

BiGGEsTS



BiclusterinG Gene Expression Time Series

Quickstart

BiGGEsTS is a software tool for time series gene expression data analysis, based on biclustering algorithms particularly suited for this kind of data. It is open source and freely available at: <http://kdbio.inesc-id.pt/software/biggests/>. The purpose of this quickstart document is to provide simple instructions on how to install and use BiGGEsTS. It is written under the assumption that the user knows what microarrays and expression data are. However, anyone who wants to try or start using BiGGEsTS without any previous knowledge on these concepts will be able to run the software using this guide. For obtaining help, sending feedback and improvements and/or reporting issues related to BiGGEsTS software, use the following e-mail: biggests.software@gmail.com.

Note that BiGGEsTS is free software: you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or any later version. BiGGEsTS is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details. A copy of the GNU General Public License is included with BiGGEsTS. For detailed information about the license, see the GNU licenses at <http://www.gnu.org/licenses/>.

Pre-Requisites

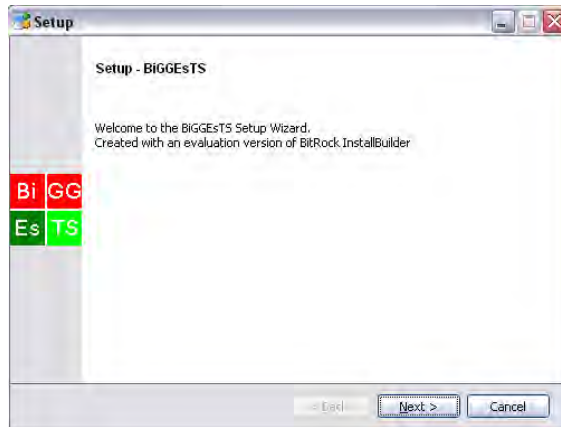
If you want to install BiGGEsTS, please check if you have a JVM (Java Virtual Machine) with JDK/JRE 1.5 or higher installed on your system.

We recommend a minimum of 1024 MB of RAM for running BiGGEsTS software.

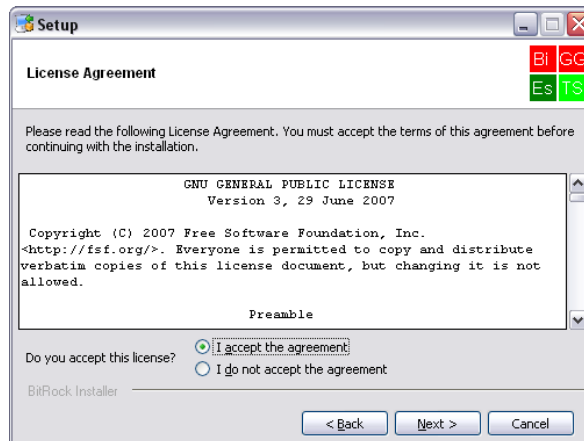
Installation

On Windows, using the installer:

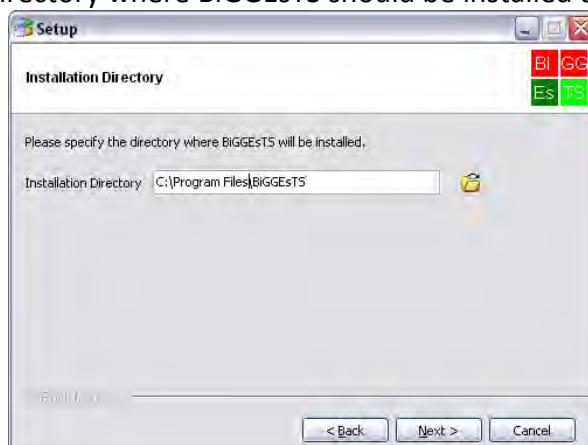
1. Double click the installer file and click **Next**.



2. After reading the terms of the license, select the **I accept the agreement** option and press **Next**.



3. Specify the directory where BiGGESTS should be installed and press **Next**.



4. After a successful installation you will be prompted to finish the process. You can do this by pressing the **Finish** button. BiGGESTS is now ready to be used. You will find shortcuts to BiGGESTS on both the **Desktop** and **Start** menu.



On Windows or Mac OS, using the multi-platform distribution:

1. After downloading the zip or tar.gz file from BiGGESTS website, decompress it to a suitable location.
2. Execute the installer file inside the resulting directory (double-click the install.bat file on Windows or run the install.sh file on Mac OS). Wait for the installation to conclude.
3. You may now execute BiGGESTS anytime by executing the biggestes script file (double-click the biggestes.bat file on Windows or run the biggestes.sh on Mac OS).

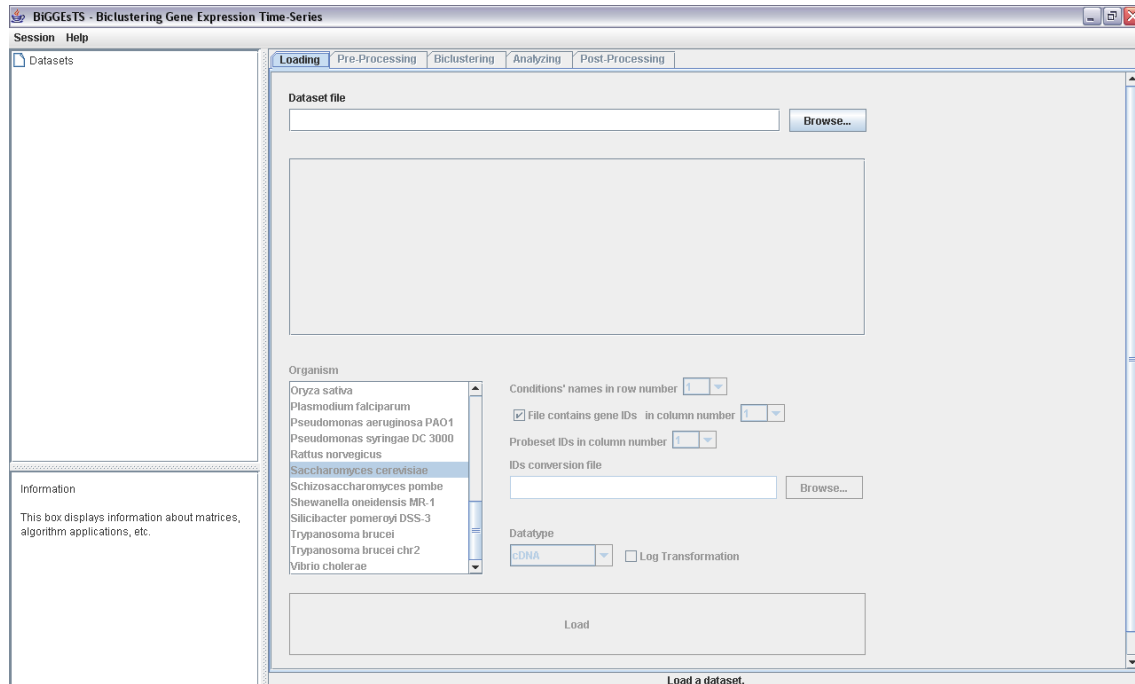
On Linux and other OSs, using the multi-platform distribution:

1. Install Graphviz on your system. You'll find the source code and binaries, as well as documentation, at <http://www.graphviz.org/>.
2. Download the zip or tar.gz file from BiGGESTS website and decompress it to a suitable location.
3. Edit BiGGESTS installer file (install.sh), appending the path to the dot binary file (usually /usr/bin/dot), preceded by a space, to the last line (e.g. "java -classpath biggestes.jar biggestes/utils/BiggestesInstall /usr/bin/dot").
4. Execute the install.sh script file.
5. You may now execute BiGGESTS whenever you want by running biggestes.sh script.

Running BiGGESTS

1. Loading a dataset

When you start BiGGESTS, the main window looks like this (program starts on the **Loading** functionality):



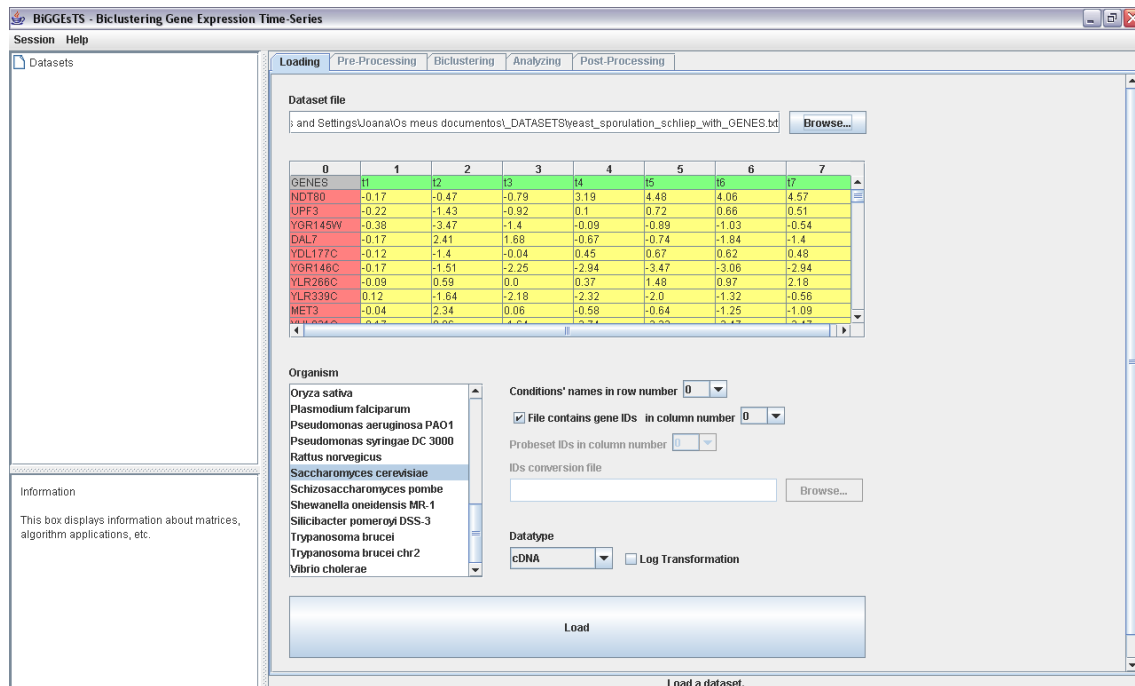
The BiGGESTS window.

You may want to **load a dataset** from a character delimited text file. Type the **path** to the file to load (or click **Browse...** button and select the dataset file in file system instead; sample files are also available in the Datasets directory, the default location for raw time series expression data, within the directory where BiGGESTS is installed; see the readme.txt file in the same directory for specific details on the contents of these sample files). You'll be presented a preview of the data in your file.

