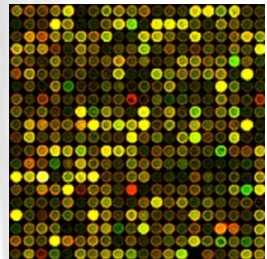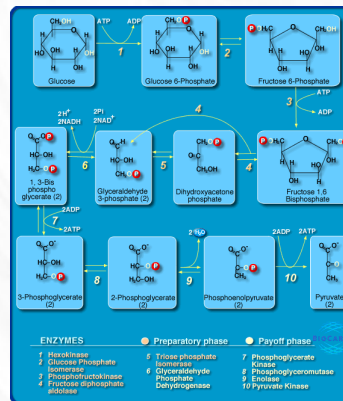# Interpreting Gene Lists

- The analysis produced 1000 hits-> Now what?
- Genome-Scale Analysis (Omics)
  - Genomics, Proteomics
- What's interesting about these genes
  - Are they enriched in known pathways, complexes, functions
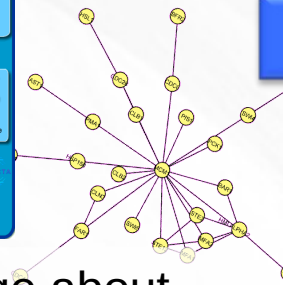


Ranking or clustering

Analysis tools

Prior knowledge about cellular processes

Eureka! New heart disease gene!

# Pathway and Network Analysis

- Any type of analysis that involves pathway or network information

- Most commonly applied to help interpret lists of genes

- Most popular type is pathway enrichment analysis, but many others are useful

- Helps gain mechanistic insight into 'omics data

# Correlation to Causation

- GWAS: find genetic markers correlated with disease – powerful approach, but:
  - genomics reduces statistical power (>multiple testing correction with >SNPs)
  - rare variants = more samples
- Associate pathways to increase power
  - Fewer pathways, organize many rare variants (damaging the system causes the disease)
- Use pathway knowledge to identify potential disease causes

# Before Analysis

- ✓ Normalization
- ✓ Background adjustment
- ✓ Quality control (garbage in, garbage out)


- ✓ Use statistics that will increase signal and reduce noise specifically for your experiment
- ✓ Other analyses you may want to use to evaluate changes
- ✓ Make sure your gene IDs are compatible with software

# Where Do Gene Lists Come From?

- Molecular profiling e.g. mRNA, protein
  - Identification → Gene list
  - Quantification → Gene list + values
  - Ranking, Clustering (biostatistics)
- Interactions: Protein interactions, microRNA targets, transcription factor binding sites (ChIP)
- Genetic screen e.g. of knock out library
- Association studies (Genome-wide)
  - Single nucleotide polymorphisms (SNPs)
  - Copy number variants (CNVs)

# **What Do Gene Lists Mean?**

- Biological system: complex, pathway, physical interactors

- Similar gene function e.g. protein kinase

- Similar cell or tissue location

- Chromosomal location (linkage, CNVs)
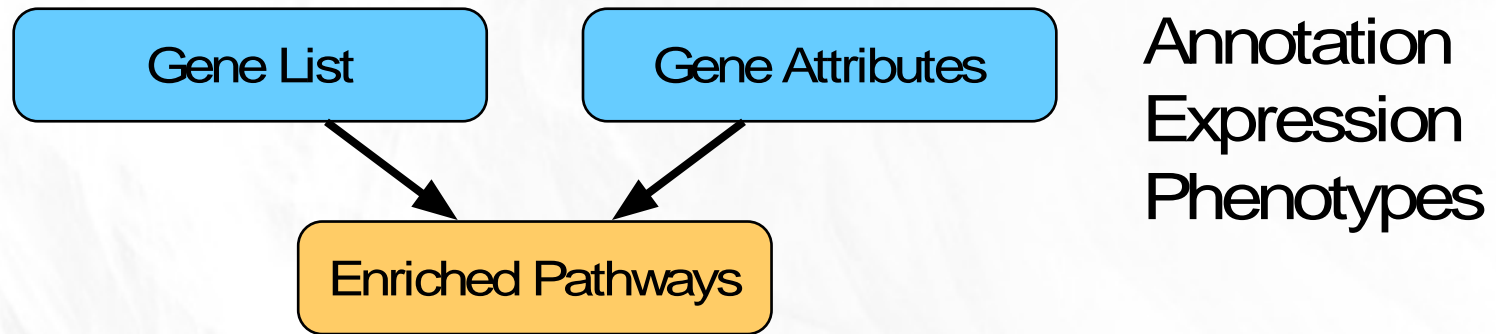
# Biological Questions

- Step 1: What do you want to accomplish with your list
  - Summarize biological processes or other aspects of gene function
  - Perform differential analysis – what pathways are different between samples?
  - Find a controller for a process (TF, miRNA)
  - Find new pathways or new pathway members
  - Discover new gene function
  - Correlate with a disease or phenotype (candidate gene prioritization)

# Biological Answers

- Computational analysis methods we will cover
  - Regulatory network analysis: find controllers
  - Pathway enrichment analysis: summarize and compare
  - Network analysis: predict gene function, find new pathway members, identify functional modules (new pathways)

# Pathway Enrichment Analysis



Annotation
Expression
Phenotypes

DAVID, GSEA, g:Profiler

- Gene identifiers
- Gene attributes/annotation
  - Gene Ontology
    - Ontology Structure
    - Annotation
  - BioMart + other sources

# Gene and Protein Identifiers

- Identifiers (IDs) are ideally unique, stable names or numbers that help track database records
  - E.g. Social Insurance Number, Entrez Gene ID 41232
- Gene and protein information stored in many databases
  - → Genes have many IDs
- Records for: Gene, DNA, RNA, Protein
  - Important to recognize the correct record type
  - E.g. Entrez Gene records don't store sequence. They link to DNA regions, RNA transcripts and proteins e.g. in RefSeq, which stores sequence.

GNAQ
GNAS
DGKZ
GUCY1A3
PDE4B
PDE4D
ATP2A2
ATP2A3
NOS1
CNN1
GSTO1
NOS3
CNN2
MYLK2
CALD1
ACTA1
MYL2

# Common Identifiers

**Gene**
Ensembl ENSG00000139618
**Entrez Gene 675**
Unigene Hs.34012

**RNA transcript**
GenBank BC026160.1
RefSeq NM_000059
Ensembl ENST00000380152

**Protein**
Ensembl ENSP00000369497
RefSeq NP_000050.2
UniProt BRCA2_HUMAN or
A1YBP1_HUMAN
IPI IPI00412408.1
EMBL AF309413
PDB 1MIU

**Species-specific**
HUGO HGNC BRCA2
MGI MGI:109337
RGD 2219
ZFIN ZDB-GENE-060510-3
FlyBase CG9097
WormBase WBGene00002299 or ZK1067.1
SGD S000002187 or YDL029W
**Annotations**
InterPro IPR015252
OMIM 600185
Pfam PF09104
Gene Ontology GO:0000724
SNPs rs28897757
**Experimental Platform**
Affymetrix 208368_3p_s_at
Agilent A_23_P99452
CodeLink GE60169
Illumina GI_4502450-S

