

# Data Integration with Cytoscape

Samad Lotia

L'institut Pasteur

6 May 2009

# The Problem

Given a biological problem:

- ▶ The “data integration” part: data from disparate sources
- ▶ Now that our data is “integrated”, what do we do with it?  
We wish to draw conclusions from this data *collectively*

# Why Integrate Data in the First Place?

- ▶ Many instruments for observing biological phenomena are noisy:
  - ▶ Gene expression profiles can very noisy if expression levels are low<sup>1</sup>
  - ▶ Up to 50% of yeast two-hybrid data are false<sup>2</sup>
- ▶ If one can observe the same hypothesis from a variety of instruments, this improves the chance of the hypothesis' validity

---

<sup>1</sup>Quackenbush, H. Weighing our measures of gene expression. *Mol. Syst. Biol.* 2, 63 (2006).

<sup>2</sup>Sprinzak, E., Sattath, S., & Margalit, H. How reliable are experimental protein-protein interaction data? *J. Mol. Biol.* 327, 919—923 (2003).

# How Are We Going to Integrate Data?

- ▶ In biology we have components that interact or influence other components, that in turn influence other components, and so on.
  - ▶ Genes
  - ▶ Proteins
  - ▶ Pathways
- ▶ The essence of this talk: *Use a network abstraction as a basis for bringing together disparate data sources and to draw conclusions from them for biological problems.*

# What Can We Do with a Network Abstraction?

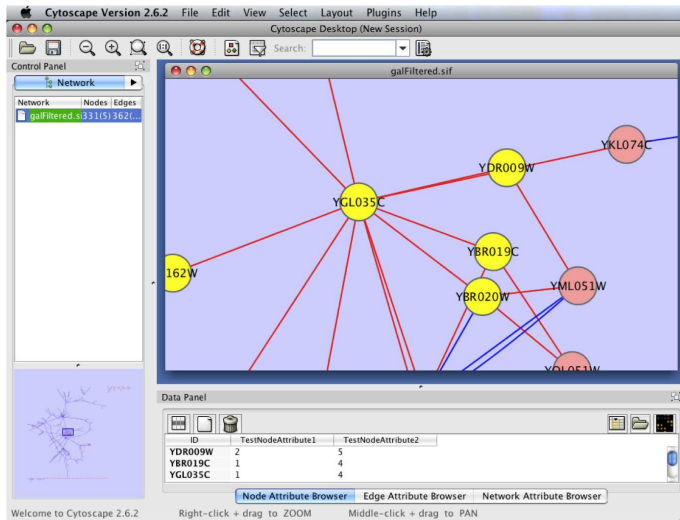
- ▶ By using networks to tackle biological problems, we can use Cytoscape to:
  - ▶ Bring together data from a variety of places
  - ▶ Visualize interactions of biological components
  - ▶ Analyze networks to determine what parts are biologically significant

# Some Background on Cytoscape

- ▶ Free & open source software application—LGPL license
- ▶ Written in Java—can run on Windows, Mac, & Linux
- ▶ Core development made by:
  - ▶ Agilent Technologies (Santa Clara, California, USA)
  - ▶ Eötvös Loránd University (Budapest, Hungary)
  - ▶ Institute for Systems Biology (Seattle, Washington, USA)
  - ▶ Memorial Sloan-Kettering Cancer Center (New York, New York, USA)
  - ▶ National Center for Integrative Biomedical Informatics (Ann Arbor, Michigan, USA)
  - ▶ L'institut Pasteur (Paris, Île-de-France, France)
  - ▶ University of Toronto (Toronto, Ontario, Canada)
  - ▶ University of California, San Diego (La Jolla, California, USA)
  - ▶ University of California, San Francisco (San Francisco, California, USA)
  - ▶ Unilever (London, UK)

