

SOFTWARE

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

Xun Zhu<sup>1,2</sup>, Thomas Wolfgruber<sup>1,2</sup>, Austin Tasato<sup>3</sup>, Lana X Garmire<sup>1, 2\*</sup>

\*Correspondence:

LGarmire@cc.hawaii.edu

<sup>1</sup>Graduate Program in Molecular Biology and Bioengineering, University of Hawaii at Manoa, Honolulu, HI 96816

<sup>2</sup>Epidemiology Program, University of Hawaii Cancer Center, Honolulu, HI 96813

<sup>3</sup>Department of Electrical Engineering, University of Hawaii at Manoa, Honolulu, HI 96816

## 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 code is freely available for research use at: <http://garmiregroup.org/granatum/code>

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

## 1 Background

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

10 Software that has been developed to address these challenges may have very limited accessibility  
11 for biologists with only general computer skills, as they typically require the ability to use a  
12 computing language like R [5,6]. Other existing workflows that can be used to analyze scRNA-seq  
13 data, such as Singular (Fluidigm, Inc., South San Francisco, CA, USA), Cell Ranger/ Loupe  
14 (Pleasanton, CA, USA), and Scater [7] all require some non-graphical interactions and they may not  
15 provide a comprehensive set of scRNA-seq analysis methods. To fill this gap, we have developed  
16 Granatum, a fully interactive graphical scRNA-seq analysis tool. Granatum is the Latin word for  
17 pomegranate, which bears many seeds, resembling single cells within the entity. This tool employs  
18 an easy-to-use web browser interface for a wide range of methods suitable for scRNA-seq analysis:  
19 removal of batch effects, removal of outlier cells, normalization of expression levels, filtering of  
20 under-informative genes, clustering of cells, identification of differentially expressed genes,  
21 identification of enriched pathways/ontologies, visualization of protein networks, and  
22 reconstruction of pseudo-time paths for cells. Our software will empower a much broader

23 audience of research communities to study single cell complexity, by allowing them to readily  
24 explore single-cell expression data from a graphical user interface.

## 25 **Implementation**

### 26 *Overview*

27 Both the front-end and the back-end of Granatum are written in the R software language, and built  
28 with the Shiny framework [8]. Multiple concurrent users are handled by Shiny and each user works  
29 on its own data space. To protect the privacy of users, the data submitted by one user is not visible  
30 to any other user. The front-end is implemented as a web page with dynamically loaded pages,  
31 and is arranged in a step-wise fashion. The default theme uses the Bootstrap framework. ShinyJS  
32 [9] is used to power some of the interactive components. To allow users to redo a task, each  
33 processing step is equipped with a reset button.

34

### 35 *Interactive widgets*

36 The package visNetwork is used for the layout and physics simulation of the network modules [10].  
37 DataTables are used to preview user submitted data and to show tabular data in various modules  
38 [11]. Plotly is used for the interactive outlier identification step [12]. The package ggplot2 is used  
39 for the scatter-plots and box-plots, which is also used by the Monocle package for the Pseudo-time  
40 construction step [3,13].

41

### 42 *Back-end variable management*

43 The expression matrix and the metadata sheet are stored separately for each user. The metadata  
44 sheet can refer to groups, batches, or other properties of the samples in the corresponding

45 expression matrix. These two types of tables are shared across all modules. Other variables shared  
46 across all modules include the log-transformed expression matrix, the filtered and normalized  
47 expression matrix, the dimensionally reduced matrix, species (human or mouse) and the primary  
48 metadata column.

49

## 50 *Deployment*

51 Granatum is deployed from a pre-configured VirtualBox Appliance (machine image), which is  
52 configured with all tool files and dependencies. VirtualBox is an open-source hypervisor developed  
53 by the Oracle Corporation – <https://www.virtualbox.org/>. The Granatum image is provided as an  
54 Open Virtual Appliance file, which is loaded by clicking through the *Import Appliance* function of  
55 VirtualBox installed on a Windows or Linux system. When the Granatum image is running, it opens  
56 an Ubuntu desktop and starts the Granatum server (Additional file 1). When the server has  
57 completely loaded a text, a welcome message will appear on the screen indicating that Granatum  
58 is ready for use. The server can also be accessed from a web browser outside of the VirtualBox  
59 instance, by navigating to the appropriate local port, e.g., <http://localhost:8028>. Accessing the  
60 server from outside of VirtualBox simplifies data transfer to/from the server, e.g., files can be  
61 loaded from the user's desktop outside of VirtualBox.

62

## 63 *Batch-effect removal*

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

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

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

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

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

### 71 *Clustering methods*

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

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

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

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

84

85 *Correlation t-SNE*

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

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

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

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

92

### 93 *Elbow-point finding algorithm in clustering*

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

102

### 103 *Differential expression analysis*

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

108

### 109 *Gene-set enrichment analysis*

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

112

### 113 *Pseudo-time construction*

114 We use Monocle (version 2.2.0) in our pseudo-time construction step. When building the  
115 *CellDataSet* required for monocle's input, we set the *expressionFamily* to *negbinomial.size()*. The  
116 dimension reduction is done using the *reduceDimension* function with *max\_components* set to be  
117 2.

