# PATHWAY STUDIO™

# Finding Major Regulators/Cellular Processes/Diseases in Experimental Data

**ELSEVIER**

Sub-Network Enrichment Analysis (SNEA) is used to identify and prioritize the targets and regulators most implicated in the experimental dataset. SNEA is based on the Gene Set Enrichment Analysis algorithm. Sub-Networks are user defined networks calculated *de novo* from the information in the database and consist of a seed/regulator and their neighbors (targets) in the database. The seeds of the sub-network whose targets are statistically enriched are implicated as important regulators (or cell processes or diseases) by the experimental data. You can specify the kind of seeds and the kind of relationships for building sub-networks from the SNEA dialog. A seed can be a protein/complex/ functional class, small molecule, cellular process or disease. Once identified, regulators can be further examined to help elucidate cellular processes, mechanisms and pathways impacted in the experiment.

Examples of results from SNEA analysis of experimental data:

- Major gene expression regulators (such as transcription factors) responsible for a differential gene expression profile
- Major miRNA regulators responsible for a differential gene expression profile
- Binding regulatory networks
- Differential gene/protein profiles that are enriched for genes/proteins known to be associated with a particular disease
- Differential gene/protein profiles that are enriched for genes/proteins known to be associated with a particular cellular process

The SNEA tool has easy-to-use preset options for defining the most useful sub-network types, and a customer menu for the experienced user who wants to perform more advanced analysis.

**Description of Sub-Network Enrichment Analysis**

The Sub-Network Enrichment Analysis (SNEA) algorithm uses existing relationships in the database to build "sub-networks" based on user specified criteria. It then uses these sub-networks with the GSEA algorithm to identify the networks that are significantly enriched. When calculating enrichment, only the targets are considered, but not the seed/regulator.

The user-defined sub-networks consist of a single "regulator" or "seed" and its nearest neighbor network. The type of relationship(s) included in the network and the directionality of the network are user determined.

Outbound (from the seed or regulator) relationships are selected when one wants to identify regulators of targets included in the gene list or experimental data. Inbound (to the seed or regulator) relationships are selected when the seed is a disease or cell process. In this instance the algorithm identifies entities known to be associated with a particular disease or cell process. Many combinations of sub-networks are possible.

Recall: in the mammal database all protein relationships to a disease (or cell process) are inbound to the disease (or cell process). Also, the only type of relationship to disease or cell process is regulation.
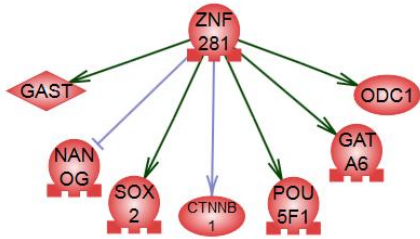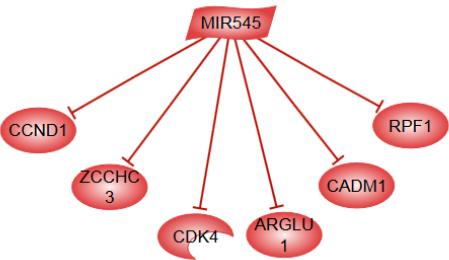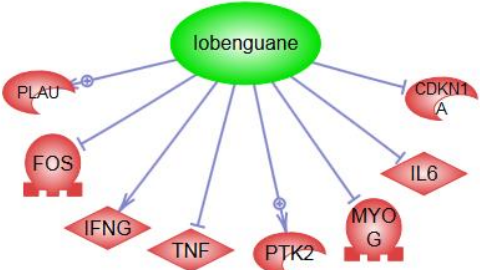
However ChemEffect® and DiseaseFx®, expand on this model.  They contain more relation types to diseases and cell processes and some relations such as GeneticChange are outbound from disease. Customer with access to these additional datasets will have additional preset options to allow utilization of the additional relations present in these databases.

Sub-networks that have as the seed a disease or cell process are useful in defining a collection of proteins known to be associated with these entities, without any implied function by the designated directionality of the relationship.


**Defining the Sub-Networks**

Selecting user-defined sub-networks involves first defining the "regulator" or "seed" type" and the nearest "neighbor" network by selecting specific relationship types and directionality.