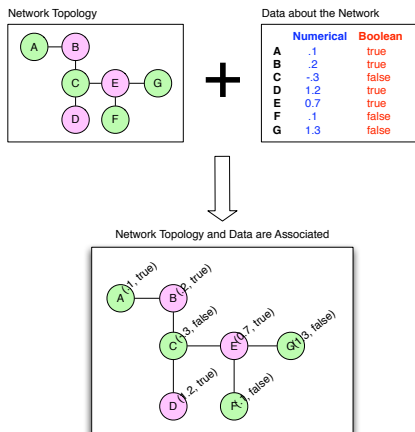
# Cytoscape's Core Functionality

## ► Visualizing networks



# Cytoscape's Core Functionality

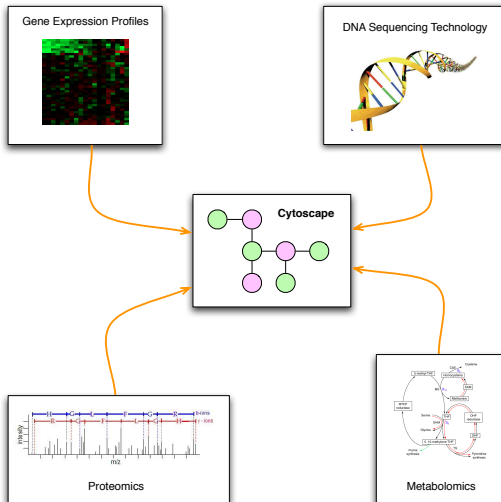
- ▶ Nodes and edges can have *attributes* associated with them
- ▶ Attributes and corresponding nodes and edges are *dynamically* bound





# Cytoscape's Core Functionality

- ▶ Agnostic semantics
- ▶ We can overlay data from a variety of instruments on top of a network

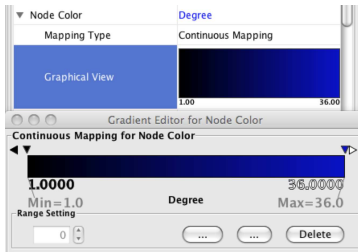






# Cytoscape's Core Functionality

- VizMapper's continuous mapping:

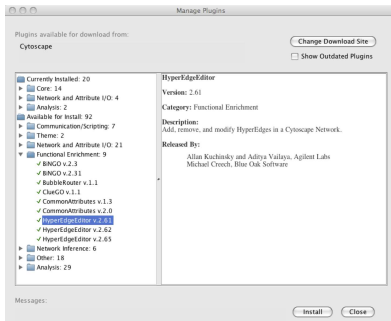


- VizMapper's discrete mapping:

Node Shape	Category
Mapping Type	Discrete Mapping
Amino Acid Metabolism	HEXAGON
Biosynthesis of Secondary ...	TRIANGLE
Carbohydrate Metabolism	ROUND_RECT
Energy Metabolism	RECT
Glycan Biosynthesis and M...	OCTAGON
Lipid Metabolism	PARALLELOGRAM
Metabolism of Cofactors a...	ELLIPSE
Metabolism of Other Amin...	DIAMOND
Nucleotide Metabolism	HEXAGON
Xenobiotics Biodegradatio...	TRIANGLE

# Cytoscape's Extended Functionality

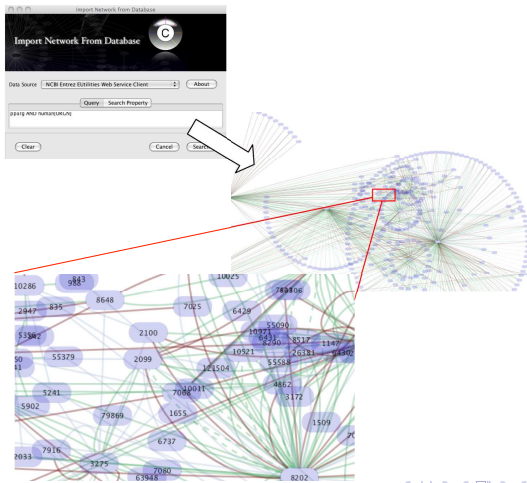
- ▶ Cytoscape extends its core functionality with *plugins*
- ▶ Developed by third parties
- ▶ Some major categories of plugins:
  - ▶ Obtain networks from online databases
  - ▶ Import node and edge attributes from online databases
  - ▶ Perform algorithmic analysis of networks



# Cytoscape's Extended Functionality: Networks from Online Databases

- ▶ We can import networks from NCBI Entrez and Pathway Commons

Obtaining Interactions of PPAR-Gamma in Humans from NCBI Entrez



# Cytoscape's Extended Functionality: Node and Edge Attributes from Online Databases

- ▶ Given a network where nodes are database IDs: overlay information about nodes onto network from a database

## Obtain Node Information from NCBI Entrez

The screenshot displays the Cytoscape interface with the NCBI Entrez Gene window open. The window is titled "NCBI Entrez Gene" and features the Entrez logo and the text "Entrez, The Life Sciences Search Engine". The "Data Source" is set to "NCBI Entrez Gene". The "Key Attribute" is set to "ID" and the "Data Type" is set to "Entrez Gene ID". The "Available Annotation Category" list includes "Summary", "Publications", "Phenotypes", "Pathways", and "General Protein Information", all of which are checked. The "Reset" and "Cancel" buttons are visible at the bottom of the window.