The dataset text file must be structured like this:

	C ₁	C ₂	C ₃	C ₄	C ₅	C ₆	...
G ₁	E ₁₁	E ₁₂	E ₁₃	E ₁₄	E ₁₅	E ₁₆	...
G ₂	E ₂₁	E ₂₂	E ₂₃	E ₂₄	E ₂₅	E ₂₆	...
G ₃	E ₃₁	E ₃₂	E ₃₃	E ₃₄	E ₃₅	E ₃₆	...
...

In which each G_x is a name of a gene, each C_y is a name of a condition (time-point) and each E_{xy} is the expression value of gene G_x in condition C_y. Values must be delimited by a specific character, which may be a **tab**, regular **space** or a **semicolon (;)**. The gene expression data may also have some missing values (blank gaps). BiGGESTS will deal with them later.



The input of time series gene expression data.

If the names of the genes in your file comply with an identification nomenclature system other than HGNC (HUGO Gene Nomenclature Committee), then you must uncheck the **File contains gene IDs in column number** checkbox, specify which column contains the names of the genes and provide an ID conversion file. This should be a character delimited text file (allowed delimiters are the same as the ones for dataset files) containing two columns: the first with the probeset IDs (the names of the genes used in the experiment) and the second with the corresponding HGNC names. You may also use BiGGEsTS with probeset IDs without converting them by leaving the **File contains gene IDs in column number** checkbox checked. Just be aware that, if you do this, function analysis won't produce valid results, since such names can't be matched with the ones used in the Gene Ontology files.

Select the number of the row which contains the **names of the experimental conditions** and the correct **data type**.

To perform a **log transformation** on the loaded data, check the **Log Transformation** checkbox. In that case, both **Original** and **Preprocessed** matrices containing the original and the log transformed data matrix, respectively, will be added to the tree.

Finally, press the **Load** button. BiGGEsTS will check if the Gene Ontology files are available for the organism that you have selected and if that is the case, it retrieves the biological functions that annotate the genes of the loaded dataset. This may take a few seconds but once it is done, one is able to check which functions annotate a given gene, by clicking on the corresponding row in the matrix.

2. Analyzing matrices

After loading the sample dataset, a new dataset is created in the dataset tree (on the left side of the window). An **Original** matrix is added to this new dataset, containing

the loaded gene expression data. By default, BiGGESTS shows the colored expression matrix (**Colors** tab), but you may also visualize the matrix with the original data by clicking on the **Values** tab (bottom of the window).

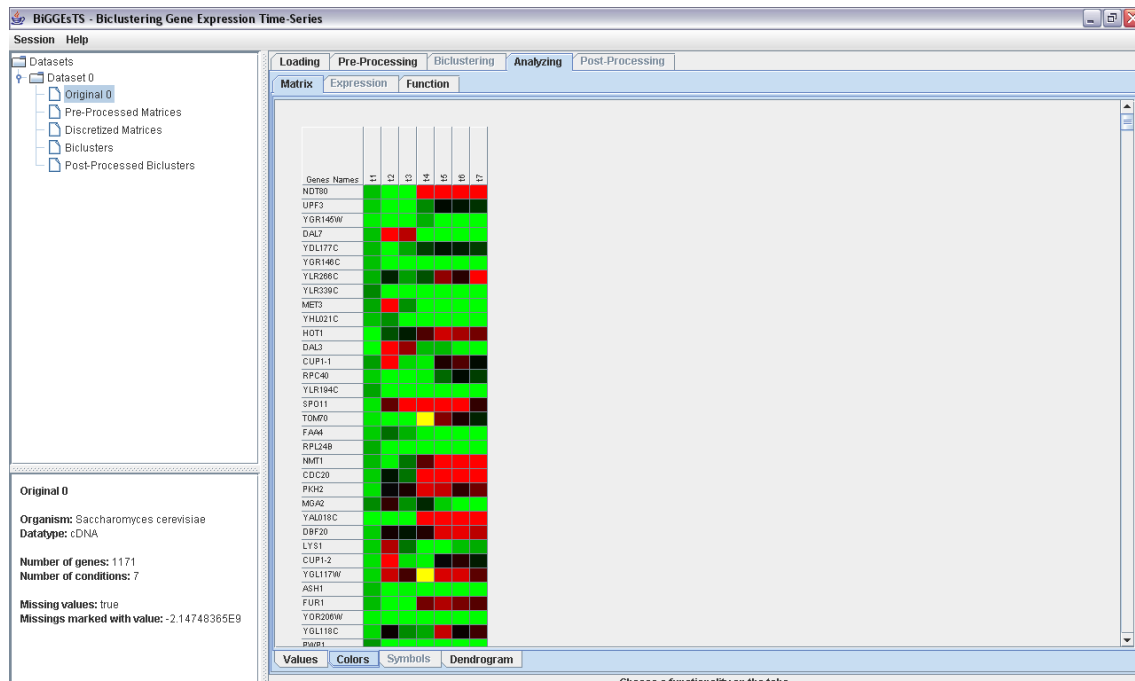


Table of colors for the original matrix. Lower values are colored green, higher values are colored red and missing values are set to yellow.

All kinds of matrices are **sortable**. This means that you can have the genes (rows) of a matrix sorted by their expression values/symbols for a given condition (column) in increasing or decreasing order. This is done by pressing the cell that corresponds to the condition, that you want the rows of the table to be sorted by, in the header of the table. The first press sorts the values in increasing order, a second one does it in decreasing order and a third causes the rows to be displayed in their original order.

Moving your mouse over one of the cells of the matrix displays a specific **text tip** according to the content of the cell.

Matrices can be **exported as image files**. The access to this functionality is provided via the right button of the mouse. When exporting a matrix as an image, you are prompted to provide a proper file name, location and format (PNG and JPG available).

As mentioned before, BiGGESTS is able to **extract biological functions that annotate the genes** of the dataset. If such information is available, it can be obtained by pressing a row of the table (left button of the mouse). The biological functions that annotate the corresponding gene are then displayed in a popup window, which additionally allows for selecting and copying the text content to paste on other documents.

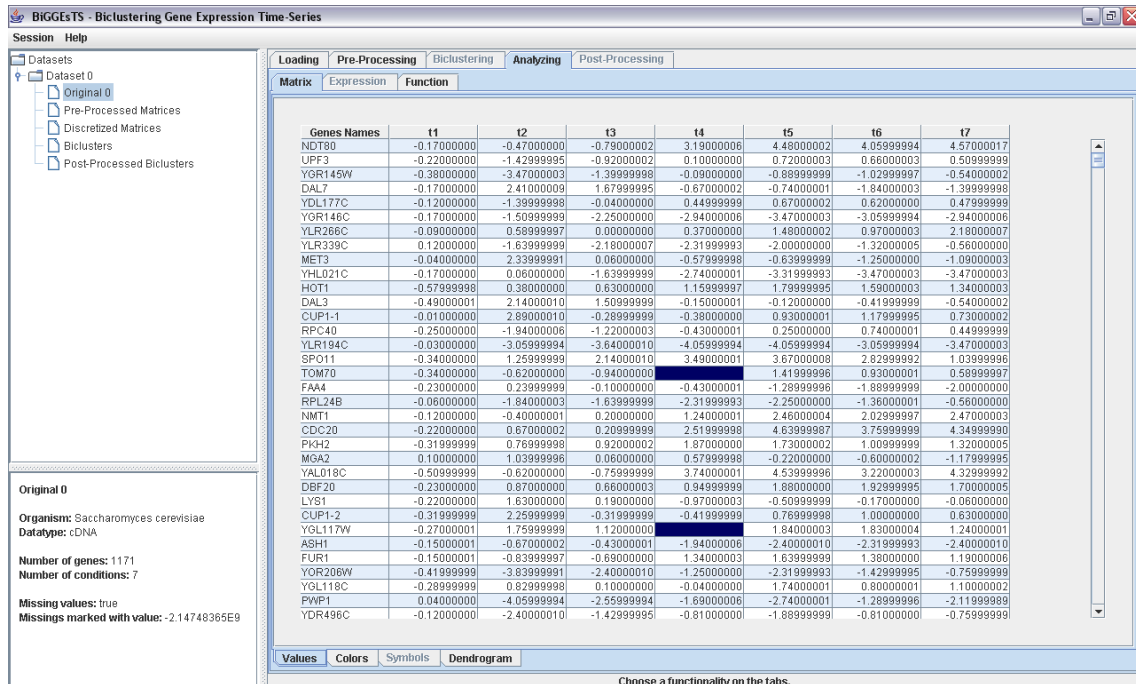


Table of values for the original matrix. Missing values are colored in dark blue.

After having loaded the time series gene expression data, BiGGESTS is quite intuitive to use. For every operation there are some basic steps to follow: select some data matrix or bicluster in the dataset tree and then choose the functionality that you want to perform on the tabs at the top (and bottom) of the window. The **input of gene expression data is always available** for every selected node in the dataset tree. Below is a summary of the additional main features, available upon the selection of a given type of node in the dataset tree together with a given set of tabs:

Original or Preprocessed matrices – (i) Preprocessing; (ii) Biclustering (CC-TSB, only available if the matrix has no missing values); (iii) Analyzing -> Matrix -> Values; (iv) Analyzing -> Matrix -> Colors; (v) Analyzing -> Matrix -> Dendrogram.

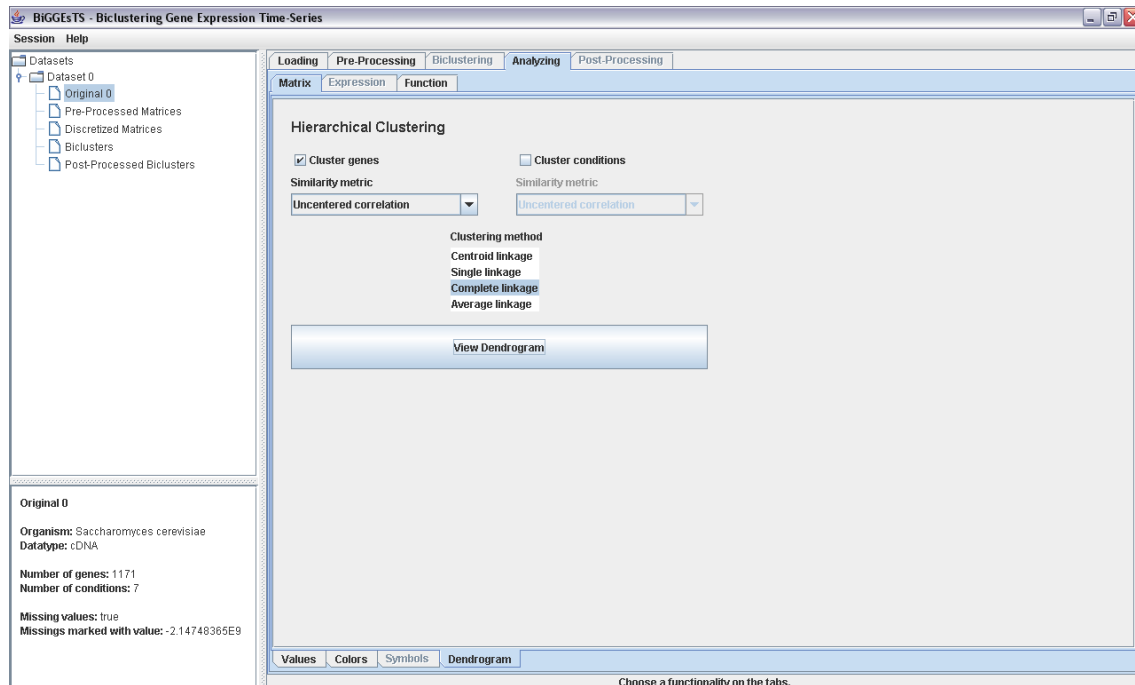
Discretized matrices – (i) Biclustering (CCC-Biclustering and e-CCC-Biclustering); (ii) Analyzing -> Matrix -> Values; (iii) Analyzing -> Matrix -> Colors; (iv) Analyzing -> Matrix -> Symbols; (v) Analyzing -> Matrix -> Dendrogram.

