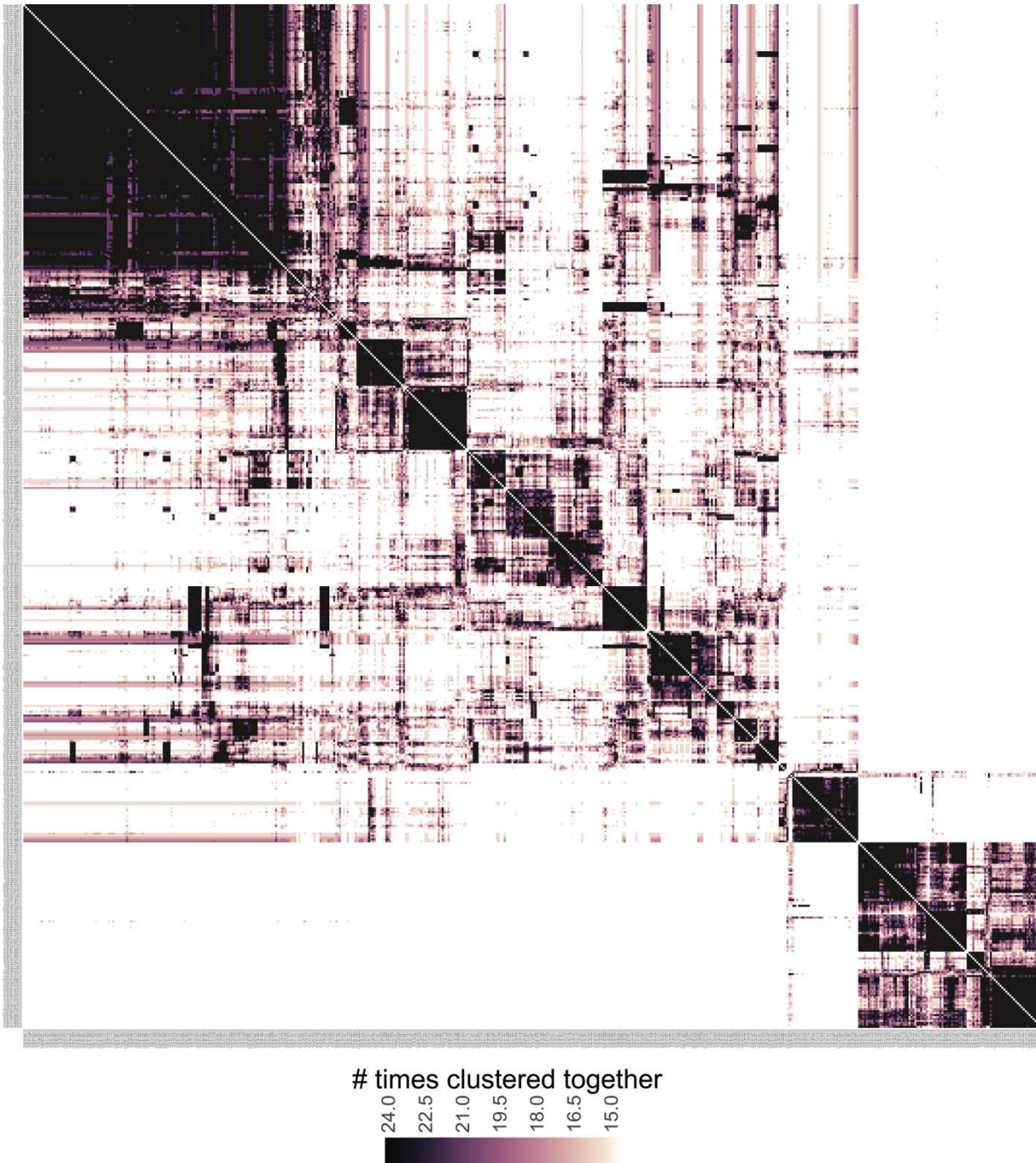
180 *Figure S1 Pairwise label comparisons*

181 Automatically generated heatmap showing pairwise comparison of labeling

182 automatically generated using hypercluster of breast cancer samples. Colors represent

183 adjusted rand index between labels.





185 *Figure S2 Pairwise sample comparisons*

186 Automatically generated pairwise comparison of breast cancer samples. Color indicates  
187 the number of times two samples were assigned the same cluster.

