

# Large-scale functional annotation

*Kris Gunsalus  
XDAS Fall 2019*

# Pathway and network analysis

- Approaches to analyze functional enrichment
- Network visualization
- Network inference

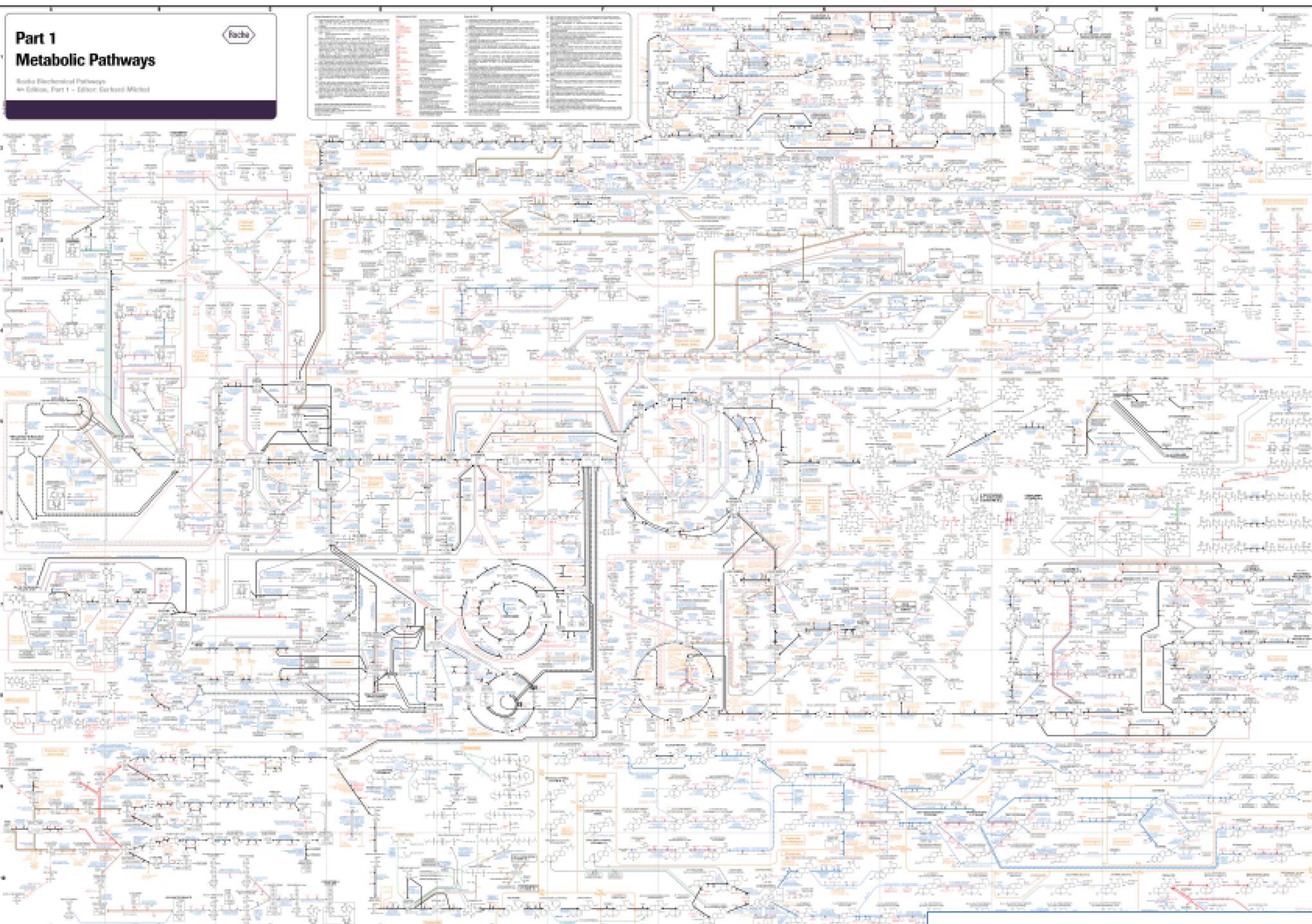
# How do we define "function"?