Biclusters Group or Post-Processed Biclusters Group – (i) Analyzing -> Matrix -> Values; (ii) Analyzing -> Matrix -> Colors; (iii) Analyzing -> Matrix -> Symbols (only for groups of biclusters obtained from discretized matrices); (iv) Analyzing -> Expression -> Bicluster time-points; (v) Analyzing -> Expression -> Bicluster pattern (only for groups of biclusters obtained from discretized matrices); (vi) Analyzing -> Function -> Table; (vii) Post-Processing.

Bicluster – (i) Analyzing -> Matrix -> Values; (ii) Analyzing -> Matrix -> Colors; (iii) Analyzing -> Matrix -> Symbols (only for biclusters obtained from discretized matrices); (iv) Analyzing -> Expression -> Bicluster time-points; (v) Analyzing -> Expression -> All time-point; (vi) Analyzing -> Expression -> Bicluster pattern (only for biclusters obtained from discretized matrices); (vii) Analyzing -> Function -> Table.

3. Visualizing hierarchical structure with dendrograms

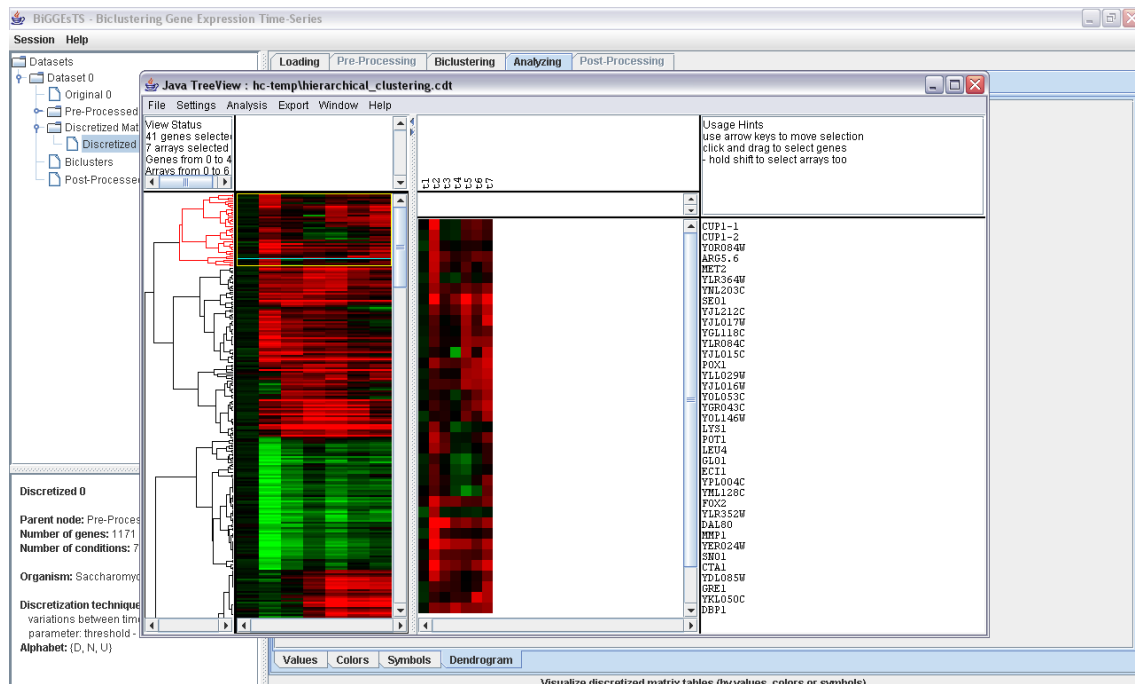
Hierarchical clustering can be useful for deriving hierarchical relationships based on the degree of similarity between the elements of the gene expression data, either genes or conditions. This technique is available via the **Analyzing, Matrix** and **Dendrogram** tabs, selected in this order.



Analyzing matrix dendrogram panel with the options for selecting and applying a hierarchical clustering algorithm to the gene expression data.

The clustering of both genes and conditions dimensions is available. We however note that when analyzing time series gene expression data, clustering conditions is not very meaningful, since they correspond to consecutive instants of time. As in the Cluster 3.0 software (<http://bonsai.ims.u-tokyo.ac.jp/~mdehoon/software/cluster/software.htm>), BiGGESTS lets you select the metric used for measuring the similarity between the elements of the gene expression matrix (either genes or conditions), either based on their correlation or distance, from the following list: (i) uncentered correlation, (ii) Pearson's correlation, (iii) uncentered absolute correlation, (iv) Pearson's absolute correlation, (v) Spearman's rank correlation, (vi) Kendall's tau correlation, (vii) Euclidean distance, and (viii) Cityblock distance (for details on these distances see <http://bonsai.ims.u-tokyo.ac.jp/~mdehoon/software/cluster/manual/Distance.html>). You may also specify the approach followed by the algorithm when computing the cluster-pairwise similarity: (i) centroid linkage, (ii) single linkage, (iii) complete linkage, and (iv) average linkage (for details on these techniques see <http://bonsai.ims.u-tokyo.ac.jp/~mdehoon/software/cluster/manual/Hierarchical.html>).

Once you have set up the correct parameters and pressed the **View Dendrogram** button, the results of the hierarchical clustering analysis are presented in a dendrogram displayed by a new Java TreeView window (for details on this application, please visit the official site of Java TreeView: <http://jtreeview.sourceforge.net/>).

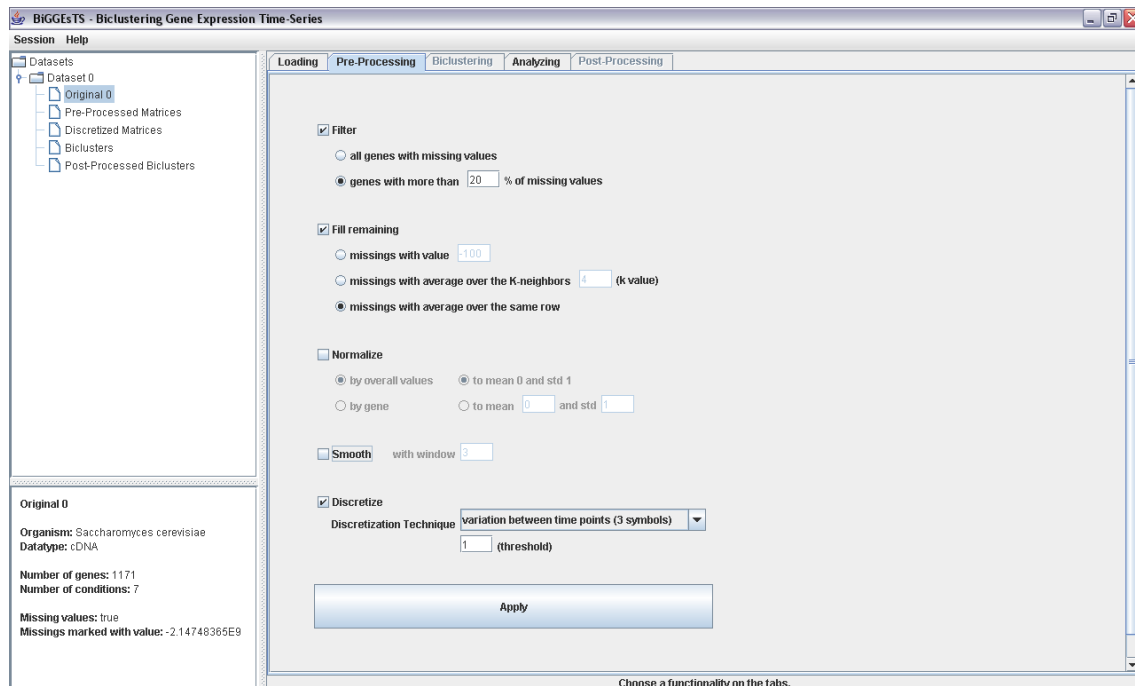


A dendrogram displayed (on the left) by the Java TreeView window superimposed on the BiGGESTS window. The figure further shows a cluster of genes selected (highlighted in red in the hierarchical structure, on the left). The cluster itself is displayed in the panel on the middle. The list of genes included in the cluster is displayed in the panel on the right.

4. Preprocessing time series gene expression data

The next step we want to exemplify involves the preprocessing of the **Original** expression matrix. Select the **Preprocessing** tab (on the top of the window), after having the matrix that you want to preprocess selected on the dataset tree. Preprocessing includes the following techniques: (i) **gene filtering**, for filtering genes with missing values and only available for matrices with missing values; (ii) **missing values filling**, for filling missing data with real values; (iii) **data normalization**, to compensate for systematical differences between data measured by the several microarrays/conditions; (iv) **smoothing**, for reducing the impact of the noise in the analysis; and (v) **discretization**, for reducing the infinite set of real gene expression values to an adequate range of discrete values. Note that the data normalization, smoothing and discretization techniques can be applied to matrices with missing values with no previous treatment, because they use an appropriate approach to minimize the impact of the missing values. Once you have selected the preprocessing techniques and set their corresponding parameters, you can press the **Apply** button.

Upon the selection of any or several of the first four preprocessing options (filter, fill, normalize, smooth) and disabling of the discretization one, only a **Preprocessed** node is added to the dataset tree. If the discretization is enabled, an additional **Discretized** node is also created and added to the dataset tree. When selecting more than one preprocessing option, the several options are applied one at a time by BiGGESTS to the gene expression data following the order of the options displayed, from top to bottom, in the preprocessing panel.



Preprocessing panel displaying the preprocessing options.

The names of the preprocessing steps are quite explicit, but we also list them here, including a brief description of their parameters:

Filter – includes two options for removing genes (rows of the matrix) with missing values: (i) **all genes with missing values** eliminates all rows which contain missing values; (ii) **genes with more than x % of missing values** eliminates all rows whose percentage of missing values exceeds the value of x.

