

SOFTWARE

Granatum: a graphical single-cell RNA-seq analysis pipeline for genomics scientists

1 Xun Zhu^{1,2}, Thomas Wolfgruber^{1,2}, Austin Tasato³, Lana X Garmire^{1, 2*}

2 _____

3 *Correspondence:

4 LGarmire@cc.hawaii.edu

5 1Graduate Program in Molecular
6 Biology and Bioengineering,
7 University of Hawaii at Manoa,
8 Honolulu, HI 96816

9 2 Epidemiology Program, University
10 of Hawaii Cancer Center, Honolulu,
11 HI 96813

12 3 Department of Electrical
13 Engineering, University of Hawaii at
14 Manoa, Honolulu, HI 96816

15

Abstract

Background: Single-cell RNA sequencing (scRNA-seq) is an increasingly popular platform to study heterogeneity at the single cell level. Computational methods to process scRNA-seq have limited accessibility to bench scientists, as they require significant amount of bioinformatics skills.

Results: We have developed Granatum, a web browser based scRNA-seq analysis pipeline to make analysis more broadly accessible to researchers. Without a single line of programming code, a user can click through the pipeline, setting parameters and visualizing results via the interactive graphical interface. The pipeline conveniently walks the users through various steps of scRNA-seq analysis. It has a comprehensive list of modules, including plate merging and batch effect removal, outlier sample removal, gene filtering, gene expression normalization, cell clustering, differential gene expression analysis, pathway/ontology enrichment analysis, protein network interaction visualization, and pseudo-time cell series construction.

Conclusions: Granatum enables much widely adoption of scRNA-seq technology by empowering the bench scientists with an easy to use graphical interface for scRNA-seq data analysis. The package is freely available for research use at: <http://garmiregroup.org/granatum/app>

Keywords: single-cell; gene expression; graphical; normalization; clustering; differential expression; pathway; pseudo-time; software

16 **Background**

17 The arrival of single-cell high-throughput RNA sequencing (scRNA-seq) has provided new
18 opportunities for researchers to identify the expression characteristics of individual cells among
19 complex tissues. This is a significant leap forward from bulk cell RNA expression analysis. In cancer,
20 for example, scRNA-seq allows tumorous cells to be separated apart from healthy cells [1] and
21 primary cells be differentiated from metastatic cells [2]. Single-cell expression data can also be
22 used to describe trajectories of cell differentiation and development [3]. However, analyzing data
23 from scRNA-seq brings new computational challenges, e.g., accounting for inherently high drop-
24 out (artificial loss of RNA expression information) [4].

25 Software that has been developed to address these challenges may have very limited accessibility
26 for biologists with only general computer skills, as they typically require the ability to use a
27 computing language like R [5,6]. Other existing workflows that can be used to analyze scRNA-seq
28 data, such as Singular (Fluidigm, Inc., South San Francisco, CA, USA), Cell Ranger/ Loupe
29 (Pleasanton, CA, USA), and Scater [7] all require some non-graphical interactions and they may not
30 provide a comprehensive set of scRNA-seq analysis methods. To fill this gap, we have developed
31 Granatum, a fully interactive graphical scRNA-seq analysis tool. Granatum is the Latin word for
32 pomegranate, which bears many seeds, resembling single cells within the entity. This tool employs
33 an easy-to-use web browser interface for a wide range of methods suitable for scRNA-seq analysis:
34 removal of batch effects, removal of outlier cells, normalization of expression levels, filtering of
35 under-informative genes, clustering of cells, identification of differentially expressed genes,

36 identification of enriched pathways/ontologies, visualization of protein networks, and
37 reconstruction of pseudo-time paths for cells. Our software will empower a much broader
38 audience of research communities to study single cell complexity, by allowing them to readily
39 explore single-cell expression data from a graphical user interface.

40 **Implementation**

41 **Overview**

42 Both the front-end and the back-end of Granatum are written in the R software language, and built
43 with the Shiny framework [8]. Multiple concurrent users are handled by Shiny and each user works
44 on its own data space. To protect the privacy of users, the data submitted by one user is not visible
45 to any other user. The front-end is implemented as a web page with dynamically loaded pages,
46 and is arranged in a step-wise fashion. The default theme uses the Bootstrap framework. ShinyJS
47 [9] is used to power some of the interactive components. To allow users to redo a task, each
48 processing step is equipped with a reset button.

49 **Interactive widgets**

50 The package visNetwork is used for the layout and physics simulation of the network modules [10].
51 DataTables are used to preview user submitted data and to show tabular data in various modules
52 [11]. Plotly is used for the interactive outlier identification step [12]. The package ggplot2 is used
53 for the scatter-plots and box-plots, which is also used by the Monocle package for the Pseudo-time
54 construction step [3,13].

55 **Back-end variable management**

56 The expression matrix and the metadata sheet are stored separately for each user. The metadata
57 sheet can refer to groups, batches, or other properties of the samples in the corresponding
58 expression matrix. These two types of tables are shared across all modules. Other variables shared
59 across all modules include the log-transformed expression matrix, the filtered and normalized
60 expression matrix, the dimensionally reduced matrix, species (human or mouse) and the primary
61 metadata column.

62 **Batch-effect removal**

63 Batch-effect removal is done using the following procedure. First, we calculate the median
64 expression of each sample, denoted as med_i for sample i . Second, we calculate the mean of med_i
65 for each batch, denoted as $batchMean_b$ for batch b ,

$$batchMean_b = geometricMean_{i \in batch_b}(med_i).$$

66 Finally, each batch will be multiplied by a factor which pulls towards the global geometric mean of
67 the sample medians, i.e., when $i \in batch_b$ and m is the number of samples,

$$sampleNew_i = sampleOld_i \cdot \frac{geometricMean_{i \in 1, \dots, m}(med_i)}{batchMean_b}.$$

68 Where $sampleNew_i$ and $sampleOld_i$ denote the expression levels (vector) for all genes within
69 sample i before (old) and after (new) batch-effect removal.

70 Clustering methods

71 The following description of clustering algorithms assumes n being the number of genes, m being
72 the number of samples, and k being the number of clusters.

73 **Non-negative matrix factorization (NMF):** the log-transformed expression matrix (n -by- m) is
74 factorized into two non-negative matrices H (n -by- k) and W (k -by- m) with k being the expected
75 number of clusters. The latter matrix is then used to determine the membership of each cluster by
76 determining, for each column in W , which of the k entries has the highest value [14,15]. The NMF
77 computation is implemented in the NMF R-package, as reported earlier [14,16].

78 **K-means:** K-means is done on either the log-transformed expression matrix or the 2-by- m
79 correlation t-SNE matrix. The algorithm is implemented by the *kmeans* function in R [17].

80 **Hierarchical clustering (Hclust):** Hclust is also done on either the log-transformed expression
81 matrix or the 2-by- m correlation t-SNE matrix. The algorithm is implemented by the *hclust* function
82 in R [18]. The heatmap with dendrograms is plotted using the *heatmap* function in R.

83 Correlation t-SNE

84 Correlation t-SNE is implemented to assess heterogeneity of the data. It is calculated using a two-
85 step process. First, a distance matrix is calculated using the correlation distance. The correlation
86 distance $D_{i,j}$ between sample i and sample j is defined as

$$D_{i,j} = 1 - \text{Correlation}(S_i, S_j),$$

87 where S_i and S_j are the i -th and j -th column (sample) of the expression matrix.

88 Next, t-SNE is performed using this distance matrix, which reduces the expression matrix to two
89 dimensions. We use the Rtsne R package for this calculation [19].