## 118 **Results**

### 119 *Overview*

120 Granatum neatly presents nine modules, arranged as steps and ordered by their dependency  
121 (Figure 1), spanning a comprehensive set of methods for single cell analysis. It starts with one or  
122 more user-supplied expression matrices and corresponding sample metadata sheet(s), followed by  
123 data-merging, batch-effect removal, outlier removal, normalization, gene filtering, clustering,  
124 differential expression, protein-protein network, and pseudo-time construction.

125 Comparing to other freely available tools, the workflow is flexible in several aspects: (1) It supports  
126 multiple dataset submission and batch effect removal; (2) at any point of the step, the user can  
127 reset the current step for re-analysis; (3) the user can bypass certain steps and still complete the  
128 workflow; (4) the user can select subsets of samples/data for their customized analysis need; (5)  
129 the user can identify outlier samples either automatically by a pre-set threshold, or manually by

130 simply clicking the samples the PCA plot or the correlation t-SNE plot; (6) multiple cores can be  
131 specified in the differential expression module for speed-up; (7) GSEA can be performed for the  
132 differentially expressed genes in all pairs of subgroups, following clustering analysis; (8) Monocle  
133 pseudo-time construction can be performed to gain insights of relationships between the cells. We  
134 elaborate the details of each step in chronological order, in the following sections.

135

### 136 *Upload data*

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

144

145 *Batch-effect removal* Samples obtained in batches can create unwanted technical variation, which  
146 confound the biological variation [22]. It is thus important to remove the expression level  
147 difference due to batches. Granatum provides a batch-effect removal step, where the batches are  
148 shown as different colors in the box-plot (Figure 2). If more than one datasets are uploaded, by  
149 default each dataset is assumed to be one batch. Alternatively, if the batch numbers are indicated  
150 in the sample metadata sheet, the user may select the column in which the batch numbers are  
151 stored (blue circled in Figure 2). For datasets with a large number of cells, to maintain legibility of



152 the box-plot a random selection of 96 sub-samples is shown in the box-plot, and can be re-  
153 sampled freely.

154

### 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

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

168

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

173

### 174 *Normalization*

175 Normalization is essential to most scRNA-seq data, except those with the UMI counts, before the  
176 down-stream functional analyses. The current version of Granatum has implemented three  
177 commonly used normalization algorithms: rescale to geometric mean, quantile normalization, and  
178 size-factor normalization [23,24]. A box-plot is shown post normalization, to help illustrate its  
179 effect to the median, mean, and extreme values across samples. As is the case in the batch-effect  
180 removal step, for a dataset with a large number of samples, 96 sub-samples are randomly chosen  
181 for the visualization purpose (Figure 3C).

182

### 183 *Gene filtering*

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

196

### 197 *Clustering*

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

207  
208 Next, the resulting cluster labels obtained above, are then super-imposed onto the two  
209 unsupervised PCA and correlation t-SNE plots (Figure 4A). The user can also represent user-defined  
210 labels in the sample metadata as different colors in these plots. By comparing the two sets of  
211 labels, the users can quickly check the concordance between the prior metadata labels and the  
212 computed clusters.

### 213 *Differential expression*

214  
215 After obtaining a set of clusters, it is intuitively important to identify genes that are differentially  
216 expressed between any two clusters. Granatum uses the state-of-the-art SCDE method for its  
217 single-cell DE analysis [20]. The DE comparison is performed in a pair-wise fashion when more than  
218 two clusters are present. This step is computationally time and memory consuming. To shorten  
219 computation time, a user can select the number of cores for parallelization on multi-core machines  
220 (Figure 5A). When SCDE is completed, tabbed tables show the genes sorted by their Z-scores,

221 along with the model coefficients (Figure 5B). As another feature to empower the users, the gene  
222 symbols are linked to their corresponding GeneCards pages ([www.genecards.org](http://www.genecards.org)) [25]. The DE  
223 results can be downloaded as a CSV file via the "Download CSV table" button.

224 To investigate the collective biological functions of these genes, the user can further perform Gene  
225 Set Enrichment Analysis (GSEA) with either KEGG pathways or Gene Ontology (GO) terms (circled  
226 blue in Figure 5B) [26–29]. We have employed a very intuitive bubble-plot to visualize the GSEA  
227 results, where the vertical position of the bubble indicates the enrichment score of the gene sets,  
228 and the size of the bubble indicates number of genes in that set (KEGG pathway or GO term)  
229 (Figure 5C).

230

### 231 *Protein network visualization*

232 Protein-protein interaction (PPI) network gives straightforward and systematic understanding of  
233 the connections between these differentially expressed genes. Granatum selects the top K (default  
234 K=200) genes in the DE results, and super impose the PPI network on them. Genes that are not  
235 connected to any other genes in the list are removed from the PPI network. We use visNetwork to  
236 enable the interactive display of the graph [10]. The user can freely rearrange the graph by  
237 dragging the nodes to the desired location, and reconfiguring the layout to achieve good visibility  
238 of the modules (via elastic-spring physics simulation) (Figure 6A). In this interactive graph, the Z-  
239 scores are mapped as colors on the nodes where red indicates up-regulation and blue indicates  
240 down-regulation.

241

242 *Pseudo-time construction* Granatum has included the Monocle algorithm, a widely-used method to  
243 reconstruct a pseudo-timeline for the samples [3]. Monocle uses the Reversed Graph Embedding

244 algorithm to learn the structure of the data, and the Principal Graph algorithm to find the time-  
245 lines and branching points of the samples. We superimpose the timeline on the samples scatter-  
246 plot projected on the two components of the learned projection matrix. The user may map any  
247 pre-defined labels or numeric assays provided in the metadata sheet on to the scatter-plot (Figure  
248 6B). The plotting functions are adapted from the visualization code in Monocle.