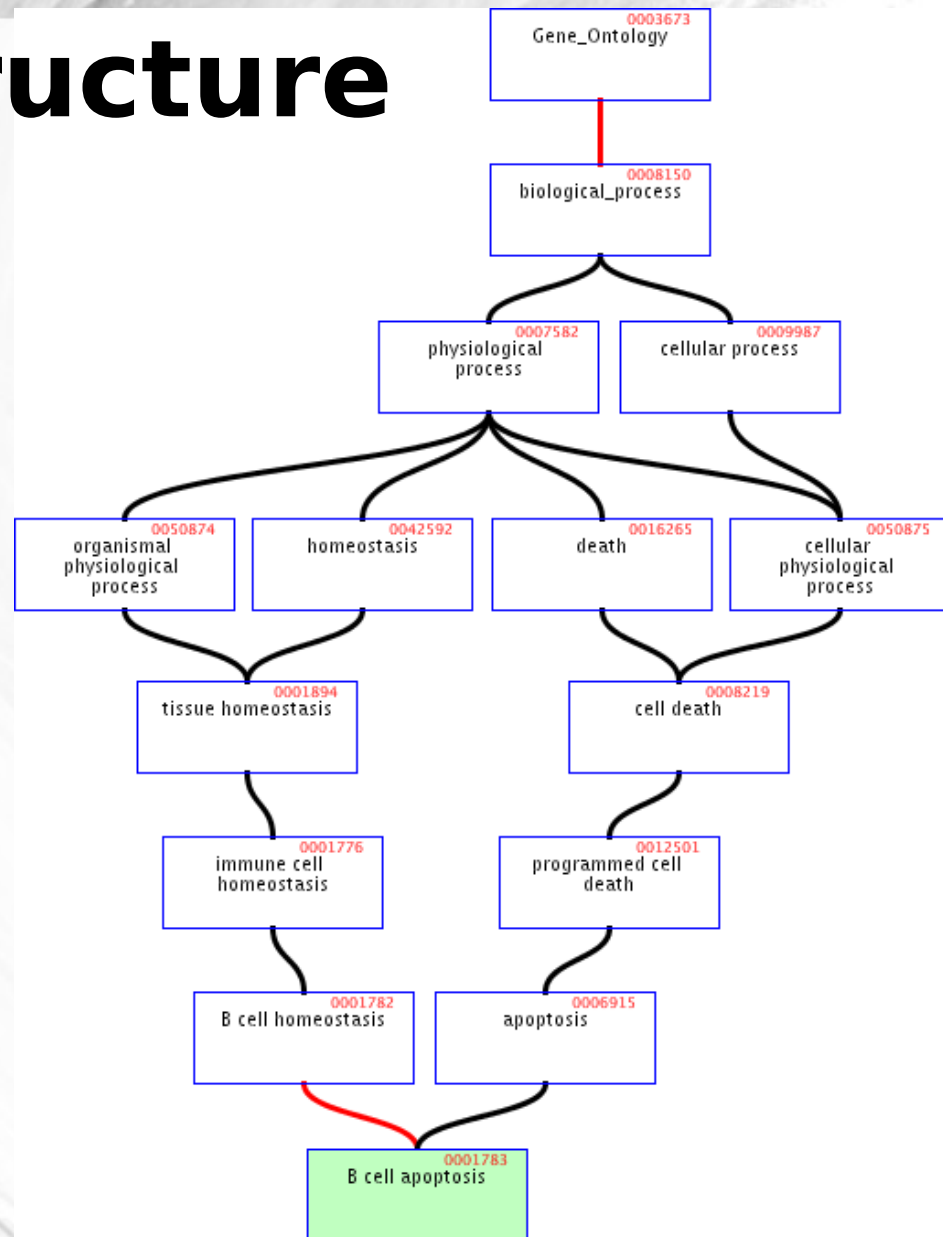Red =
Recommended

# Identifier Mapping

- So many IDs!
  - Software tools recognize only a handful
  - May need to map from your gene list IDs to standard IDs

- Four main uses
  - Searching for a favorite gene name
  - Link to related resources
  - Identifier translation
    - E.g. Proteins to genes, Affy ID to Entrez Gene
  - Merging data from different sources
    - Find equivalent records

# What is the Gene Ontology (GO)?

- Set of biological phrases (terms) which are applied to genes:
  - protein kinase
  - apoptosis
  - membrane
- Dictionary: term definitions
- Ontology: A formal system for describing knowledge
- www.geneontology.org

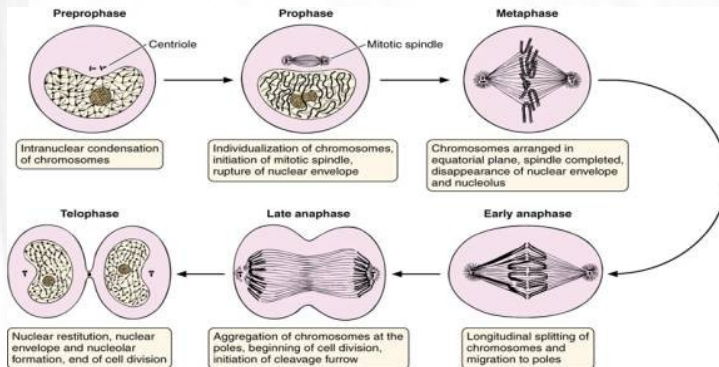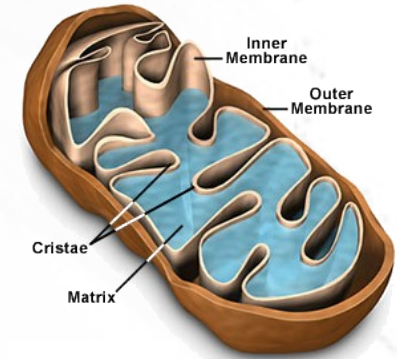Jane Lomax @ EBI

**www.geneontology.org**

# GO Structure

- Terms are related within a hierarchy
  - is-a
  - part-of
- Describes multiple levels of detail of gene function
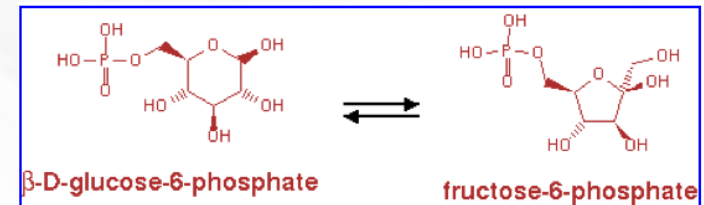- Terms can have more than one parent or child

# What GO Covers?

- GO terms divided into three aspects:
  – cellular component
  – molecular function
  – biological process

glucose-6-phosphate isomerase activity

Cell division

# Part 1/2: Terms

- Where do GO terms come from?
  - GO terms are added by editors at EBI and gene annotation database groups
  - Terms added by request
  - Experts help with major development
  - 37104 terms, with definitions
    - 23074 biological_process
    - 2994 cellular_component
    - 9392 molecular_function
    - As of June 2012

# Part 2/2: Annotations

- Genes are linked, or associated, with GO terms by trained curators at genome databases
  - Known as 'gene associations' or GO annotations
  - Multiple annotations per gene
- Some GO annotations created automatically (without human review)

# Species Coverage

- All major eukaryotic model organism species and human

- Several bacterial and parasite species through TIGR and GeneDB at Sanger

- New species annotations in development

- Current list:
  - http://www.geneontology.org/GO.downloads.annotations.shtml

# Gene Attributes

- Function annotation
  - Biological process, molecular function, cell location
- Chromosome position
- Disease association
- DNA properties
  - TF binding sites, gene structure (intron/exon), SNPs
- Transcript properties
  - Splicing, 3' UTR, microRNA binding sites
- Protein properties
  - Domains, secondary and tertiary structure, PTM sites
- Interactions with other genes

# Sources of Gene Attributes

- Ensembl BioMart (general)
  - http://www.ensembl.org
- Entrez Gene (general)
  - http://www.ncbi.nlm.nih.gov/sites/entrez?db=gene
- Model organism databases
  - E.g. SGD: http://www.yeastgenome.org/
- Many others.....

# Ensembl BioMart

- Convenient access to gene list annotation



Select genome

Select filters

Select attributes
to download

www.ensembl.org

# **Enrichment analysis**

- Introduction to enrichment analysis

- Hypergeometric Test, Fisher's Exact Test

- GSEA enrichment analysis for ranked lists.

- Multiple test corrections:

– Bonferroni correction

– False Discovery Rate computation using Benjamini-Hochberg procedure

# The "result"

| Probe Set ID | log.ratio | pvalue | adj.p |
|---|---|---|---|
| 73554_at | 1.4971 | 0.0000 | 0.0004 |
| 91279_at | 0.8667 | 0.0000 | 0.0017 |
| 74099_at | 1.0787 | 0.0000 | 0.0104 |
| 83118_at | -1.2142 | 0.0000 | 0.0139 |
| 81647_at | 1.0362 | 0.0000 | 0.0139 |
| 84412_at | 1.3124 | 0.0000 | 0.0222 |
| 90585_at | 1.9859 | 0.0000 | 0.0258 |
| 84618_at | -1.6713 | 0.0000 | 0.0258 |
| 91790_at | 1.7293 | 0.0000 | 0.0350 |
| 80755_at | 1.5238 | 0.0000 | 0.0351 |
| 85539_at | 0.9303 | 0.0000 | 0.0351 |
| 90749_at | 1.7093 | 0.0000 | 0.0351 |
| 74038_at | -1.6451 | 0.0000 | 0.0351 |
| 79299_at | 1.7156 | 0.0000 | 0.0351 |
| 72962_at | 2.1059 | 0.0000 | 0.0351 |
| 88719_at | -3.1829 | 0.0000 | 0.0351 |
| 72943_at | -2.0520 | 0.0000 | 0.0351 |
| 91797_at | 1.4676 | 0.0000 | 0.0351 |
| 78356_at | 2.1140 | 0.0001 | 0.0359 |