90 **Elbow-point finding algorithm in clustering**

91 In the clustering module with automatic determination of the number of clusters, the
92 identification of the optimum number of clusters is done prior to presenting the clustering results.
93 First, we calculate the k-means clusters from $k = 2$ to $k = 10$. For each k , we calculate the
94 percentage of the explained variance (EV). To find the elbow-point $k = m$ where the EV plateaus,
95 we fit the k -EV data points with a linear elbow function. This function consists of a linearly
96 increasing piece from 0 to m , and a constant piece from m to 10. We iterate from $m = 1$ to 10 and
97 identify m which gives the best coefficient of determination (R^2) of linear regression as the "elbow
98 point".

99 **Differential expression analysis**

100 We use SCDE (version 1.99.4) in our Differential expression (DE) analysis step. The minimum size
101 entries parameter of the *scde.error.models* function is set to be the lesser of 2000 or the number
102 of genes after filtering [20]. When more than two clusters are present, a pair-wise DE analysis is
103 performed.

104 **Gene-set enrichment analysis**

105 The GSEA algorithm is implemented in the *fgsea* R-package which uses an optimized algorithm for
106 fast calculation speed [21].

107 **Pseudo-time construction**

108 We use Monocle (version 2.2.0) in our pseudo-time construction step. When building the
109 *CellDataSet* required for monocle's input, we set the *expressionFamily* to *negbinomial.size()*. The
110 dimension reduction is done using the *reduceDimension* function with *max_components* set to be
111 2.

112 **Results**

113 **Overview and comparison with scRNA-seq pipelines**

114 Granatum is by far the most comprehensive tool web browser based scRNA-seq analysis pipeline
115 without any programming requirement (Table 1). We have systematically compared Granatum
116 with 12 other existing tools, to demonstrate its versatile functions. Among other tools, methods
117 such as SCDE / PAGODA and Flotilla, are developed for programmers and requires expertise in a
118 particular programming language. In contrast, Granatum with its simple graphical interface
119 requires no programming knowledge, and is very easy to navigate through. Current version of
120 Granatum neatly presents nine modules, arranged as steps and ordered by their dependency
121 (Figure 1). It starts with one or more user-supplied expression matrices and corresponding sample
122 metadata sheet(s), followed by data merging, batch-effect removal, outlier removal,
123 normalization, gene filtering, clustering, differential expression, protein-protein network, and
124 pseudo-time construction.

125 Comparing to other freely available tools, Granatum workflow has many superior functionalities
126 that make it flexible (Table 1). Below we enlist some of them. (1) Unlike tools such as SCRAT
127 (<https://zhiji.shinyapps.io/scrat/>), ASAP [22] and Sake (<http://sake.mhammell.tools/>), it is the only
128 GUI pipeline that supports multiple dataset submission as well as batch effect removal; (2) at any
129 point of the step, the user can reset the current step for re-analysis; (3) the user can bypass certain
130 steps and still complete the workflow; (4) the user can select subsets of samples/data for their
131 customized analysis need; (5) the user can identify outlier samples either automatically by a pre-

132 set threshold, or manually by simply clicking the samples the PCA plot or the correlation t-SNE
133 plot; (6) the user can specify multiple cores in the differential expression module for speed-up; (7)
134 Both GSEA and network analysis can be performed for the differentially expressed genes in all
135 pairs of subgroups, following clustering analysis; (8) Monocle pseudo-time construction can be
136 performed to gain insights of relationships between the cells. In the following sections, we
137 elaborate the details of each step in Granatum in chronological order.

138 **Upload data**

139 Granatum accepts one or multiple expression matrices as the input. Each expression matrix can be
140 accompanied by a table describing the groups, batches, or other properties of the samples in the
141 corresponding matrix. This accompanying table is called the metadata sheet. Multiple matrices
142 may be uploaded sequentially. The user also specifies the species of the data, either human or
143 mouse, for downstream functional analysis. After the input files are uploaded, preview tables for
144 the matrix and metadata are displayed, providing the user an opportunity check that the data they
145 have input is as expected.

146 **Batch-effect removal**

147 Samples obtained in batches can create unwanted technical variation, which confound the
148 biological variation [23]. It is thus important to remove the expression level difference due to
149 batches. Granatum provides a batch-effect removal step, where the batches are shown as
150 different colors in the box-plot (Figure 2). If more than one datasets are uploaded, by default each
151 dataset is assumed to be one batch. Alternatively, if the batch numbers are indicated in the sample

152 metadata sheet, the user may select the column in which the batch numbers are stored (blue
153 circled in Figure 2). For datasets with a large number of cells, to maintain legibility of the box-plot a
154 random selection of 96 sub-samples is shown in the box-plot, and can be re-sampled freely.

155 **Outlier identification**

156 Computationally abnormal samples pose serious problems for many down-stream analysis
157 procedures. It is thus crucial to identify and remove them in the early stage. Granatum's outlier
158 identification step features PCA plot and t-SNE plot, two connected interactive scatter-plots that
159 have different computational characteristics. A PCA plot illustrates the Euclidean distance between
160 the samples, and a correlation t-SNE plot shows the associative distances between the samples.
161 The interactive mode of these plots is realized by the Plotly library [12] (Figure 3A).

162 Outliers can be identified automatically by either using a z score threshold or setting a fixed
163 number of outliers. In addition, the user can select or de-select each sample, by clicking, boxing or
164 drawing a lasso on its corresponding points on either PCA or t-SNE plot (Figure 3A and 3B). This
165 level of interaction from users is one of the many examples of thoughtful tool design, in order to
166 empower them.

167 To help users select sample of a particular property, Granatum also allows for mapping any of the
168 columns in the metadata sheet onto the scatter-plots (circled blue in Figure 3A). The complete
169 metadata information of the selected samples can be found in a table at the bottom of the page
170 (circled red in Figure 3A).

171 **Normalization**

172 Normalization is essential to most scRNA-seq data, except those with the UMI counts, before the
173 down-stream functional analyses. The current version of Granatum has implemented three
174 commonly used normalization algorithms: rescale to geometric mean, quantile normalization, and
175 size-factor normalization [24,25]. A box-plot is shown post normalization, to help illustrate its
176 effect to the median, mean, and extreme values across samples. As is the case in the batch-effect
177 removal step, for a dataset with a large number of samples, 96 sub-samples are randomly chosen
178 for the visualization purpose (Figure 3C).

179 **Gene filtering**

180 Due to scRNA-seq's relative high level of noise, it has been recommended to remove lowly
181 expressed genes as well as lowly dispersed genes [4]. To this end, Granatum has a step to remove
182 these genes. The user can interactively select both the average expression level threshold and the
183 dispersion threshold (Figure 3D). The dispersion calculation and negative binomial model fitting
184 are calculated by modifying the output of the Monocle package [3]. We have customized the
185 visualization code to enhance integration with the other components, by setting up the threshold
186 selection sliders and number of genes statistics message on the Granatum web page (Figure 3D).
187 On the mean-dispersion plot, each gene is represented by a point, where the x-axis is the mean of
188 the expression levels after log transformation, and the y-axis is the dispersion factor calculated
189 from a negative binomial model. The preserved genes are highlighted as black and the genes to be

190 removed are labeled as gray colors. The number of genes before and after filtering are also
191 displayed.