188

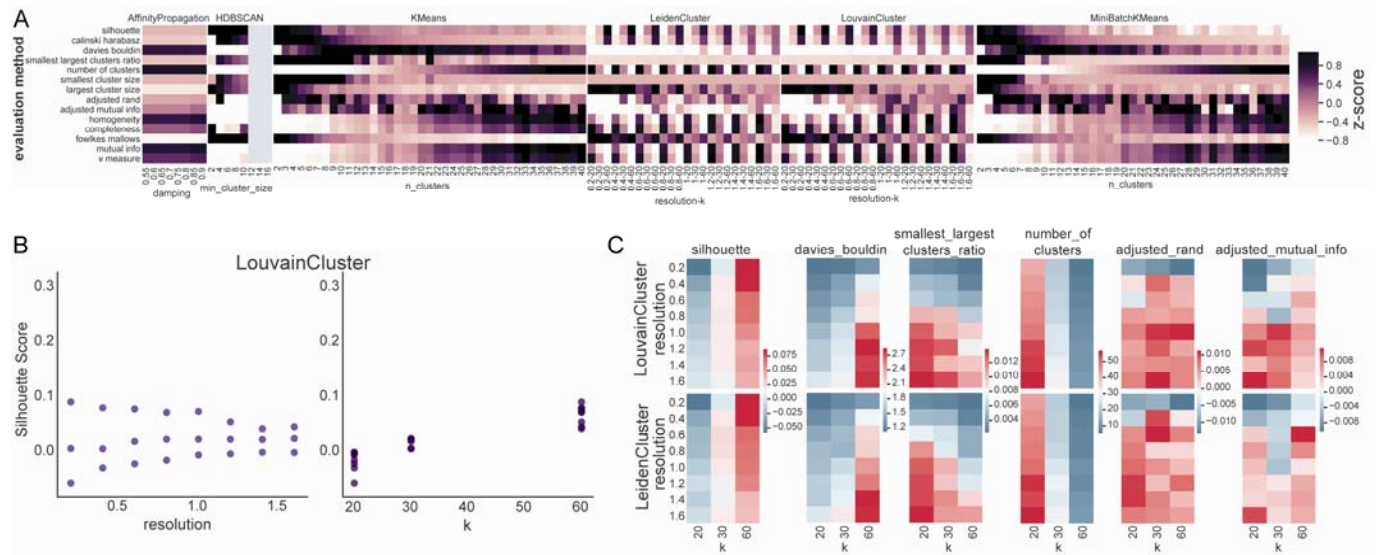
189 Results

190 *Unsupervised clustering of RNA-seq on breast cancer patient samples*

191 To illustrate the utility of hypercluster in a disease-relevant context, we applied  
192 our method to RNA-seq data from 526 breast cancer patient samples from the Cancer  
193 Genome Atlas (TCGA) (32), a dataset that has been previously used for benchmarking  
194 clustering algorithms (33). As demonstrated, RNA-seq can be used to classify breast  
195 cancer patients into four major PAM50 subtypes (Basal-like, LuminalA, LuminalB, and  
196 Her2-enriched), which are based on the expression of 50 specific genes (7,34,35). We  
197 removed genes with any missing values and subset to the 500 most variable genes as  
198 input for all available algorithms with ranges of hyperparameter conditions. We then  
199 compared the sample clustering results from our 500 gene clustering compared with  
200 subtypes defined by the PAM50 classifier. This workflow is available on the github  
201 examples folder (<https://github.com/liliblu/hypercluster/tree/dev/examples>).

202 Hypercluster automatically outputs a visualization of evaluation metrics for all  
203 hyperparameter combinations (Fig. 2A), which allows users to quickly see how  
204 changing hyperparameters affects clustering result quality. These results highlight how  
205 evaluation metrics are not generally convex over ranges of hyperparameters (e.g.  
206 silhouette score as n\_clusters changes with the KMeans algorithm (Fig. 2A)  
207 demonstrating the utility of the exhaustive grid search approach. In addition, our pipeline  
208 optionally creates a pairwise comparison of labeling, with a specified user metric (Figure  
209 S1) to make it easier to understand how robust and consistent labeling is across  
210 algorithms and parameters.

211



212

213 *Fig. 2 Visualizations of clustering metrics from breast cancer RNA-seq*

214 a) Example automatic output from hypercluster, showing z-scored evaluation values of

215 evaluation metrics for each clustering algorithm and hyperparameter set, including

216 those in 2E. Evaluations applied to clustering from TCGA breast cancer samples.

217 b) Effect of varying resolution (left) and k for shared nearest neighbor matrix (right) on

218 silhouette score, an inherent metric measuring clustering quality, for Louvain clustering.

219 c) Effect of resolution and choice of k on various evaluation metrics for both Louvain

220 (top) and Leiden (bottom) clustering.

221

222 Labels and evaluation results are easily accessible for further custom analyses.

223 To demonstrate a possible downstream workflow that hypercluster facilitates, we

224 investigated results from Louvain and Leiden clustering, which are commonly used in

225 scRNA-seq analysis, on the same breast cancer RNA-Seq dataset (36)). Louvain and