## What about the Biology???

# Slightly more informative results

| Probe Set ID | Gene Symbol | Gene Title | go biological process term | go molecular function term | log.ratio | pvalue | adj.p |
|---|---|---|---|---|---|---|---|
| 73554_at | CCDC80 | coiled-coil domain contain | --- | --- | 1.4971 | 0.0000 | 0.0004 |
| 91279_at | C1QTNF5 /// | C1q and tumor necrosis fa | visual perception /// embry | --- | 0.8667 | 0.0000 | 0.0017 |
| 74099_at | --- | --- | --- | --- | 1.0787 | 0.0000 | 0.0104 |
| 83118_at | RNF125 | ring finger protein 125 | immune response /// modi | protein binding /// zinc ion | -1.2142 | 0.0000 | 0.0139 |
| 81647_at | --- | --- | --- | --- | 1.0362 | 0.0000 | 0.0139 |
| 84412_at | SYNPO2 | synaptopodin 2 | --- | actin binding /// protein bir | 1.3124 | 0.0000 | 0.0222 |
| 90585_at | C15orf59 | chromosome 15 open read | --- | --- | 1.9859 | 0.0000 | 0.0258 |
| 84618_at | C12orf39 | chromosome 12 open read | --- | --- | -1.6713 | 0.0000 | 0.0258 |
| 91790_at | MYEOV | myeloma overexpressed ( | --- | --- | 1.7293 | 0.0000 | 0.0350 |
| 80755_at | MYOF | myoferlin | muscle contraction /// bloc | protein binding | 1.5238 | 0.0000 | 0.0351 |
| 85539_at | PLEKHH1 | pleckstrin homology doma | --- | binding | 0.9303 | 0.0000 | 0.0351 |
| 90749_at | SERPINB9 | serpin peptidase inhibitor, | anti-apoptosis /// signal tra | endopeptidase inhibitor ac | 1.7093 | 0.0000 | 0.0351 |
| 74038_at | --- | --- | --- | --- | -1.6451 | 0.0000 | 0.0351 |
| 79299_at | --- | --- | --- | --- | 1.7156 | 0.0000 | 0.0351 |
| 72962_at | BCAT1 | branched chain aminotrans | G1/S transition of mitotic | catalytic activity /// branch | 2.1059 | 0.0000 | 0.0351 |
| 88719_at | C12orf39 | chromosome 12 open read | --- | --- | -3.1829 | 0.0000 | 0.0351 |
| 72943_at | --- | --- | --- | --- | -2.0520 | 0.0000 | 0.0351 |
| 91797_at | LRRC16A | leucine rich repeat contain | --- | --- | 1.4676 | 0.0000 | 0.0351 |
| 78356_at | TRDN | triadin | muscle contraction | receptor binding | 2.1140 | 0.0001 | 0.0359 |

If we are lucky, some of the top genes mean something to us

But what if they don't?

And how what are the results for other genes with similar biological functions

Apply some methods to incorporate biological knowledge into microarray analysis

The type of knowledge to deal with is rather simple: We know groups/sets of genes that for example

•Belong to the same pathway

•Have a similar function

•Are located on the same chromosome, etc…

We will assume these groupings to be given, i.e we will not discuss methods how to detect pathways, networks, gene clusters

# What is a pathway?

- No clear definition
  - Wikipedia: "In biochemistry, **metabolic pathways** are series of chemical reactions occurring within a cell. In each pathway, a principal chemical is modified by chemical reactions."
  - These pathways describe enzymes and metabolites
- But often the word "pathway" is also used to describe gene regulatory networks or protein interaction networks
- In all cases a pathway describes a biological function very specifically

# What is a Gene Set?

- All genes involved in a pathway are an example of a Gene Set

- All genes corresponding to a Gene Ontology term are a Gene Set

- All genes mentioned in a paper of Smith et al might form a Gene Set


A Gene Set is a much more general and less specific concept than a pathway

# What is Gene Set/Pathway analysis?

The aim is to give one number (score, p-value) to a Gene Set/Pathway to answer questions like:

•Are many genes in the pathway differentially expressed (up-regulated/downregulated)?

•Can we give a number (p-value) to the probability of observing these changes just by chance?

# Pathway and Gene Set data resources

The Gene Ontology (GO) database

http://www.geneontology.org/

GO offers a relational/hierarchical database

•Parent nodes: more general terms

•Child nodes: more specific terms

At the end of the hierarchy there are genes/proteins

At the top there are 3 parent nodes: biological process, molecular function and cellular component

Example: we search the database for the term "inflammation"

The genes on our array that code for one of the 44 gene products would form the corresponding "inflammation" gene set

# KEGG pathway database

KEGG = Kyoto Encyclopedia of Genes and Genomes

http://www.genome.jp/kegg/pathway.html

The pathway database gives far more detailed information than GO

- Relationships between genes and gene products

But: this detailed information is only available for selected organisms and processes

# Types of enrichment analysis

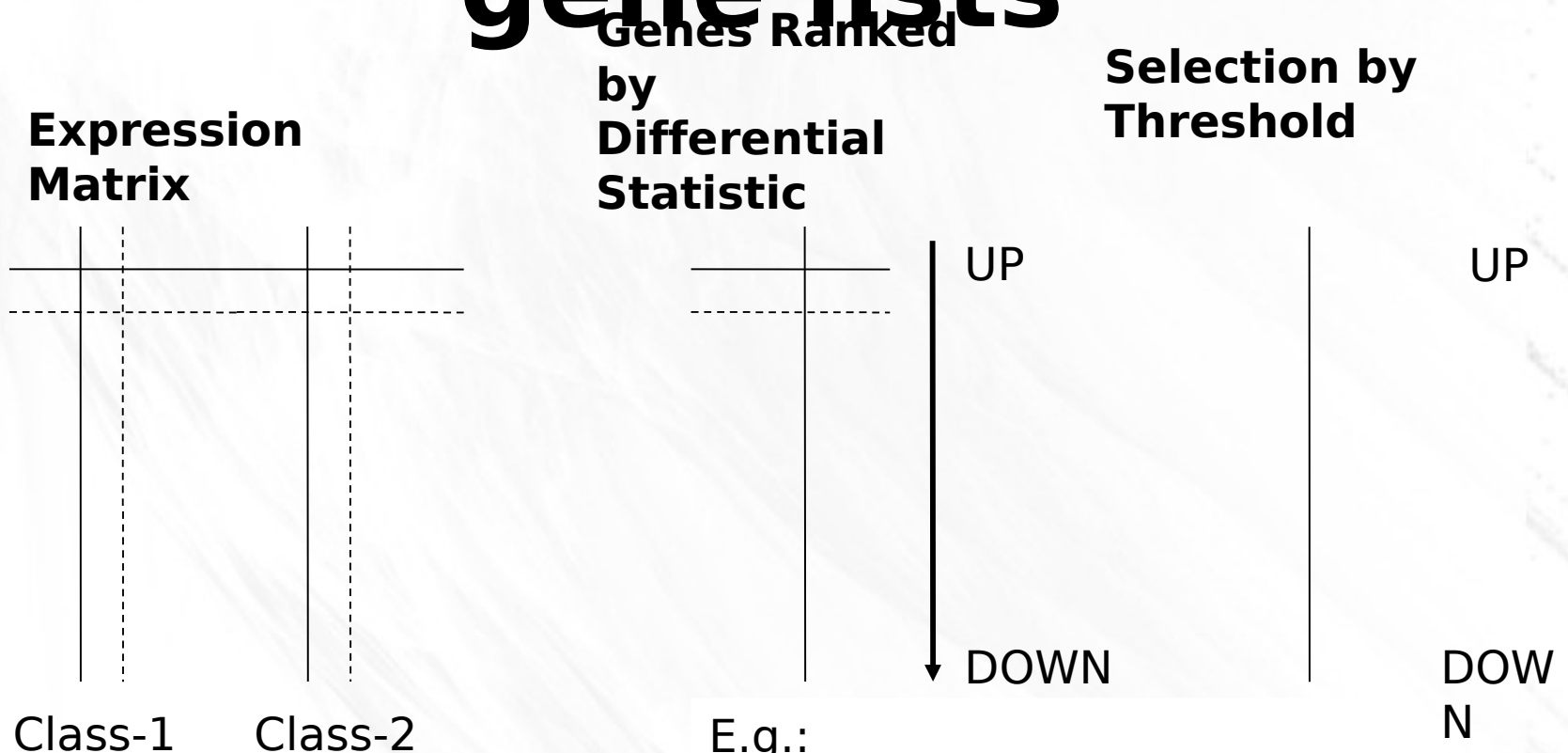- **<u>Gene list</u>** (e.g. expression change > 2-fold)

  - Answers the question: **Are any gene sets surprisingly enriched (or depleted) in my gene list?**

  - Statistical test: Fisher's Exact Test (Hypergeometric test)

- **<u>Ranked list</u>** (e.g. by differential expression)

  - Answers the question: **Are any gene set ranked surprisingly high or low in my ranked list of genes?**

  - Statistical test: GSEA

# Gene list enrichment analysis

- Given:

  1. Gene list: e.g. RRP6, MRD1, RRP7, RRP43, RRP42 (yeast)

  2. Gene sets or annotations: e.g. Gene ontology, transcription factor binding sites in promoter

- Question: *Are any of the gene annotations <u>surprisingly</u> enriched in the gene list?*

# Two-class design for gene lists

**Expression Matrix**

**Genes Ranked by Differential Statistic**

**Selection by Threshold**

UP

UP

DOWN

DOW N

Class-1     Class-2

E.g.:
- Fold change
- Log (ratio)
- t-test
- Significance analysis of microarrays

# Recipe for gene list enrichment test

- **Step 1:** Rank your gene list,

- **Step 2:** Select your gene sets to test for enrichment,

- **Step 3:** Run enrichment tests and correct for multiple testing, if necessary,