192 **Clustering**

193 Clustering is a routine heuristic analysis for scRNA-seq data. Granatum selects five commonly used
194 algorithms: non-negative matrix factorization [14], k-means, k-means combined with correlation t-
195 SNE, hierarchical clustering (hclust), and hclust combined with correlation t-SNE. The number of
196 clusters may be set manually, or automatically determined using an elbow-point finding algorithm
197 (Methods, Figure 4A). For the latter approach, the algorithm will attempt to cluster samples with
198 number of clusters (k) ranging from 2 to 10, and determine the best number by finding the elbow-
199 point k . k indicates the starting point of plateau for explained variance (EV), above which EV
200 creases only minimally. If hclust is selected, a heatmap with hierarchical grouping and
201 dendrograms be shown in a pop-up window (Figure 4B).

202 Next, the resulting cluster labels obtained above, are then super-imposed onto the two
203 unsupervised PCA and correlation t-SNE plots (Figure 4A). The user can also represent user-defined
204 labels in the sample metadata as different colors in these plots. By comparing the two sets of
205 labels, the users can quickly check the concordance between the prior metadata labels and the
206 computed clusters.

207 **Differential expression**

208 After obtaining a set of clusters, it is intuitively important to identify genes that are differentially
209 expressed between any two clusters. Granatum uses the state-of-the-art SCDE method for its

210 single-cell DE analysis [20]. The DE comparison is performed in a pair-wise fashion when more than
211 two clusters are present. This step is computationally time and memory consuming. To shorten
212 computation time, a user can select the number of cores for parallelization on multi-core machines
213 (Figure 5A). When SCDE is completed, tabbed tables show the genes sorted by their Z-scores, along
214 with the model coefficients (Figure 5B). As another feature to empower the users, the gene
215 symbols are linked to their corresponding GeneCards pages (www.genecards.org) [26]. The DE
216 results can be downloaded as a CSV file via the "Download CSV table" button.

217 To investigate the collective biological functions of these genes, the user can further perform Gene
218 Set Enrichment Analysis (GSEA) with either KEGG pathways or Gene Ontology (GO) terms (circled
219 blue in Figure 5B) [27–30]. We have employed a very intuitive bubble-plot to visualize the GSEA
220 results, where the vertical position of the bubble indicates the enrichment score of the gene sets,
221 and the size of the bubble indicates number of genes in that set (KEGG pathway or GO term)
222 (Figure 5C).

223 **Protein network visualization**

224 Protein-protein interaction (PPI) network gives straightforward and systematic understanding of
225 the connections between these differentially expressed genes. Granatum selects the top K (default
226 K=200) genes in the DE results, and super impose the PPI network on them. Genes that are not
227 connected to any other genes in the list are removed from the PPI network. We use visNetwork to
228 enable the interactive display of the graph [10]. The user can freely rearrange the graph by
229 dragging the nodes to the desired location, and reconfiguring the layout to achieve good visibility

230 of the modules (via elastic-spring physics simulation) (Figure 6A). In this interactive graph, the Z-
231 scores are mapped as colors on the nodes where red indicates up-regulation and blue indicates
232 down-regulation.

233 **Pseudo-time construction**

234 Granatum has included the Monocle algorithm, a widely-used method to reconstruct a pseudo-
235 timeline for the samples [3]. Monocle uses the Reversed Graph Embedding algorithm to learn the
236 structure of the data, and the Principal Graph algorithm to find the time-lines and branching points
237 of the samples. We superimpose the timeline on the samples scatter-plot projected on the two
238 components of the learned projection matrix. The user may map any pre-defined labels or numeric
239 assays provided in the metadata sheet on to the scatter-plot (Figure 6B). The plotting functions are
240 adapted from the visualization code in Monocle.

241 **Discussion**

242 The field of scRNA-seq is fast-evolving both in terms of the development of instrumentation and
243 the innovation of computational methods. However, it becomes exceedingly hard for a wet-lab
244 researcher without formal bioinformatics training to catch up with the latest iterations of
245 algorithms [5]. This poses major barriers to them and many resort to sending their generated data
246 to third-party bioinformaticians, before they are able to visualize the data themselves. This
247 segregation often prolongs the research cycle time, as it often takes significant effort to maintain
248 effective communications between the two sides (sometimes even more complicated with a third

249 party of the genomics core). Also, issues with the experimentations do not get the chance to be
250 spotted early enough, to avoid significance loss of time and cost in the projects. It is thus very
251 attractive to have a non-programming graphical application which includes state-of-the-art
252 algorithms as routine procedures, in the hands of the bench-scientist who generate the scRNA-seq
253 data.

254 Granatum is our attempt to fill this void. It is to our knowledge the first solution that aims to cover
255 the entire scRNA-seq workflow with an intuitive, step-wise graphical user interface. Throughout
256 the development process our priority has been to make sure that it is fully accessible to
257 researchers with no programming experiments. We have strived to achieve that the plots and
258 tables are self-explanatory, interactive and visually pleasant. We have sought inputs from our
259 single-cell bench-side collaborators, to ensure that the terminologies are easy to understand by
260 them. We also supplement Granatum with a manual and video that guide the users through the
261 entire workflow, using example datasets. Currently Granatum targets users who have their
262 expression matrices and metadata sheets ready. However, we are developing the next version of
263 Granatum, which will handle the entire scRNA-seq data processing and analysis pipeline including
264 FASTQ quality control, alignment, and expression quantification. In the future, we will enrich
265 Granatum with capacities to analyze and integrate other types of genomics data in single cells,
266 such as exome-seq and methylation data.

267 **Conclusions**

268 We have developed a graphical web application called Granatum, which enables bench
269 researchers with no programming expertise to analyze state-of-the-art scRNA-Seq data. This tool
270 offers many interactive features to allow routine computational procedures with a great amount
271 of flexibility. We expect that this platform will empower the bench-side researchers with more
272 independence in the fast-evolving single cell genomics field.

273 **Figure legends**

274 **Figure 1: Granatum workflow.** Granatum is built with the Shiny framework, which supports both
275 front-end and the back-end. The user uploads one or more expression matrices with
276 corresponding metadata for samples. The back-end stores data separately for each individual user,
277 and invokes third-party libraries on demand.

278 **Figure 2: The batch-effect removal steps.** A box-plot is shown for the samples. The colors indicate
279 the batch labels, which can be selected using the batch factor selection box circled in blue. In cases
280 where more than 96 cells are present in the data, only a random sample of 96 cells are shown. The
281 user can re-sample the data by clicking the “Re-plot random 96 cells” button.

282 **Figure 3: The outlier removal, normalization and gene filtering steps.** A) The main interface of the
283 outlier removal step. The two scatter-plots are the PCA and correlation t-SNE plots, with colors
284 indicate the cell labels (box circled in blue). The metadata table (circled in red) shows the labels for

285 the selected cells. B) The pop-up window for automatic outlier detection options after the “auto-
286 identify” button is clicked. C) The normalization step. The box-plot shows the expression levels of
287 each cell in log-scale. In cases where more than 96 cells are present in the data, only a random
288 sample of 96 cells are shown. D) The Gene filtering step. The y-axis of the scatter-plot is the
289 empirical dispersion, estimated by a negative binomial model. The x-axis is the log mean
290 expression of each gene. The user can change the threshold by dragging the two sliders circled in
291 blue.

292 **Figure 4: The Clustering step.** A) Main interface. PCA and t-SNE plots are shown with colors
293 mapped to user-selected sample labels. After clustering, samples are marked with their assigned
294 cluster numbers. The user can either choose a specific number of clusters or let Granatum
295 compute the best number of clusters. B) When Hclust (Euclidean) is selected, a pop-up window will
296 show a heatmap of the expression matrix with dendrograms.