226 Leiden clustering are community detection algorithms for networks, usually generated

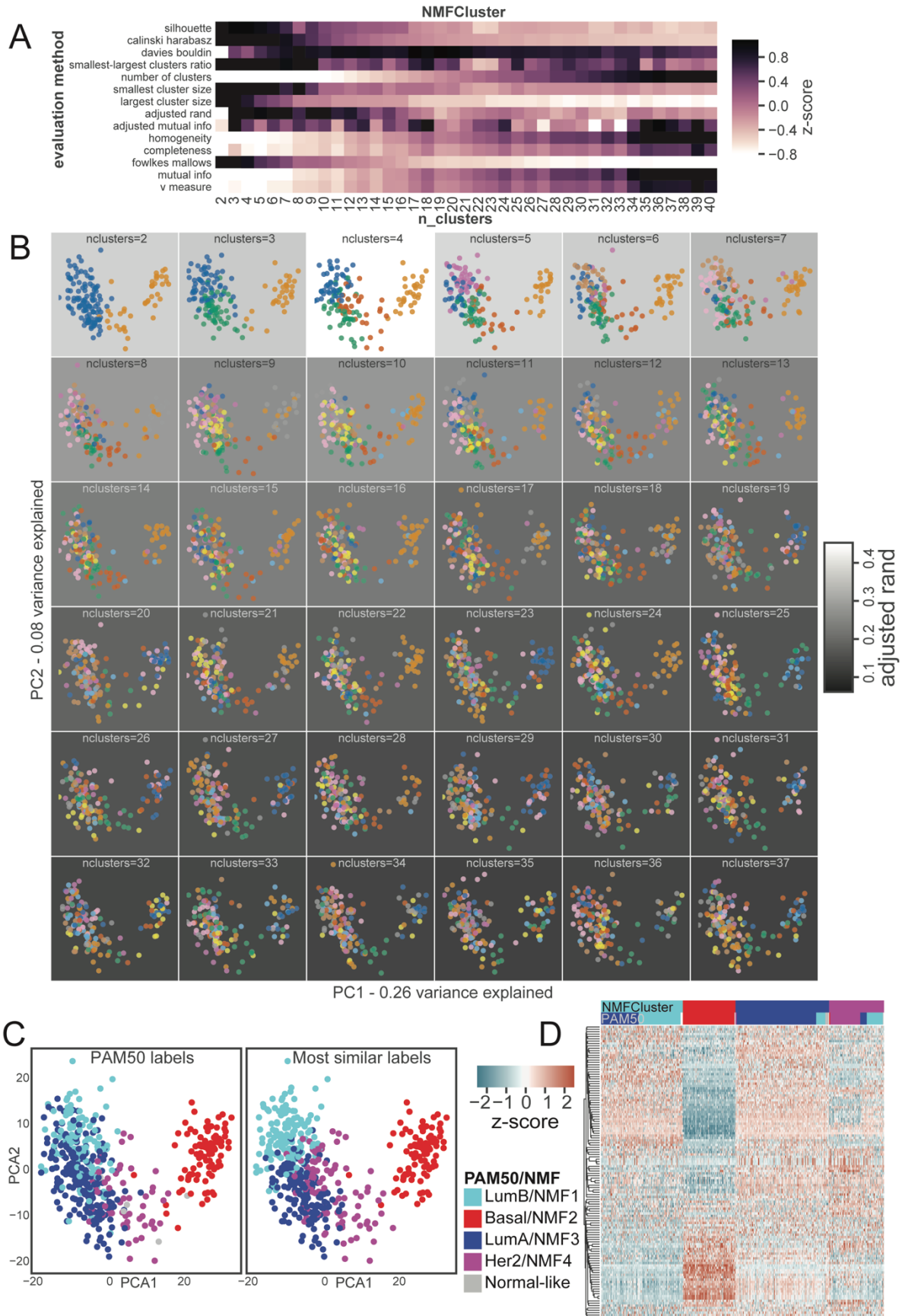
227 from shared K-nearest-neighbor adjacency matrices. We varied resolution, which

228 affects the number of members in final communities, and the k defining how many

229 nearest neighbors are measured for constructing the adjacency graph (Fig. 2B, C).  
230 Resolution and k have significant effects on labeling results and their corresponding  
231 evaluations. Interestingly, increasing resolution appears to have opposite effects on  
232 clustering quality (e.g. as measured by silhouette score) depending on k, with a large  
233 spread of silhouette scores dependent on k at low resolution, converging to similar  
234 silhouette scores at higher resolution (Fig. 2B, C). These results highlight the  
235 importance of simultaneous tuning of multiple hyperparameters. Plots like those in Fig.  
236 2B, showing the effect of varying each parameter individually on evaluation metrics, can  
237 be automatically generated by the `visualize_for_picking_best_labels` function or listing  
238 evaluations in the “`screepplot_evals`” section of a `config.yml` file.

239         To observe if clustering on 500 variable genes can recapitulate PAM50  
240 classification, we identified results that best match PAM50 subtypes according to the  
241 adjusted rand score while labeling all samples (Fig. 3). By this metric, the best labels  
242 were generated by NMF clustering (37) with `n_clusters=4` (Fig. 3A-C). These labels that  
243 do diverge from the PAM50 classification correspond to a subset of Luminal A samples  
244 that cluster with Luminal B samples (Fig. 3D). Hypercluster allows researchers to  
245 compare different algorithms and hyperparameter combinations in a reproducible and  
246 convenient way.





248 *Fig. 3 Exploration of NMF clustering results on breast cancer RNA-seq*

249 a) Automatically generated heatmap of evaluation metrics for NMF clustering results.

250 b) PCA projection of 200 random samples colored by labels assigned by NMF

251 clustering. Background color indicates similarity to PAM50 labeling calculated by

252 adjusted rand index.

253 c) PCA projection of samples colored by PAM50 subtypes and most similar NMF

254 clustering labels.