- **Step 4:** Interpret your enrichments

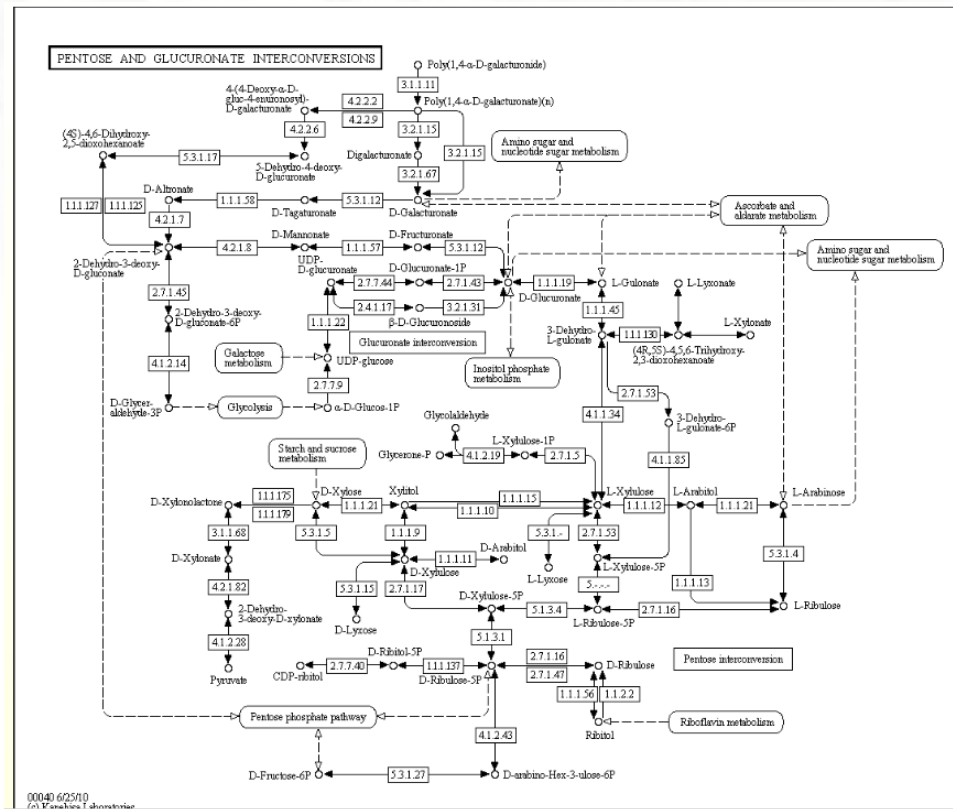- **Step 5:** Make conclusions

# **Theory component**

- Hypergeometric test for calculating enrichment P-values for gene lists

- GSEA for computing enrichment P-values for ranked lists

- Multiple test corrections:

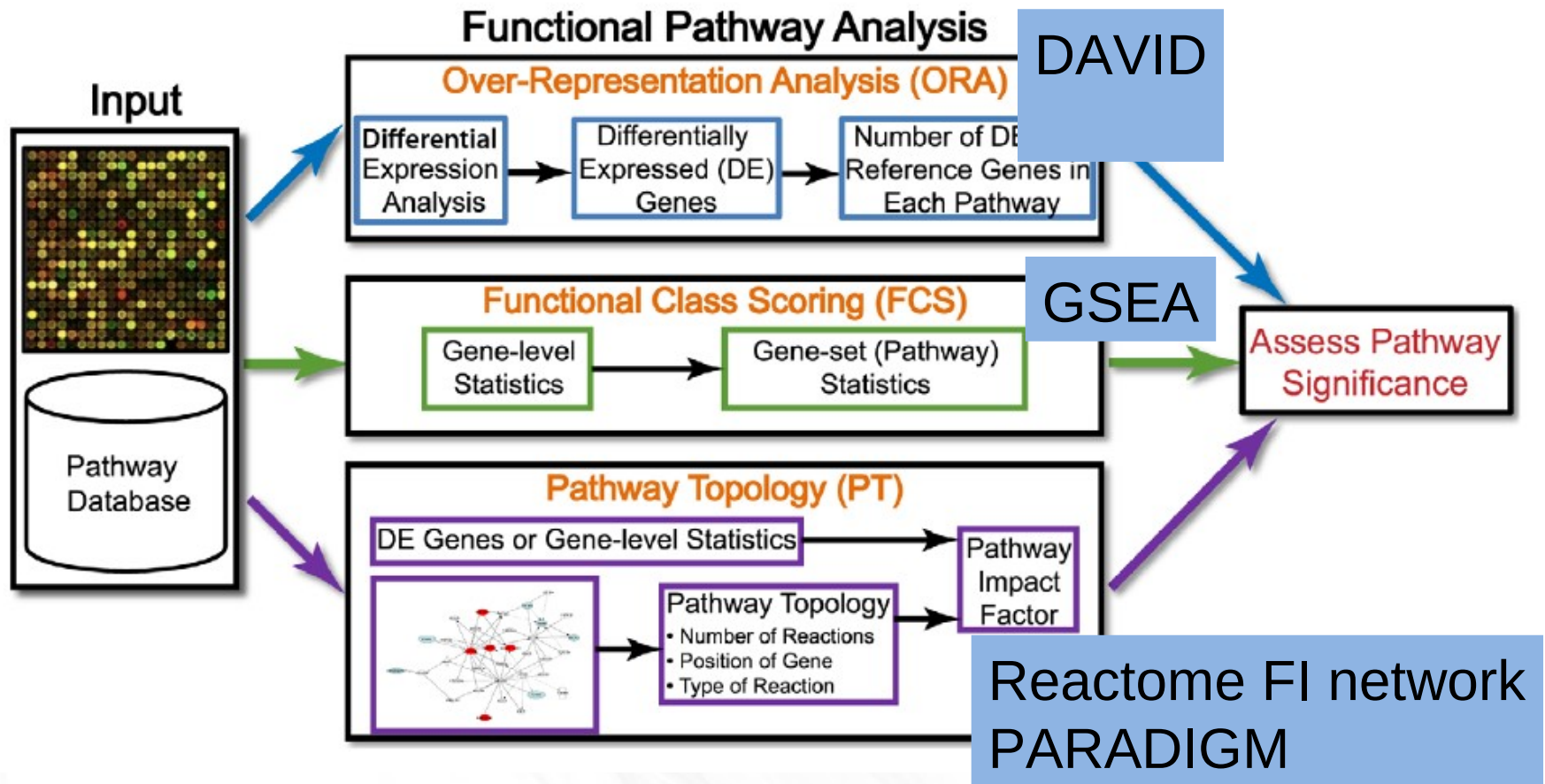  - Bonferroni

  - Benjamini-Hochberg FDR

# Important details

- To test for *under-enrichment* of "black", test for *over-enrichment* of "red".

- Need to choose "background population" appropriately, e.g., if only portion of the total gene complement is queried (or available for annotation), only use that population as background.

- To test for enrichment of more than one independent types of annotation (red vs black and circle vs square), apply Fisher's exact test separately for each type.

# Pathway and Network Analysis of –omics Data

# Classes of Gene Set Analysis



Khatri *et al.* PLOS Comp Bio. 8:1 2012

# Limitations of Gene Set Enrichment Analysis

- Many possible gene sets – diseases, molecular function, biological process, cellular compartment, pathways...

- Gene sets are heavily overlapping; need to sort through lists of enriched gene sets!

