# Accepted Manuscript
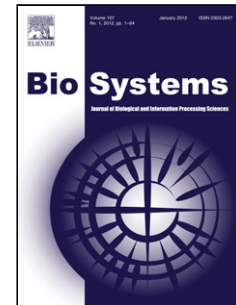
Title: GeRNet: A Gene Regulatory Network Tool

Author: J.S. Dussaut C.A. Gallo F. Cravero M.J. Martínez
J.A. Carballido I. Ponzoni

Please cite this article as: Dussaut, J.S., Gallo, C.A., Cravero, F., Martínez, M.J., Carballido, J.A., Ponzoni, I.,GeRNet: A Gene Regulatory Network Tool, *BioSystems* (2017), http://dx.doi.org/10.1016/j.biosystems.2017.08.006

# GeRNet: A Gene Regulatory Network Tool

J.S. Dussaut[1], C.A. Gallo[1], F. Cravero[2], M.J. Martínez[1], J.A. Carballido[1], I. Ponzoni[1,*]

[1] Instituto de Ciencias e Ingeniería de la Computación Universidad Nacional del Sur - CONICET - Bahía Blanca – Argentina

[2] Planta Piloto de Ingeniería Química Universidad Nacional del Sur - CONICET - Bahía Blanca - Argentina

*Correspondence author e-mail: ip@cs.uns.edu.ar

**Abstract.** Gene regulatory networks (GRNs) are crucial in every process of life since they govern the majority of the molecular processes. Therefore, the task of assembling these networks is highly important. In particular, the so called model-free approaches have an advantage modeling the complexities of dynamic molecular networks, since most of the gene networks are hard to be mapped with accuracy by any other mathematical model. A highly abstract model-free approach, called rule-based approach, offers several advantages performing data-driven analysis; such as the requirement of the least amount of data. They also have an important ability to perform inferences: its simplicity allows the inference of large size models with a higher speed of analysis. However, regarding these techniques, the reconstruction of the relational structure of the network is partial, hence incomplete, for an effective biological analysis. This situation motivated us to explore the possibility of hybridizing with other approaches, such as biclustering techniques. This led to incorporate a biclustering tool that finds new relations between the nodes of the GRN. In this work we present a new software, called GeRNeT that integrates the algorithms of GRNCOP2 and BiHEA along a set of tools for interactive visualization, statistical analysis and ontological enrichment of the resulting GRNs. In this regard, results associated with Alzheimer disease datasets are presented that show the usefulness of integrating both bioinformatics tools.

**Keywords:** Gene Regulatory Network, Biclustering, Microarray, Machine-Learning.

## 1 Introduction

A gene regulatory network (GRN) is a collection of molecular regulators that interact among them and with other substances in the cell to govern the gene expression levels of mRNA and proteins (Karlebach and Shamir, 2008; Zamani et al., 2010). These networks are crucial in every process of life, and several cellular

processes are affected by them. The nodes of this network are genes, proteins, their corresponding mRNAs, or protein/protein complexes, while the edges indicate a chain of dependencies, with cycles corresponding to feedback loops.

The microarray technology enabled large scale studies that allow the parallel measure of gene expression under a given set of conditions. This technology introduces a large amount of information to analyze, hence several data analysis issues (Alves et al., 2010). There are several techniques that carry out the reverse engineering of GRNs from analyzing gene expression data (Alves et al., 2010; De Jong, 2002; Espanés et al., 2016; Gallo et al., 2011; Karlebach and Shamir, 2008; Lakshmanan et al., 2013; Li et al., 2008; Ponzoni et al., 2007; Pridgeon and Corne, 2004; Zamani et al., 2010). They vary from Boolean models, model-free approaches up to data mining techniques, each of them tries to overcome the disadvantages of the previous one. Particularly, one machine-learning approach for the inference of time-lagged rules, called GRNCOP2 (Gallo et al., 2011), performs this task. These time-delayed gene regulation rules are a common phenomenon (Bulashevska and Elis, 2005; Li et al., 2006; Ponzoni et al., 2007; Soinov et al., 2003; van Someren et al., 2000) and they add complexity and computational cost. The resulting potential interactions between genes from the algorithm are used to predict the gene expression states of a gene in terms of gene expression values of other genes and to assemble a GRN by applying and combining these rules.

Model-free approaches have an advantage over any mathematical model because of the complexities of dynamic molecular networks that are hard to be mapped (Li et al., 2006) since this kind of approaches does not need any underlying model. In particular, rule-based approaches, a highly abstract model-free technique, offer several advantages when data-driven analysis is performed since they require the least amount of data, with an important ability to perform inferences, and allow the inference of large size models with a higher speed of analysis (Karlebach and Shamir, 2008). The main disadvantage of this approach is that the reconstruction of the relational structure of the network is partial, hence incomplete, for an effective biological analysis. This open the doors to the idea of exploring the hybridization of these methods with other techniques that enhance the relational structure of the network. In this regard, other techniques focus on finding a behavioral pattern for the genes on a subset of conditions for the microarray expression data. These techniques are called biclustering and aim to infer the biological role of the genes. Particularly BiHEA (Gallo et al., 2009) is a Hybrid Biclustering Evolutionary Algorithm that performs this task.

In this work a software suit for the GRNCOP2 algorithm and its integration with the BIHEA algorithm is presented. The software also provides several features for handling, pre-processing and visualizing microarray data. In the following sections the GeRNeT software is presented. First, its features are pronounced along with the results of the Alzheimer disease datasets experiments. Then, a comparison with other authors work is performed and the biological inferences found by the tool are described.

## 2 Material and Methods

GeRNet is a software suit for the GRNCOP2 algorithm (Gallo et al., 2011) which infers time-lagged rules for the reconstruction of GRNs. The tool incorporates the BiHEA biclustering algorithm (Gallo et al., 2009) to the GRNCOP2 in order to enhance the network with putative relations that GRNCOP2 alone was not able to recover from the data. Additionally, the suit provides several features for handling, pre-processing and visualizing microarray data. In the following sections, all the features of GeRNet are described. The tool is available at: http://lidecc.cs.uns.edu.ar/index.php/sw/gernet.