255 d) Heatmap of 125 most variable genes with PAM50 and NMF;n\_cluster=4 labels

256 indicated on the top.

257

258 *Exploration of bone marrow microenvironment scRNA-seq*

259 To demonstrate hypercluster's utility for analysis of single cell data sets, we

260 analyzed scRNA-seq from a study investigating the hematopoietic stem cell

261 microenvironment (38) and performed comparative analysis of several clustering

262 algorithms in parallel on a high performance computing cluster utilizing a Slurm

263 scheduler (39). We used normalized expression data from untreated cells sorted for

264 mesenchymal stromal and vascular endothelial, and osteoblast markers, subset to the

265 2000 most variable genes from the seurat object containing the data (36,38). We then

266 used hypercluster to explore the labeling results from all available clustering algorithms

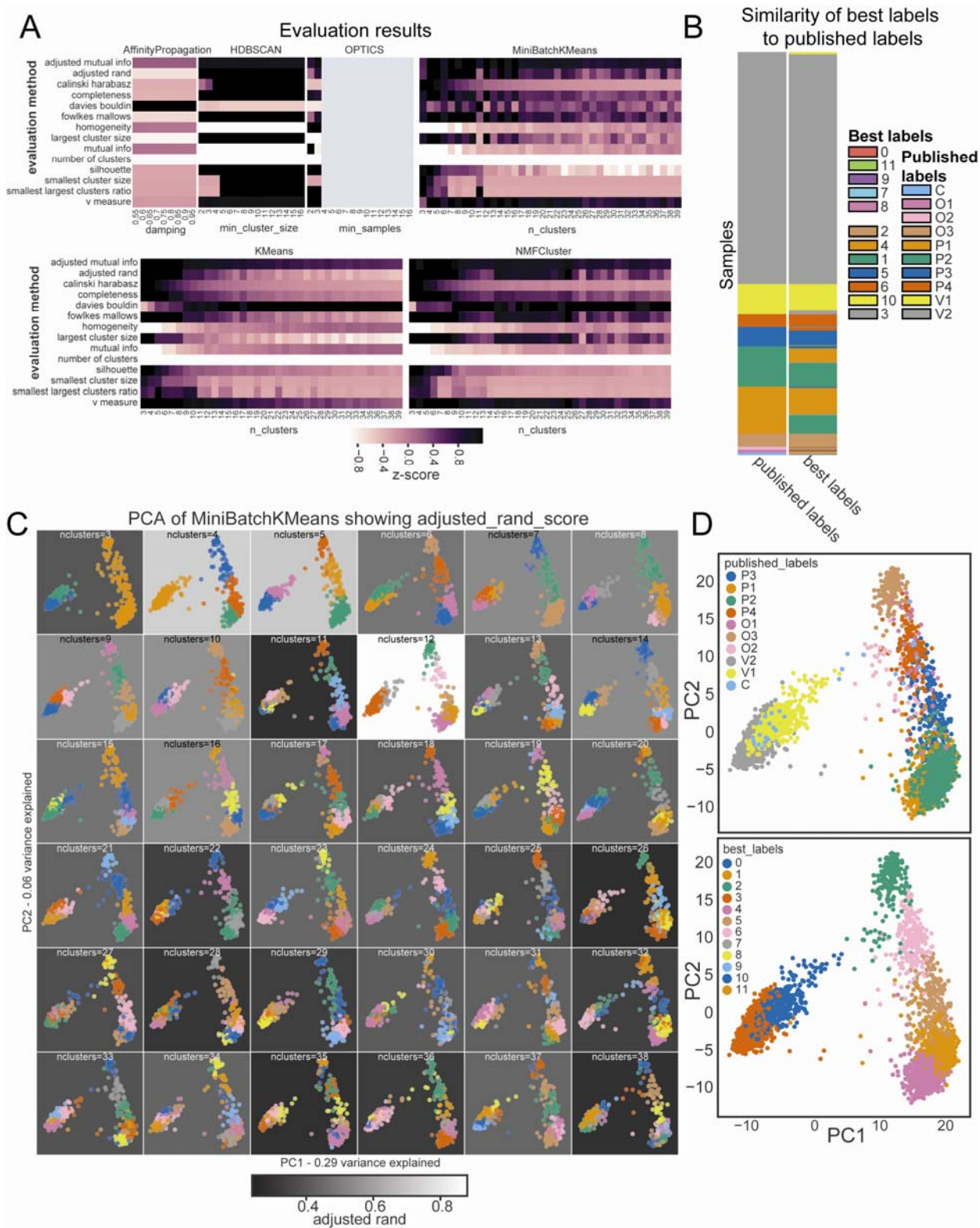
267 and ranges of relevant hyperparameters. Hypercluster was then used to evaluate labels

268 with every available metric, including metrics that measure inherent labeling quality, as

269 well as comparing new labels to cell types identified in the original study (Fig. 4A). The

270 approach that best recapitulated the published labels was clustering with

271 MiniBatchKMeans with 12 clusters (Fig. 4B-D). These labels differed from published  
272 labels largely from swapping cells in the P1 and P2 groups (Fig. 4B), which are both  
273 LEPR<sup>+</sup> subgroups, that were shown to be very similar in the original paper (38). While  
274 the original labels were generated using community detection methods like Louvain and  
275 Leiden clustering, those methods performed poorly compared to others (Figure S3),  
276 likely due to differences in data pre-processing. Varying the number of clusters has  
277 variable effects on evaluation results (Fig. 4A, 4C, Figure S3), again highlighting the  
278 importance of an exhaustive approach.