- "Bags of genes" obscure regulatory relationships among them.
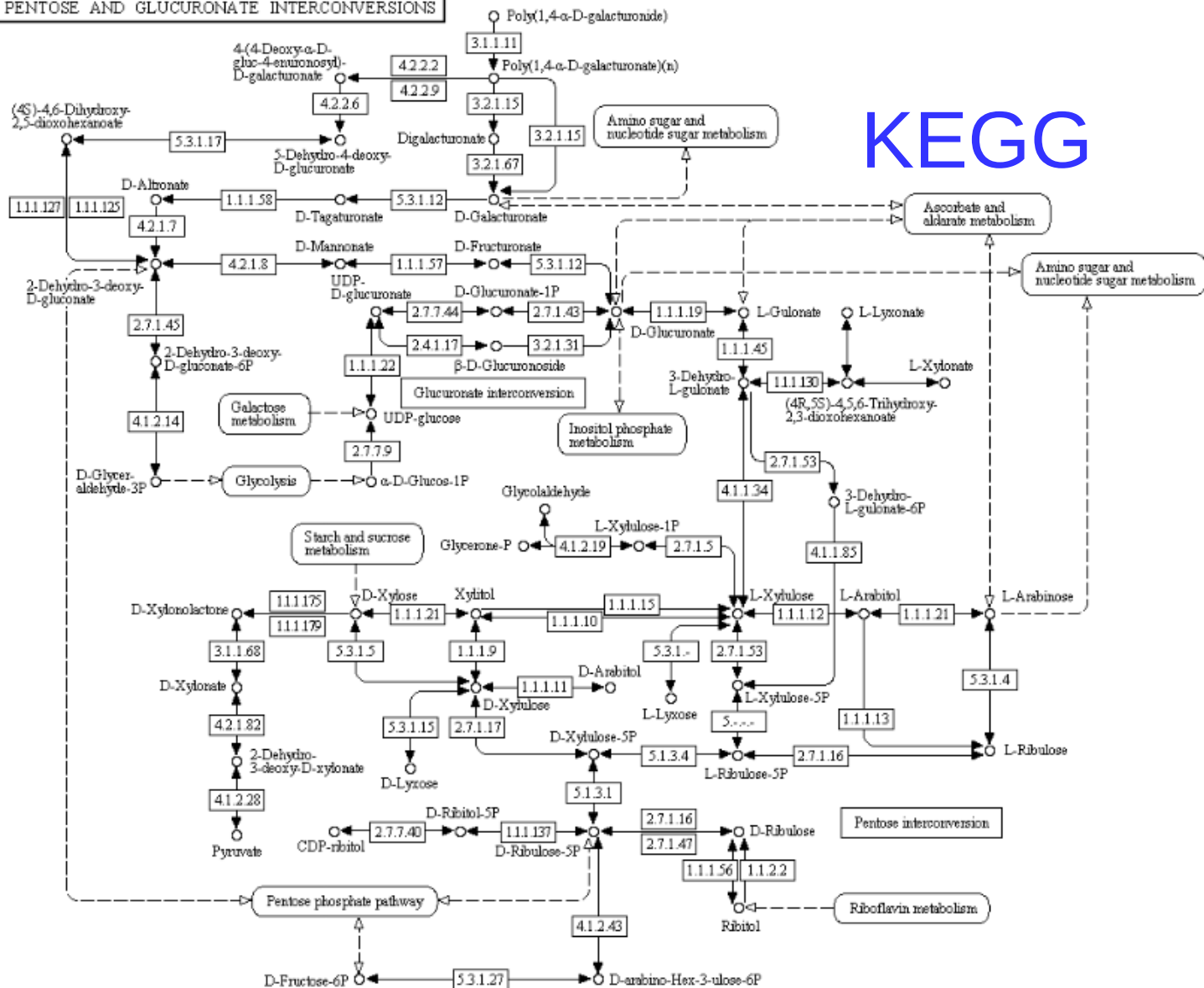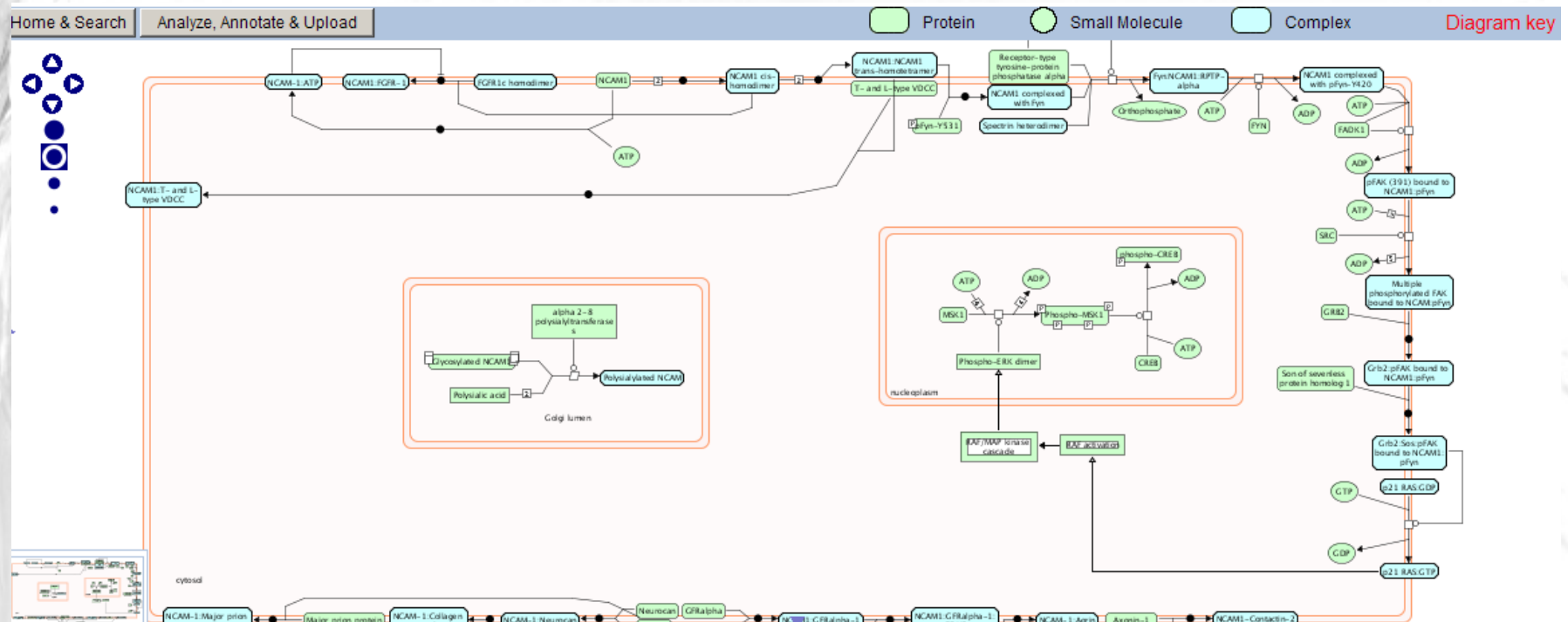
# Pathway Databases

- Advantages:
  - Usually curated.
  - Biochemical view of biological processes.
  - Cause and effect captured.
  - Human-interpretable visualizations.
- Disadvantages:
  - Sparse coverage of genome.
  - Different databases disagree on boundaries of pathways.

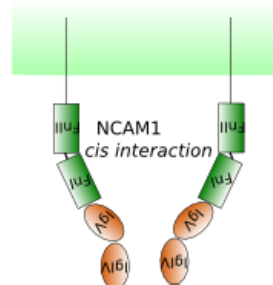PENTOSE AND GLUCURONATE INTERCONVERSIONS

KEGG

# Reactome



NCAM1 mediated intracellular signal transduction is represented in the figure below. The Ig domains in NCAM1 are represented in orange ovals and Fn domains in green squares. The tyrosine residues susceptible to phosphorylation are represented in red circles and their positions are numbered. Phosphorylation is represented by red arrows and dephosphorylation by yellow. Ig, Immunoglobulin domain; Fn, Fibronectin domain; Fyn, Proto-oncogene tyrosine-protein kinase Fyn; FAK, focal adhesion kinase; RPTPalpha, Receptor-type tyrosine-protein phosphatase; Grb2, Growth factor receptor-bound protein 2; SOS, Son of sevenless homolog; Raf, RAF proto-oncogene serine/threonine-protein kinase; MEK, MAPK and ERK kinase; ERK, Extracellular signal-regulated kinase; MSK1, Mitogen and stress activated protein kinase 1; CREB, Cyclic AMP-responsive element-binding protein; CRE, cAMP response elements. [Ditlevsen *et al* 2008]
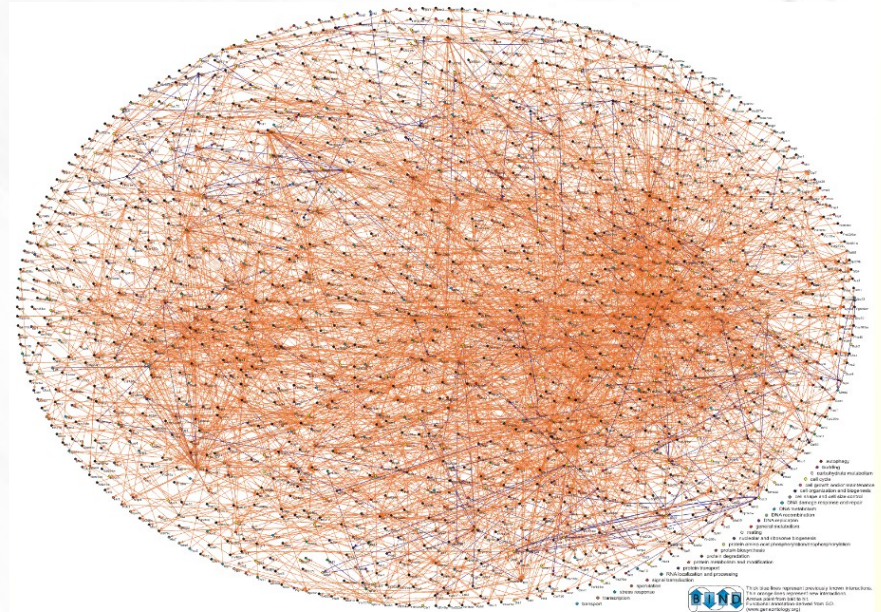
# Pathway Colorization

- Main feature offered by all pathway databases.

- Upload a gene list

- Database calculates an enrichment score on each pathway and displays ranked list.