## 249 Discussion

250 The field of scRNA-seq is fast-evolving both in terms of the development of instrumentation and  
251 the innovation of computational methods. However, it becomes exceedingly hard for a wet-lab  
252 researcher without formal bioinformatics training to catch up with the latest iterations of  
253 algorithms [5]. This poses major barriers to them and many resort to sending their generated data  
254 to third-party bioinformaticians, before they are able to visualize the data themselves. This  
255 segregation often prolongs the research cycle time, as it often takes significant effort to maintain  
256 effective communications between the two sides (sometimes even more complicated with a third  
257 party of the genomics core). Also, issues with the experimentations do not get the chance to be  
258 spotted early enough, to avoid significance loss of time and cost in the projects. It is thus very  
259 attractive to have a non-programming graphical application which includes state-of-the-art  
260 algorithms as routine procedures, in the hands of the bench-scientist who generate the scRNA-seq  
261 data.

262 Granatum is our attempt to fill this void. It is to our knowledge the first solution that aims to cover  
263 the entire scRNA-seq workflow with an intuitive, step-wise graphical user interface. Throughout  
264 the development process our priority has been to make sure that it is fully accessible to

265 researchers with no programming experiments. We have strived to achieve that the plots and  
266 tables are self-explanatory, interactive and visually pleasant. We have sought inputs from our  
267 single-cell bench-side collaborators, to ensure that the terminologies are easy to understand by  
268 them. We also supplement Granatum with a manual and video that guide the users through the  
269 entire workflow, using example datasets. Currently Granatum targets users who have their  
270 expression matrices and metadata sheets ready. However, we are developing the next version of  
271 Granatum, which will handle the entire scRNA-seq data processing and analysis pipeline including  
272 FASTQ quality control, alignment, and expression quantification. In the future, we will enrich  
273 Granatum with capacities to analyze and integrate other types of genomics data in single cells,  
274 such as exome-seq and methylation data.

## 275 **Conclusions**

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

281

## 282 **Figure legends**

283 **Figure 1: Granatum workflow.** Granatum is built with the Shiny framework, which supports both  
284 front-end and the back-end. The user uploads one or more expression matrices with

285 corresponding metadata for samples. The back-end stores data separately for each individual user,  
286 and invokes third-party libraries on demand.

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

291 **Figure 3: The outlier removal, normalization and gene filtering steps.** A) The main interface of the  
292 outlier removal step. The two scatter-plots are the PCA and correlation t-SNE plots, with colors  
293 indicate the cell labels (box circled in blue). The metadata table (circled in red) shows the labels for  
294 the selected cells. B) The pop-up window for automatic outlier detection options after the “auto-  
295 identify” button is clicked. C) The normalization step. The box-plot shows the expression levels of  
296 each cell in log-scale. In cases where more than 96 cells are present in the data, only a random  
297 sample of 96 cells are shown. D) The Gene filtering step. The y-axis of the scatter-plot is the  
298 empirical dispersion, estimated by a negative binomial model. The x-axis is the log mean  
299 expression of each gene. The user can change the threshold by dragging the two sliders circled in  
300 blue.

301 **Figure 4: The Clustering step.** A) Main interface. PCA and t-SNE plots are shown with colors  
302 mapped to user-selected sample labels. After clustering, samples are marked with their assigned  
303 cluster numbers. The user can either choose a specific number of clusters or let Granatum

304 compute the best number of clusters. B) When Hclust (Euclidean) is selected, a pop-up window will  
305 show a heatmap of the expression matrix with dendrograms.

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

312 **Figure 6: The Protein network and Pseudo-time construction steps.** A) The Protein network step.  
313 The A tabbed panel shows the connected gene modules on the PPI network between each pair of  
314 clusters. The color on each node (gene) indicates its Z-score in the differential expression test. Red  
315 and blue colors indicates up- and down- regulation. B) The Pseudo-time construction step.  
316 Monocle algorithm is customized to visualize the paths among individual cells. The user can  
317 represent sample labels from the metadata as colors in the plot.

318

## 319 **Supplementary files**

320 **Additional file 1: Granatum deployment.** A screenshot of an activated VirtualBox Appliance  
321 running the Granatum server is shown behind a web browser outside of the Appliance, which is  
322 accessing the server with the URL <http://localhost:8028/>. The server can be started by double-  
323 clicking the Granatum desktop icon within the Appliance and stopped by closing the Terminal



324 window, which pops up when the server is activated. All data to/from the server can be handled  
325 outside of the Appliance from the external browser.

326

## 327 **Availability of data and material**

328 The package Granatum is freely available for research use, and can be downloaded at:

329 <http://garmiregroup.org/granatum/code>

330 A demonstration video can be found at

331 <http://garmiregroup.org/granatum/video>

332

## 333 **Declarations**

334 NA

335

## 336 **Competing interests**

337 The authors declared no conflict of interest.

338

## 339 **Funding**

340 This research was supported by grants K01ES025434 awarded by NIEHS through funds provided by  
341 the trans-NIH Big Data to Knowledge (BD2K) initiative (<http://datascience.nih.gov/bd2k>), P20  
342 COBRE GM103457 awarded by NIH/NIGMS, NICHD R01 HD084633 and NLM R01LM012373 and  
343 Hawaii Community Foundation Medical Research Grant 14ADVC-64566 to LX Garmire.