279

280 *Fig. 4 Clustering and evaluation of scRNA-seq data*

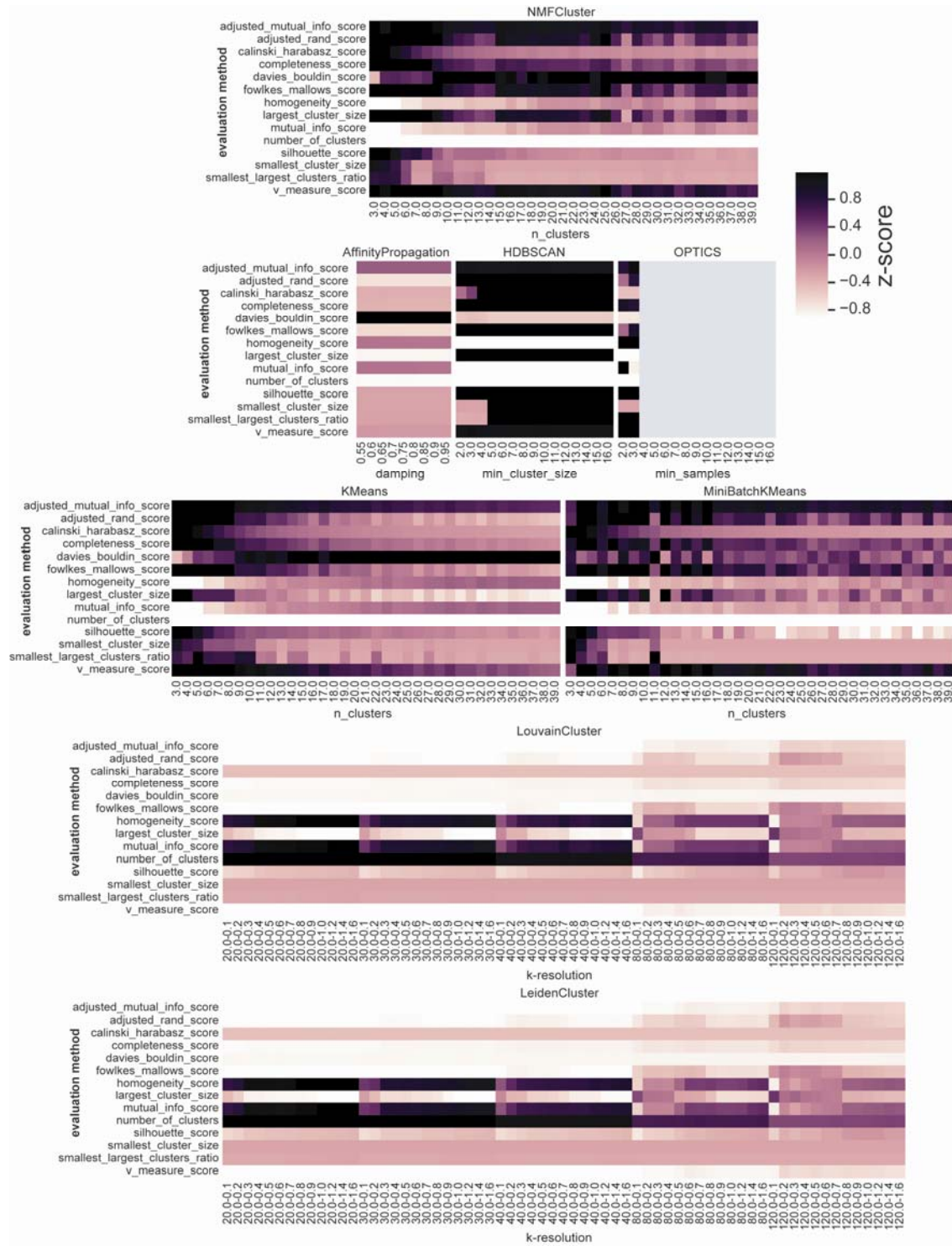
281 a) Evaluation metrics for clustering conditions, automatically generated by hypercluster  
282 for single cell RNA-seq data.

283 b) Comparison of published labels with best matching calculated labels,  
284 MiniBatchKMeans;n\_clusters=12. Legend shows mismatched clusters for the best  
285 labels on top and clusters with high correspondence to published clusters in the bottom  
286 section.

287 c) PCA projection of 700 random cells labeled by MiniBatchKMeans across  
288 hyperparameters.

289 d) PCA projection of cells colored by published labels.





290

291 *Figure S3 Full evaluations of scRNA-seq clustering*

292 Automatically generated full evaluation metric table from clustering of scRNA-seq stem

293 cell niche cells.