297 **Figure 5: The Differential expression (DE) step.** A) Before running DE, the user may select the
298 number of cores to use for speed. B) After DE, top differentially expressed genes for each pair of
299 clusters are shown. Gene Set Enrichment Analysis (GSEA) can be performed, using either KEGG
300 pathways or GO terms (circled in blue). C) The results of GSEA. The pathways on the x-axis are
301 sorted top 20 enriched gene sets. The height of the bubble indicates the absolute normalized
302 enrichment score, and the size of the bubble indicates the number of genes in the set.

303 **Figure 6: The Protein network and Pseudo-time construction steps.** A) The Protein network step.
304 The A tabbed panel shows the connected gene modules on the PPI network between each pair of

305 clusters. The color on each node (gene) indicates its Z-score in the differential expression test. Red
306 and blue colors indicates up- and down- regulation. B) The Pseudo-time construction step.
307 Monocle algorithm is customized to visualize the paths among individual cells. The user can
308 represent sample labels from the metadata as colors in the plot.

309 **Tables**

310 **Table 1: Comparison of existing single-cell analysis pipelines.**

311 **Declarations**

312 **Ethics approval and consent to participate**

313 Not Applicable.

314 **Consent for publication**

315 Not Applicable.

316 **Availability of data and material**

317 Granatum can be visited at: <http://garmiregroup.org/granatum/app>

318 A demonstration video can be found at: <http://garmiregroup.org/granatum/video>

319 **Competing interests**

320 The authors declared no conflict of interest.

321 **Funding**

322 This research was supported by grants K01ES025434 awarded by NIEHS through funds provided by
323 the trans-NIH Big Data to Knowledge (BD2K) initiative (<http://datascience.nih.gov/bd2k>), P20
324 COBRE GM103457 awarded by NIH/NIGMS, NICHD R01HD084633 and NLM R01LM012373 and
325 Hawaii Community Foundation Medical Research Grant 14ADVC-64566 to LX Garmire.

326 **Authors' contributions**

327 LXG envisioned the project. XZ developed the majority of the pipeline. TW and AT assisted in
328 developing the pipeline. TW documented the user manual and performed packaging. XZ, TW and
329 LXG wrote the manuscript. All authors have read, revised, and approved the final manuscript.

330 **Acknowledgements**

331 We thank Drs. Michael Ortega and Paula Benny for providing valuable feedback during testing the
332 tool. We also thank other group members in Garmire group for suggestions in the tool
333 development.

334 **List of abbreviations**

335 **scRNA-seq:** Single-cell high-throughput RNA sequencing

336 **DE:** differential expression

337 **GSEA:** Gene-set enrichment analysis

338 **KEGG:** Kyoto Encyclopedia of Genes and Genomes

- 339 **GO:** Gene ontology
- 340 **PCA:** Principal component analysis
- 341 **t-SNE:** t-Distributed Stochastic Neighbor Embedding
- 342 **NMF:** Non-negative matrix factorization
- 343 **Hclust:** Hierarchical clustering
- 344 **PPI:** Protein-protein interaction

345 **References**

- 346 1. Patel AP, Tirosh I, Trombetta JJ, Shalek AK, Gillespie SM, Wakimoto H, et al. Single-cell RNA-seq
347 highlights intratumoral heterogeneity in primary glioblastoma. *Science* (80-.). [Internet].
348 Department of Neurosurgery, Massachusetts General Hospital and Harvard Medical School,
349 Boston, MA 02114, USA. Department of Pathology and Center for Cancer Research, Massachusetts
350 General Hospital and Harvard Medical School, Boston, MA 02114, USA. *Broad I*; 2014;344:1396–
351 401. Available from: <http://dx.doi.org/10.1126/science.1254257>
- 352 2. Lewis BP, Burge CB, Bartel DP. Conserved seed pairing, often flanked by adenosines, indicates
353 that thousands of human genes are microRNA targets. *Cell*. Elsevier; 2005;120:15–20.
- 354 3. Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, et al. The dynamics and
355 regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat.*

- 356 Biotechnol. Nature Research; 2014;32:381–6.
- 357 4. Brennecke P, Anders S, Kim JK, Kołodziejczyk AA, Zhang X, Proserpio V, et al. Accounting for
358 technical noise in single-cell RNA-seq experiments. Nat. Methods. Nature Publishing Group; 2013;
- 359 5. Poirion OB, Zhu X, Ching T, Garmire L. Single-Cell Transcriptomics Bioinformatics and
360 Computational Challenges [Internet]. Front. Genet. . 2016. p. 163. Available from:
361 <http://journal.frontiersin.org/article/10.3389/fgene.2016.00163>
- 362 6. Team RC. R: A language and environment for statistical computing. R Foundation for Statistical
363 Computing, Vienna, Austria. 2015, URL [http. www. R-project. org](http://www.R-project.org). 2016;
- 364 7. McCarthy DJ, Campbell KR, Lun ATL, Wills QF. scater: pre-processing, quality control,
365 normalisation and visualisation of single-cell RNA-seq data in R. bioRxiv [Internet]. Cold Spring
366 Harbor Labs Journals; 2016; Available from: <http://biorxiv.org/content/early/2016/08/15/069633>
- 367 8. RStudio, Inc. Easy web applications in R. 2013.
- 368 9. Attali D. shinyjs: Easily Improve the User Experience of Your Shiny Apps in Seconds [Internet].
369 2016. Available from: <https://cran.r-project.org/package=shinyjs>
- 370 10. Almende B.V., Thieurmel B. visNetwork: Network Visualization using “vis.js” Library [Internet].
371 2016. Available from: <https://cran.r-project.org/package=visNetwork>
- 372 11. Xie Y. DT: A Wrapper of the JavaScript Library “DataTables” [Internet]. 2016. Available from:
373 <https://cran.r-project.org/package=DT>

- 374 12. Sievert C, Parmer C, Hocking T, Chamberlain S, Ram K, Corvellec M, et al. plotly: Create
375 Interactive Web Graphics via “plotly.js” [Internet]. 2016. Available from: [https://cran.r-](https://cran.r-project.org/package=plotly)
376 [project.org/package=plotly](https://cran.r-project.org/package=plotly)
- 377 13. Wickham H. ggplot2: Elegant Graphics for Data Analysis [Internet]. Springer-Verlag New York;
378 2009. Available from: <http://ggplot2.org>
- 379 14. Zhu X, Ching T, Pan X, Weissman S, Garmire L. Detecting heterogeneity in single-cell RNA-Seq
380 data by non-negative matrix factorization. PeerJ Prepr. PeerJ Inc. San Francisco, USA;
381 2016;4:e1839v1.
- 382 15. Brunet J-P, Tamayo P, Golub TR, Mesirov JP. Metagenes and molecular pattern discovery using
383 matrix factorization. Proc. Natl. Acad. Sci. [Internet]. 2004;101:4164–9. Available from:
384 <http://www.pnas.org/content/101/12/4164.abstract>
- 385 16. Gaujoux R, Seoighe C. Algorithms and framework for nonnegative matrix factorization (NMF).
386 2010.
- 387 17. Lloyd S. Least squares quantization in PCM. IEEE Trans. Inf. theory. IEEE; 1982;28:129–37.
- 388 18. Murtagh F, Contreras P. Methods of hierarchical clustering. arXiv Prepr. arXiv1105.0121. 2011;
- 389 19. Krijthe J. Rtsne: T-Distributed Stochastic Neighbor Embedding using Barnes-Hut
390 Implementation. R Packag. version 0.10, URL [http://CRAN.R-project.org/package= Rtsne](http://CRAN.R-project.org/package=Rtsne). 2015;
- 391 20. Kharchenko P V, Silberstein L, Scadden DT. Bayesian approach to single-cell differential