344

## 345 **Authors' contributions**

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

## 350 **Acknowledgements**

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

## 354 **List of abbreviations**

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

356 **DE:** differential expression

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

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

359 **GO:** Gene ontology

360 **PCA:** Principal component analysis

361 **t-SNE:** t-Distributed Stochastic Neighbor Embedding

362 **NMF:** Non-negative matrix factorization

363 **Hclust:** Hierarchical clustering

364 **PPI:** Protein-protein interaction

365

## References

366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388

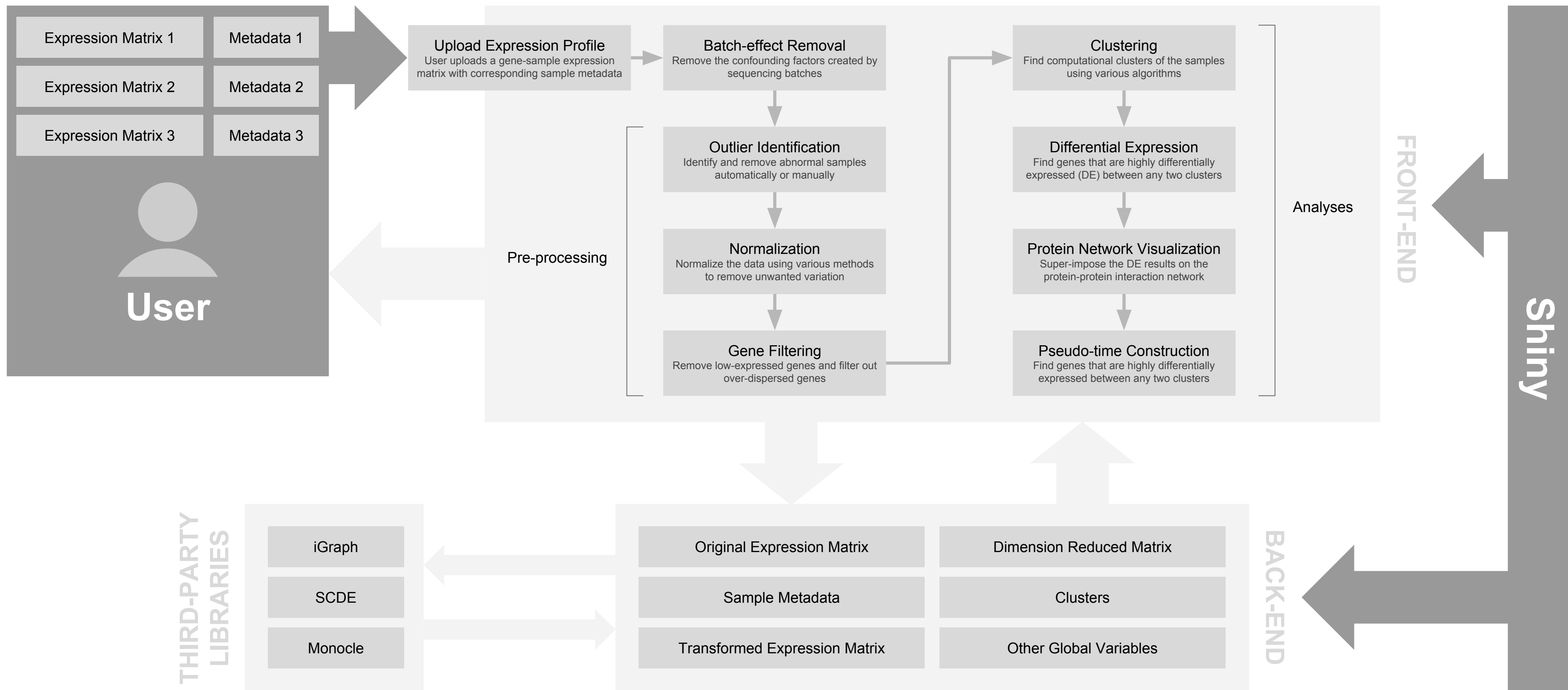
1. Patel AP, Tirosh I, Trombetta JJ, Shalek AK, Gillespie SM, Wakimoto H, et al. Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science* (80-. ). [Internet]. Department of Neurosurgery, Massachusetts General Hospital and Harvard Medical School, Boston, MA 02114, USA. Department of Pathology and Center for Cancer Research, Massachusetts General Hospital and Harvard Medical School, Boston, MA 02114, USA. *Broad J*; 2014;344:1396–401. Available from: <http://dx.doi.org/10.1126/science.1254257>
2. Lewis BP, Burge CB, Bartel DP. Conserved seed pairing, often flanked by adenosines, indicates that thousands of human genes are microRNA targets. *Cell*. Elsevier; 2005;120:15–20.
3. Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, et al. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat. Biotechnol. Nature Research*; 2014;32:381–6.
4. Brennecke P, Anders S, Kim JK, Kołodziejczyk AA, Zhang X, Proserpio V, et al. Accounting for technical noise in single-cell RNA-seq experiments. *Nat. Methods*. Nature Publishing Group; 2013;
5. Poirion OB, Zhu X, Ching T, Garmire L. Single-Cell Transcriptomics Bioinformatics and Computational Challenges [Internet]. *Front. Genet.* . 2016. p. 163. Available from: <http://journal.frontiersin.org/article/10.3389/fgene.2016.00163>
6. Team RC. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. 2015, URL [http. www. R-project. org](http://www.R-project.org). 2016;
7. McCarthy DJ, Campbell KR, Lun ATL, Wills QF. scater: pre-processing, quality