294

295 Discussion

296 Defining groups of molecularly similar patient samples is key to personalizing  
297 medical prognosis, diagnosis and treatment strategies, making unsupervised clustering  
298 a workhorse for researchers advancing personalized medicine. It is therefore essential  
299 that unsupervised clustering is rigorous and not biased by arbitrary hyperparameter  
300 selection. While extremely high quality open-source tools such as scikit-learn make  
301 unsupervised clustering accessible to many, exhaustively and reproducibly comparing  
302 hyperparameters is still challenging; hypercluster solves these issues.

303 Nearly every step in data analysis pipelines require hyperparameter selection,  
304 during which biased or arbitrary parameter selection can greatly impact results. Further,  
305 data preprocessing, involving the filtering of datasets to remove low quality or low  
306 coverage samples or features (e.g. removing genes with very few reads in RNA-seq),  
307 also greatly impacts downstream clustering results. Hypercluster provides a workflow to  
308 address the former issue, allowing for comprehensive evaluation of multiple  
309 hyperparameters and clustering algorithms simultaneously. The package auto-sklearn  
310 (14) provides functionality for automating pre-processing of data tables, which could  
311 easily be incorporated upstream of hypercluster to automate the latter. In addition to the  
312 simple command line functions, we have also employed SnakeMake for parallelization,  
313 a workflow management system already widely used for pipeline optimization (40–46).

314 If unsupervised clustering is a downstream analytic method of interest,  
315 determining which parameters to select can be cumbersome, and possibly inaccurate,  
316 without a clustering optimization tool like hypercluster. While it is not always clear how

317 to choose hyperparameters or algorithms in a consistent way (e.g. when two different  
318 conditions optimize for different metrics), it is essential to at least understand if the  
319 labels one obtains are robust to small changes in algorithm or hyperparameter choice  
320 (e.g. as shown in Figure S1). Our package greatly improves the ability of researchers to  
321 gain this understanding. In addition to assisting researchers in choosing  
322 hyperparameters, hypercluster aids computational biologists who are benchmarking  
323 new clustering algorithms, evaluation metrics and pre- or post-processing steps (3). In  
324 conclusion, hypercluster streamlines the use of unsupervised clustering to derive  
325 biologically relevant structure within data. Most importantly, it eases the prioritization of  
326 rigor and reproducibility for researchers using these techniques.

327

#### 328 Acknowledgements

329 We thank the members of Ruggles and Fenyö labs for their helpful discussions and  
330 input. We would like to thank MacIntosh Cornwell for his advice with the SnakeMake  
331 pipeline. We would also like to thank Joseph Copper Devlin for his help and advice with  
332 implementing Louvain and Leiden clustering.

#### 333 Availability of data and materials

334 Hypercluster is released on pip (pip install hypercluster) and conda (conda install -c  
335 bioconda hypercluster). Development versions and installation instructions can be found  
336 at our github (<https://github.com/liliblu/hypercluster/>), tutorials and examples, including  
337 all of the code used to create the figures in this paper, can be found here:  
338 <https://github.com/ruggleslab/hypercluster/tree/master/examples>, and documentation  
339 can be found here: <https://hypercluster.readthedocs.io/en/latest/>. Hypercluster is written



340 in python and was developed and tested on MacOS and Linux. Requirements are listed  
341 on the github and in the documentation. Hypercluster is available with the MIT licence.

342

#### 343 Financial Disclosure

344 This work has been supported by the National Cancer Institute (NCI) through CPTAC  
345 award U24 CA210972. The funders had no role in study design, data collection and  
346 analysis, decision to publish, or preparation of the manuscript.

347

#### 348 Authors' contributions

349 Conceptualization, Project Administration, Writing: LB and KVR. Data Curation, Formal  
350 analysis, Investigation, Methodology, Software, Validation, Visualization: LB. Funding  
351 acquisition, Resources, Supervision: KVR.

352

#### 353 Competing interests

354 The authors declare no competing interests.

355

#### 356 Related manuscripts

357 The authors do not have other related or duplicate manuscripts.

358

#### 359 References

360 1. Kiselev VY, Andrews TS, Hemberg M. Publisher Correction: Challenges in  
361 unsupervised clustering of single-cell RNA-seq data. Nat Rev Genet. 2019  
362 May;20(5):310.