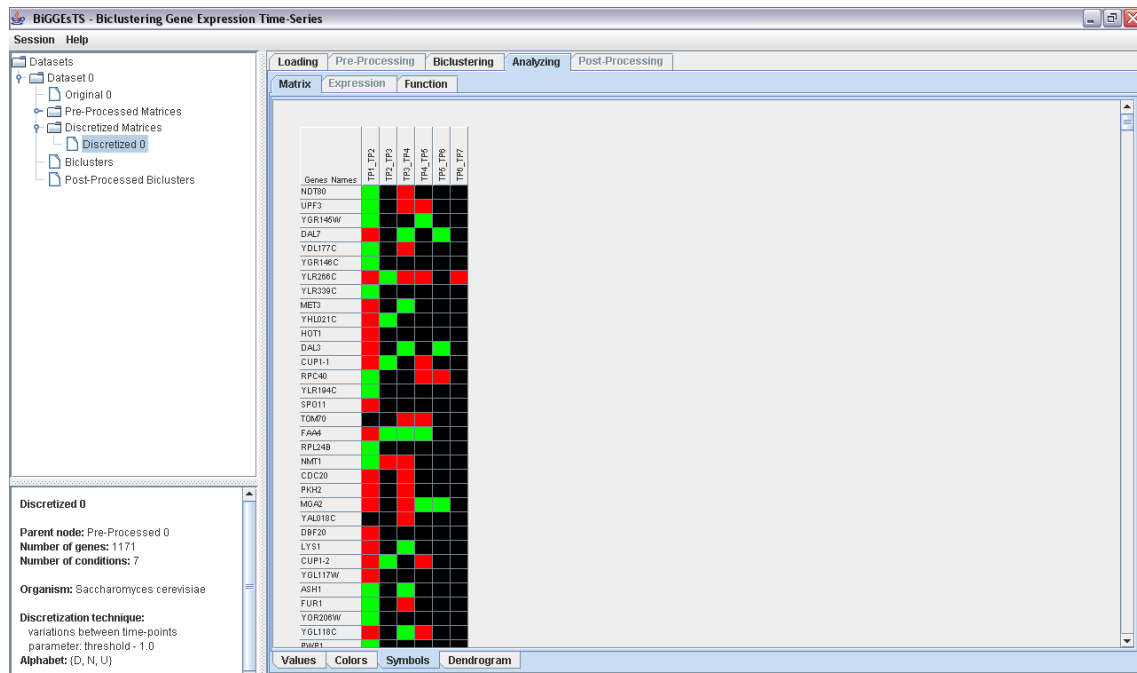
Fill remaining – provides three alternatives to fill the values missing from the expression matrix with: (i) a given **value**, which has to be typed in the corresponding text field; (ii) the **average of the values of the k-neighbor cells** of the same gene (row); (iii) the **average over all the values of the same gene** (row).

Normalize – normalizes the expression values; they can be normalized altogether using the (i) **by overall values** option to mean 0 and standard deviation 1 or to a given mean and standard deviation, which have to additionally be specified in the corresponding text fields; or by row using the (ii) **by gene** option to mean 0 and standard deviation 1.

Smooth – acts like a low-pass filter for attenuating the negative effect of outliers; requires a parameter: the **length of a window** of neighbors to consider when computing the new value for substituting an outlier (note that you must provide an odd value, since the window includes the outlier value).

Discretize – provides a number of different discretization techniques; a first group of methods may be applied to the overall values of the matrix or by gene and computes the corresponding discrete value for each real element: (i) **expression average**; (ii) **mid-range**; (iii) **max-minus percent-max**; (iv) **equal frequency**; (v) **equal width**; (vi) **expression mean and standard deviation**; the second group of methods computes the

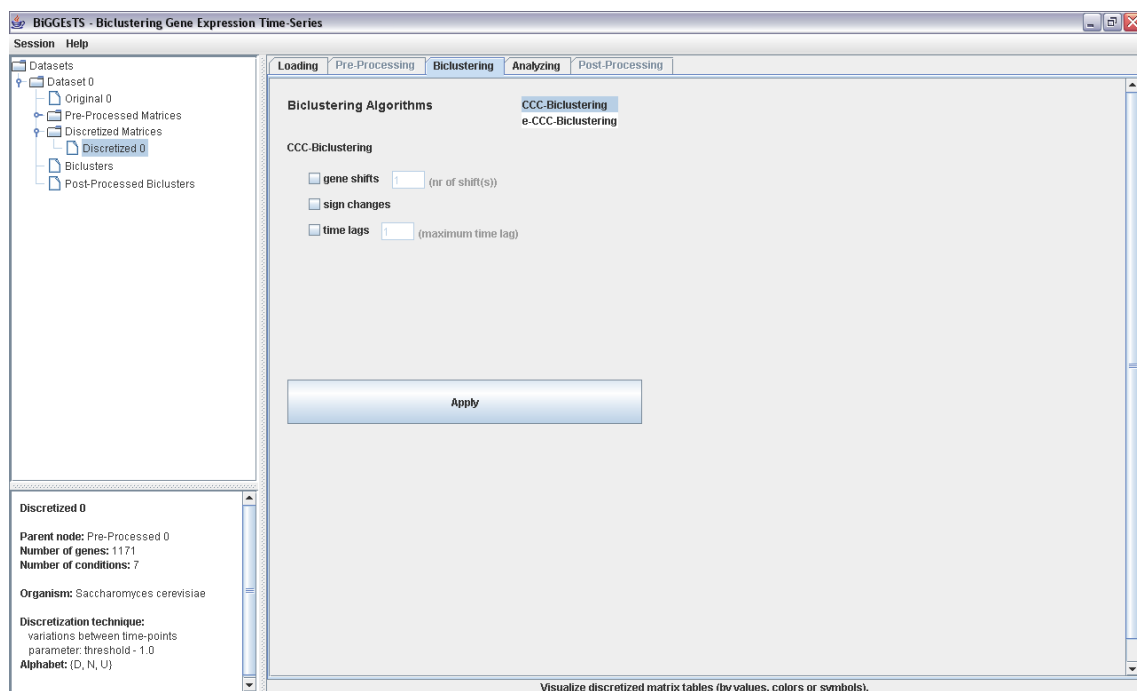
discrete values based on the variation of the expression level between pair of consecutive conditions: (vii) **transitional state discriminations**; (viii) **variation between time points**.



Heatmap of a discretized matrix (obtained from the discrete symbols instead of the real values). You may check which symbol is actually contained in each cell by moving the mouse over it (text tip).

5. Biclustering time series expression data

Accessing the **Biclustering** tab, you are able to choose the biclustering algorithm and parameterize it.



CCC-Biclustering with the default parameterization.

BiGGEsTS integrates three biclustering algorithms: one for real data with no missing values, CC-TSB-Biclustering; and two for discretized matrices that are also able to compute in the presence of missing values, CCC-Biclustering and e-CCC-Biclustering. In our example, we will apply the CCC-Biclustering algorithm to the previously mentioned **Discretized** matrix. We select the CCC-Biclustering algorithm, take the default parameterization and press the **Apply** button. Next, we also apply the CCC-Biclustering algorithm, but this time enabling the sign changes option.

The available biclustering algorithms and parameterizations are:

1. **CCC-Biclustering** – for finding biclusters with exact expression patterns in discrete data; when in the presence of missing values, the algorithm uses an appropriate approach for disregarding them from the analysis of the gene expression values.
2. **e-CCC-Biclustering** – for finding biclusters with approximate expression patterns (maximum of e errors) in discrete data; required parameters: the **maximum number of errors** allowed per pattern; when activated, the **restricted errors** variation considers as valid errors only the substitutions of symbols which are on a given neighborhood in the alphabet of discretization. When in the presence of missing values, the algorithm may follow one of two approaches: ignore them, as in the case of the CCC-Biclustering algorithm, or consider them as errors. The default behavior is to ignore the missing elements.

Common variations of both CCC-Biclustering and e-CCC-Biclustering:

- (i) **gene shifts**, for grouping in biclusters genes with similar expression evolutions, but at different expression levels; parameter: the **number of shifts**, that is, of different expression levels to consider;
- (ii) **sign changes**, for including in biclusters genes with symmetric expression patterns;
- (iii) **time lags**, for considering in the same biclusters genes with similar expression evolutions, but starting at different points in time, in a predefined order as generated by temporal programs; parameter: the **maximum time lag** allowed between expression patterns;

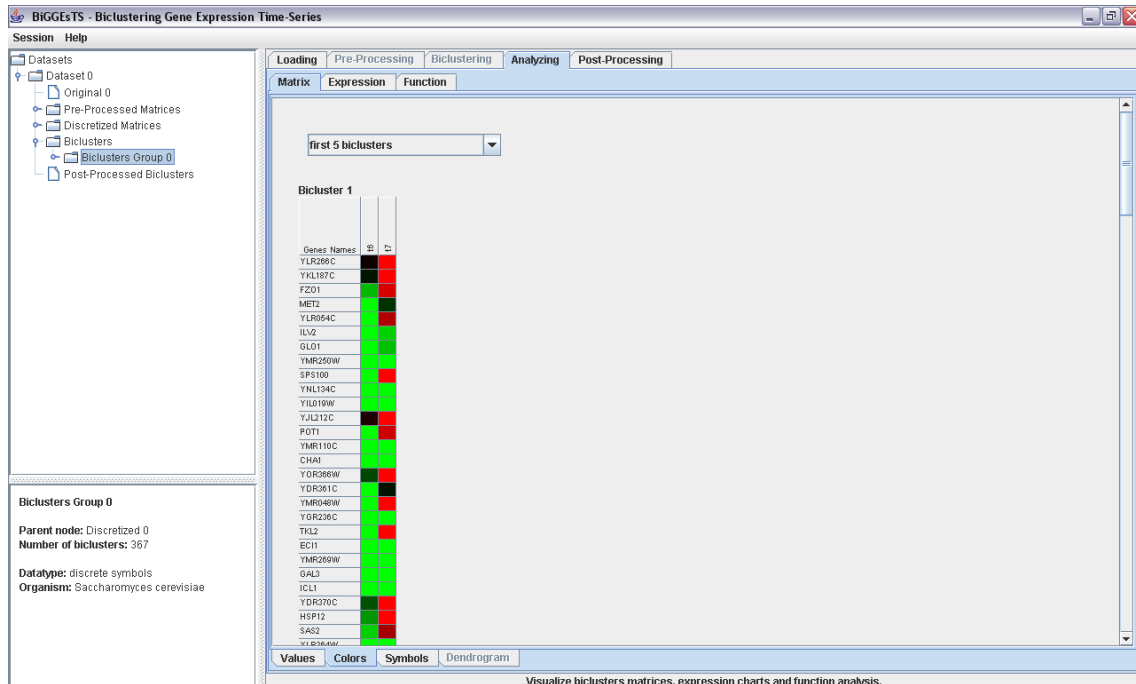
Note: gene shifts cannot be combined with sign changes or time lags.