Examples of commonly used sub-networks available through presets:

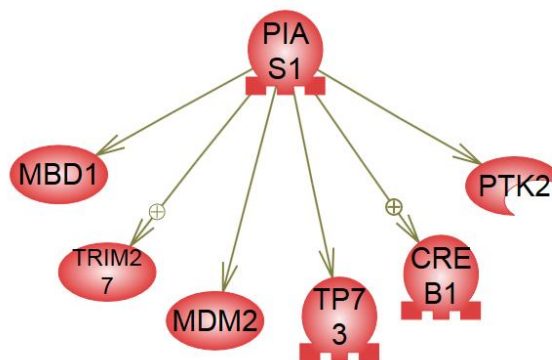| **Preset Name:  Expression Targets**<br><br>seed/regulator = protein/complex/functional class;<br>relations = promoter binding, expression<br>neighbors/targets = proteins<br>input data set = gene expression (most common),<br>miRNA array<br>results:  identifies major expression regulators active in experiment |  |
|---|---|
| **Preset Name: miRNA Targets**<br><br>seed/regulator = miRNA<br>relations = miRNAEffect<br>neighbors/targets = proteins<br>input data set =  gene expression (most common),<br>miRNA array<br>results:  identifies major miRNAs regulating gene expression in experiment |  |
| **Preset Name:  Chemical Expression Targets**<br><br>seed/regulator = small molecules<br>relations =  expression<br>neighbors/targets = proteins<br>input data set = gene expression<br>results:  identifies small molecules/drugs that regulate expression of proteins |  |

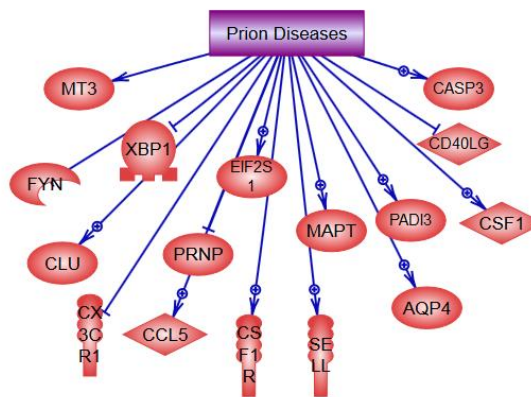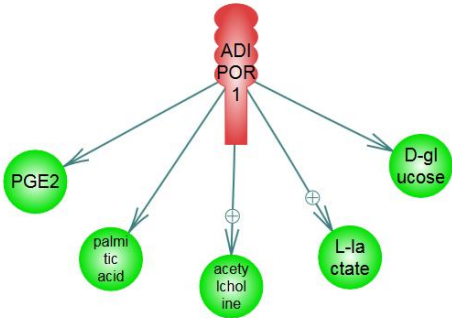| **Preset Name: Binding Proteins**<br><br>seed/regulator = protein<br>relations = binding<br>neighbors/targets = protein<br>input data set = proteomics<br>results: identifies enriched binding sub-networks mediated by an individual protein |  |
|---|---|
| **Preset Name: Protein Modification Targets**<br><br>seed/regulator = proteins<br>relations = protein modification<br>neighbors/targets = proteins<br>input data set = proteomics (most common)<br>results: identifies binding protein modification sub-networks through activities such as: acetylation, cleavage, deacetylation, demethylation, dephosphorylation, direct interaction, methylation, phosphorylation,posttrascriptional inhibition, proteolysis, ubiquitination |  |
| **Preset Name: Disease Biomarkers: Quantity**<br><br>seed/regulator = disease<br>relations = quantitative change<br>neighbors/targets = proteins<br>input data set = gene expression (most common), proteomics<br>results: identifies enrichment of proteins that are associated with specific diseases through changes in the proteins abundance, expression or activity |  |
| **Preset Name: Disease Biomarkers: Mutation**<br><br>seed/regulator = disease<br>relations = genetic change<br>neighbors/targets = proteins<br>input data set = gene expression (most common), proteomics<br>results: identifies enrichment of proteins that are associated with a specific disease through genetic changes in the genes such as: gene amplification, epigenic methylation |  |

| | |
|---|---|
| **Preset Name: Proteins/Chemicals Regulating Diseases**<br><br>seed/regulator = disease<br>relations = regulation<br>neighbors/targets = proteins (shown) or small molecules<br>input data set = gene expression (most common), proteomics or metabolomics<br>results: identifies proteins (gene expression or proteomics data) or small molecules (metabolomics data) enriched for a specific disease |  |
| **Preset Name: Proteins/Chemicals Regulating Cell Processes**<br><br>seed/regulator = disease<br>relations = regulation<br>neighbors/targets = proteins (shown) or small molecules<br>input data set = gene expression (most common), proteomics or metabolomics<br>results: identifies proteins (gene expression or proteomics data) or small molecules (metabolomics data) enriched for a specific cellular process |  |
| **Preset Name: Metabolomics Targets**<br><br>seed/regulator = protein<br>relations = molsynthesis<br>neighbors/targets = small molecules<br>input data set = metabolomics<br>results: identifies enrichment of small molecules where changes in abundance regulated by a specific protein |  |