- 363 2. Sun S, Zhu J, Ma Y, Zhou X. Accuracy, robustness and scalability of dimensionality  
364 reduction methods for single-cell RNA-seq analysis [Internet]. Vol. 20, Genome  
365 Biology. 2019. Available from: <http://dx.doi.org/10.1186/s13059-019-1898-6>
- 366 3. Liu X, Song W, Wong BY, Zhang T, Yu S, Lin GN, et al. A comparison framework  
367 and guideline of clustering methods for mass cytometry data. *Genome Biol.* 2019  
368 Dec 23;20(1):297.
- 369 4. Parker JS, Mullins M, Cheang MCU, Leung S, Voduc D, Vickery T, et al.  
370 Supervised risk predictor of breast cancer based on intrinsic subtypes. *J Clin Oncol.*  
371 2009 Mar 10;27(8):1160–7.
- 372 5. Ohnstad HO, Borgen E, Falk RS, Lien TG, Aaserud M, Sveli MAT, et al. Prognostic  
373 value of PAM50 and risk of recurrence score in patients with early-stage breast  
374 cancer with long-term follow-up. *Breast Cancer Res.* 2017 Nov 14;19(1):120.
- 375 6. Ali HR, Rueda OM, Chin S-F, Curtis C, Dunning MJ, Aparicio SA, et al. Genome-  
376 driven integrated classification of breast cancer validated in over 7,500 samples.  
377 *Genome Biol.* 2014 Aug 28;15(8):431.
- 378 7. Perou CM, Sørlie T, Eisen MB, van de Rijn M, Jeffrey SS, Rees CA, et al. Molecular  
379 portraits of human breast tumours. *Nature.* 2000 Aug 17;406(6797):747–52.
- 380 8. Capper D, Jones DTW, Sill M, Hovestadt V, Schrimpf D, Sturm D, et al. DNA  
381 methylation-based classification of central nervous system tumours. *Nature.* 2018  
382 Mar 22;555(7697):469–74.
- 383 9. Sturm D, Orr BA, Toprak UH, Hovestadt V, Jones DTW, Capper D, et al. New Brain  
384 Tumor Entities Emerge from Molecular Classification of CNS-PNETs. *Cell.* 2016  
385 Feb 25;164(5):1060–72.

- 386 10. Hoadley KA, Yau C, Hinoue T, Wolf DM, Lazar AJ, Drill E, et al. Cell-of-Origin  
387 Patterns Dominate the Molecular Classification of 10,000 Tumors from 33 Types of  
388 Cancer. *Cell*. 2018 Apr 5;173(2):291–304.e6.
- 389 11. Aure MR, Vitelli V, Jernström S, Kumar S, Krohn M, Due EU, et al. Integrative  
390 clustering reveals a novel split in the luminal A subtype of breast cancer with impact  
391 on outcome. *Breast Cancer Res*. 2017 Mar 29;19(1):44.
- 392 12. Curtis C, Shah SP, Chin S-F, Turashvili G, Rueda OM, Dunning MJ, et al. The  
393 genomic and transcriptomic architecture of 2,000 breast tumours reveals novel  
394 subgroups. *Nature*. 2012 Apr 18;486(7403):346–52.
- 395 13. Jaskowiak PA, Costa IG, Campello RJGB. Clustering of RNA-Seq samples:  
396 Comparison study on cancer data. *Methods*. 2018 Jan 1;132:42–9.
- 397 14. Feurer M, Klein A, Eggensperger K, Springenberg J, Blum M, Hutter F. Efficient and  
398 Robust Automated Machine Learning. In: Cortes C, Lawrence ND, Lee DD,  
399 Sugiyama M, Garnett R, editors. *Advances in Neural Information Processing*  
400 *Systems 28*. Curran Associates, Inc.; 2015. p. 2962–70.
- 401 15. Barber RF, Ha W. Gradient descent with non-convex constraints: local concavity  
402 determines convergence. *Inf Inference*. 2018 Dec 11;7(4):755–806.
- 403 16. Van Craenendonck T, Blockeel H. Using internal validity measures to compare  
404 clustering algorithms. *Benelearn 2015 Poster presentations (online)*. 2015;1–8.
- 405 17. Köster J, Rahmann S. Snakemake--a scalable bioinformatics workflow engine.  
406 *Bioinformatics*. 2012 Oct 1;28(19):2520–2.
- 407 18. Cluster and cloud execution — Snakemake 5.9.1+0.g138720f.dirty documentation  
408 [Internet]. [cited 2020 Jan 5]. Available from:

- 409 <https://snakemake.readthedocs.io/en/stable/executing/cluster-cloud.html>
- 410 19. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-  
411 learn: Machine Learning in Python. *J Mach Learn Res.* 2011;12(Oct):2825–30.
- 412 20. Kim H, Park H. Sparse non-negative matrix factorizations via alternating non-  
413 negativity-constrained least squares for microarray data analysis [Internet]. Vol. 23,  
414 *Bioinformatics.* 2007. p. 1495–502. Available from:  
415 <http://dx.doi.org/10.1093/bioinformatics/btm134>
- 416 21. Traag VA, Waltman L, van Eck NJ. From Louvain to Leiden: guaranteeing well-  
417 connected communities. *Sci Rep.* 2019 Mar 26;9(1):5233.
- 418 22. Traag VA, Krings G, Van Dooren P. Significant scales in community structure. *Sci*  
419 *Rep.* 2013 Oct 14;3:2930.
- 420 23. Csardi G, Nepusz T, Others. The igraph software package for complex network  
421 research. *InterJournal, complex systems.* 2006;1695(5):1–9.
- 422 24. Traag V. *leidenalg* [Internet]. Github; [cited 2020 Jan 27]. Available from:  
423 <https://github.com/vtraag/leidenalg>
- 424 25. Traag V. *louvain-igraph* [Internet]. Github; [cited 2020 Jan 27]. Available from:  
425 <https://github.com/vtraag/louvain-igraph>
- 426 26. McKinney W, Others. Data structures for statistical computing in python. In:  
427 *Proceedings of the 9th Python in Science Conference.* Austin, TX; 2010. p. 51–6.
- 428 27. Walt S van der, Colbert SC, Varoquaux G. The NumPy Array: A Structure for  
429 Efficient Numerical Computation. *Comput Sci Eng.* 2011 Mar 1;13(2):22–30.
- 430 28. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al.  
431 *SciPy 1.0--Fundamental Algorithms for Scientific Computing in Python* [Internet].