3. **CC-TSB-Biclustering** – for finding biclusters with approximate expression patterns in real data; parameters: (i) **delta**, the threshold for the MSR of each bicluster; (ii) **alpha**, the ratio between the MSR of a row and the MSR of the matrix; (iii) the **number of biclusters to extract**; (iv) the **maximum number of iterations**.

Note that the computation may take time, especially when a large matrix is involved. Sometimes, if the dimension of the matrix is really demanding, computation may also end up aborting, usually due to memory issues.

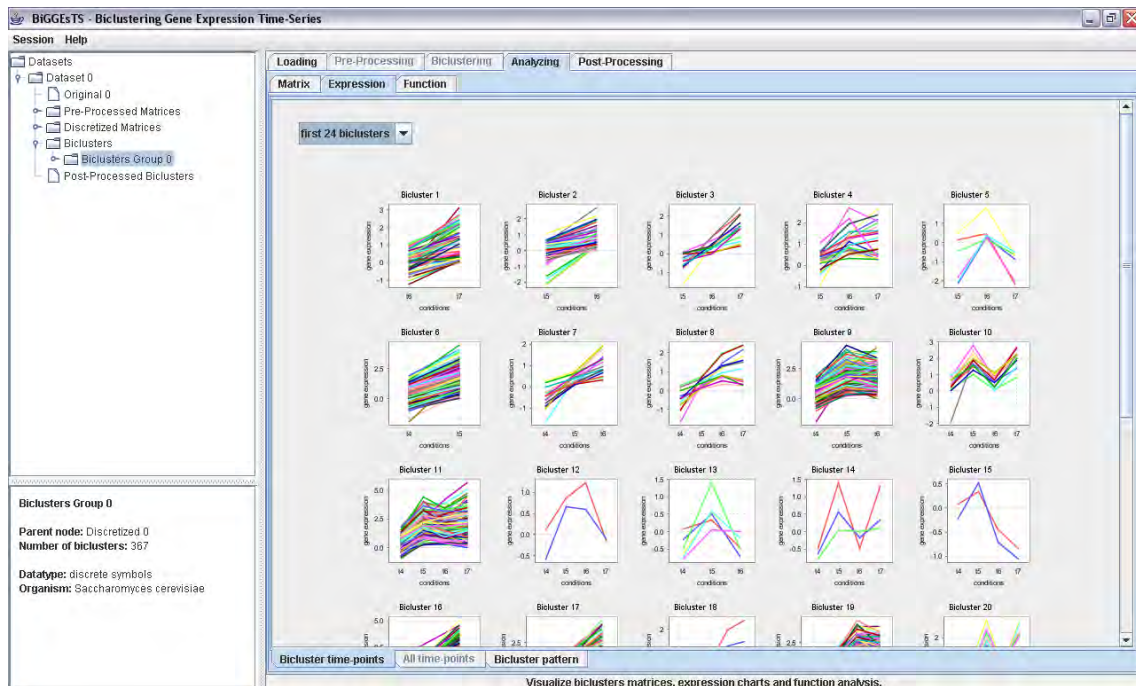
6. Visualizing biclustering results

Upon a successful application of a biclustering algorithm, a group of biclusters is added to the tree, the **Analyzing**, **Matrix** and **Colors** tabs are selected and the colored matrices of the first 5 biclusters of the group are displayed by default. This number can be changed using the combo box available in the **Colors** panel. The matrices of values and symbols are accessed by selecting the **Values** or **Symbols** tab, respectively.



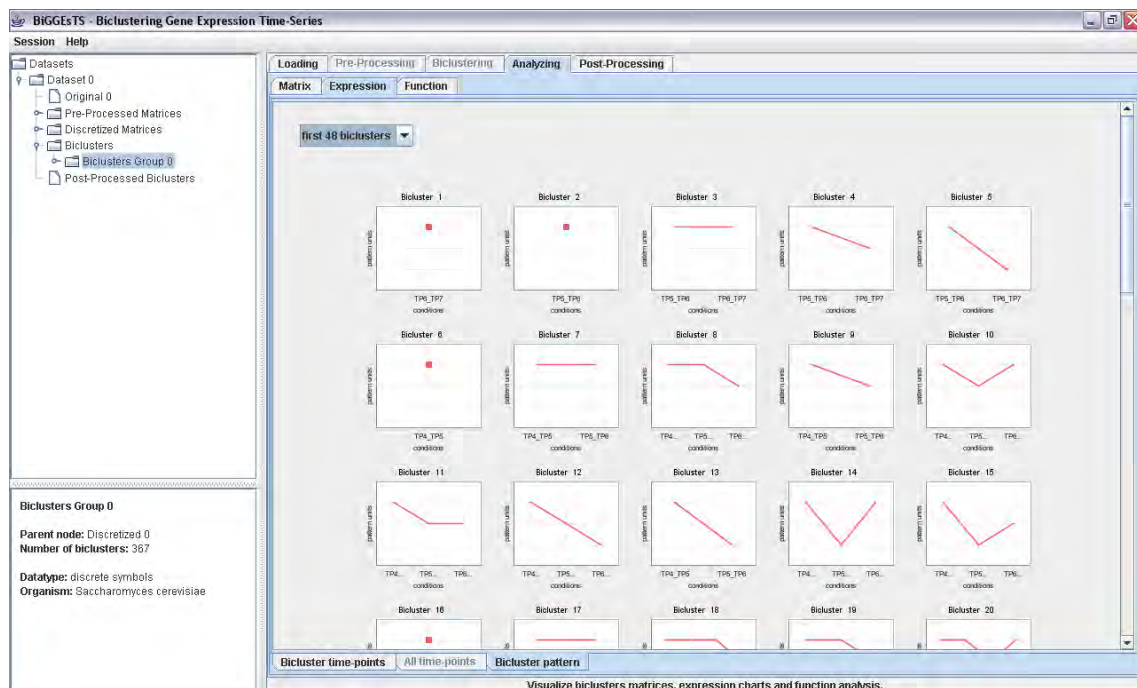
Heatmap (matrix of colors) of the first 5 CCC-biclusters (only the first matrix is visible).

BiGGEsTS enables the display of miniaturized expression charts of biclusters in a group.

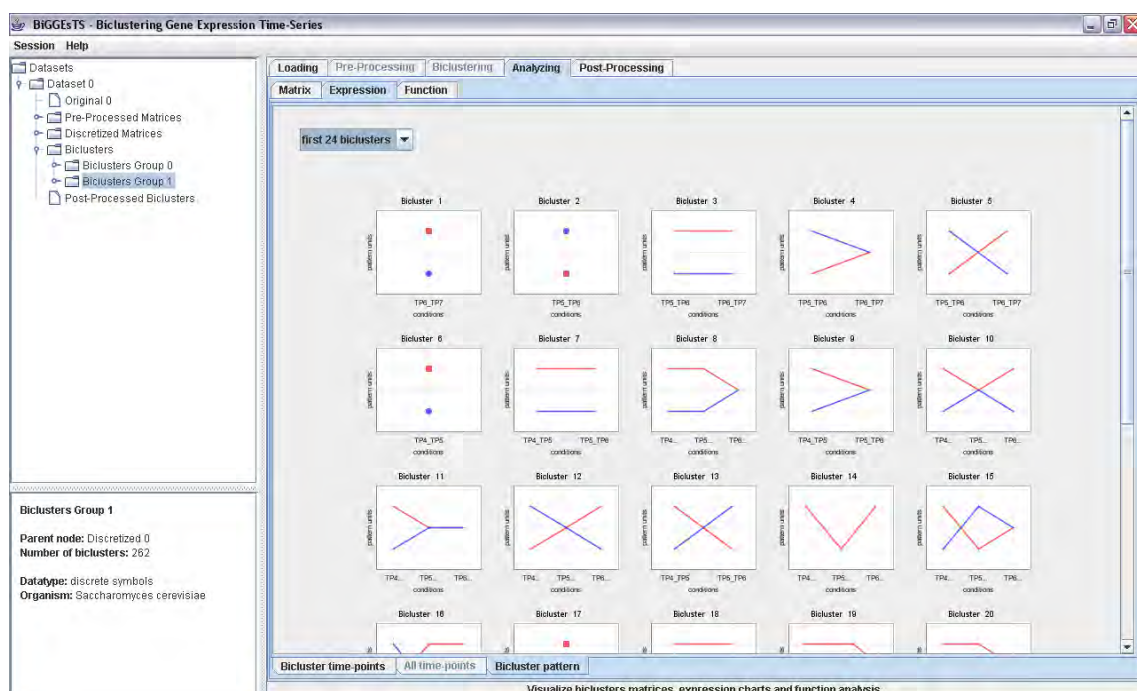


Miniatures of the time-points charts of the first 20 CCC-biclusters identified using the default CCC-Biclustering algorithm.

The expression charts and pattern charts are available upon selection of the **Analyzing, Expression, Bicluster time-points** tabs (in this order) and the **Analyzing, Expression, Bicluster pattern** tabs (in this order), respectively. Note that when biclusters are large, that is, composed of many genes and/or conditions, or when you are trying to display the charts of a considerable number of biclusters in a group, BiGGESTS may take some time to draw all these data. By default, only the first 12 biclusters are displayed. You may change this number in the corresponding combo box on the top of the panel.



Miniatures of the pattern charts of the first 20 biclusters identified using the default CCC-Biclustering algorithm.



Miniatures of the pattern charts of the first 20 biclusters identified by the sign changes CCC-Biclustering.

To access the individual biclusters you must open the folder of the group of biclusters in the dataset tree (try pressing the key on the left of the group folder). Once the group is opened you may select a specific bicluster and view its corresponding information. Matrices of values, colors and symbols are displayed in a similar fashion to the one used for original, preprocessed and discretized matrices. Bicluster expression charts, all time-points expression charts and bicluster pattern charts are upon selection of the **Bicluster time-points**, **All time-points** and **Bicluster pattern** tabs, respectively (after selecting **Analyzing** and **Expression** tabs).