- 392 expression analysis. *Nat. Methods*. Nature Publishing Group; 2014;11:740–2.
- 393 21. Sergushichev A. An algorithm for fast preranked gene set enrichment analysis using cumulative
394 statistic calculation. *bioRxiv* [Internet]. Cold Spring Harbor Labs Journals; 2016; Available from:
395 <http://biorxiv.org/content/early/2016/06/20/060012>
- 396 22. Gardeux V, David F, Shajkofci A, Schwalie PC, Deplancke B. ASAP: a Web-based platform for the
397 analysis and inter-active visualization of single-cell RNA-seq data. *bioRxiv*. Cold Spring Harbor Labs
398 Journals; 2016;96222.
- 399 23. Hicks SC, Teng M, Irizarry RA. On the widespread and critical impact of systematic bias and
400 batch effects in single-cell RNA-Seq data. *bioRxiv*. Cold Spring Harbor Labs Journals; 2015;25528.
- 401 24. Bolstad BM, Irizarry RA, Åstrand M, Speed TP. A comparison of normalization methods for high
402 density oligonucleotide array data based on variance and bias. *Bioinformatics*. Oxford Univ Press;
403 2003;19:185–93.
- 404 25. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-Seq
405 data with DESeq2. *bioRxiv*. Cold Spring Harbor Labs Journals; 2014;
- 406 26. Rebhan M, Chalifa-Caspi V, Prilusky J, Lancet D. GeneCards: integrating information about
407 genes, proteins and diseases. *Trends Genet*. Elsevier Current Trends; 1997;13:163.
- 408 27. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set
409 enrichment analysis: a knowledge-based approach for interpreting genome-wide expression
410 profiles. *Proc. Natl. Acad. Sci. National Acad Sciences*; 2005;102:15545–50.

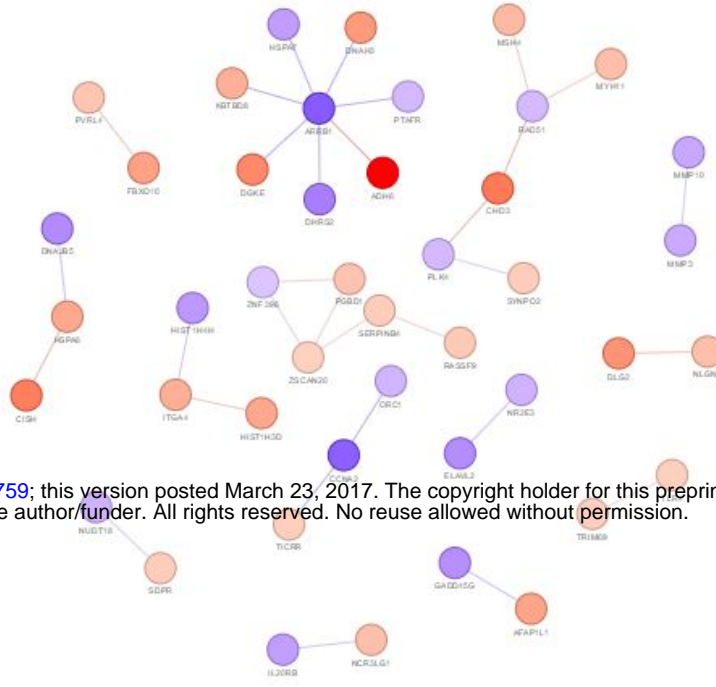
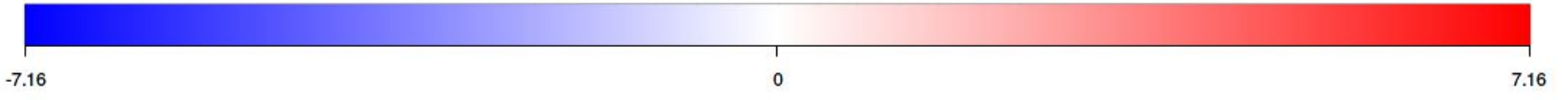
- 411 28. Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K. KEGG: new perspectives on
412 genomes, pathways, diseases and drugs. *Nucleic Acids Res. Oxford Univ Press*; 2017;45:D353--
413 D361.
- 414 29. Consortium GO, others. Gene ontology consortium: going forward. *Nucleic Acids Res. Oxford*
415 *Univ Press*; 2015;43:D1049--D1056.
- 416 30. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene Ontology: tool for
417 the unification of biology. *Nat. Genet. Nature Publishing Group*; 2000;25:25–9.
- 418
- 419

A

Protein network

1 vs. 2 1 vs. 3 2 vs. 3

Z-scores (blue = Down-regulation, red = Up-regulation)



bioRxiv preprint doi: <https://doi.org/10.1101/110759>; this version posted March 23, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