432 arXiv [cs.MS]. 2019. Available from: <http://arxiv.org/abs/1907.10121>

433 29. Hunter JD. Matplotlib: A 2D Graphics Environment. *Comput Sci Eng*. 2007 May  
434 1;9(3):90–5.

435 30. Waskom M, Botvinnik O, O’Kane D, Hobson P, Lukauskas S, Gemperline DC, et al.  
436 mwaskom/seaborn: v0.8.1 (September 2017) [Internet]. 2017. Available from:  
437 <https://zenodo.org/record/883859>

438 31. 2.3. Clustering — scikit-learn 0.22 documentation [Internet]. [cited 2019 Dec 23].  
439 Available from: <https://scikit-learn.org/stable/modules/clustering.html>

440 32. Cancer Genome Atlas Network. Comprehensive molecular portraits of human  
441 breast tumours. *Nature*. 2012 Oct 4;490(7418):61–70.

442 33. Chalise P, Fridley BL. Integrative clustering of multi-level ’omic data based on non-  
443 negative matrix factorization algorithm. *PLoS One*. 2017 May 1;12(5):e0176278.

444 34. Sorlie T, Tibshirani R, Parker J, Hastie T, Marron JS, Nobel A, et al. Repeated  
445 observation of breast tumor subtypes in independent gene expression data sets.  
446 *Proc Natl Acad Sci U S A*. 2003 Jul 8;100(14):8418–23.

447 35. Sørliie T, Perou CM, Tibshirani R, Aas T, Geisler S, Johnsen H, et al. Gene  
448 expression patterns of breast carcinomas distinguish tumor subclasses with clinical  
449 implications. *Proc Natl Acad Sci U S A*. 2001 Sep 11;98(19):10869–74.

450 36. Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM 3rd, et al.  
451 Comprehensive Integration of Single-Cell Data. *Cell*. 2019 Jun 13;177(7):1888–  
452 902.e21.

453 37. Tang J, Ceng X, Peng B. New Methods of Data Clustering and Classification Based  
454 on NMF [Internet]. 2011 International Conference on Business Computing and

- 455 Global Informatization. 2011. Available from:  
456 <http://dx.doi.org/10.1109/bcgin.2011.114>
- 457 38. Tikhonova AN, Dolgalev I, Hu H, Sivaraj KK, Hoxha E, Cuesta-Domínguez Á, et al.  
458 The bone marrow microenvironment at single-cell resolution. *Nature*. 2019  
459 May;569(7755):222–8.
- 460 39. Yoo AB, Jette MA, Grondona M. SLURM: Simple Linux Utility for Resource  
461 Management. In: *Job Scheduling Strategies for Parallel Processing*. Springer Berlin  
462 Heidelberg; 2003. p. 44–60.
- 463 40. Wang D. hppRNA—a Snakemake-based handy parameter-free pipeline for RNA-  
464 Seq analysis of numerous samples. *Brief Bioinform*. 2018 Jul 20;19(4):622–6.
- 465 41. Pranzatelli TJF, Michael DG, Chiorini JA. ATAC2GRN: optimized ATAC-seq and  
466 DNase1-seq pipelines for rapid and accurate genome regulatory network inference.  
467 *BMC Genomics*. 2018 Jul 31;19(1):563.
- 468 42. Abdelaal T, Michielsen L, Cats D, Hoogduin D, Mei H, Reinders MJT, et al. A  
469 comparison of automatic cell identification methods for single-cell RNA sequencing  
470 data [Internet]. Vol. 20, *Genome Biology*. 2019. Available from:  
471 <http://dx.doi.org/10.1186/s13059-019-1795-z>
- 472 43. Dirmeier S, Emmenlauer M, Dehio C, Beerenwinkel N. PyBDA: a command line tool  
473 for automated analysis of big biological data sets. *BMC Bioinformatics*. 2019 Nov  
474 12;20(1):564.
- 475 44. single-cell-rna-seq [Internet]. Github; [cited 2020 Jan 8]. Available from:  
476 <https://github.com/snakemake-workflows/single-cell-rna-seq>
- 477 45. Lun ATL, McCarthy DJ, Marioni JC. A step-by-step workflow for low-level analysis

478 of single-cell RNA-seq data with Bioconductor [Internet]. Vol. 5, F1000Research.  
479 2016. p. 2122. Available from: <http://dx.doi.org/10.12688/f1000research.9501.2>  
480 46. Sonesson C, Robinson MD. Bias, robustness and scalability in single-cell differential  
481 expression analysis. Nat Methods. 2018 Apr;15(4):255–61.