### 2.1 Project handling

GeRNet is based on the concept of project to manipulate microarray data. In this regard, the tool provides with a dialog box were the datasets can be added to a new project by pressing the "Add datasets..." button. The application allows the selection of multiple files, so it is possible to select all the data at once or one by one by pressing the button and choosing each file iteratively. The datasets are displayed in the dialog box to select if they correspond to Time Series data or Steady State data. It is also possible to delete the data from the list individually or by pressing the "Clear list" button.

Additionally, the application requires the specification of the particular CSV format of the data, by selecting whether the data has the gene names in the first column and/or if it has the sample names in the first row. Also, the data may contain gene annotations in the second column of the datasets. If none of these options are selected, the application assigns anonymous names in each case. In this regard, all the data must be consistent with the selected "Gene and column names" options. As an example, if the gene symbol ID in the first column option is selected, then all the datasets must meet this convention. This allows the automatic resize and reordering of the data rows in the case that they were provided with a different number of genes and/or with a different gene ordering. In the case that the data does not have the gene names in the first column, they must be provided with the same number of genes and in the same order.

Finally, the software allows to save and to open previously saved projects. The saved project contains all the datasets, as well as the gene regulatory network (in case that one was inferred) and the parameters for the algorithms.

### 2.2 Datasets views and manipulation

The framework allows different views of the datasets and provides useful tools to manipulate them. The datasets viewing window (Fig. 1) shows each data file in a different tab and it also indicates whether they correspond to a Time Series or to a Steady State data. If the data contains gene annotations, the corresponding information is displayed in a second column. It is also possible to choose between a Heat Map view and a Numeric view of the data. In the Heat Map view, the framework

allows the selection of the colors representing the above of average and the below of average values of the heat map by pressing the "Heat Map Colors" button.

Whenever the button "Heat Map View" is unselected, the view of the datasets turns into the Numeric view mode (Fig. 2). In this mode, the values of the datasets on each sample can be modified by double clicking with the left mouse button into the desired cell. Also, it is possible to estimate the missing values of all the datasets by pressing the "Estimate Missing Values" button in the secondary option bar. This function replaces the cells marked as missed (represented with a 999 value) with an estimation of the values obtained by a Bayesian Principal Component Analysis method (BPCA).

Both views allow for zooming of the datasets with several scales, thus improving the visualization. Also, it is possible to transform the data by means of log2, loge, log10, z-score, translation (add) and escalation (mult).
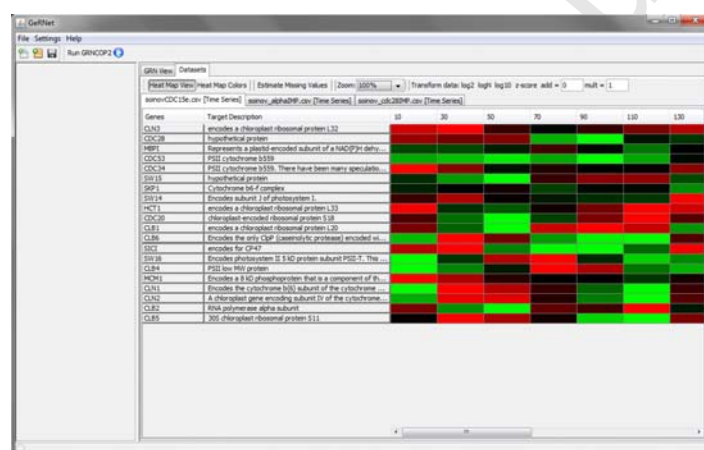


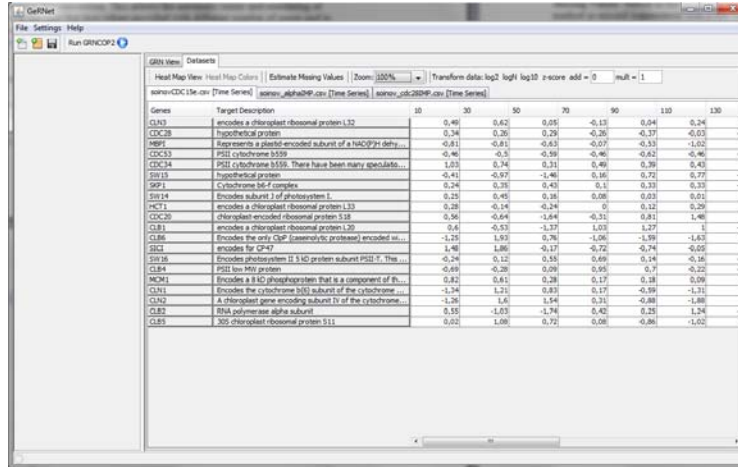**Fig. 1.** GeRNeT, Datasets view window, Heat Map View.

**Fig. 2.** GeRNeT, Datasets view window, Numeric view (Heat Map unselected).

### 2.3 The GRNCOP2 algorithm

In order to run the GRNCOP2 algorithm, it is necessary to have at least one dataset loaded in the framework, either by creating a new project or by opening an existing one. The algorithm runs by pressing the "Run GRNCOP2" ⬤ button in the main option bar. As a result, a progress bar appears that shows the movement of the search.

The results of the search depend on the type of datasets loaded and also on the parameters selected for the algorithm. The type of the datasets determines primarily the delay of the rules that can be inferred. If all the datasets are Time Series type, then the maximum time delay searchable is determined by the dataset with less samples minus 4 (this is due to restrictions in the GNRCOP2 algorithm regarding the minimum number of samples of the datasets which is 4). If all the datasets are Steady State type, then only simultaneous rules can be inferred. In the case when there are mixed data types (Time Series and Steady State), the simultaneous rules are inferred with all the datasets, whereas the time delay rules are inferred only with the Time Delay datasets.