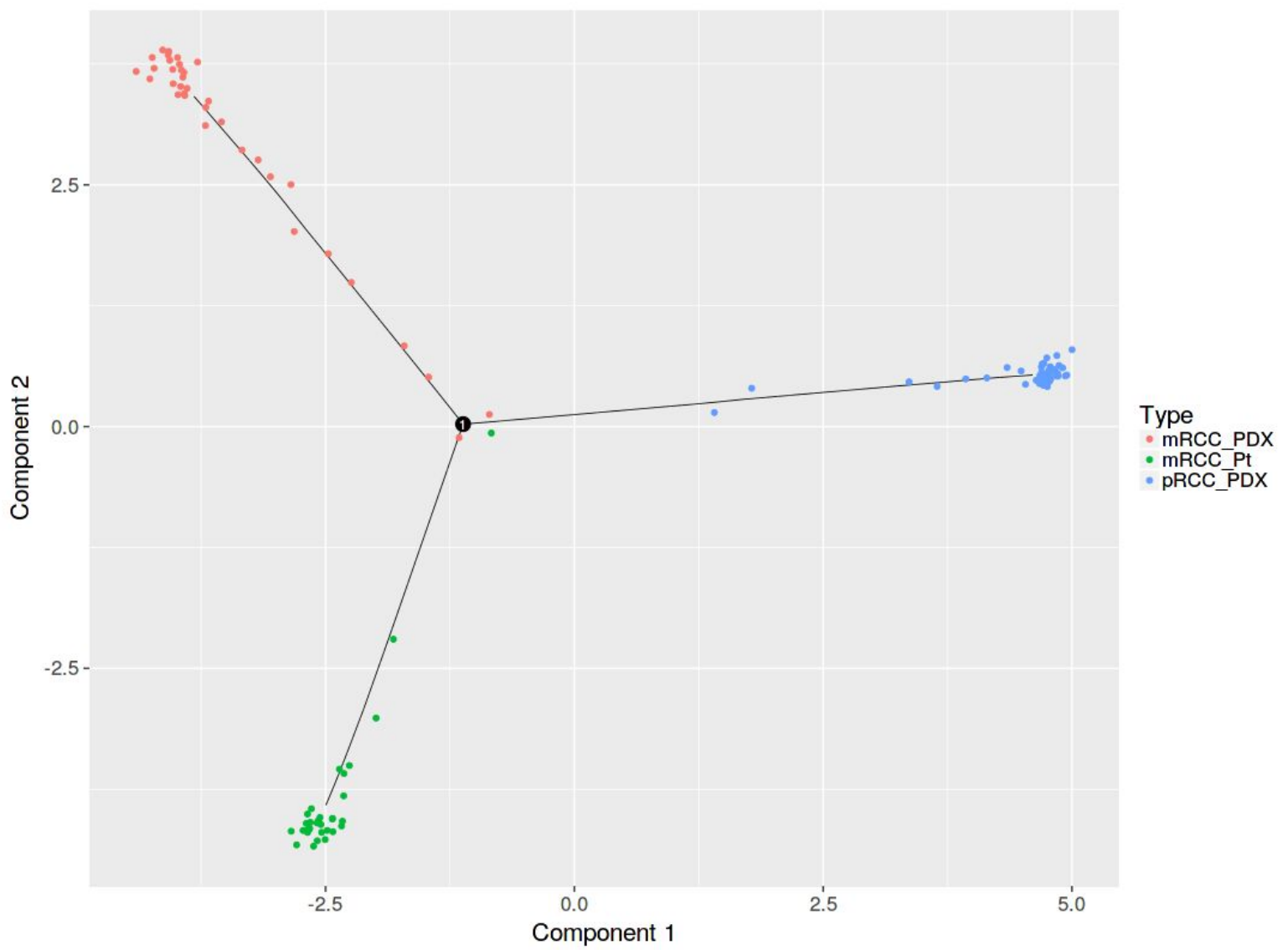
Proceed

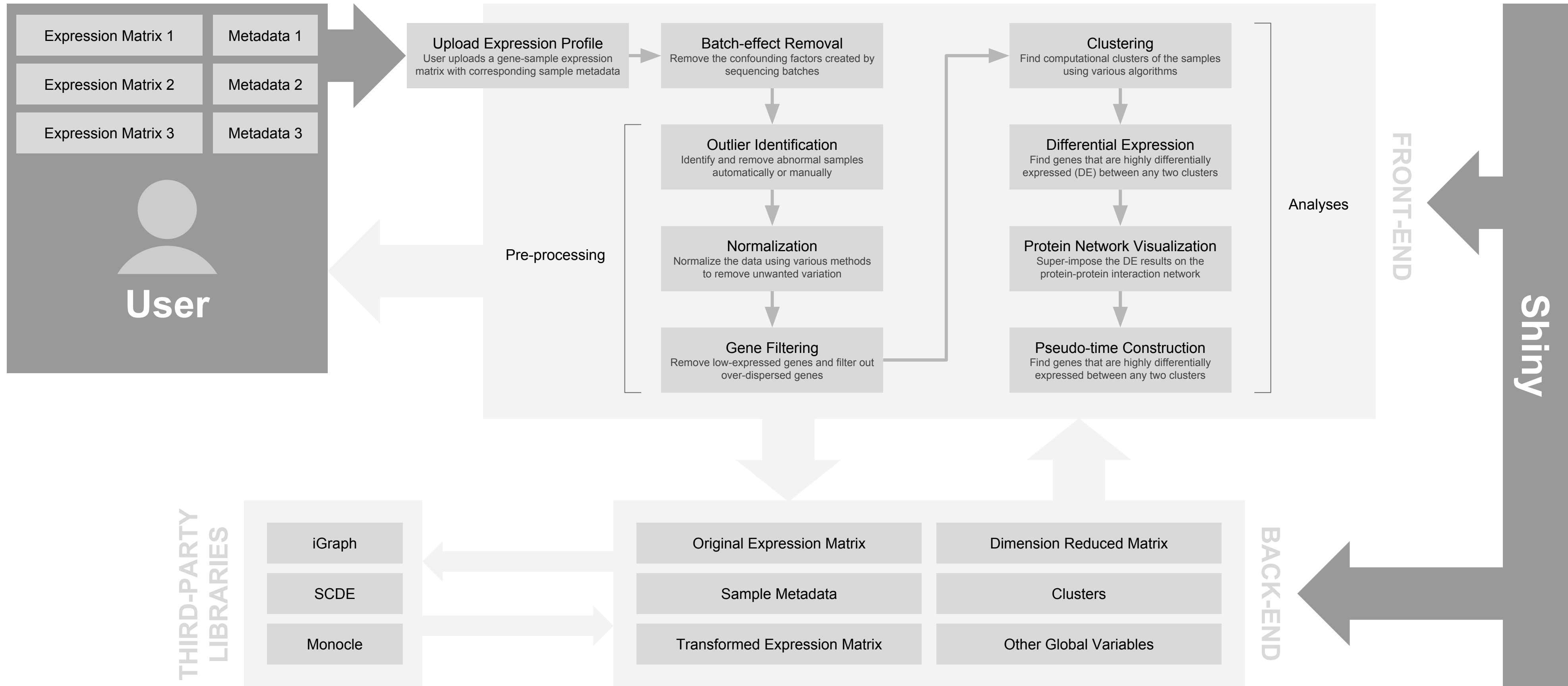
B

Pseudo-time construction

Cell labels

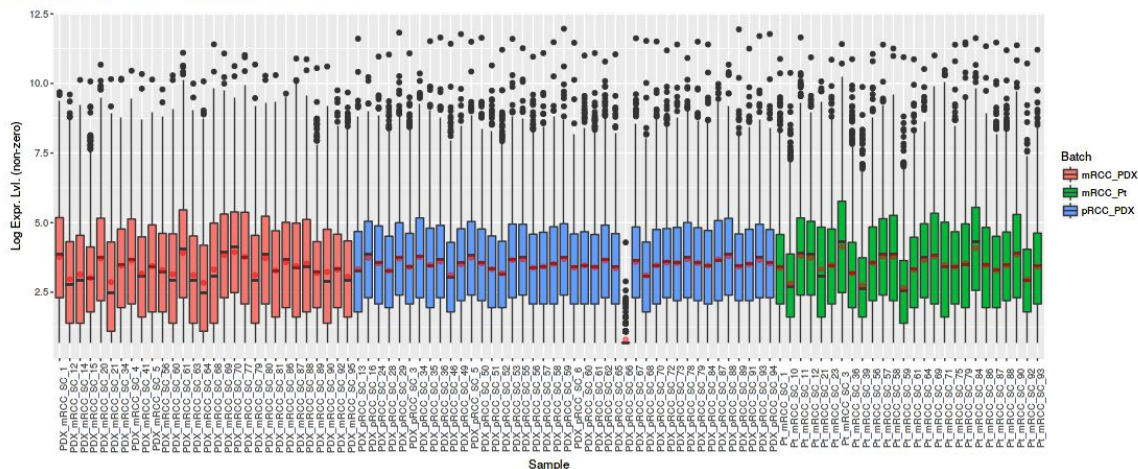
Type





Batch-effect removal

Data generated in batches may have confounding effects on results. To address this, select the factor that distinguishes cells in different batches, e.g., "Dataset", and check the underlying box before clicking a normalization button.