- **KEGG: Kyoto Encyclopedia of Genes and Genomes**
- **Reactome: expert-curated pathways**
- **GO: Gene Ontology**
  - Directed acyclic graph
  - Three ontologies:
    - Biological Process (BP)
    - Cellular Component (CC)
    - Molecular Function (MF)
  - Multiple sources of annotation
    - Biochemistry, genetics, electronically inferred (e.g. from homologs)

**Part 1**  
**Metabolic Pathways**

Rocke Biochemical Pathways  
4th Edition, Part 1 - Editor: Gerhard Michel

Rocke



## Part 1 Metabolic Pathways

Reed's Biochemical Pathways  
4th Edition, Part 1 – Editor: Gerhard Michael

Facts

### Carbohydrate Metabolism Acidic Carbohydrate Derivatives

### Carbohydrate Metabolism Inositol

### Carbohydrate Metabolism Di- and Polysaccharides

### Carbohydrate Metabolism Nucleotide Sugars

### Carbohydrate Metabolism Glycolysis and Gluconeogenesis

### C1-Metabolism

### Bacterial Metabolism Methanogenesis

### Amino Acid Metabolism Leucine, Isoleucine, Valine

### Cofactors and Vitamins Coenzyme A

### Lipid Metabolism Glyco- and Phospholipids

### Lipid Metabolism Sphingolipids

### Lipid Metabolism Fatty Acids

### Bacterial Metabolism Alkane Oxidation

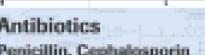
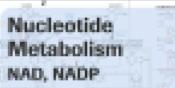
### Lipid Metabolism Carotenoids and Isoprenoids

### Steroid Metabolism Phytosterols

### Steroid Metabolism Cholesterol Synthesis

### Steroid Metabolism Mineralocorticoids and Glucocorticoids

### Steroid Metabolism Androgens and Estrogens



### Amino Acid Metabolism Histidine

### Carbohydrate Metabolism Amino Sugar Derivatives

### Bacterial Metabolism Methane Oxidation

### Carbohydrate Metabolism Pyruvate Turnover

### Citrate and Glyoxalate Cycle

### Amino Acid Metabolism Serine, Threonine, Cysteine, Methionine

### Amino Acid Metabolism Lysine

### Amino Acid Metabolism Tryptophan, Tyrosine, Phenylalanine

### Tetrapyrrole Metabolism Porphyrins, Cobalamin

### Tetrapyrrole Metabolism Heme, Cytochromes, Chlorophyll

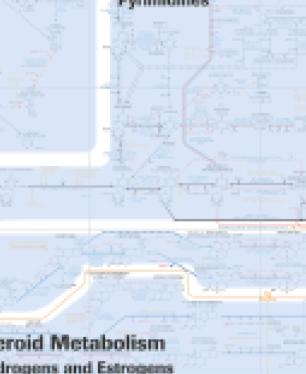
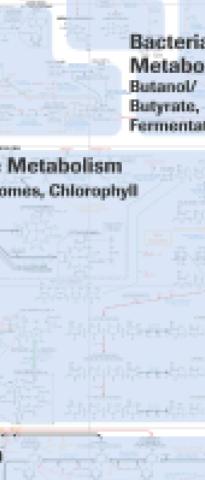
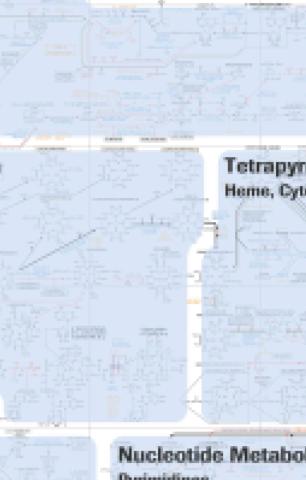
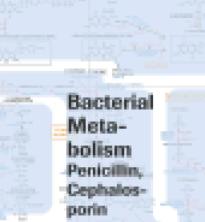
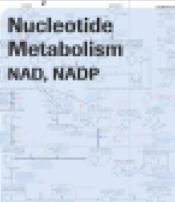
### Bacterial Metabolism Penicillin, Cephalosporin