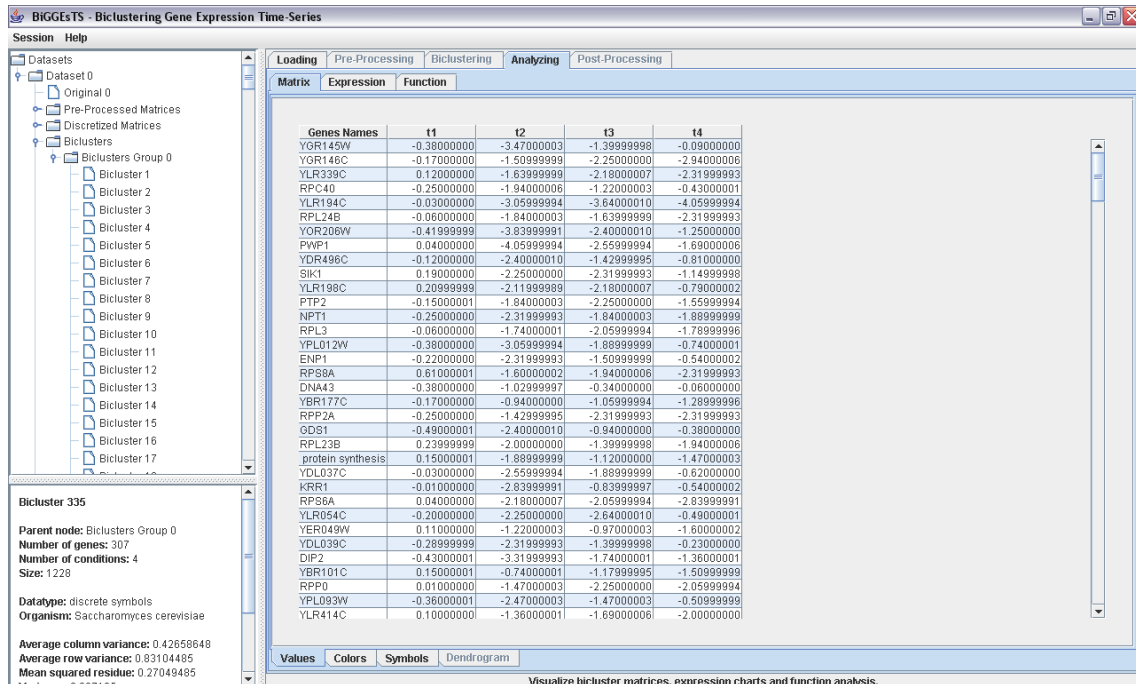


Table of values of the bicluster 335 identified using the default CCC-Biclustering algorithm.

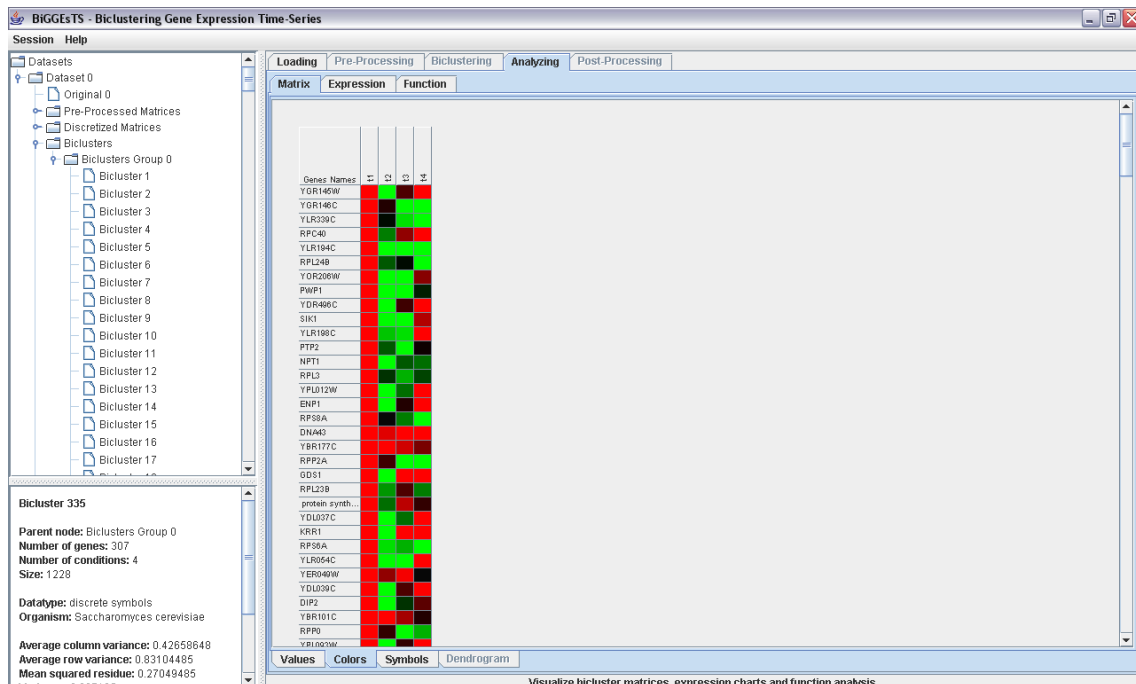


Table of colors of the bicluster 335 identified by the default CCC-Biclustering algorithm.

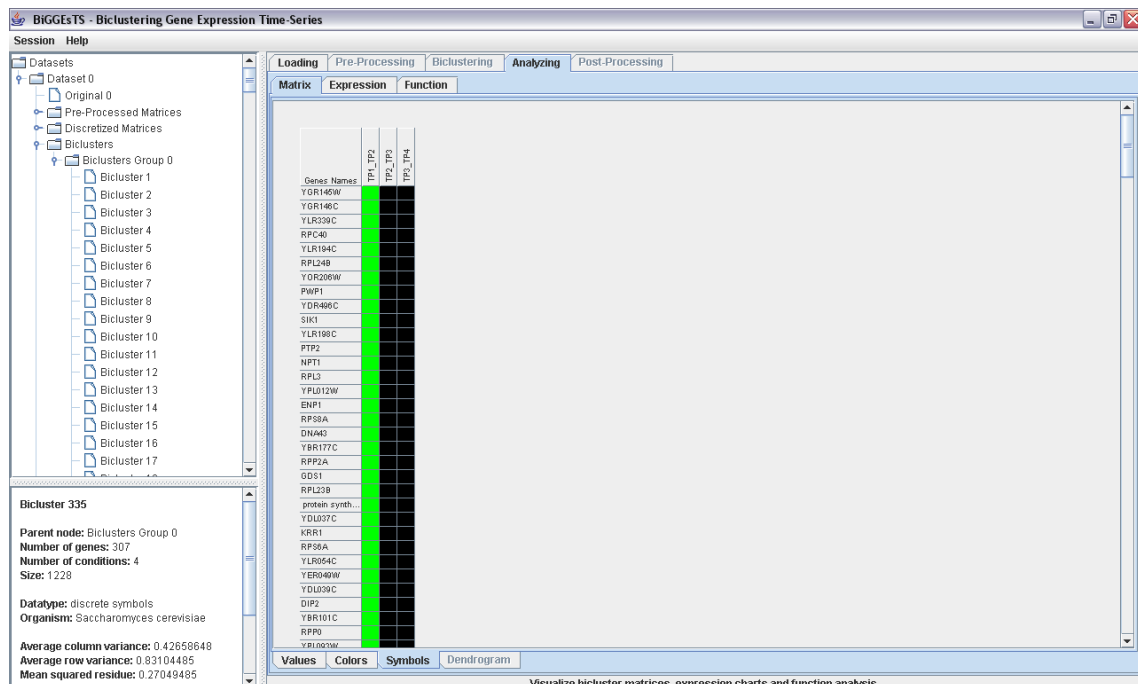


Table of symbols of the bicluster 335 identified using the default CCC-Biclustering algorithm.

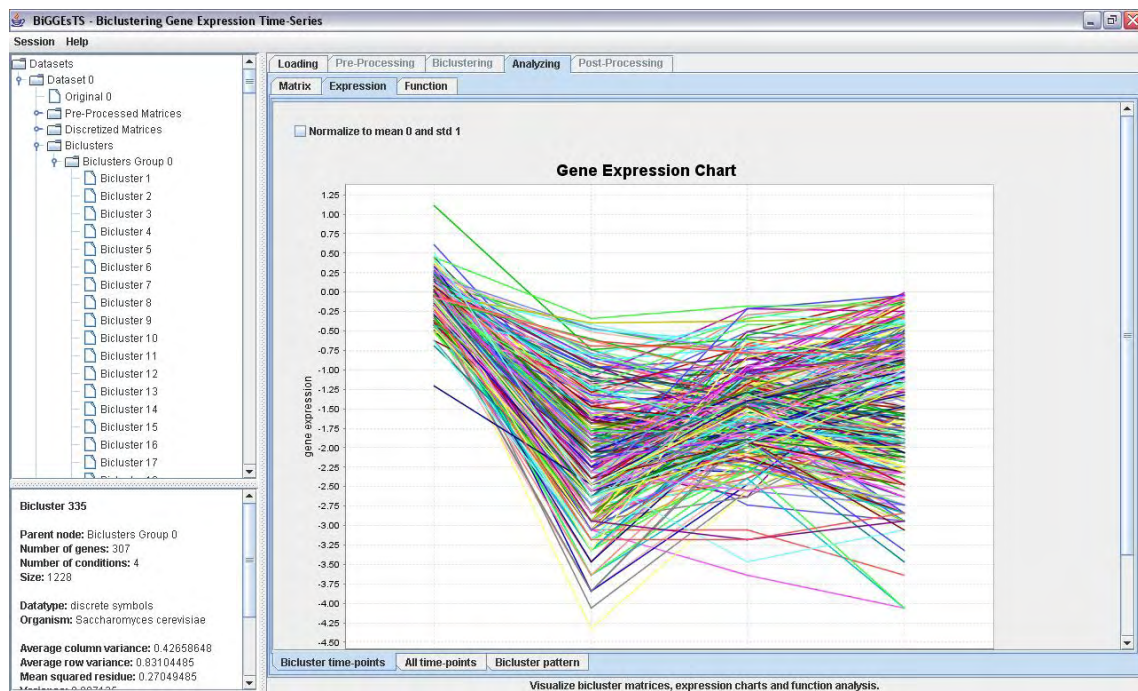


Chart of the evolution of the expression level of the genes in the CCC-bicluster 335 along the corresponding conditions of the bicluster. It is possible to normalize the expression levels by checking the **Normalize to mean 0 and std 1** checkbox.

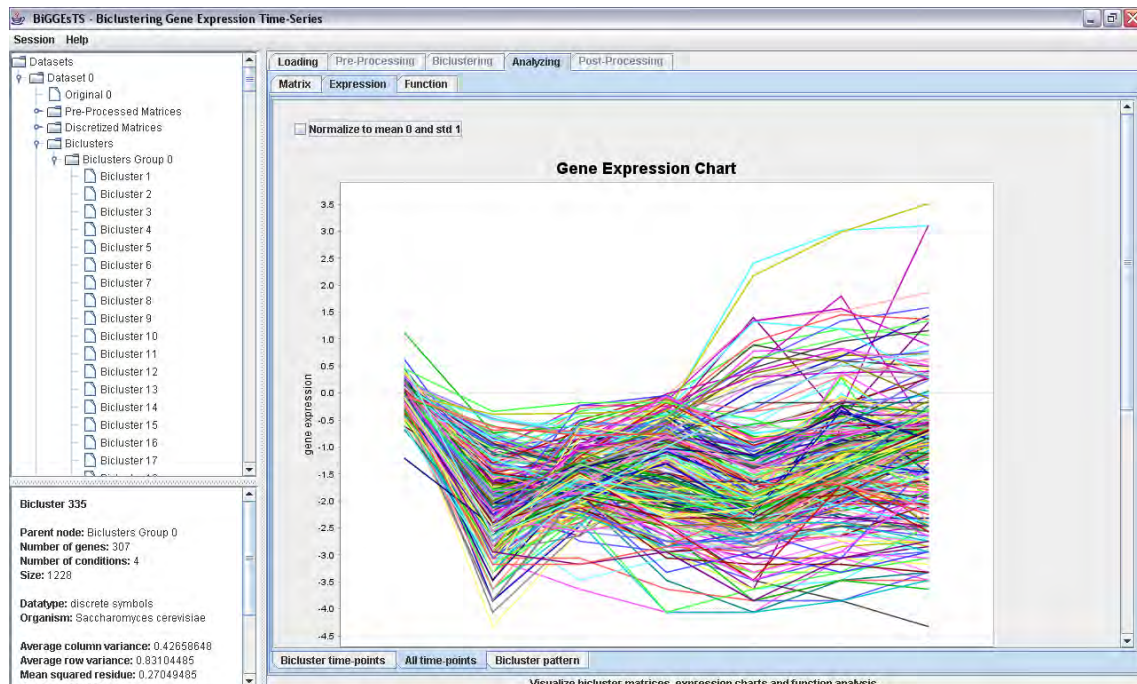


Chart displaying the evolution of the expression level of the genes in the CCC-bicluster 335 along all the conditions of the dataset. It is possible to normalize the expression level by checking the **Normalize to mean 0 and std 1** checkbox.