Re-plot random 96 cells

Batch factor:

Type

Remove batch effect

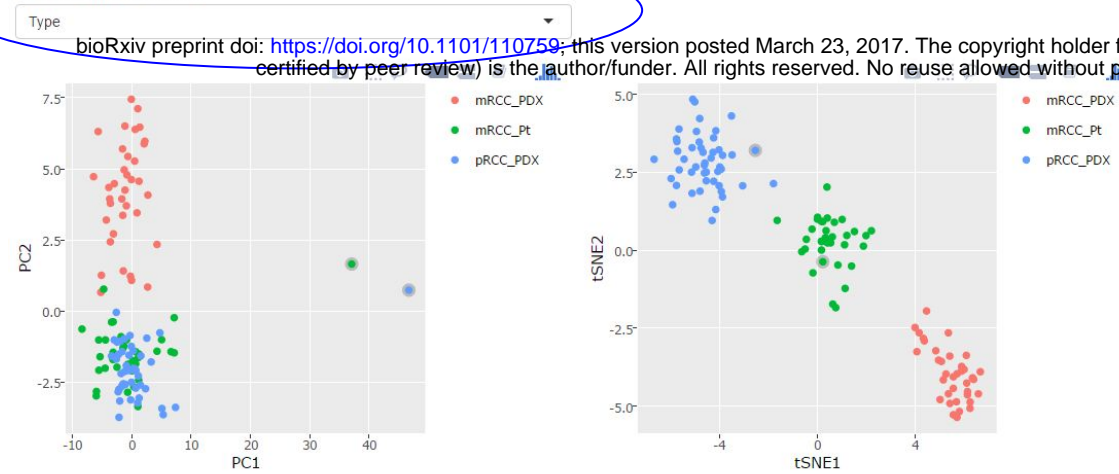
Reset

Submit

A

Outlier removal

Cell labels (from metadata)

 Cluster using only top expressed genes

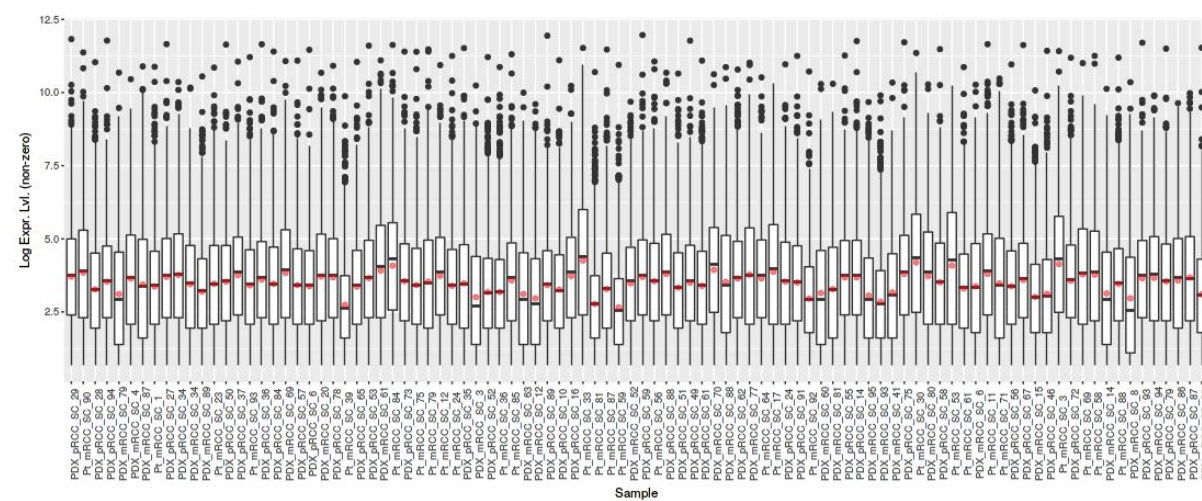
Selected cells:

id	Type	State	Pt_PDX	Mapped_reads	GSM	SRX	SRR
Pt_mRCC_SC_5	mRCC_Pt	mRCC	Pt	36234	GSM1887310	SRX1253756	SRR2431431
PDX_pRCC_SC_66	pRCC_PDX	pRCC	PDX	7603	GSM1887283	SRX1253736	SRR2431411

Showing 1 to 2 of 2 entries

C

Normalization



B

Outlier removal

Z-score threshold

4

Number of Outliers

1

Using

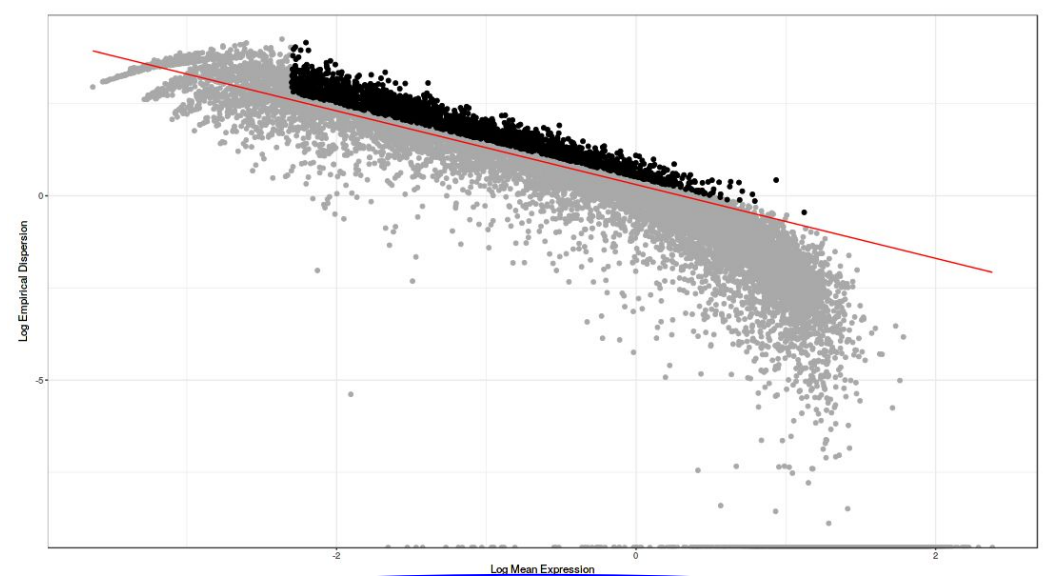
- Z-score threshold
- Fixed number of samples

According to

- PCA
- Correlation t-SNE

D

Gene filtering



Log Mean Expression Threshold

-5.63

-2.3

Dispersion Fit Threshold

0

1.22

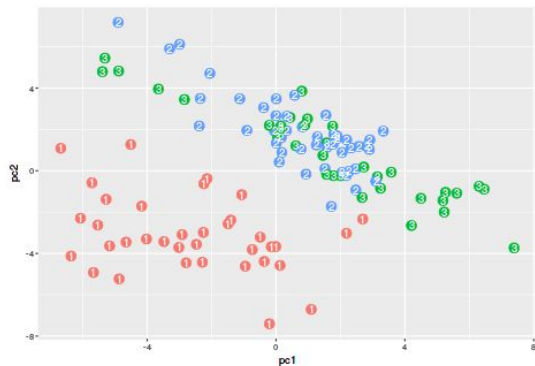
Starting number of genes:
19924Post-filtering number of genes:
2252

A

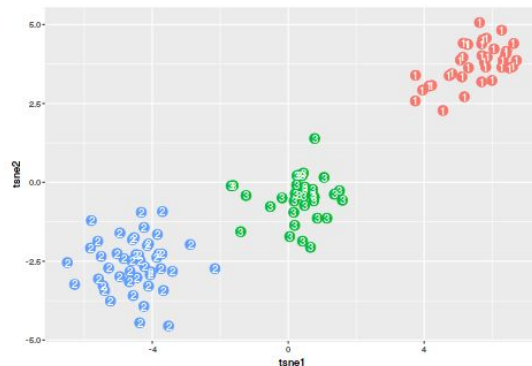
Clustering

Cell labels

Type



● mROC_PDX
● mROC_Pt
● pROC_PDX



● mROC_PDX
● mROC_Pt
● pROC_PDX

Clustering method

- Non-negative matrix factorization
- K-means (Euclidean)
- K-means (correlation t-SNE)
- Hierarchical clustering (Euclidean) with heatmap
- Hierarchical clustering (correlation t-SNE)

Automatically choose the number of clusters (might take long time)