- Browse into pathways of interest; download colorized pictures.
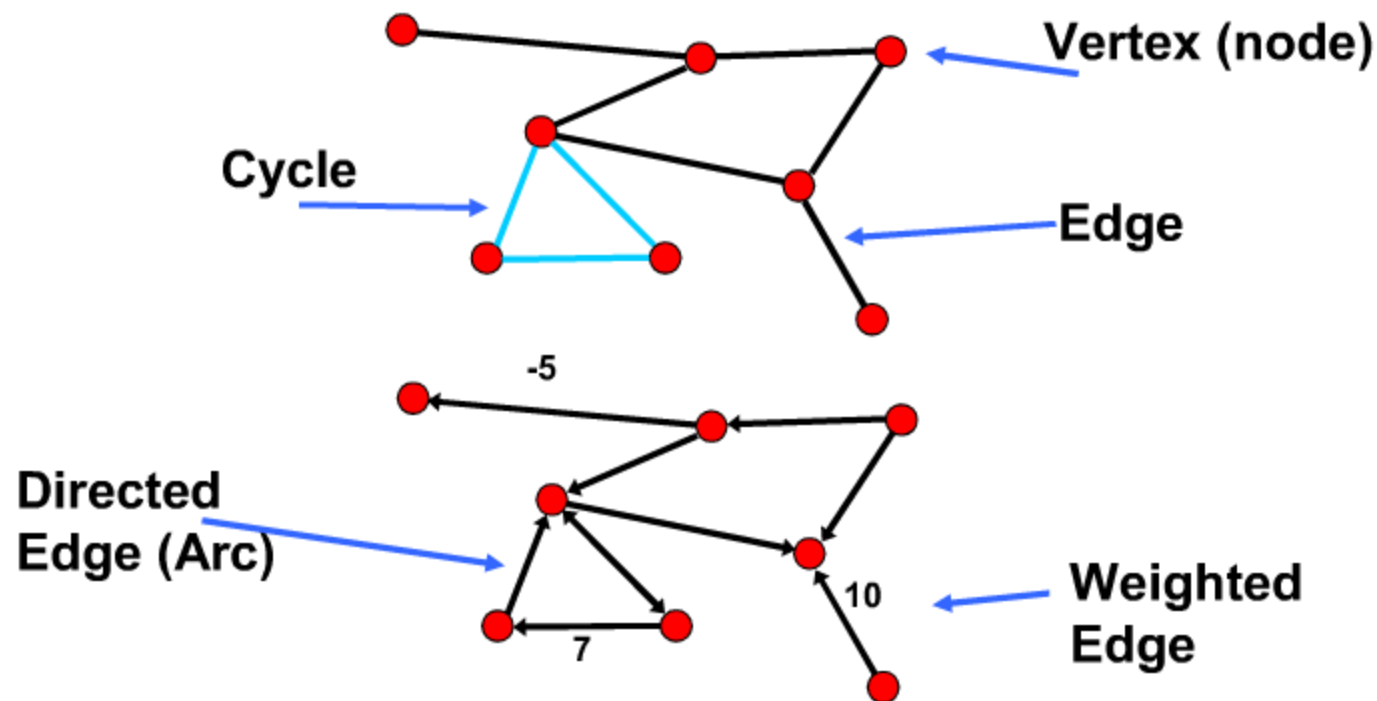
# Networks

- Pathways capture only the "well understood" portion of biology.

- Networks cover less well understood relationships:
  - Genetic interactions
  - Physical interaction
  - Coexpression
  - GO term sharing
  - Adjacency in pathways

# Networks

- E.g. Protein-protein interaction networks
- Useful if we don't know pathways
  - Could discover new pathways
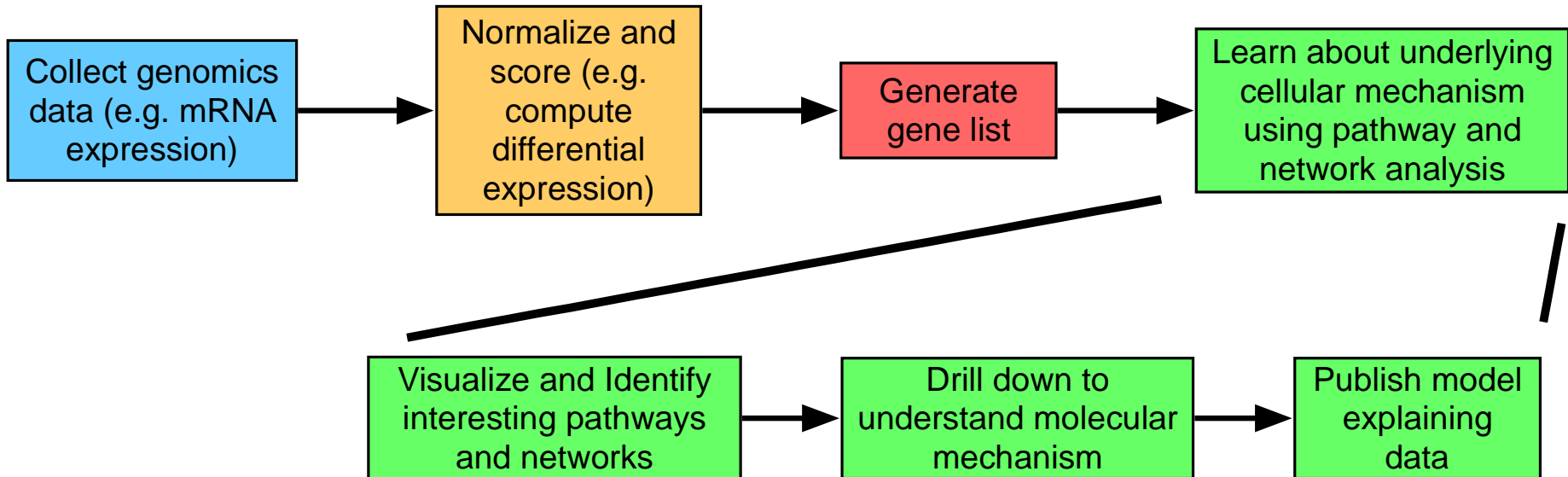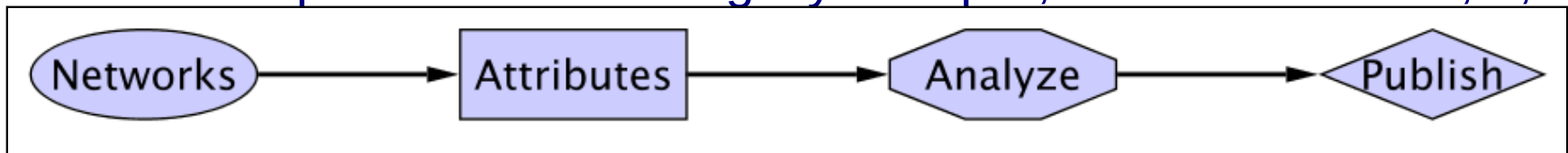
# Mapping Biology to a Network

- A simple mapping: Protein-protein interactions
  - one protein/node, one interaction/edge
- Edges can represent other relationships
  - Physical e.g. protein-protein interaction
  - Regulatory e.g. kinase activates target
  - Genetic e.g. epistasis
  - Similarity e.g. protein sequence similarity
- Critical: understand the mapping for network analysis