The default parameters of the GRNCOP2 algorithm are shown in the dialog box of Fig. 3. All the parameters of the algorithm are described as follow:
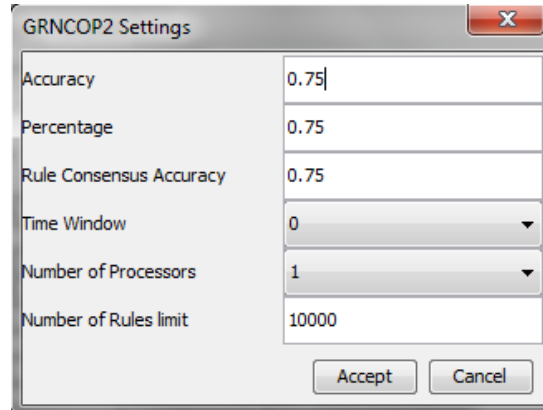
**Fig. 3.** Parameters of GRNCOP2.

*Accuracy.* A real value between 0 and 1 that represents the minimum acceptable confidence of the rules that will be reported by the algorithm. All the rules that doesn't meet this threshold will be discarded and not reported.

*Percentage.* A real value between 0 and 1 that represents the minimum percentage of samples that the rules +g1 ━━➤ +g2, -g1 ┅┅┅➤ -g2, +g1 ━━━┃ -g2 and -g1 ┅┅┅┃ +g2 need to meet in order to be reported. All the rules that don't meet this threshold will be discarded and will not be reported. This allows a fine tuning of the search, thus avoiding too many false positives of these kinds of rules since they are more likely to be inferred by chance.

*Rule Consensus Accuracy.* A real value between 0 and 1 that represents the minimum acceptable percentage of datasets that the rules must meet with the Accuracy parameter. All the rules that don't meet this threshold will be discarded and will not be reported.

*Time Window.* An integer value greater or equal to 0 that sets the maximum delay of the rules to be inferred. A zero value states that only simultaneous rules will be inferred. An $x$ value greater than zero states that the delays of the rules to be inferred will be between 0 (simultaneous) and $x$, including both ends

*Number of Processors.* An integer value greater or equal to 1 that sets the maximum number of processors that the algorithm can use in order to perform the search. If more processors are assigned to the algorithm, then less time will take the search. The maximum value is limited to the maximum number of cores available in the host machine.

*Number of Rules limit.* An integer value greater or equal to 1 that sets the maximum number of rules to be displayed by the framework. If the GRNCOP2 algorithm

exceeds this limit when it finishes the search, a dialog box will appear asking for two options: save the rules in a csv file or continue. If continue is selected, the responsiveness of the interface it is not guaranteed since it will depend on the hardware over which the software is running.

Once the search is complete, the GRN View tab will be made available with a graph representation of the inferred gene regulatory network along with a list of the obtained rules.

### 2.4 Gene Regulatory Network view and manipulation

This section describes the features of this framework to view, modify and manipulate the gene regulatory networks inferred by GRNCOP2. Once the GRNCOP2 algorithm finishes the search, the gene regulatory network associated with the rules inferred is automatically drawn in the "GRN View" tab (Fig. 4). Also, a table that contains all the rules appears in the left side of the main window. Each gene in the GRN View is represented as a green circle vertex, whereas the rules that relate them are directed edges that are drawn accordingly to the following scheme:

- **+/-g1** *d*==> **+/- g2** rules are represented by a directed edge ⟶ from g1 to g2 and state that when g1 is over expressed (under expressed), g2 is activated (inhibited) with a delay of *d* units of time.

- **+ g1** *d* ==> **+ g2** rules are represented by a directed edge ‐‐‐➔ from g1 to g2 and state that when g1 is over expressed, g2 is activated with a delay of *d* units of time.

- **- g1** *d* ==> **- g2** rules are represented by a directed edge ⋯⋯➔ from g1 to g2 and state that when g1 is under expressed, g2 is inhibited with a delay of *d* units of time.

- **+/- g1** *d* ==> **-/+ g2** rules are represented by a directed edge ⟶| from g1 to g2 and state that when g1 is over expressed (under expressed), g2 is inhibited (activated) with a delay of *d* units of time.

- **+ g1** *d* ==> **- g2** rules are represented by a directed edge ‐‐‐| from g1 to g2 and state that when g1 is over expressed, g2 is inhibited with a delay of *d* units of time.

- **- g1** *d* ==> **+ g2** rules are represented by a directed edge ⋯⋯| from g1 to g2 and state that when g1 is under expressed, g2 is activated with a delay of *d* units of time.

Only the genes that are associated by at least one rule are drawn in the gene regulatory network.

**Fig. 4.** Rules inferred by GeRNeT in the GRN View tab.



**Fig. 5.** A different view of the GRN View tab.

The GRN View offers several views of the network (Fig. 5). It contains one tab called "All delays" that shows the entire network with all the time delays, and also contains one tab per time delay inferred that shows only the rules of that specific time delay. Below the network graph view (Fig. 4) there is the Reference panel that contains the indication of the meaning of each edge in the graph. It also allows to hide any type of interaction by clicking on the button that contains the edge. Above the network graph view, there is an option bar aimed at customizing the current view of the network. It allows the selection of several layouts that automatically rearrange the genes according to some specific criteria. Also, it is possible to add or hide several types of labels for the edges and vertexes of the graph, like the confidence ("Show Accuracy / BiHEA Score" button) and time delay ("Show Time Delay" button) of the rules, and the names of the genes ("Show Gene Names" button). It also allows to hide genes accordingly to a specific selection in the graph or by a previously created gene list file. The "Filter by Selected Genes" button hides all the genes except the ones that

are selected in the graph, whereas the "Filter by Gene List File" button hides all the genes except the ones that are in the gene list of the selected file. The format of this file is a text archive with each gene name separated by a new line.

On the left side of the GRN View there is the tool bar used to interact with the network. It contains three tools: the arrow ▸ tool, the hand ✊ tool and the pen ✎ tool.

*Arrow* ▸ *tool.* This tool turns the view in picking mode:

- Left mouse button press on a gene or association to select it.

- Left mouse button + shift press on a gene or association to add or toggle selection.

- Left mouse button + drag on a gene to move all selected genes.

- Left mouse button + drag to select genes in a rectangle.

- Left mouse button + shift + drag to add to selection with genes in a rectangle.

*Hand* ✊ *tool.* This tool turns the view in the transforming mode:

- Left mouse button + drag to translate the view.

- Left mouse button + shift + drag to rotate the view.