- 389 control, normalisation and visualisation of single-cell RNA-seq data in R. bioRxiv  
390 [Internet]. Cold Spring Harbor Labs Journals; 2016; Available from:  
391 <http://biorxiv.org/content/early/2016/08/15/069633>
- 392 8. RStudio, Inc. Easy web applications in R. 2013.
- 393 9. Attali D. shinyjs: Easily Improve the User Experience of Your Shiny Apps in Seconds  
394 [Internet]. 2016. Available from: <https://cran.r-project.org/package=shinyjs>
- 395 10. Almende B.V., Thieurmel B. visNetwork: Network Visualization using “vis.js”  
396 Library [Internet]. 2016. Available from: [https://cran.r-](https://cran.r-project.org/package=visNetwork)  
397 [project.org/package=visNetwork](https://cran.r-project.org/package=visNetwork)
- 398 11. Xie Y. DT: A Wrapper of the JavaScript Library “DataTables” [Internet]. 2016.  
399 Available from: <https://cran.r-project.org/package=DT>
- 400 12. Sievert C, Parmer C, Hocking T, Chamberlain S, Ram K, Corvellec M, et al. plotly:  
401 Create Interactive Web Graphics via “plotly.js” [Internet]. 2016. Available from:  
402 <https://cran.r-project.org/package=plotly>
- 403 13. Wickham H. ggplot2: Elegant Graphics for Data Analysis [Internet]. Springer-Verlag  
404 New York; 2009. Available from: <http://ggplot2.org>
- 405 14. Zhu X, Ching T, Pan X, Weissman S, Garmire L. Detecting heterogeneity in single-  
406 cell RNA-Seq data by non-negative matrix factorization. PeerJ Prepr. PeerJ Inc. San  
407 Francisco, USA; 2016;4:e1839v1.
- 408 15. Brunet J-P, Tamayo P, Golub TR, Mesirov JP. Metagenes and molecular pattern  
409 discovery using matrix factorization. Proc. Natl. Acad. Sci. [Internet]. 2004;101:4164–  
410 9. Available from: <http://www.pnas.org/content/101/12/4164.abstract>
- 411 16. Gaujoux R, Seoighe C. Algorithms and framework for nonnegative matrix

- 412 factorization (NMF). 2010.
- 413 17. Lloyd S. Least squares quantization in PCM. IEEE Trans. Inf. theory. IEEE;
- 414 1982;28:129–37.
- 415 18. Murtagh F, Contreras P. Methods of hierarchical clustering. arXiv Prepr.
- 416 arXiv1105.0121. 2011;
- 417 19. Krijthe J. Rtsne: T-Distributed Stochastic Neighbor Embedding using Barnes-Hut
- 418 Implementation. R Packag. version 0.10, URL [http://CRAN.R-project.org/package=](http://CRAN.R-project.org/package=Rtsne)
- 419 Rtsne. 2015;
- 420 20. Kharchenko P V, Silberstein L, Scadden DT. Bayesian approach to single-cell
- 421 differential expression analysis. Nat. Methods. Nature Publishing Group;
- 422 2014;11:740–2.
- 423 21. Sergushichev A. An algorithm for fast preranked gene set enrichment analysis
- 424 using cumulative statistic calculation. bioRxiv [Internet]. Cold Spring Harbor Labs
- 425 Journals; 2016; Available from: <http://biorxiv.org/content/early/2016/06/20/060012>
- 426 22. Hicks SC, Teng M, Irizarry RA. On the widespread and critical impact of systematic
- 427 bias and batch effects in single-cell RNA-Seq data. bioRxiv. Cold Spring Harbor Labs
- 428 Journals; 2015;25528.
- 429 23. Bolstad BM, Irizarry RA, Åstrand M, Speed TP. A comparison of normalization
- 430 methods for high density oligonucleotide array data based on variance and bias.
- 431 Bioinformatics. Oxford Univ Press; 2003;19:185–93.
- 432 24. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion
- 433 for RNA-Seq data with DESeq2. bioRxiv. Cold Spring Harbor Labs Journals; 2014;
- 434 25. Rebhan M, Chalifa-Caspi V, Prilusky J, Lancet D. GeneCards: integrating