# Network Analysis Workflow

```
Collect genomics          Normalize and                              Learn about underlying
data (e.g. mRNA    →       score (e.g.      →    Generate     →       cellular mechanism
expression)                compute               gene list           using pathway and
                           differential                               network analysis
                           expression)
```

```
Visualize and Identify         Drill down to              Publish model
interesting pathways    →   understand molecular   →       explaining
and networks                   mechanism                      data
```
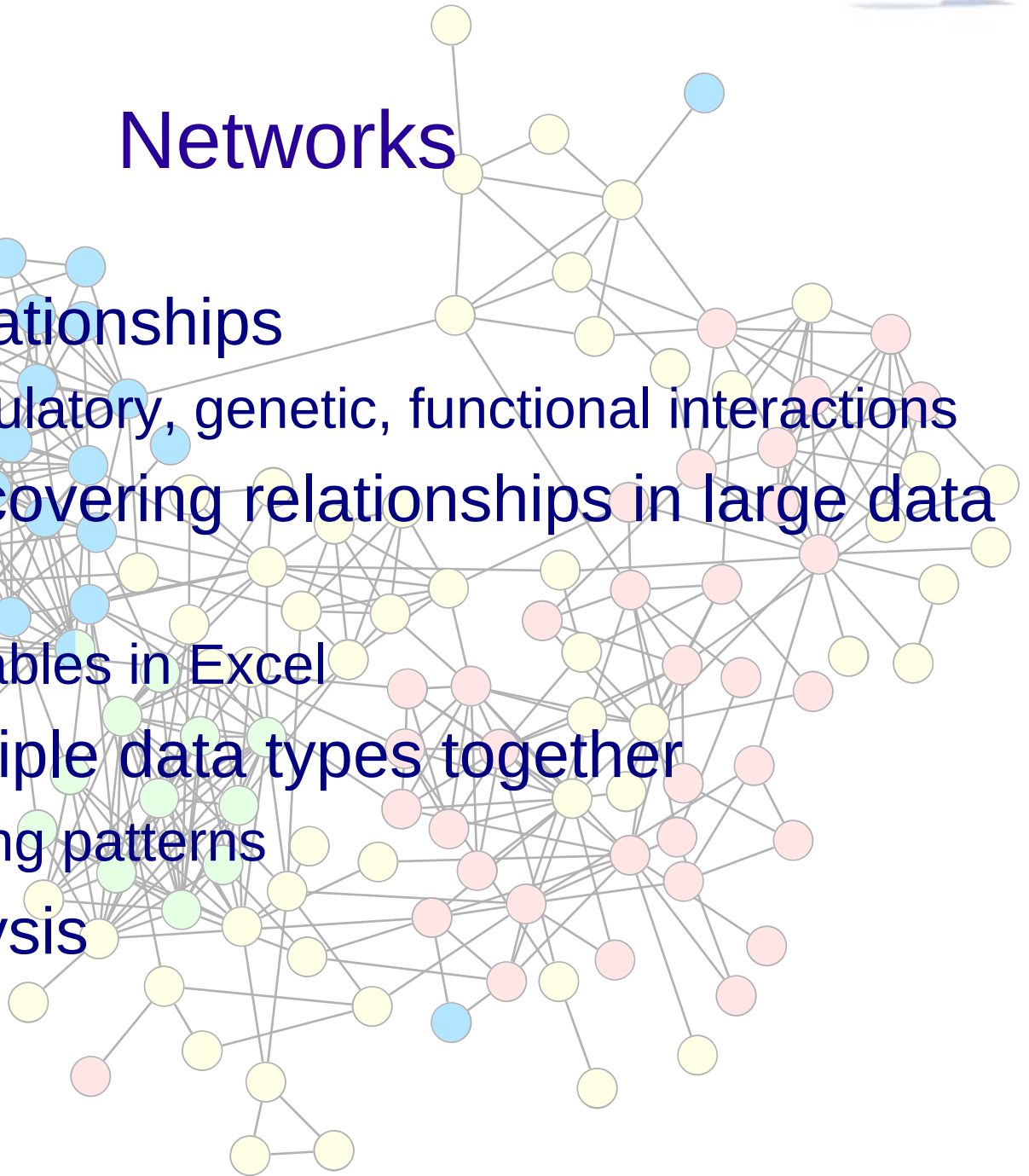
- A specific example of this workflow:
    - Cline, et al. "Integration of biological networks and gene expression data using Cytoscape", Nature Protocols, 2,

```
Networks   →   Attributes   →   Analyze   →   Publish
```
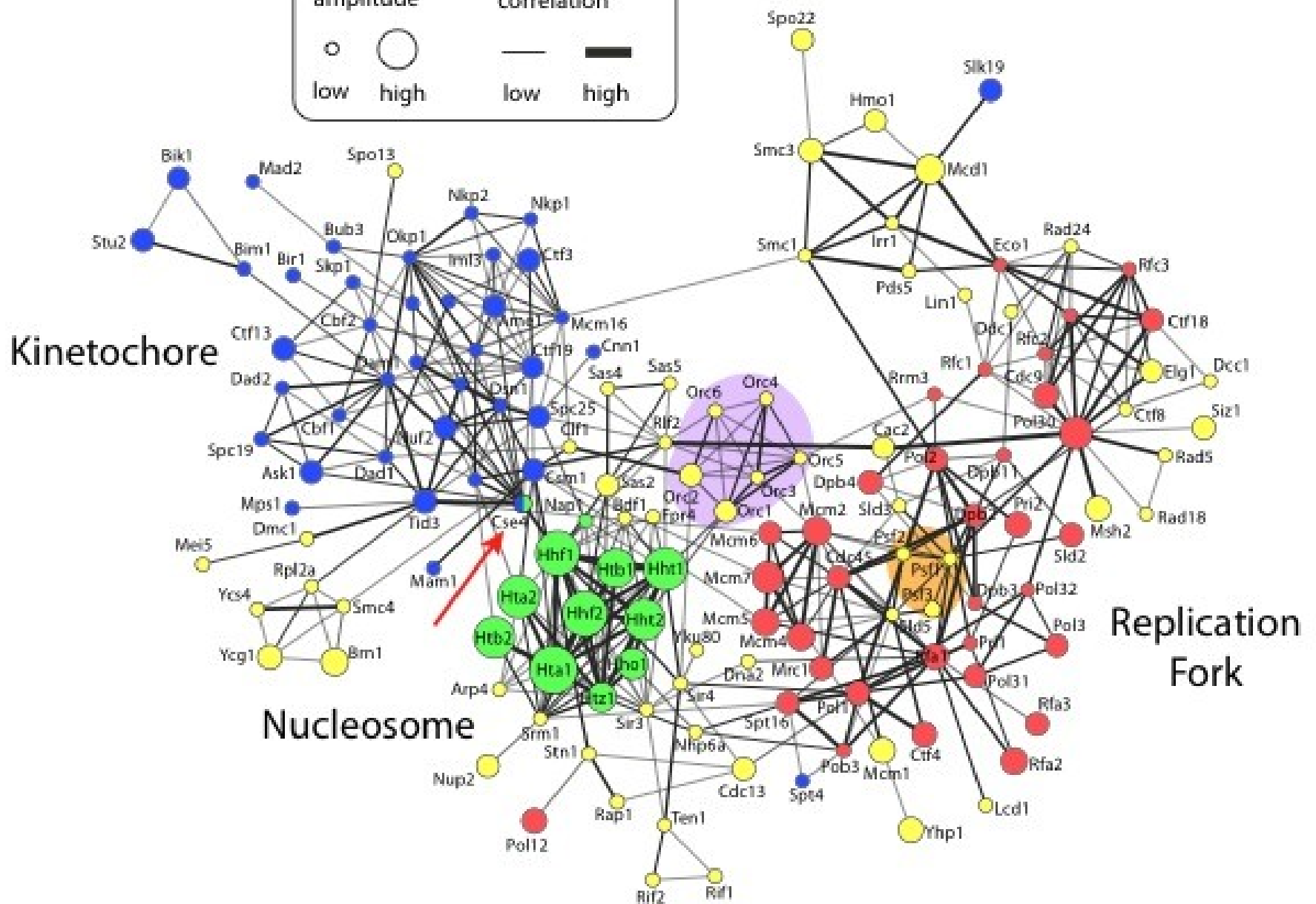
# Networks

- Represent relationships
  - Physical, regulatory, genetic, functional interactions
- Useful for discovering relationships in large data sets
  - Better than tables in Excel
- Visualize multiple data types together
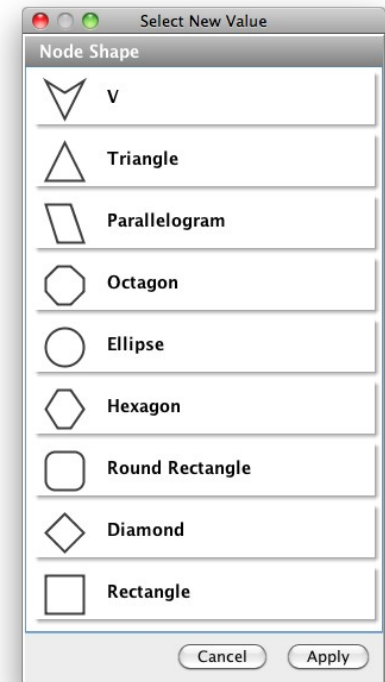  - See interesting patterns
- Network analysis