### Bacterial Metabolism, Butanol/ Butyrate, Fermentation

### Amino Acid Metabolism Urea Cycle

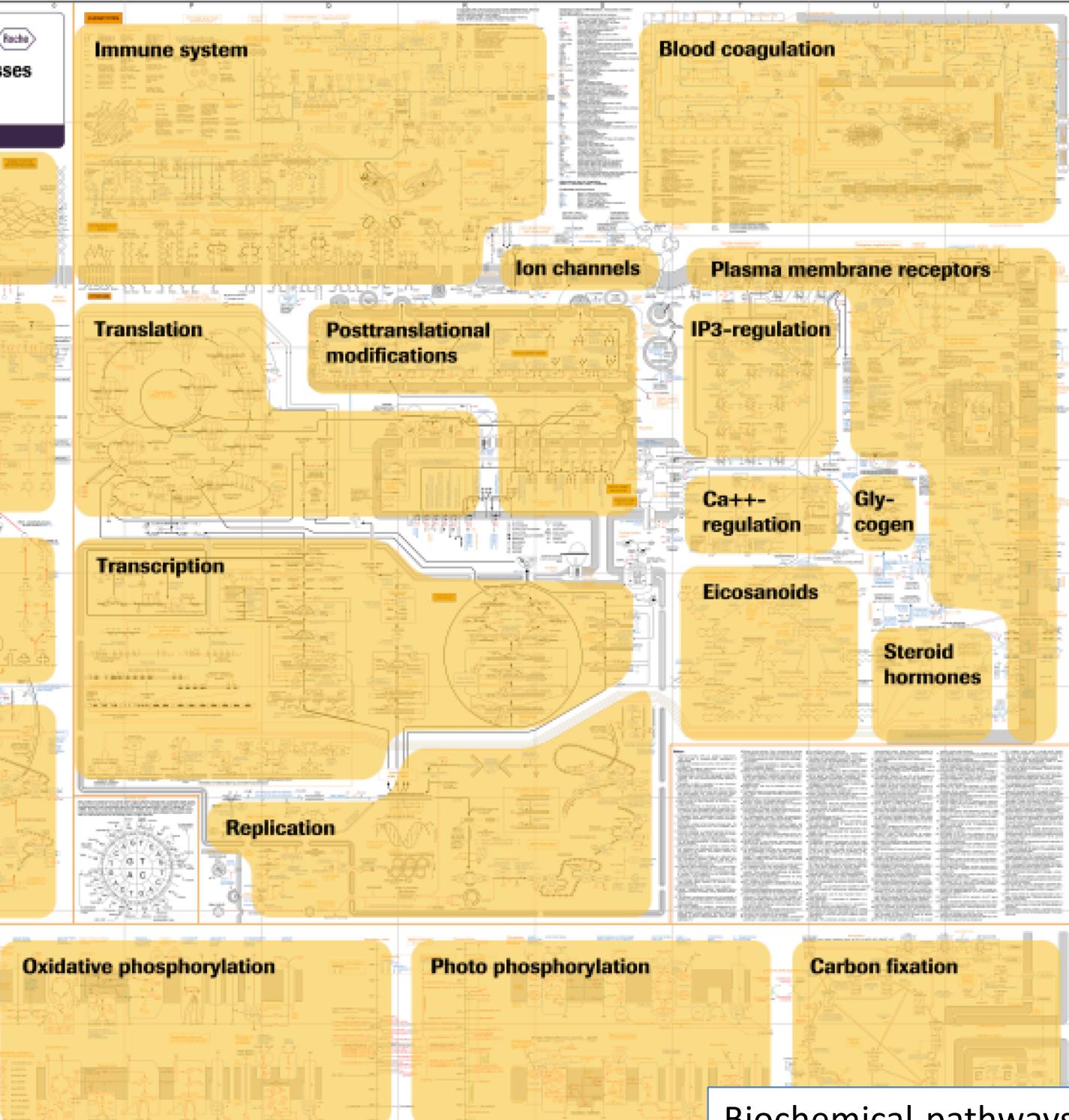
### Amino Acid Metabolism Glutamate, Proline, Hydroxyproline