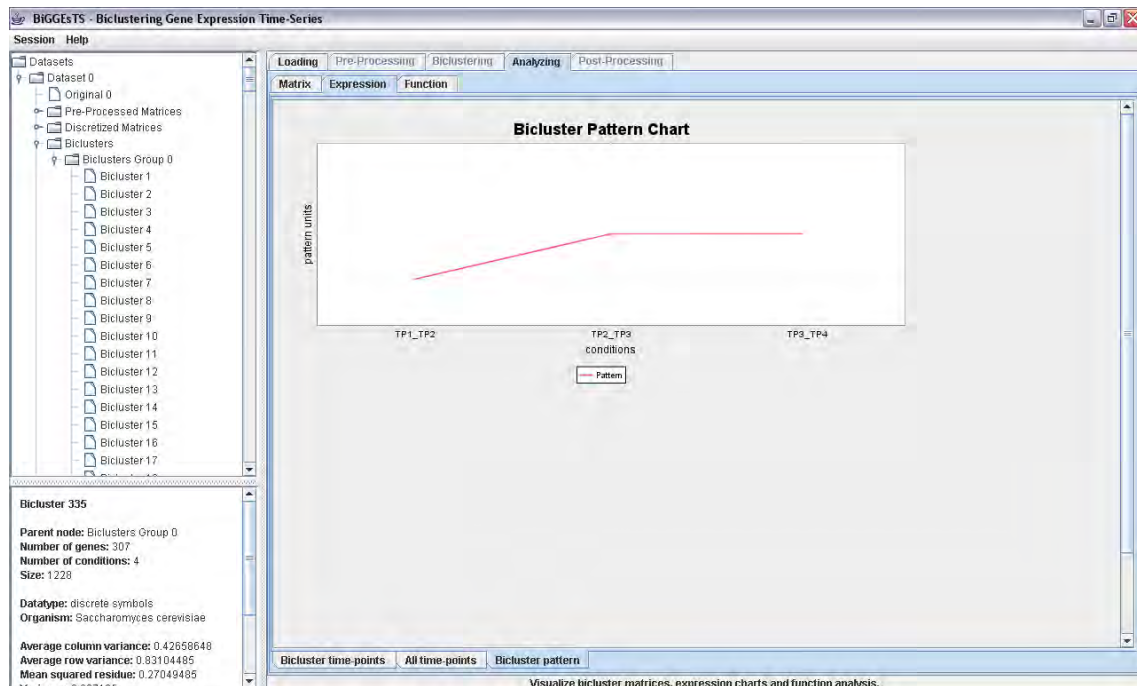


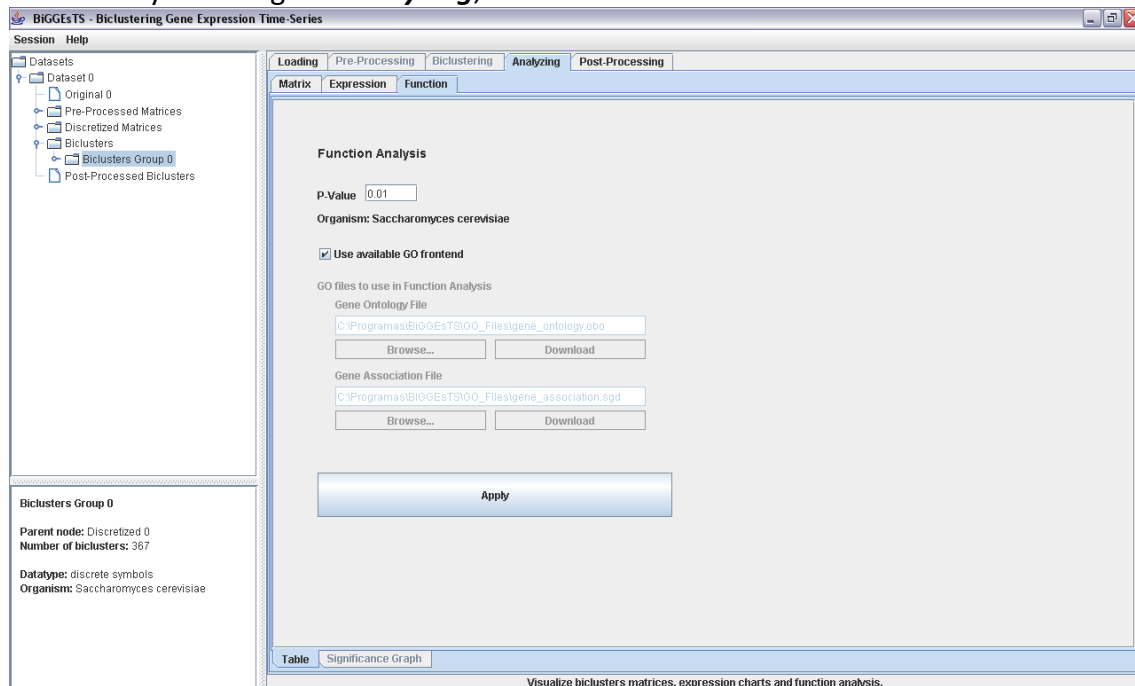
Chart displaying the expression pattern of the CCC-bicluster 335.

All charts, except the miniature ones, allow **saving as image**, **printing** and **zooming**. These options can be accessed by pressing the right button of the mouse on the chart.

7. Analyzing biological functions

BiGGESTS automatically extracts the biological functions (terms) that annotate the genes of the dataset when the Gene Ontology files are available. For biclusters and

groups of biclusters, it is additionally possible to perform a functional enrichment using the term-for-term analysis, which computes the statistical significance of each biological term that annotates the genes in the biclusters. The term-for-term analysis is available by selecting the **Analyzing, Function** and **Table** tabs in this order.



Parameters for the term-for-term analysis.

The required parameters for term-for-term analysis are: a **p-value**, the threshold below which the Bonferroni corrected p-values computed for the biological functions are considered statistically significant; the general **ontology** and specific organism **annotation files**. If these files are available at the proper location (the GO_Files directory within the BiGGESTS installation directory), their corresponding file paths are displayed in the text boxes. Also note that when this is the case, most likely BiGGESTS has already used these files to extract the biological functions that annotate the genes in a previous step. When such happens, the **Use available GO frontend** check box appears selected, which means that the annotations have already been retrieved and the term-for-term analysis can use them, avoiding to repeat the parsing of the Gene Ontology files. This accelerates the computation process of the term-for-term analysis.

Additionally, BiGGESTS enables downloading the files from the Gene Ontology repository by pressing the **Download** button below the corresponding file path text box.

The results of the term-for-term analysis are displayed in a table. Significant and highly significant terms are highlighted in green. The threshold for statistical high significance is always 0.01. The one for statistical significance is also 0.01 by default, but can be changed to a different value.



Values calculated using p-value = 0.01 Recalculate

Bicluster ID	# Genes	# Conditions	Best p-value	Best corrected p-value	# Significant terms	# Highly sig terms	Significance threshold	# Sig terms (th)
1	68	2	0.00004522	0.02613557	59.0	0.0	0.01000000	0
2	44	2	0.00019324	0.08792201	57.0	0.0	0.01000000	0
3	15	3	0.00000677	0.00133282	42.0	1.0	0.01000000	1
4	23	3	0.01823254	1.00000000	49.0	0.0	0.01000000	0
5	6	3	0.00512383	0.73783094	37.0	0.0	0.01000000	0
6	213	2	0.00108959	1.00000000	66.0	0.0	0.01000000	0
7	16	3	0.00896492	1.00000000	58.0	0.0	0.01000000	0
8	13	4	0.00591784	1.00000000	43.0	0.0	0.01000000	0
9	191	3	0.00069930	0.71755838	69.0	0.0	0.01000000	0
10	12	4	0.01024765	1.00000000	29.0	0.0	0.01000000	0
11	177	4	0.00060708	0.50859731	82.0	0.0	0.01000000	0
12	2	4	0.86003637	1.00000000	0.0	0.0	0.01000000	0
13	6	3	0.00032546	0.04198460	32.0	0.0	0.01000000	0
14	3	4	0.00256191	0.18445773	23.0	0.0	0.01000000	0
15	2	4	0.00002190	0.00129190	25.0	5.0	0.01000000	5
16	305	2	0.00000000	0.00000001	123.0	12.0	0.01000000	12
17	118	3	0.00009579	0.00476688	118.0	1.0	0.01000000	1
18	2	5	0.03389148	1.00000000	5.0	0.0	0.01000000	0
19	116	4	0.00004473	0.03836694	120.0	0.0	0.01000000	0
20	7	5	0.00597780	0.69342440	12.0	0.0	0.01000000	0
21	108	5	0.00002470	0.01966528	138.0	0.0	0.01000000	0
22	175	3	0.00000000	0.00000021	106.0	13.0	0.01000000	13
23	5	5	0.00426985	0.51865241	35.0	0.0	0.01000000	0
24	170	4	0.00000000	0.00000010	102.0	14.0	0.01000000	14
25	3	5	0.00767260	0.27621377	12.0	0.0	0.01000000	0

Results of the term-for-term analysis applied to the genes in the group of biclusters. It's possible to recalculate the significance of the biological terms using different Gene Ontology files or p-value threshold by pressing the **Recalculate** button. Pressing a row of this table will select the corresponding bicluster in the dataset tree and redirect the content of the panel to the results of the following selection of tabs: Analyzing, Matrix, Colors.