- Left mouse button + ctrl (or command) + drag to shear the view.

*Pen* ✎ *tool.* This tool turns the view in the editing mode:

- Left mouse button press on empty space creates a new gene.

- Left mouse button on a gene, followed by a drag to another gene creates a directed edge between them.

*With any tool.*

- A right mouse button click on an empty space with no gene selection shows a popup menu (Fig. 6) to create a new gene, and allows to export the current view (see the Exporting the Results section). If at least one gene is picked, then the popup menu shows options for creating a new gene, deleting the selected genes, finding new associations that involve the selected genes through BiHEA (see the Integration with BiHEA section) and exporting the current view.

- A right mouse button click on a gene shows a popup with several actions. If no genes are picked, then the only option that is shown is for exporting the current view. Otherwise, you can create new associations between the selected genes and the clicked gene as a target, delete the selected genes, find new associations that involve the selected genes through BiHEA, and export the current view.

- The mouse wheel allows to scale the view. When the scale is equal or less than one, the view is scaled. When the scale is more than one, the layout is scaled.

- Hover a gene with the mouse and the corresponding gene annotation (if any) will pop up on screen.



**Fig. 6.** Right click popup menu on the GRN View tab.

## 2.5    Integration with BiHEA

The framework incorporates the BiHEA algorithm to allow the inference of new interactions that were not recovered by the GRNCOP2 algorithm. The BiHEA is an evolutionary biclustering algorithm aimed at finding co-expression (similar or opposed) between genes that may only occur in a subset of the samples. In this regard, since BiHEA only find subsets of genes that are related in a subset of samples, there is the need of additional mechanisms to allow the recovery of interactions that can be compared with the rules inferred by GRNCOP2, and then, added to the gene regulatory network.

The integration of the BiHEA algorithm involves several steps. First, the user must select in the GRN View tab what genes are of interest in order to find new interactions. These genes will act as potential regulators for the other genes in the datasets. Then, by clicking with the right mouse button on the view, a popup menu will be displayed showing the option "Find new associations with BiHEA" (ref. Fig. 6). The search begins once that option is clicked, the results are displayed in Fig. 8.

The search is related to the current time delay displayed in the network, i.e., if the time delay 1 is displayed, then the algorithm will search for rules with a delay of one-time unit. If all the delays are shown (i.e., the tab "All rules" is selected) then a dialog will show up to allow the user to select one specific time delay. This time delay affects the way that the data is considered by the algorithm, since all the genes that will act as potential regulators will be displaced in the selected time delay in the

datasets with respect to the other genes. As an example, suppose a dataset D1 with five genes and ten samples (Fig. 7). If the genes that are selected are g2 and g4, and the time delay is 1, then the actual matrix in which the search will be performed is the D1'.



**Fig. 7.** Example for the matrix D1' in which the algorithm will perform the search if D1 is the original dataset, the genes selected are g2 and g4 and the time delay is 1.

The next step performs a search of biclusters with the BiHEA algorithm in all the datasets separately. Once all the searches finish, a pairwise analysis is performed for each one of the selected genes. To make the explanation simpler, let us use the previous example (Fig. 7). If g2 and g4 are selected, then the pairwise analysis computes for each pair of genes g2-g1, g2-g3, g2-g4, g2-g5 and g4-g1, g4-g2, g4-g3, g4-g5, the number of biclusters of D1' for which both genes of the pair are members. Then those values are averaged with the values of the other datasets in order to obtain the final BiHEA Score for each interaction.

Finally, when all this process is done, a dialog box is shown (Fig. 8) with the highest BiHEA Scores along with the corresponding pair of genes of the interaction, allowing the user to select the desired associations and discarding the rest. The selected interactions will be added automatically to the gene regulatory network only if they were not present before.

**Fig. 8.** Dialog box with the highest BiHEA Scores and the corresponding pair of genes in the interaction. The user can select the desired rules and discard the others.

If the search was performed for a time delay 0, i.e. for simultaneous rules, then each pair of genes is displayed twice, one for each gene acting as a regulator. This is because a priori it cannot be assumed any causality in the relation since the only information that is known is the co-expression of both genes. On the other hand, if the search was performed with a delay greater than 0, i.e., for time delayed rules, then it can be assumed that the expression (or inhibition) of the target gene is due to the expression (or inhibition) of the regulator gene in the previous samples. The opposite regulation is also inferred since the BiHEA algorithm builds the biclusters considering both the similar and the opposed expression profiles of each gene. In order to differentiate the BiHEA Score of the Accuracy of GRNCOP2, given that both values have different interpretations, the BiHEA Score number is visualized in the network with the letters BS at the end of the value.

Finally, there are several parameters for the BiHEA algorithm that can be tuned (Fig. 9), as described below:



**Fig. 9.** BiHEA parameters.

*Population size.* An integer value greater than 1 that represents the size of the main population maintained by the evolutionary algorithm. This is the number of biclusters that BiHEA will return.

*Max Generations.* An integer value greater than 1 that refers to the number of iterations that the algorithm will perform improving the solutions until it stops. Greater values yield to solutions that are more optimized, but this also entails more execution time.

*Mutation Probability.* A real value between 0 and 1 that sets the probability of the mutation operator in the evolutionary process.

*Alpha.* A real value that sets a threshold for the local search in the evolutionary algorithm. A greater value improves the execution time but also deteriorates the quality of the solutions; and vice versa.

*Seed.* An integer value that sets the seed for the random number generator of the evolutionary algorithm.

*Max Number of Rows.* An integer positive value that sets the size of the table in the dialog box that shows the rules found by BiHEA.

## 2.6    Exporting the results

There are two kinds of results that can be exported from the framework: the rules inferred by the algorithms and the graphical view of the network. In order to export the rules, the software suit opens a dialog box to save the file in a coma separated value format. On the other hand, in order to export the graphical view of the network, the framework shows a dialog box with several options (Fig. 10). There are several formats available, such as PDF, EPS, JPG, PNG and many other along with options to tune it.