# Summary

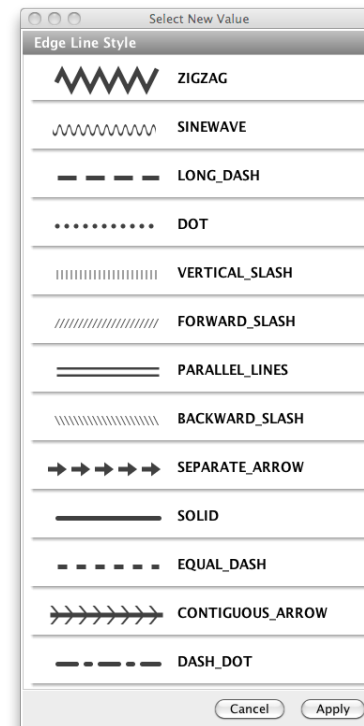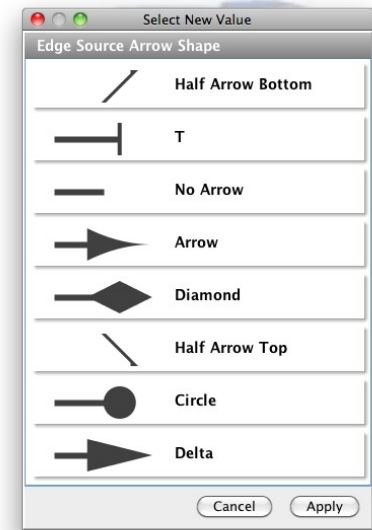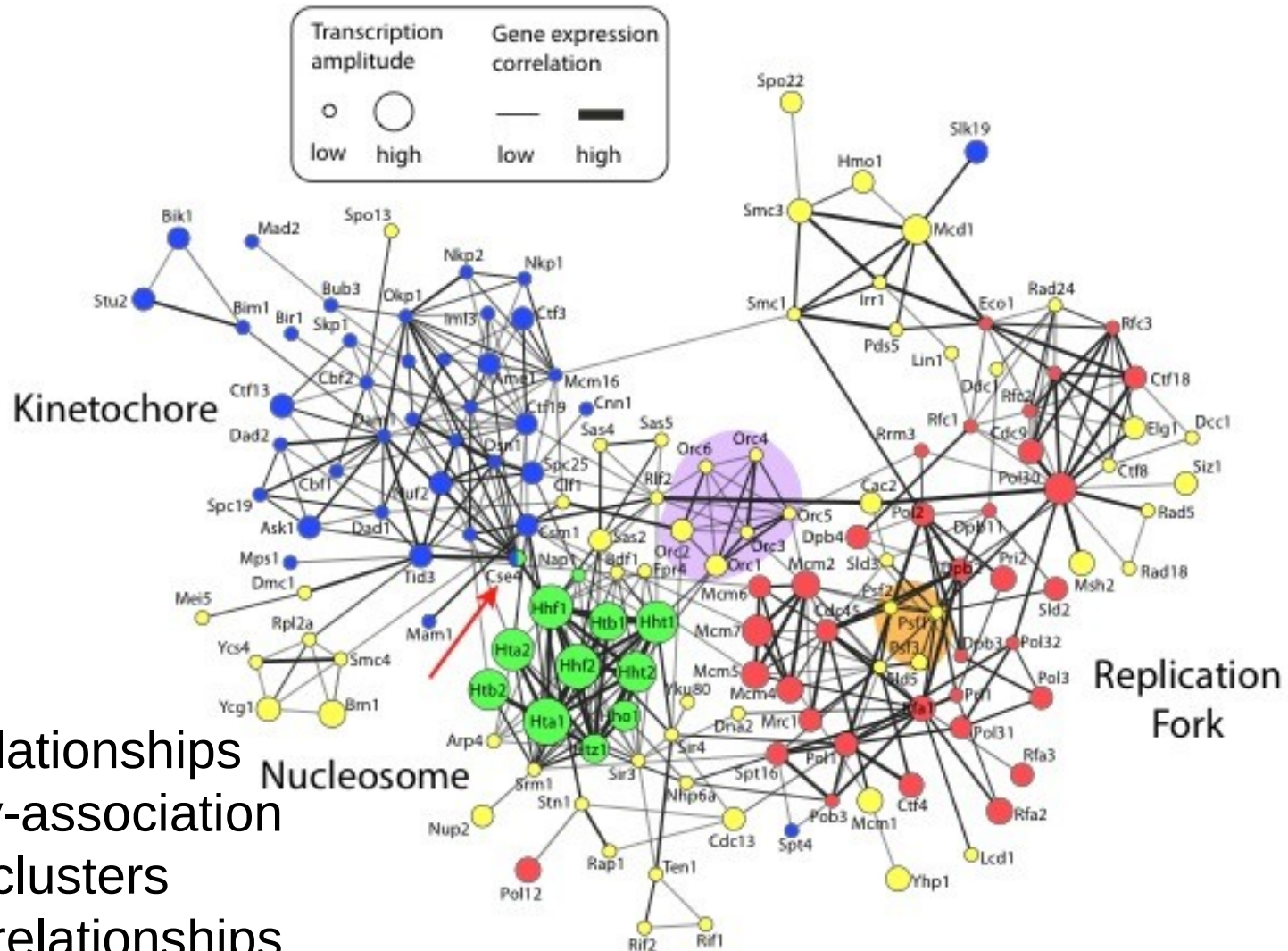- Networks are useful for seeing relationships in large data sets
- Important to understand what the nodes and edges mean
- Important to define the biological question - know what you want to do with your gene list or network
- Many methods available for gene list and network analysis
  - Good to determine your question and search for a solution
  - Or get to know many methods and see how they can be applied to your data
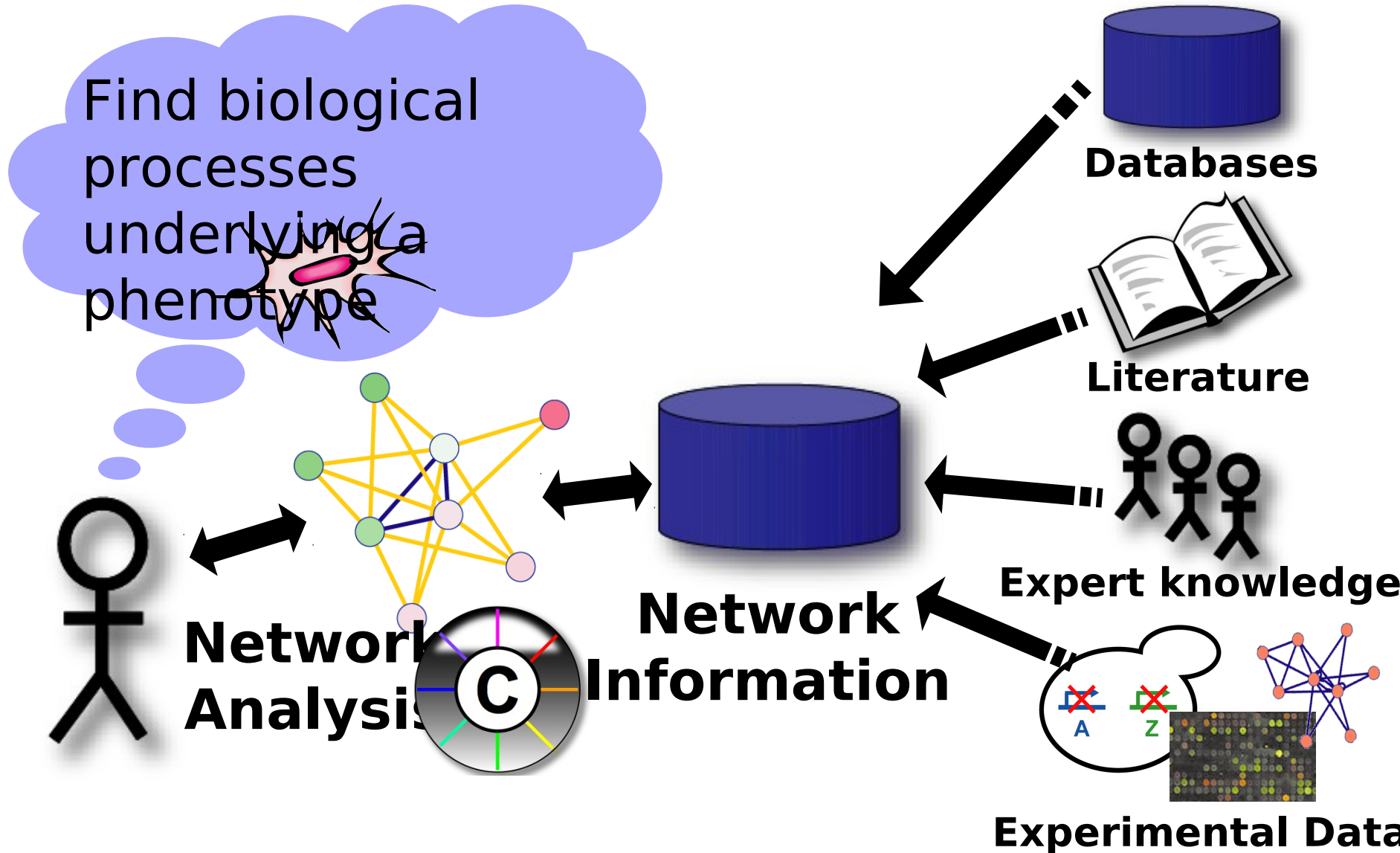
# Visual Features

- Node and edge attributes
  - Text (string), integer, float, Boolean, list
  - E.g. represent gene, interaction attributes
- Visual attributes
  - Node, edge visual properties
  - Colour, shape, size, borders, opacity...



Select New Value — Edge Source Arrow Shape

| | |
|---|---|
| | Half Arrow Bottom |
| | T |
| | No Arrow |
| | Arrow |
| | Diamond |
| | Half Arrow Top |
| | Circle |
| | Delta |

Cancel    Apply

Select New Value — Edge Line Style

| | |
|---|---|
| | ZIGZAG |
| | SINEWAVE |
| | LONG_DASH |
| | DOT |
| | VERTICAL_SLASH |
| | FORWARD_SLASH |
| | PARALLEL_LINES |
| | BACKWARD_SLASH |
| | SEPARATE_ARROW |
| | SOLID |
| | EQUAL_DASH |
| | CONTIGUOUS_ARROW |
| | DASH_DOT |

Cancel    Apply

Select New Value — Node Shape

| | |
|---|---|
| | V |
| | Triangle |
| | Parallelogram |
| | Octagon |
| | Ellipse |
| | Hexagon |
| | Round Rectangle |
| | Diamond |
| | Rectangle |

Cancel    Apply

# Visually Interpreting a Network



Data relationships
Guilt-by-association
Dense clusters
Global relationships

# Active Community

http://www.cytoscape.org

- 10,000s users, >5000 downloads/month
- Help

  Cline MS et al. Integration of biological networks and gene expression data using Cytoscape Nat Protoc. 2007;2(10):2366-82

  - Documentation, data sets
  - Mailing lists
  - **http://tutorials.cytoscape.org**
- Annual Conference: TBD, North America 2014
- >200 Apps Extend Functionality
  - Build your own, requires programming

- Cytoscape is a useful, free software tool for network visualization and analysis
- Provides basic network manipulation features
- Apps are available to extend the functionality

# Gene List and Network Analysis Overview