### Nucleotide Metabolism Pyrimidines



**Part 2**  
**Cellular and Molecular Processes**

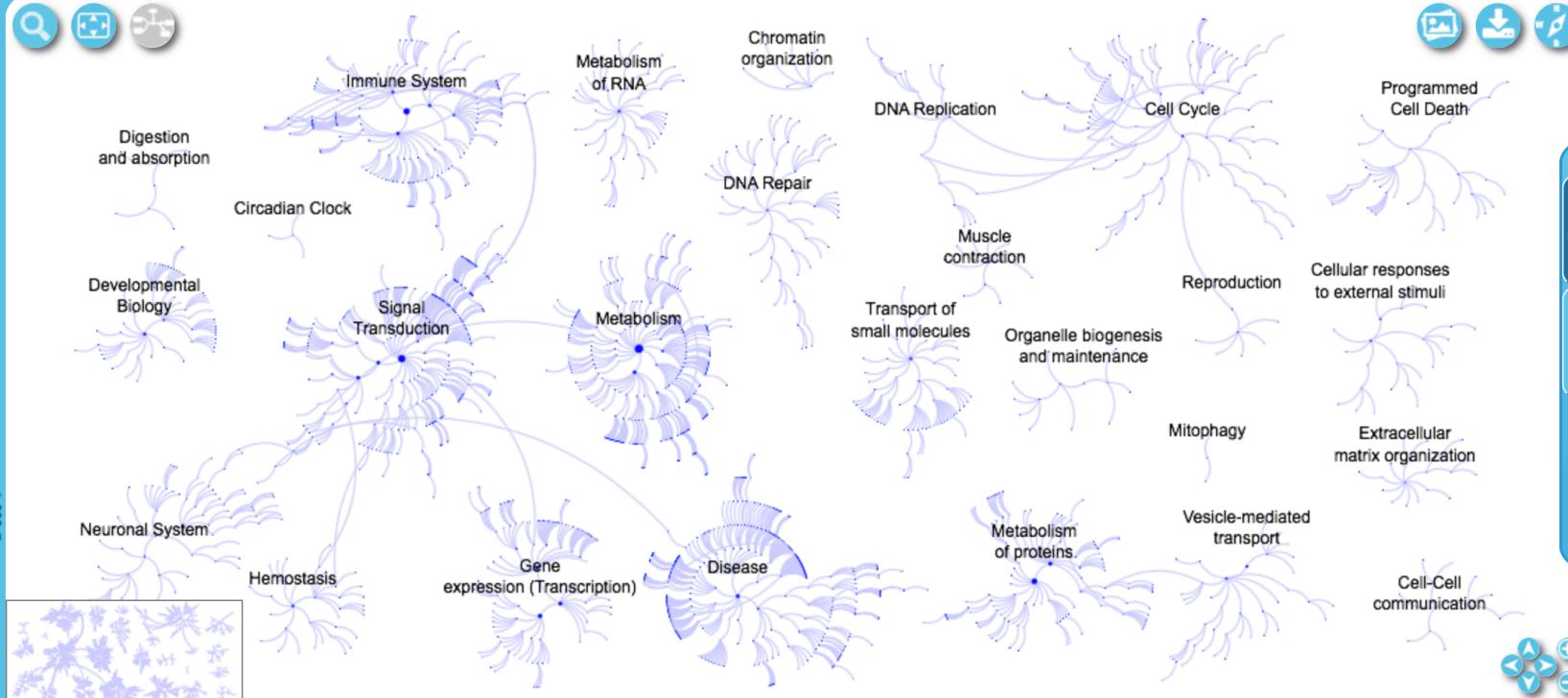
Reka Biochemical Pathways  
4th Edition, Part 2 – Editor: Gerhard Michel





## Event Hierarchy:

- + Cell Cycle
- + Cell-Cell communication
- + Cellular responses to external stimuli
- + Chromatin organization
- + Circadian Clock
- + Developmental Biology
- + Digestion and absorption
- + Disease
- + DNA Repair
- + DNA Replication
- + Extracellular matrix organization
- + Gene expression (Transcription)
- + Hemostasis
- + Immune System
- + Metabolism
- + Metabolism of proteins
- + Metabolism of RNA
- + Mitophagy
- + Muscle contraction
- + Neuronal System
- + Organelle biogenesis and maintenance
- + Programmed Cell Death
- + Reproduction
- + Signal Transduction
- + Transport of small molecules
- + Vesicle-mediated transport



## Description

Displays details when you select an item in the Pathway Browser. For example, when a reaction is selected, shows details including the input and output molecules, summary and references containing supporting evidence. When relevant, shows details of the catalyst, regulators, preceding and following events.