**Fig. 10.** Several options to export the graphical view of the network.

## 3 Results and Discussion

In order to test the tool's performance, a microarray dataset obtained from GEO (Gene Expression Omnibus: http://www.ncbi.nlm.nih.gov/geo/) was used. The dataset, called GSE1297 (Blalock et al., 2004), analyzes hippocampal gene expression of nine control and 22 AD subjects of varying severity; it shows the progression of the Alzheimer Disease (AD) in terms of patients with no disease present, called *control* samples, patients with the disease on their first stages, called *incipient* samples, patients with moderate symptoms of the disease, called *moderate* samples, and patients on their final stages, called *severe* samples.

The GeRNet Tool was used on each matrix from the dataset GSE1297 (control, incipient, moderate and severe) to obtain a list of rules from the gene regulatory network obtained. This list of rules was calculated for each matrix from the dataset, from now on called list R. Notice that there is an R list for each of the four matrix of the dataset, each of them indicates the progression of the disease. The next step was to run the BiHEA integration, as stated in section 2.5, in order to obtain a new set of extra rules, called list B for each matrix. This list contains the rules or interactions not detected by the algorithm of GRNCOP2 by using the Hybrid MOEA explained in the previous section.

In order to compare the rules found by GeRNeT with the ones reported in Kong et al. (Kong et al., 2014), the first ones were reduced using the same Gene Regulators and Gene Targets as those reported by Kong et al. The results show that with the inclusion of BiHEA, several new rules were found. Kong et al. method uses an independent component analysis (ICA) instead of clustering methods, in order to select the differentially expressed genes on the same dataset that was introduced previously. The results show an increase of rules found by GeRNeT, only after the application of the BiHEA integration, some of them are found by Kong et al. and several others are not presented in their work. These rules are explained in the following section.

Figures 13 to 16 contrast the association rules obtained by GeRNeT and Kong et al. for each AD dataset. The interactions highlighted in blue corresponds to the rules found by Kong et al. that are also inferred by GRNCOP2, without using BiHEA tools. The interactions highlighted in orange represent those found by Kong et al. that are also inferred by hybridizing GRNCOP2 with BiHEA. In other words, the orange rules cannot be detected using GRNCOP2 alone. Highlighted in purple are those rules found by hybridizing GRNCOP2 with BiHEA, but are not reported by Kong et al. Finally, the rules in grey are found for GRNCOP2 alone, but are not detected by Kong et al. Regarding the nodes, the ones that are filled in blue are the Regulator genes, and the ones that are empty are the Target ones. We focus our discussion on the purple highlighted rules, which are those that are not found in the paper of Kong et al. while are found by the GERNET tool using the BiHEA integration. This set of rules are explained on the following section in detail and additional bibliography supporting them is cited.

**Fig. 11.** All the rules or interactions found by GeRNeT using GSE1297 **control** dataset. Highlighted in blue are the rules found also by Kong et al., highlighted in orange are the rules found both by Kong et al. and by the GeRNeT tool using the BiHEA integration, and finally the rules that are not present in Kong's work while are found by GeRNeT using the BiHEA integration are highlighted in purple. The nodes filled in blue are the Regulator Genes.

**Fig. 12.** All the rules or interactions found by GeRNeT using GSE1297 **incipient** dataset. The highlighted colors are as stated in Fig. 11



**Fig. 13.** All the rules or interactions found by GeRNeT using GSE1297 **moderate** dataset. The highlighted colors are as stated in Fig. 11.

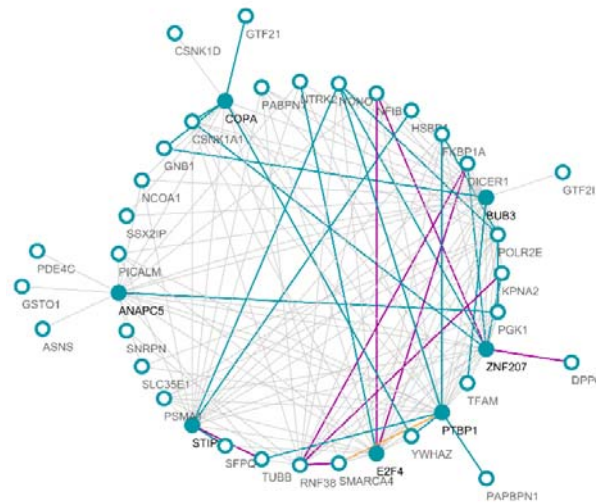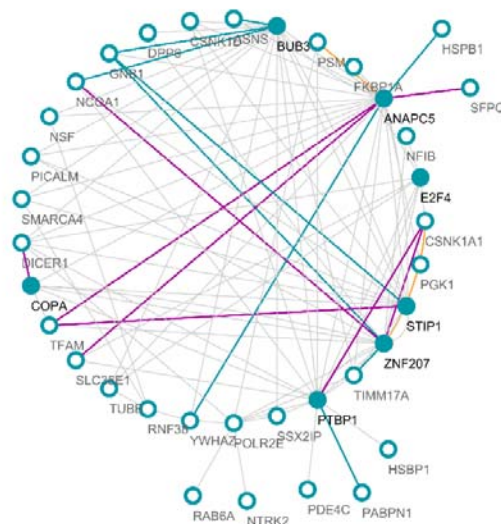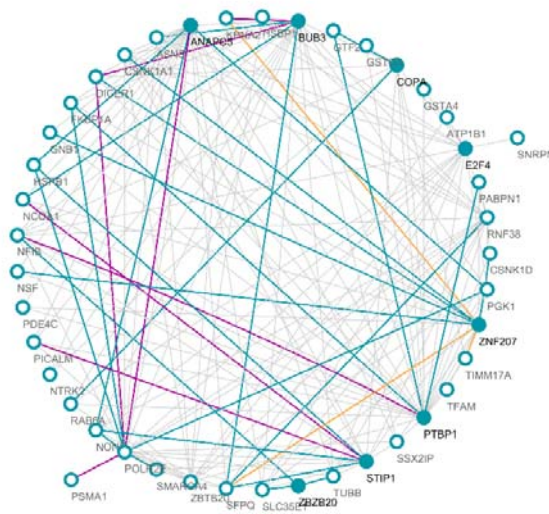**Fig. 14.** All the rules or interactions found by GeRNeT using GSE1297 **severe** dataset. The highlighted colors are as stated in Fig. 11.