An arrow points from the "ID" attribute field to a network diagram. The network diagram is titled "PPAR-Gamma from NCBI" and shows a complex network of nodes and edges. The nodes are labeled with IDs, and the edges represent interactions. A yellow node labeled "5469" is highlighted. The "Data Panel" is open, showing the details for node "5469". The "GO Term: Biological Process" is selected, and the "Node Attribute" is set to "ID". The "Network Attribute Browser" is also visible, showing a list of attributes.

GO Term: Biological Process  
GO Term: Cellular Component  
GO Term: Molecular Function

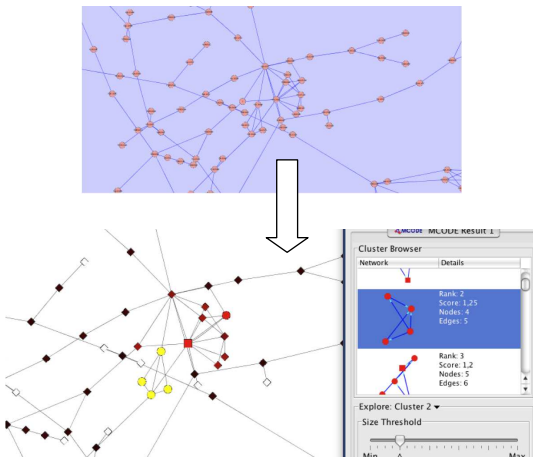
Node Attribute  
Network Attribute Browser

Right-click + d  
to PAN

androgen receptor signaling pathway  
fat cell differentiation  
positive regulation of transcription from RNA polymerase II promoter  
regulation of transcription, DNA-dependent  
transcription  
transcription initiation from RNA polymerase II promoter

# Cytoscape's Extended Functionality: Network Analysis

- ▶ Extract subnetworks and score them; the means to do this could be done by looking at:
  - ▶ Network topology
  - ▶ Values of node or edge attributes





# Problem of Glioblastoma

- ▶ Study of glioblastoma, a common form of brain cancer<sup>3</sup>
- ▶ Biological problem: People newly diagnosed with glioblastoma have a variety of genetic aberrations.
- ▶ How can one develop treatments if there are a large number of genetic aberrations?

---

<sup>3</sup>Cancer Genome Atlas Research Network. Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature* vol. 455 (23 Oct 2008).

# Data of Glioblastoma Patients

- ▶ Analyzed biospecimens from people in early stages of glioblastoma
- ▶ Looked at copy number alterations (CNAs)  $\Rightarrow$  about  $\frac{3}{4}$  of affected genes had expression levels proportional to number of copies
- ▶ Determined what genes had non-silent and silent mutations
  - ▶ p53 (regulates cell growth cycle) was mutated
  - ▶ NF1 (tumor suppressor gene) was inactivated or deleted
  - ▶ EGFR (signals the cell to grow) was activated
  - ▶ PI(3)K complex (signals a transition in cell cycle) was mutated
  - ▶ MGMT (DNA repair enzyme) was methylated, thus reducing its expression



# What Can We Do with the Data?

- ▶ Problem: data from a variety of sources: CNAs, mutations of many genes
- ▶ How can one draw conclusions from disparate data sources?
- ▶ Solution they used: project this data onto an established glioblastoma pathway network<sup>4</sup>

---

<sup>4</sup>Furnari, F. B. *et al.* Malignant astrocytic glioma: genetics, biology, and paths to treatment. *Genes Dev.* 21, 2683-2710 (2007).

# What Can We Learn from the Data?

- ▶ Almost all patients had mutations in any of these four pathways
- ▶ These pathways are *highly* interconnected
- ▶ While a variety of these genes were affected by different types of mutations, *these genes are strongly related to each other*

# Concluding Remarks

We can use Cytoscape for bringing data together under a network paradigm; we could potentially improve our understanding of a network by:

- ▶ obtaining metadata from a variety of online databases,
- ▶ visualizing it,
- ▶ and performing algorithmic analysis of it.