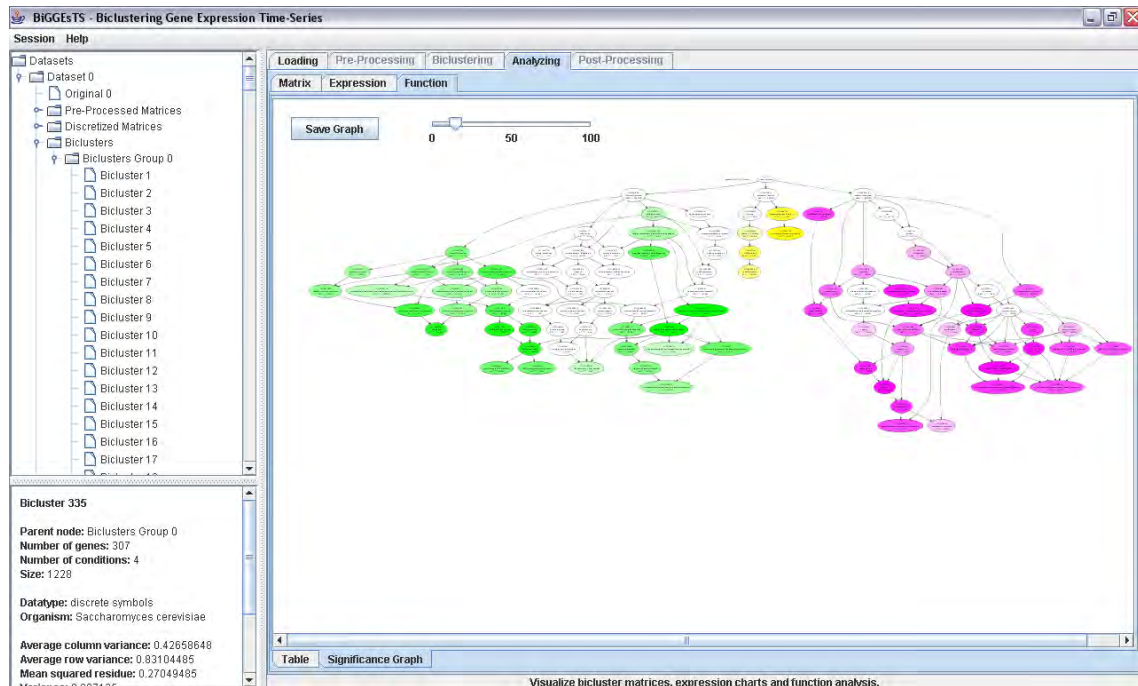
Values calculated using p-value = 0.01 Recalculate

Number of Genes in Population: 1171
Number of Genes in Bicluster: 307

GO Term ID	GO Term Name	#A. Population G.	#A. Set G.	P-Value	Corrected P-Value
GO:000607	ethanol metabolic process	6	3	0.18865435	1.00000000
GO:000415	negative regulation of histone H3-K36 methylation	1	1	0.26216909	1.00000000
GO:0051321	meiotic cell cycle	55	2	0.99999923	1.00000000
GO:0007093	mitotic checkpoint	5	1	0.78199703	1.00000000
GO:0008793	phosphorus metabolic process	32	5	0.95048198	1.00000000
GO:0007049	cell cycle	119	7	1.00000000	1.00000000
GO:0007154	cell communication	27	10	0.14242116	1.00000000
GO:0044262	cellular carbohydrate metabolic process	60	13	0.83463967	1.00000000
GO:0009068	aspartate family amino acid catabolic process	4	2	0.28225508	1.00000000
GO:0009260	ribonucleotide biosynthetic process	11	4	0.32023612	1.00000000
GO:0006090	pyruvate metabolic process	14	8	0.01339011	1.00000000
GO:0042594	response to starvation	4	2	0.28225508	1.00000000
GO:0006344	RNA processing	70	66	0.00000000	0.00000000
GO:0031060	regulation of histone methylation	1	1	0.26216909	1.00000000
GO:0006658	phosphatidylserine metabolic process	1	1	0.26216909	1.00000000
GO:0043038	amino acid activation	3	1	0.59869504	1.00000000
GO:0000082	G1/S transition of mitotic cell cycle	9	3	0.43351528	1.00000000
GO:0000375	RNA splicing, via transesterification reactions	7	1	0.88171804	1.00000000
GO:0006873	cell ion homeostasis	14	5	0.29422271	1.00000000
GO:0043624	cellular protein complex disassembly	3	2	0.16892256	1.00000000
GO:0006094	gluconogenesis	11	5	0.13407262	1.00000000
GO:0006886	intracellular protein transport	35	11	0.29550899	1.00000000
GO:0048164	alcohol catabolic process	20	11	0.00546348	1.00000000
GO:0045091	nitrogen fixation	3	1	0.59869504	1.00000000

Results of the term-for-term analysis applied to the genes in the bicluster 335.

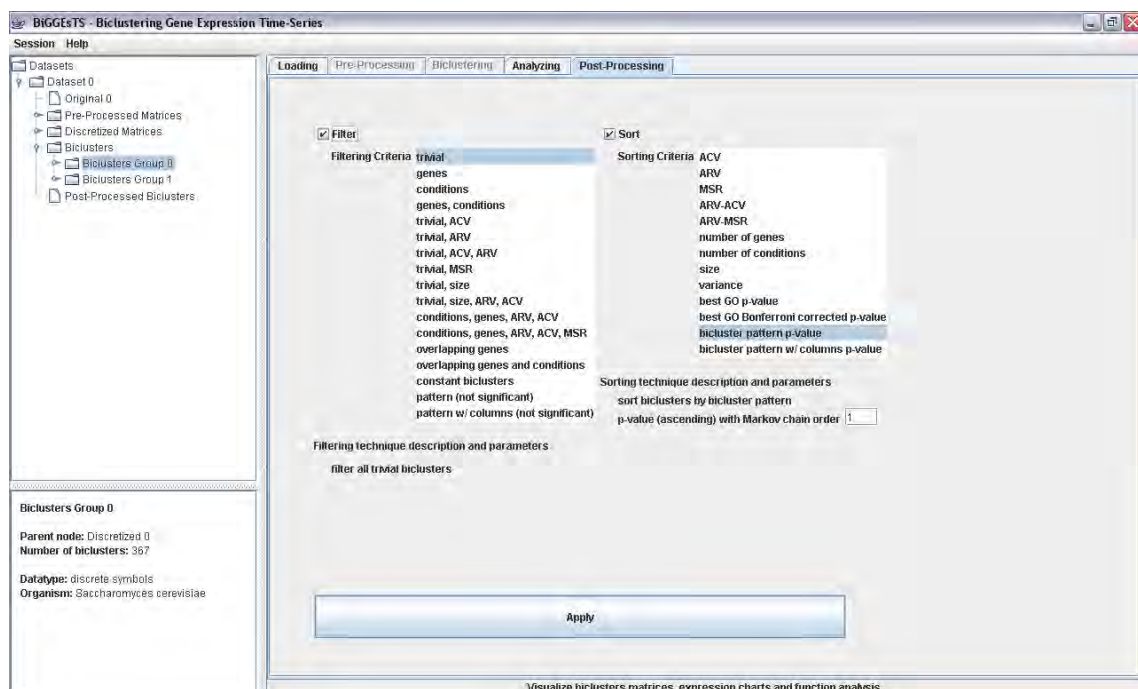
When the term-for-term analysis reveals significant functions (highlighted in green), the **Significance Graph** tab is enabled and the corresponding panel displays a graph with the ontology of the biological terms. The significant terms are highlighted in green, yellow or purple depending on which of the main biological function they specialize from. You may further **zoom** the graph or **save** it as a **PNG** or **SVG** image.



Graph of significant biological functions extracted by the term-for-term analysis for the genes in the bicluster 335.

8. Post-Processing groups of biclusters

The groups of biclusters can be post-processed (**Post-Processing** tab). This operation consists in filtering and/or sorting the biclusters of a group of biclusters based on specific criteria, available for selection in the post-processing panel. A short description of each filtering/sorting criterion and parameters is provided below the corresponding list in the panel, upon the selection of a given item. To apply the post-processing operations to the data press the **Apply** button.



Post-Processing panel: available filtering and sorting options.

Below is a list of the available post-processing criteria for filtering and/or sorting the biclusters in a group of biclusters:

Filtering

Filtering processes the group of biclusters and removes biclusters which: (i) are **trivial**, that is, are composed by a single row or a single column; contain (ii) a **number of genes** less than a given threshold, (iii) a **number of conditions** less than a given threshold, (iv) a **number of genes and a number of conditions** less than two corresponding thresholds; have (v) an **average column variance (ACV)** greater than a given threshold, (vi) an **average row variance (ARV)** less than a given threshold, (vii) an ACV greater than and an ARV less than two given thresholds, (viii) a **MSR** greater than a given threshold.

BiGGESTS can also eliminate biclusters: whose (ix) **size**, the number of genes times the number of conditions, is less than a given threshold; which satisfy a combination of the previously mentioned criteria, such as the thresholds for their (x) **size, ARV and ACV**, (xi) **numbers of conditions and genes, ARV and ACV**, (xii) **numbers of conditions and genes, ARV, ACV and MSR**; which are very similar to other biclusters in the group, according to a given **percentage of similarity** on the dimension(s) of (xiii) the **genes** or (xiv) both **genes and conditions**; which (xv) are **constant**, that is, have no variation of the level of expression from one condition to another.

Two additional options for removing biclusters with **non significant expression patterns** are available for biclusters with discrete data, both using Markov chains to assess the statistical significance of the pattern and computing a p-value based on the (xvi) **overall** or (xvii) **column-wise** background probability of the occurrence of the pattern.

Sorting

Sorting the biclusters in a group reorganizes the biclusters according to: their values of (i) **ACV**, (ii) **ARV**, (iii) **MSR**; the **difference between their** (iv) **ARV and ACV**, (v) **ARV and MSR**. You can also sort biclusters by their **number of** (vi) **genes** or (vii) **conditions**, (viii) **size** or (ix) **variance**. These criteria allow biclusters to be sorted either in decreasing or decreasing order of the considered value(s).

Additional **sorting by the best p-value** obtained for the GO terms that annotate the genes of each bicluster in the group, both (x) **standard** and (xi) **corrected** for multiple testing, is also available. Groups of biclusters with discrete data can further be sorted based on the **significance of their expression pattern**, measured by a p-value computed using Markov models and based on the (xii) **overall** or (xiii) **column-wise** background probability of the occurrence of the pattern.

After post-processing a group of biclusters, a new group of biclusters is generated, which we call a **Post-Processed Biclusters Group**, and added to the dataset tree. This group is similar to a group of biclusters and allows for the very same functionalities. The resulting biclusters in a post-processed group can also be found in the original group. Because biclusters are just filtered and/or sorted, their data does not change in



BiGGESTS Quickstart

relation to the original data. However, as a consequence of the filtering operation, the post-processed group may not contain all the original biclusters. Moreover, as a consequence of the sorting operation, the post-processed biclusters may also not be in the same order as in the original group. To address this, the post-processed biclusters are given new identifiers, following their order in the post-processed group, but they also maintain the original identifiers in parenthesis, (), for easier tracking the bicluster in the original group.