### 3.1    Biological relevance of the results

As stated before, the rules found by GeRNeT were highlighted in different colors according to the source of its finding (see figures 13 to 16). In blue are the rules reported by Kong (Kong et al., 2014) and found by the GRNCOP2 without the integration of the biclustering technique. In orange are the rules also found by Kong and GRNCOP2 hybridized with BiHEA. The biological relevance of these sets of rules, blue and orange, are explained in Kong's work and, therefore, we decide do not focus on them. Nevertheless, it is important to mention that combination of GRNCOP2 with the biclustering technique detects several rules that were unable to be reported for GRNCOP2 alone.

The other highlighted rules, the purple ones, are those that were found by using the hybridized method integrating BiHEA, and they were not reported by Kong's work. These rules are explained in more detail in Tables 1 and 2. In the first table, a brief explanation of each gene involved is presented. In the second table, information related to the biological relevance of each *purple* rule is presented. Notably, most of these relations are associated with the GO term *quinolinate catabolic process* that is strongly associated with the Alzheimer disease (AD) in the bibliography. The chemical reactions and pathways contained in the quinolinate catabolic process result in the breakdown of quinolinate, the anion of quinolinic acid. The quinolinate synthase is an enzyme that catalyzes glycerone phosphate and iminosuccinate and affords quinolinic acid (Ashihara et al., 2011). Studies have found in the post-mortem

brains of AD patients higher neuronal quinolinic acid levels and that quinolinic acid can associate with tau protein (Schwarcz et al., 2012; Severino et al., 2011). Furthermore, it has been demonstrated that quinolinic acid increases tau phosphorylation *in vitro* in human fetal neurons and induces ten neuronal genes, including some known to correlate with AD. In immunoreactivity studies, it has been found that quinolinic acid immunoreactivity is strongest in glial cells that are located close to amyloid plaques and that there is immunoreactivity with neurofibrillary tangles.

**Table 1.** In the following table a list of the genes involved in the GRN is shown in the column Gene. In the column Type the gene role is explained, if it is Target, Regulator or both. Finally in the Function column the gene function is briefly explained.

| Type | Gene | Function |
|---|---|---|
| Target & Regulator genes | ANAPC5 | This is a protein-coding gene. Anaphase Promoting Complex Subunit 5 protein, is involved in the pathway protein ubiquitination, which is part of protein modification. It is a cell cycle-regulated E3 ubiquitin ligase that regulates progression through mitosis and the G1 phase of the cell cycle, as member of the anaphase promoting complex/cyclosome (APC/C). |
| Target genes | CSNK1A1 | Casein Kinase 1 Alpha 1 is protein-coding gene. Casein kinases are operationally defined by their preferential utilization of acidic. It can phosphorylate several proteins. Contributes in Wnt signaling. Might be relevant in segregating chromosomes during mitosis, and in keratin cytoskeleton disassembly and in that way, it can regulate epithelial cell migration |
| | DICER1 | This gene encodes a protein possessing an RNA helicase motif containing a DEXH box in its amino terminus and an RNA motif in the carboxyl terminus. The encoded protein functions as a ribonuclease and is necessary by the RNA interference and small temporal RNA (stRNA) pathways to generate the active small RNA component that inhibits gene expression. Alternative splicing results in multiple transcript variants. |
| | DPP6 | This gene encodes a single-pass type II membrane protein, which belongs to the peptidase S9B family of serine proteases. This protein does bind particular voltage-gated potassium channels and modifies their expression and biophysical properties. |
| | FKBP1A | The protein encoded by this gene belongs to the immunophilin protein family, which is relevant in immunoregulation and basic cellular processes comprising protein folding and trafficking. |
| | KPNA2 | In the energy-dependent translocation through the nuclear pore complex, in the process of import of proteins into the nucleus, is required a nuclear localization sequence (NLS). The protein encoded by KPNA2 gene interacts with the NLSs of DNA helicase Q1 and SV40 T antigen and may be involved in the nuclear transport of proteins. |

| | | | |
|---|---|---|---|
| | NCOA1 | This gene encodes a protein that is a transcriptional coactivator for nuclear hormone and steroid receptors. The Nuclear Receptor Coactivator 1 protein binds nuclear receptors directly and stimulates the transcriptional activities in a hormone-dependent fashion. | |
| | NFIB | The Nuclear factor 1 B-type protein-coding by NFIB gene, recognizes and binds a palindromic sequence existing in cellular and viral promoters and in the begining of replication of adenovirus type 2. These proteins are individually capable of activating transcription and replication. | |
| | PDE4C | The protein encoded by this gene hydrolyzes the second messenger cAMP, which is a regulator and mediator of a number of cellular responses to extracellular signals. Thus, this protein plays a key role in many important physiological processes. | |
| | PICALM | This gene encodes a Phosphatidylinositol-binding clathrin assembly protein that collects clathrin and adaptor protein complex 2 (AP2) to cell membranes at sites of coated pit formation and clathrin vesicle assembly. The polymorphisms of this gene are associated with the risk of Alzheimer disease. | |
| | PSMA1 | The proteasome is a multi-catalytic proteinase complex, whose main role is to cleave peptides with Arg, Phe, Tyr, Leu, and Glu. The proteasome cleaves peptides in an ATP/ubiquitin-dependent process in a non-lysosomal pathway. This gene encodes a represent of the peptidase T1A family, which is a 20S core alpha subunit. | |
| | SFPQ | RNA and DNA binding protein, which participates in many nuclear processes. Essential pre-mRNA splicing factor necessary in spliceosome formation, at early stages, and for splicing catalytic step II, probably as a heteromer with NONO. The SFPQ-NONO heteromer associated with MATR3 may play a role in nuclear retention of defective RNAs, and may be involved in DNA unwinding by modulating the function of topoisomerase I/TOP1 | |
| | SLC35E1 | The solute carrier family 35 member E1 protein coding by this gene has a putative transporter function. | |
| | SMARCA4 | The transcription activator BRG1 protein-coding by this gene is a transcriptional coactivator collaborating with nuclear hormone receptors to increase transcriptional activation. In resting neurons, transcription of the c-FOS promoter is inhibited by BRG1-dependent recruitment of a phospho-RB1-HDAC repressor complex. | |
| | TFAM | This gene encodes a relevant mitochondrial transcription factor comprising two high mobility group motifs. The encoded protein also functions in mitochondrial DNA replication and repair. | |
| | TUBB | This gene encodes a beta tubulin protein. This protein and alpha tubulin conforms a dimer, working as a structural component of microtubules. Mutations in this gene are responsible for cortical dysplasia, complex, together with several brain malformations. | |