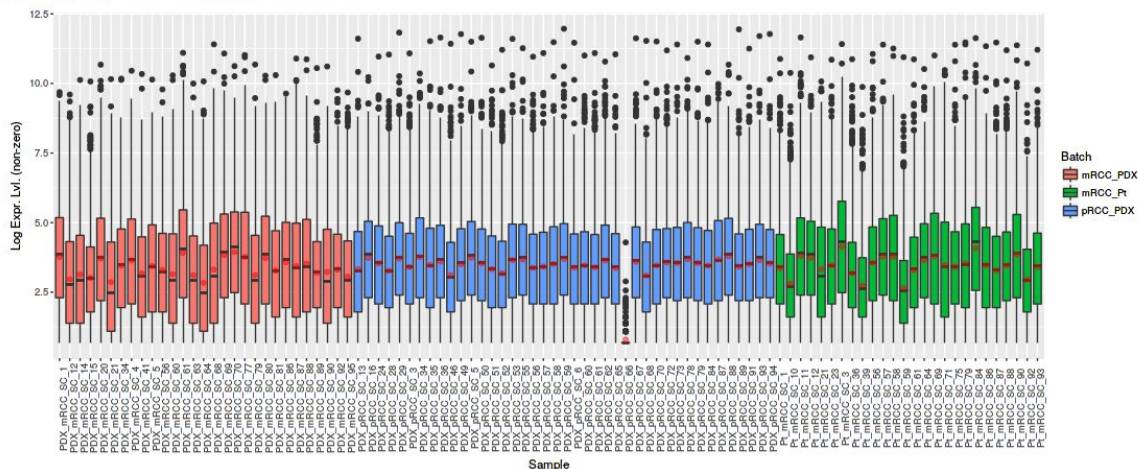
- 435 information about genes, proteins and diseases. *Trends Genet. Elsevier Current*  
436 *Trends*; 1997;13:163.
- 437 26. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al.  
438 Gene set enrichment analysis: a knowledge-based approach for interpreting genome-  
439 wide expression profiles. *Proc. Natl. Acad. Sci. National Acad Sciences*;  
440 2005;102:15545–50.
- 441 27. Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K. KEGG: new  
442 perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res. Oxford*  
443 *Univ Press*; 2017;45:D353--D361.
- 444 28. Consortium GO, others. Gene ontology consortium: going forward. *Nucleic Acids*  
445 *Res. Oxford Univ Press*; 2015;43:D1049--D1056.
- 446 29. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. *Gene*  
447 *Ontology: tool for the unification of biology. Nat. Genet. Nature Publishing Group*;  
448 2000;25:25–9.
- 449
- 450

# Granatum



# 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

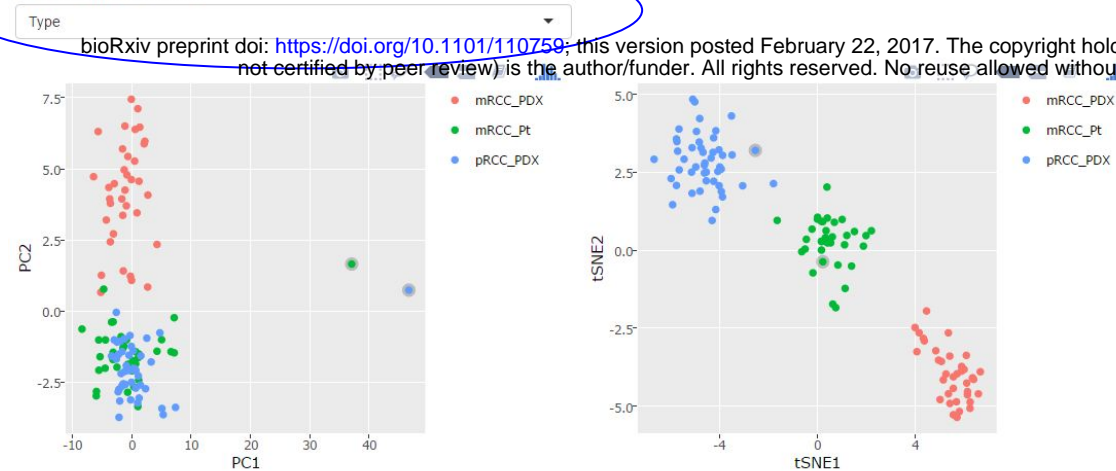
Submit



A

## Outlier removal

Cell labels (from metadata)

 Cluster using only top expressed genes

Auto-identify Remove selected De-select all

Reset Submit

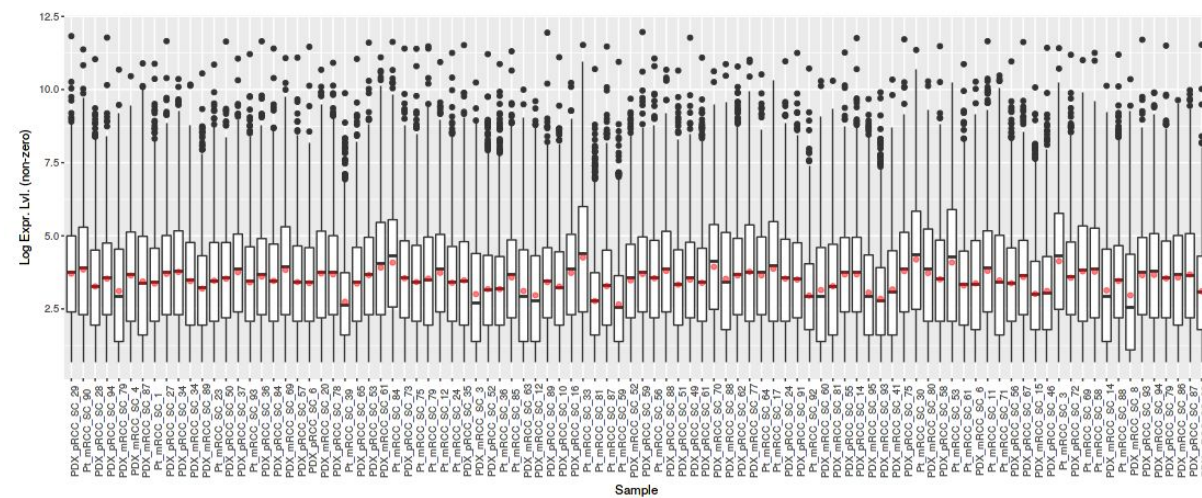
Selected cells:

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



Re-plot random 96 cells

Rescale to genomic mean Quantile normalization Size-factor normalization

Reset Submit

B

## Outlier removal

Z-score threshold

4

Number of Outliers

1

Using

- Z-score threshold
- Fixed number of samples

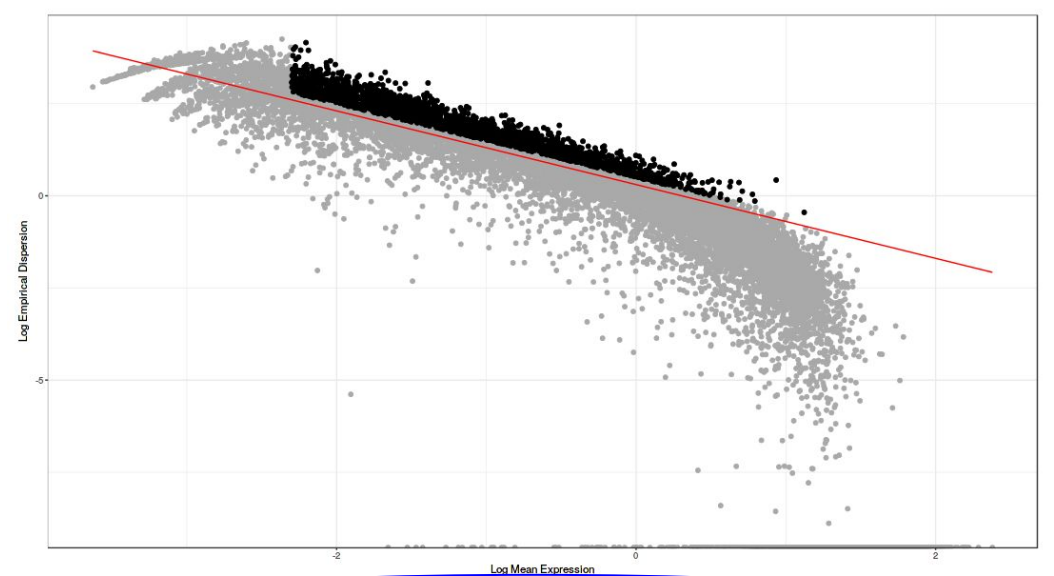
According to

- PCA
- Correlation t-SNE

Cancel OK

D

## Gene filtering



Log Mean Expression Threshold

-5.63

-2.3

Dispersion Fit Threshold

2.98

0

1.22

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

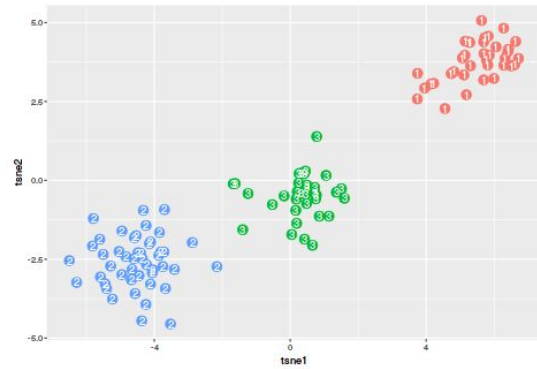
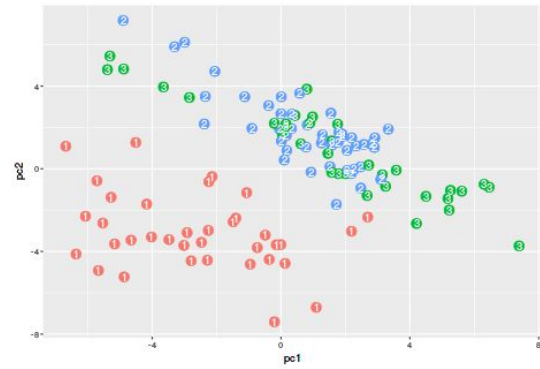
Submit

# A

## Clustering

Cell labels

Type



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)