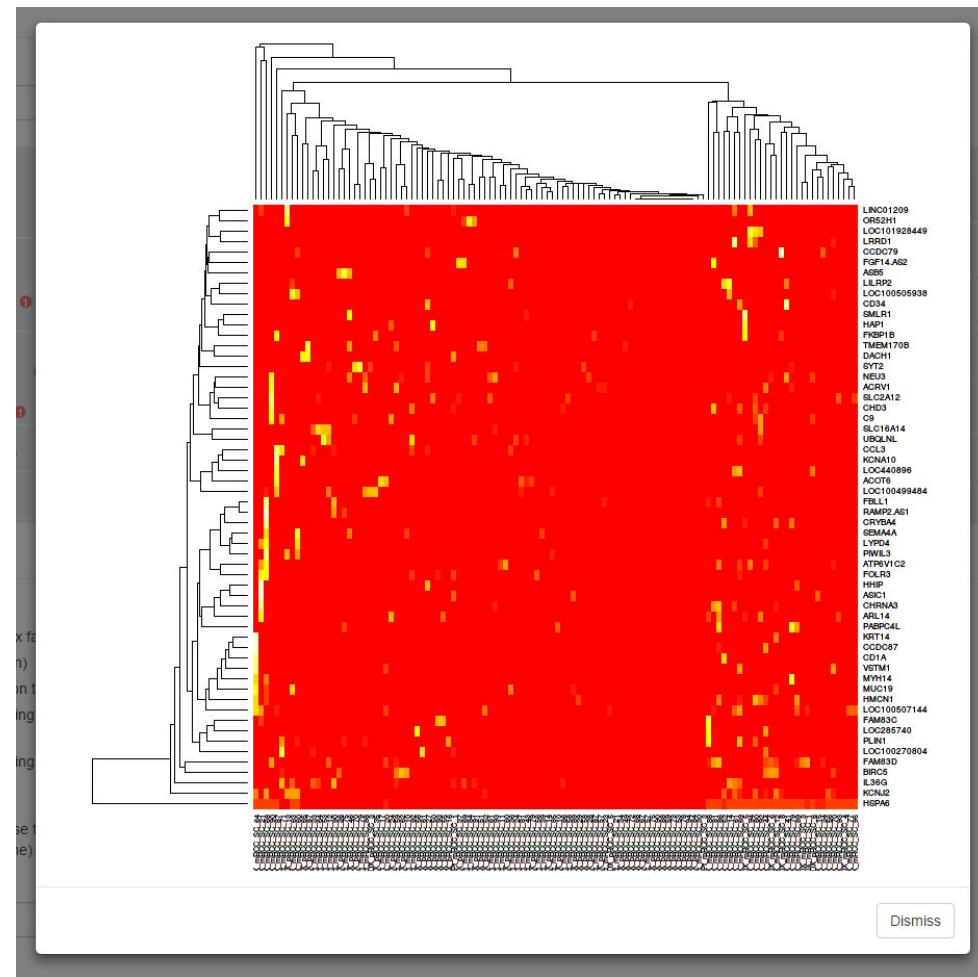
Number of clusters

3

Run clustering

Submit

B



Dismiss

A

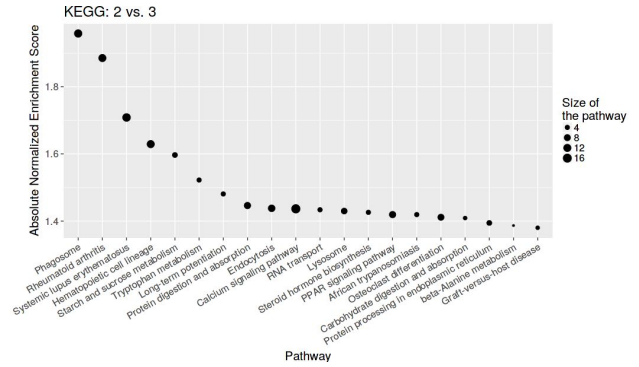
Differential expression

Number of processor cores

1

Start analysis

C



B

Differential expression

Cell labels

Type

Numbers in tabs below indicate which clusters have been compared. Genes are sorted most to least differentially expressed by absolute Z-score value.

1 vs. 2 1 vs. 3 2 vs. 3

Show 10 entries

Search:

gene	lb	mle	ub	ce	Z	cZ
CDH6	-14.916794	-5.746980	-4.563778	-4.563778	-7.160847	-6.337979
HSPA6	5.070865	6.338581	14.536479	5.070865	7.160813	6.337979
KRT81	12.001047	12.930705	13.437792	12.001047	7.160813	6.337979
CSF2	11.493960	12.465876	13.141991	11.493960	7.160809	6.337979
TCN1	-13.353277	-12.803934	-8.028869	-8.028869	-7.157471	-6.337979
DKK1	10.986874	12.043304	12.677162	10.986874	7.155977	6.337979
SLC15A1	-13.564563	-13.057477	-6.296324	-6.296324	-7.155594	-6.337979
SAMD5	-12.592648	-12.043304	-11.324932	-11.324932	-7.146775	-6.337979
MEG3	-12.592648	-11.874275	-11.155903	-11.155903	-7.140434	-6.337979
DCAF4L1	6.761153	12.423619	13.015220	6.761153	6.836788	6.028128

gene

lb

mle

ub

ce

Z

cZ

Showing 1 to 10 of 2,252 entries

Previous

1

2

3

4

5

...

226

Next

KEGG enrichment

Gene Ontology enrichment

Download CSV table

Submit

Software	GUI driven workflow	Live web site	Video tutorial	Interactive plots	Batch-effect removal	Outlier removal	Normalization	Over-dispersed genes identification	Clustering analysis	Differential expression analysis	Gene-set enrichment analysis	Network analysis	Pseudo-time construction	Citation
Granatum	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	
SCRAT / TSCAN / GSCA	✓	(*)	✓	✓	✗	✓	✓	✗	✓	✓	✓	✗	✓	Ji et al. 2016
ASAP	✓	✓	✗	✓	✗	✗	✓	✓	✓	✓	✓	✗	✗	Gardeux et al. 2016
Sake	✓	✓	✗	✓	✗	✗	✓	✓	✓	✓	✓	✓	✗	NA
Singular	✗	✗	✗	✓	✗	✓	✗	✗	✓	✓	✗	✗	✗	Fluidigm Corp. 2015
Cell Ranger / Loupe	✗	✗	✓	✓	✓	✓	✓	✗	✓	✓	✗	✗	✗	Zheng et al. 2017
Seurat	✗	✗	✗	✗	✗	✓	✓	✗	✓	✓	✗	✗	✗	Satija et al. 2016
Scater	✗	✗	✗	✓	✓	✓	✓	✗	✓	✗	✗	✗	✗	McCarthy et al. 2017
Monocle	✗	✗	✗	✗	✗	✓	✓	✓	✓	✓	✗	✗	✓	Trapnell et al. 2014
SCDE / PAGODA	✗	(**)	(***)	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	Kharchenko et al. 2014
Flotilla	✗	✗	✗	✓	✓	✓	✗	✗	✓	✓	✓	✓	✗	NA
Sincell	✗	✗	✗	✗	✗	✗	✗	✗	✓	✗	✓	✗	✓	Juliá et al. 2015
Sincera	✗	✗	✗	✗	✗	✗	✓	✗	✓	✓	✓	✓	✗	Guo et al. 2015

(*) The three components (SCRAT, TSCAN and GSCA) are not integrated.

(**) Results can be shown interactively using a web interface. However, the results themselves have to be pre-computed in R.

(***) For the interactive interface only