| | | | |
|---|---|---|---|
| **Regulator genes** | BUB3 | This gene encodes a mitotic checkpoint protein involved in spindle-assembly checkpoint function and a function in promoting the establishment of correct kinetochore-microtubule (K-MT) attachments. The BUB1/BUB3 complex plays a role in the inhibition of anaphase-promoting complex or cyclosome (APC/C) when a spindle-assembly checkpoint is activated. | |
| | COPA | Protein transport from the endoplasmic reticulum to Golgi compartments is facilitated in part by non-clathrin-coated vesicular coat proteins (COPs). Seven coat proteins have been identified, and they represent subunits of a complex known as coatomer, encoded by COPA gene. | |
| | E2F4 | The protein encoded by this gene is a member of the E2F family of transcription factors. The E2F family has a central role in the regulation of cell cycle and action of tumor suppressor proteins. Besides, it is a target for small DNA tumor viruses. | |
| | POLR2E | The POLR2E gene encodes the fifth largest subunit of RNA polymerase II, the polymerase responsible for the synthesis of messenger RNA (mRNA) in eukaryotes. | |
| | PTBP1 | The PTBP1 gene belongs to the subfamily of ubiquitously expressed heterogeneous nuclear ribonucleoproteins (hnRNPs). The hnRNPs are RNA-binding proteins and they complicate with heterogeneous nuclear RNA (hnRNA). These proteins are linked with pre-mRNAs in the nucleus and occur to influence pre-mRNA processing and other features of mRNA metabolism and transport. | |
| | RNF38 | The protein coding by RNF38 gene plays as an E3 ubiquitin-protein ligase capable to ubiquitinate p53/TP53 that promotes its relocalization to discrete foci associated with PML nuclear bodies. The RING motif is a zinc-binding domain found in a large set of proteins, which is relevant in several cellular processes such as development, signal transduction, oncogenesis, and apoptosis. Multiple transcript variants encoding different isoforms have been found for this gene. | |
| | STIP1 | The protein encoded by this gene is an adaptor protein that coordinates some functions in protein folding. It is also involved in pre-mRNA splicing, specifically in spliceosome disassembly during late-stage splicing events. | |
| | ZNF207 | ZNF207/BuGZ is principally integrated of disordered low-complexity regions and undergoes phase transition or coacervation to define temperature-dependent liquid droplets. Coacervation activates microtubule bundling and concentrates tubulin, activating microtubule polymerization and assembly of spindle by increasing the concentration of its building blocks. Besides, acts as a regulator of mitotic chromosome alignment by mediating the kinetochore loading and stability of BUB3. | |

**Table 2.** In the following table each *purple* rule (see figures 13 to 16) is explained using the Enrichment Rule taken from GO.

| Regulator | Target | Enrichment Rule |
|---|---|---|
| ANAPC5 | NFIB | Regulation of ubiquitin-protein ligase activity involved in mitotic cell cycle |
| ANAPC5 | SFPQ | Quinolinate catabolic process |
| ANAPC5 | SLC35E1 | Regulation of ubiquitin-protein ligase activity involved in mitotic cell cycle |
| ANAPC5 | TFAM | Quinolinate catabolic process |
| BUB3 | DICER1 | Quinolinate catabolic process |
| BUB3 | FKBP1A | Quinolinate catabolic process |
| BUB3 | KPNA2 | Quinolinate catabolic process |
| BUB3 | PSMA1 | Quinolinate catabolic process |
| COPA | DICER1 | Quinolinate catabolic process |
| COPA | KPNA2 | Quinolinate catabolic process |
| COPA | PDE4C | Quinolinate catabolic process |
| E2F4 | ANAPC5 | Quinolinate catabolic process |
| E2F4 | DICER1 | Quinolinate catabolic process |
| E2F4 | NFIB | Quinolinate catabolic process |
| E2F4 | SMARCA4 | Quinolinate catabolic process |
| POLR2E | ANAPC5 | Quinolinate catabolic process |
| POLR2E | DICER1 | Quinolinate catabolic process |
| POLR2E | PSMA1 | Quinolinate catabolic process |
| PTBP1 | CSNK1A1 | Quinolinate catabolic process |
| PTBP1 | NCOA1 | Quinolinate catabolic process |
| PTBP1 | NFIB | Quinolinate catabolic process |
| PTBP1 | SLC35E1 | Quinolinate catabolic process |
| RNF38 | DICER1 | Quinolinate catabolic process |
| RNF38 | KPNA2 | Quinolinate catabolic process |
| RNF38 | SMARCA4 | Quinolinate catabolic process |
| STIP1 | CSNK1A1 | Quinolinate catabolic process |
| STIP1 | NCOA1 | Quinolinate catabolic process |
| STIP1 | PICALM | Quinolinate catabolic process |

| STIP1 | TFAM | Quinolinate catabolic process |
|-------|------|-------------------------------|
| STIP1 | TUBB | Quinolinate catabolic process |
| ZNF207 | CSNK1A1 | Quinolinate catabolic process |
| ZNF207 | DPP6 | Quinolinate catabolic process |
| ZNF207 | NCOA1 | Protein stabilization |
| ZNF207 | NFIB | Protein stabilization |

## 4 Conclusions

In this work the importance and complexity of the inverse engineering of GRNs was assessed, along with the benefits and advantages of using model-free methods above the mathematical ones. The main disadvantage of model-free approaches was also addressed along with the importance of hybridizing these rule-based methods with other strategies that improve the inference of biological meaningful regulatory networks. This led to the presentation of a new software tool GeRNeT that incorporated this hybridization between the algorithm GRNCOP2 (Gallo et al., 2011) and the biclustering technique BiHEA (Gallo et al., 2009), together with several analysis facilities.