**Running Sub-Network Enrichment Analysis in Pathway Studio**

With an experiment open, go to the experiment view Tools menu and select "Analyze Experiment." From the Analyze Experiment window select "Sub-Network Enrichment Analysis" in the Analysis Type drop down menu.

The box "clean up resulting sub-networks by removing neighbors not present in the experiment" is checked by default. This will limit the sub-networks to only entities included in the experiment being analyzed.

4

If the experiment has multiple comparisons, select the desired comparison for analysis. The default *p*-value is < 0.05 and the maximum number of networks is set to 100. These values can be changed at the user's discretion.

The list of preset options is dependent on the type of experiment. For example, a metabolomics data set will have a different list of preset options than will a gene expression experiment. For most users' needs, the list of preset options will suffice. However, there is an option to define custom sub-networks for more advanced users.

For the example below the experiment type is gene expression and the preset selected is "Expression Targets." This analysis will identify the top gene expression regulators for this experiment.

The results of the analysis are displayed in the list pane at the bottom of the screen.

Each *de novo* calculated sub-network is named based on its seed/regulator. The total number of neighbors in the sub-network may be higher than the # of measured neighbors. If the box was checked to "clean up" the resulting sub-networks, then only the # of measured neighbors will be displayed in the graph view when selected.

The Gene Set Seed column provides the list of identified regulators. In this example these are the most important expression regulators for this experiment. You can copy this column by selecting the entire table and then going to Edit > Copy Gene Seeds to Clipboard.



The Median Change value indicates the median expression value for all the targets of the sub-network.

The Activation Score indicates how closely the changes in the expression of the targets closely match the predicted effect (positive or negative) that the regulator has on the target, where a positive number indicates concordance and a negative number indicated discordance.

**The Activation Score is a measure of whether the regulator is "active" or "repressed" in the experimental conditions.**

| Regulator Effect on Target | Target Experimental Results | Concordant/ Disconcordant | *Implied activation/repression of seed/regulator** |
|---|---|---|---|
| Positive | Down regulated | Disconcordant | Repression |
| Positive | Up regulated | Concordant | Activated |
| Negative | Down regulated | Concordant | Activated |
| Negative | Up regulated | Disconcordant | Repressed |

\* The user must examine if the actual state of the seed/regulator matches or does not match the implied state of activation or repression.

Activation Score calculation considerations:
- Any genes with absolute log change below 0.5 are disregarded from calculation
- Any genes that are connected with relationships without effect are disregarded from calculation

Thus, differential genes (that may or may not been filtered by a user) with log changes above 0.5 threshold AND connected with signed relationships to the seed are analyzed. The ones that differentially change consistent with the sign of the relationship are considered "concordant." The activation score = (N_concordant – N_disconcordant)/sqrt(N_total)

Further insight into the experiment may be gained by examining the list of top regulators individually or as a group and examining the concordance/disconconrdance of the sub-networks.

If you have any questions about Pathway Studio please contact Customer Care:

**USA, Canada and Latin America:**

(8am-8pm CET - St.Louis)
Tel: US toll-free: +1 (888) 615 4500
Tel: Non toll-free: +1 (314) 523 4900
Email: usinfo@elsevier.com
Email Brazil: brinfo@elsevier.com

**Europe, Middle East and Africa:**

(9am-6pm GMT+1, Amsterdam office)
Tel: +31 20 485 3767
Email: nlinfo@elsevier.com

**Japan:**

(9,30am-5,30pm JST, Tokyo office)
Tel: +81 (3) 5561 5035
Email: jpinfo@elsevier.com
Website: japan.elsevier.com

**Asia and Australasia:**

(9am-6pm SST, Singapore office)
Tel: +65 6349 0222
Email: sginfo@elsevier.com