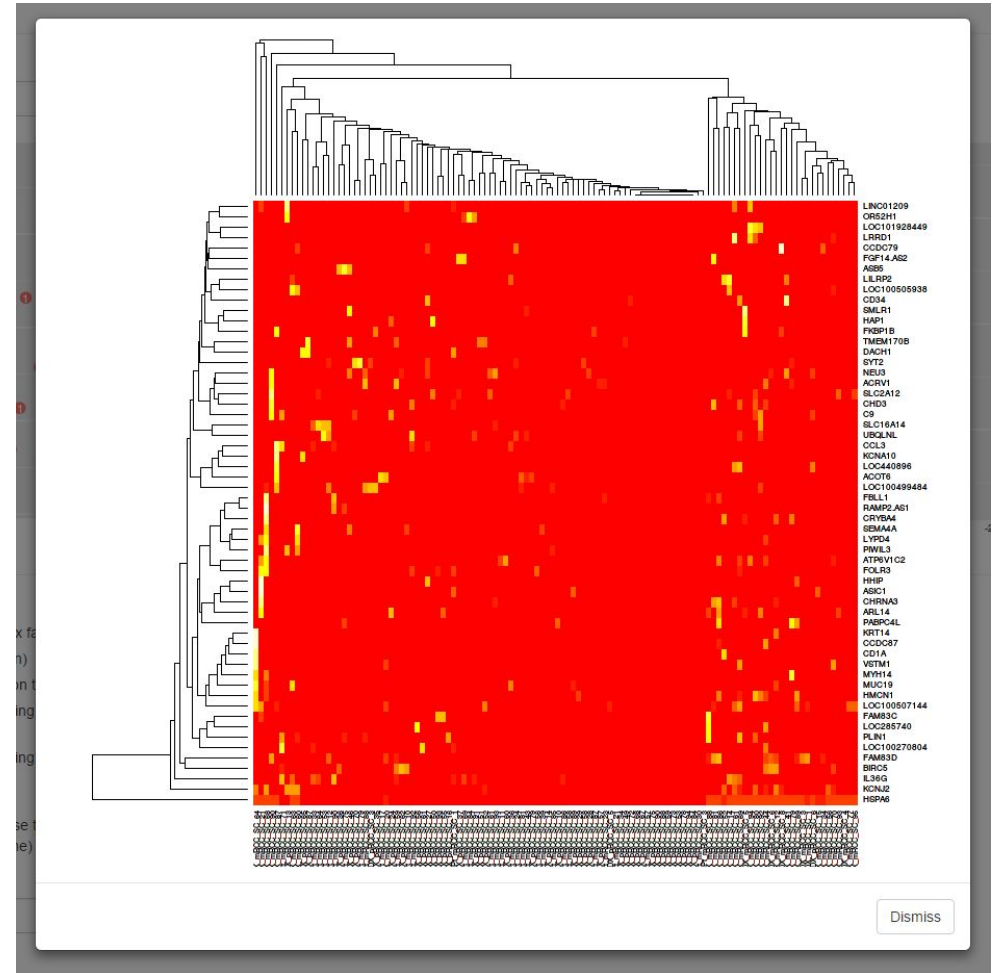
Number of clusters

3

Run clustering

Submit

# B



Dismiss

A

## Differential expression

Number of processor cores



B

## Differential expression

Cell labels

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








Showing 1 to 10 of 2,252 entries



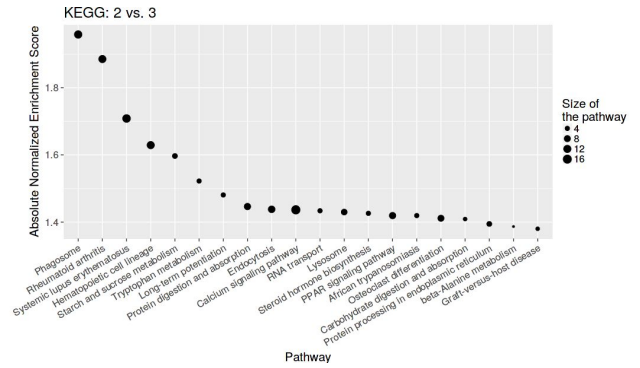









C

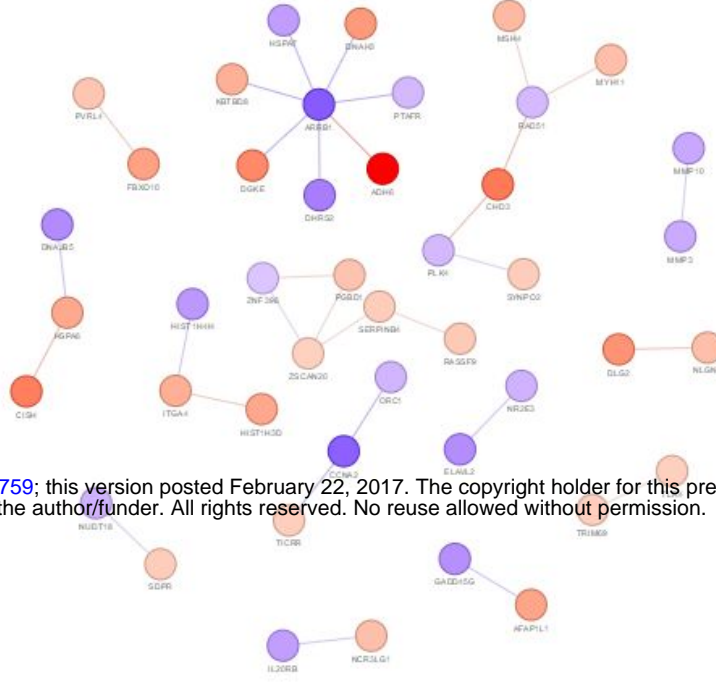
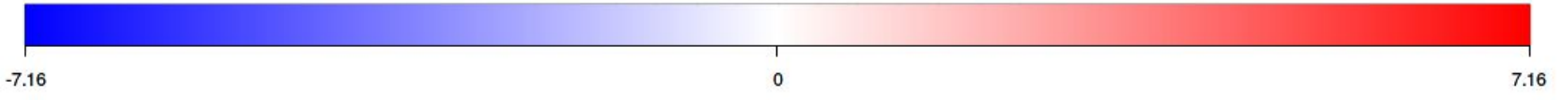


# A

## 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 February 22, 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