In order to evaluate the integration with BiHEA, a case study analysis was performed, using an AD progression dataset (Blalock et al., 2004). The analysis involved a contrast between the results obtained with the algorithm GRNCOP2 and the results added by the integration of BiHEA using the GeRNeT tool. These results were compared with the ones reported by Kong et al. (Kong et al., 2014) using the same datasets. From these studies, it was possible to conclude that the knowledge extracted by BiHEA provides relevant associations that cannot be detected by GRNCOP2 nor by Kong et al. In this way, it is evident that the collaboration between two model-free strategies of different characteristics allows to reach GRNs with more information. In the future, we hope to deepen in the evaluation and improvement of this tool with other analysis scenarios and also consider novel additional hybridization strategies.

## 5 Acknowledgements

## 6       Author contributions

JSD, CAG, and IP conceived and designed the tool facilities and experiments, and also analyzed the data. CAG coded the integration of the software package. JSD and IP worked in the statistical analysis of the data obtained for the software package. JSD, CAG, and IP wrote the first draft of the manuscript. FC and JSD analyzed the biological relevance of the Alzheimer Disease experiment results. MJM contributed to the design of the data visualizations. JSD, CAG, FC, JAC, and IP contributed to the writing of the manuscript. JSD, CAG, FC, MJM, JAC, and IP jointly developed the structure and arguments for the paper and made critical revisions. All authors reviewed and approved the final manuscript.

## 7       Disclosures and ethics

As a requirement of publication author(s) have provided to the publisher a signed confirmation of compliance with legal and ethical obligations including but not limited to the following: authorship and contributorship, conflicts of interest, privacy and confidentiality and (where applicable) protection of human and animal research subjects. The authors have read and confirmed their agreement with the ICMJE authorship and conflict of interest criteria. The authors have also confirmed that this article is unique and not under consideration or published in any other publication and that they have permission from rights holders to reproduce any copyrighted material. Any disclosures are made in this section. The external blind peer reviewers report no conflicts of interest.

## 8       References

Alves, R., Rodriguez-Baena, D.S., Aguilar-Ruiz, J.S., 2010. Gene association analysis: a survey of frequent pattern mining from gene expression data. Brief. Bioinform. 11, 210–224.

Ashihara, H., Crozier, A., Komamine, A., 2011. Plant metabolism and biotechnology.

Blalock, E.M., Geddes, J.M., Chen, K.C., Porter, N.M., Markesbery, W.R., Landfield, P.W., 2004. Incipient Alzheimer's disease: microarray correlation analyses reveal major transcriptional and tumor suppressor responses, in: Proceedings of the National Academy of Sciences. pp. 2173–2178.

Bulashevska, S., Elis, R., 2005. Inferring Genetic Regulatory Logic from Expression Data. Bioinformatics 21, 2706–2713.

De Jong, H., 2002. Modeling and Simulation of Genetic Regulatory Systems: A Literature Review. J. Comput. Biol. 9, 67–103.

Espanés, P. de, Osses, A., Rapaport, I., 2016. Fixed-points in random Boolean networks: The impact of parallelism in the Barabási–Albert scale-free topology case. Biosystems.

Gallo, C.A., Carballido, J.A., Ponzoni, I., 2011. Discovering time-lagged rules from microarray data using gene profile classifiers. BMC Bioinformatics 12, 1.

Gallo, C.A., Carballido, J.A., Ponzoni, I., 2009. BiHEA: A Hybrid Evolutionary Approach for Microarray Biclustering. Lect. Notes Comput. Sci. 5676, 36–47.

Karlebach, G., Shamir, R., 2008. Modeling and analysis of gene regulatory networks. Nat. Rev. Mol. Cell Biol. 9, 770–780.

Kong, W., Mou, X., Zhi, X., Xin, Z., Yang, Y., 2014. Dynamic Regulatory Network Reconstruction for Alzheimer's Disease Based on Matrix Decomposition Techniques. Comput. Math. Methods Med.

Lakshmanan, S., Park, J., Jung, H., 2013. Design of state estimator for genetic regulatory networks with time-varying delays and randomly occurring uncertainties. Biosystems.

Li, H., Xuan, J., Wang, Y., Zhan, M., 2008. Inferring regulatory networks. Front. Biosci. 13, 263–275.

Li, X., Rao, S., Jiang, W., Li, C., Xiao, Y., Guo, Z., Zhang, Q., Wang, L., Du, L., Li, J., Li, L., Zhang, T., Wang, Q.K., 2006. Discovery of Time-Delayed Gene Regulatory Networks based on temporal gene expression profiling. BMC Bioinformatics 7, 26.

Ponzoni, I., Azuaje, F., Augusto, J., Glass, D., 2007. Inferring Adaptive Regulation Thresholds and Association Rules from Gene Expression Data through Combinatorial Optimization Learning. IEEE/ACM Trans. Comput. Biol. Bioinforma. 4, 624–634.

Pridgeon, C., Corne, D., 2004. Genetic Network Reverse-Engineering and Network Size; Can We Identify Large GRNs?, in: IEEE Symp Computational Intelligence in Bioinformatics and Computational Biology 2004. pp. 32–36.

Schwarcz, R., Bruno, J., Muchowski, P., 2012. Kynurenines in the mammalian brain: when physiology meets pathology. Nat. Rev.

Severino, P., Muller, G., Vandresen-Filho, S., Tasca, C., 2011. Cell signaling in NMDA preconditioning and neuroprotection in convulsions induced by quinolinic acid. Life Sci.

Soinov, L.A., Krestyaninova, M.A., Brazma, A., 2003. Towards reconstruction of gene networks from expression data by supervised learning. Genome Biol. 4, 1.

van Someren, E.P., Wessels, L.F., Reinders, M.J., 2000. Linear modeling of genetic networks from experimental data, in: Ismb. pp. 355–366.

Zamani, Z., Hajihosseini, A., Masoudi-Nejad, A., 2010. Computational Methodologies for Analyzing, Modeling and Controlling Gene Regulatory Networks. Biomed. Eng. Comput. Biol. 2, 47–62.