## Molecules

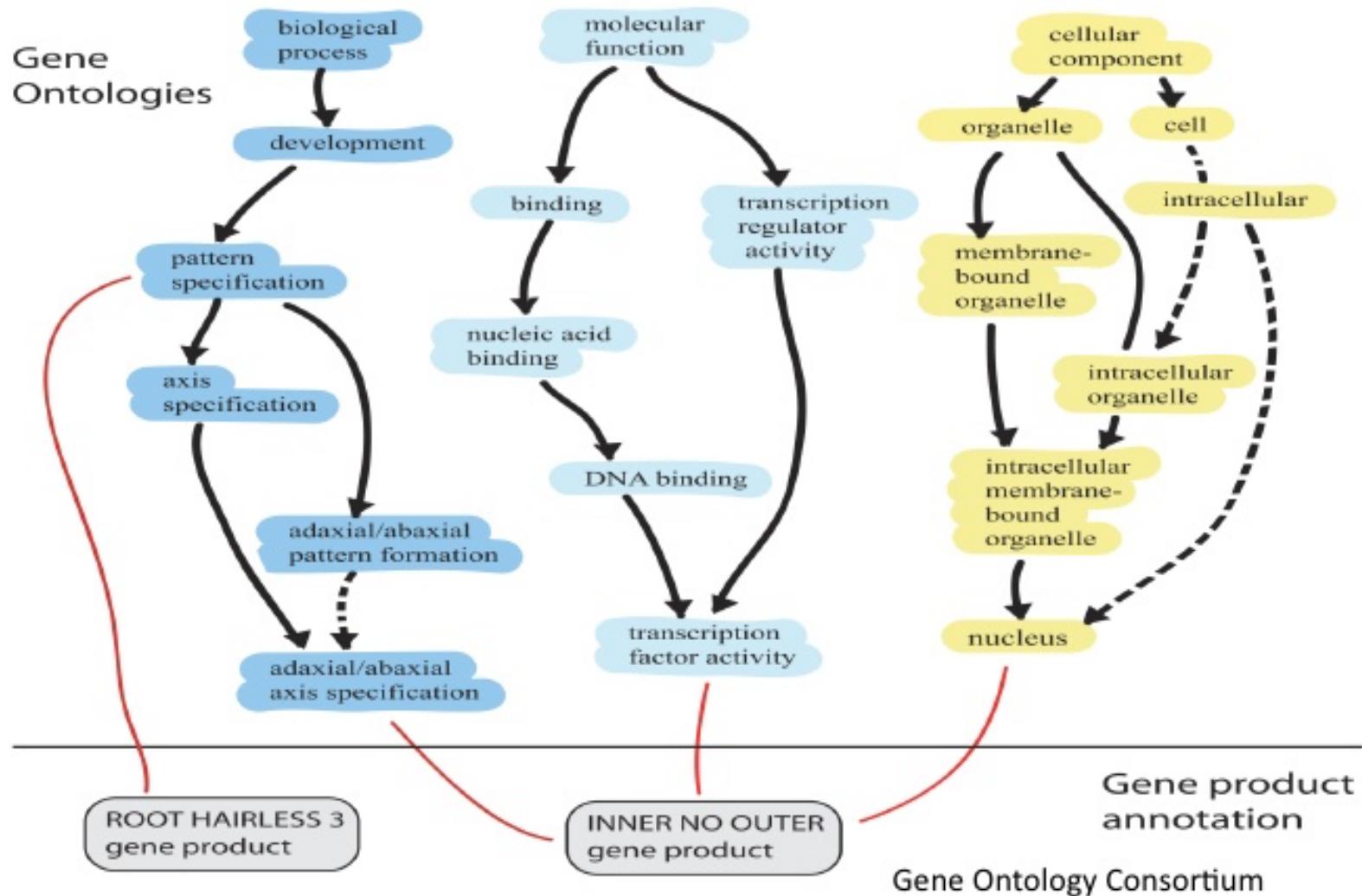
## Structures

## Expression

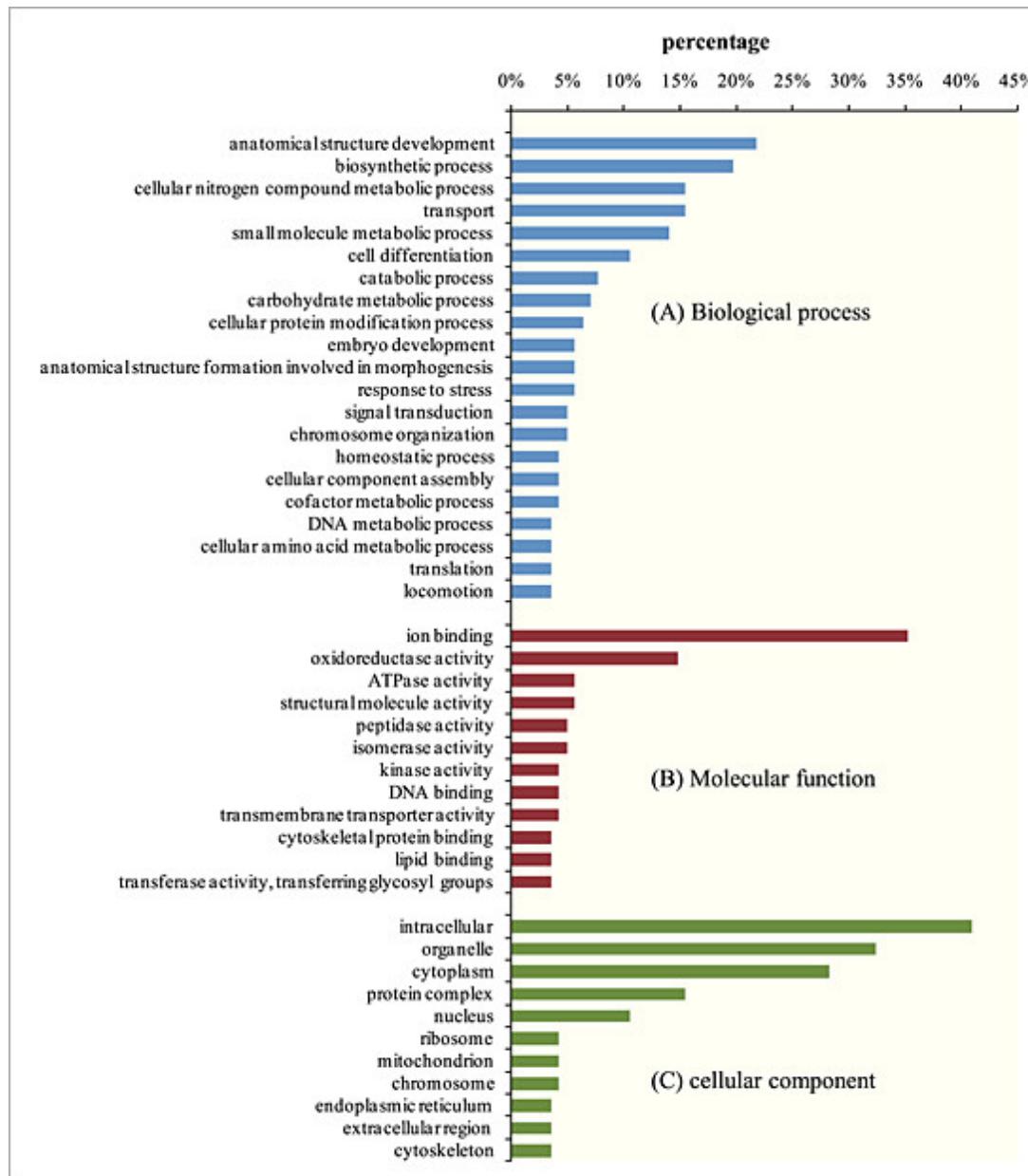
## Analysis

## Downloads

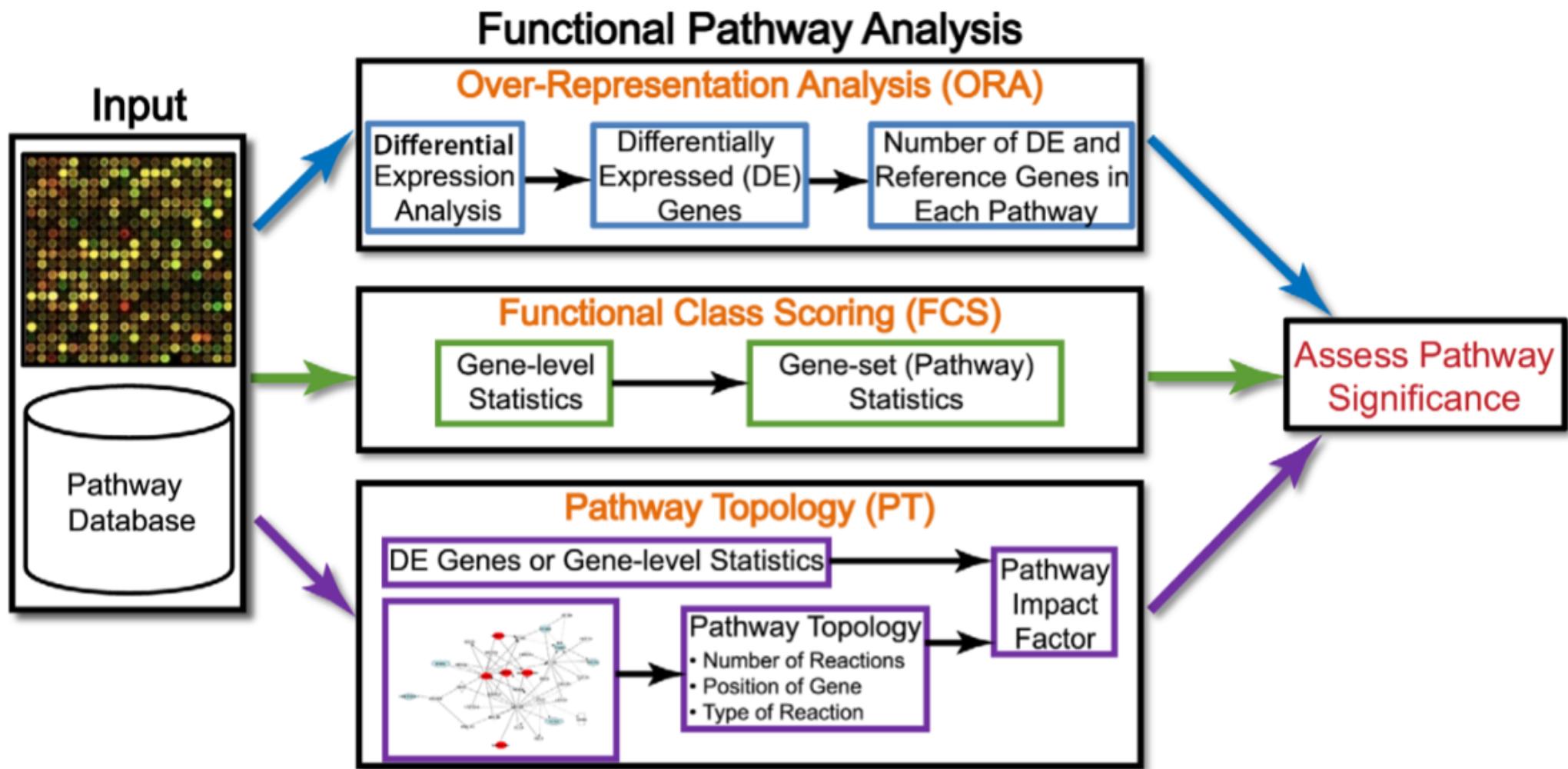
# Gene Ontology: example



# Gene Ontology: example



# Three generations of pathway analysis



- ⇒ Large-scale studies often produce **lists of genes** that are perturbed in some way
- ⇒ Pathway analysis **reduces complexity** and **increases explanatory power**

# Over-representation Analysis (ORA)

- Approach
  - Input list
    - up-or down-regulated genes, RNAi hits, etc.
  - Assess proportion of genes with a particular annotation
    - GO term, pathway membership
  - Compare against universe of all possible genes
  - Assess over- or under-representation
  - "2x2" method: Tabular Statistics
    - hypergeometric, chi-squared, binomial distributions
- Limitations
  - Statistical tests independent of measured changes
  - Uses only most significant genes
  - Assumes genes (and annotations) are independent of each other

# Functional Class Scoring (FCS)

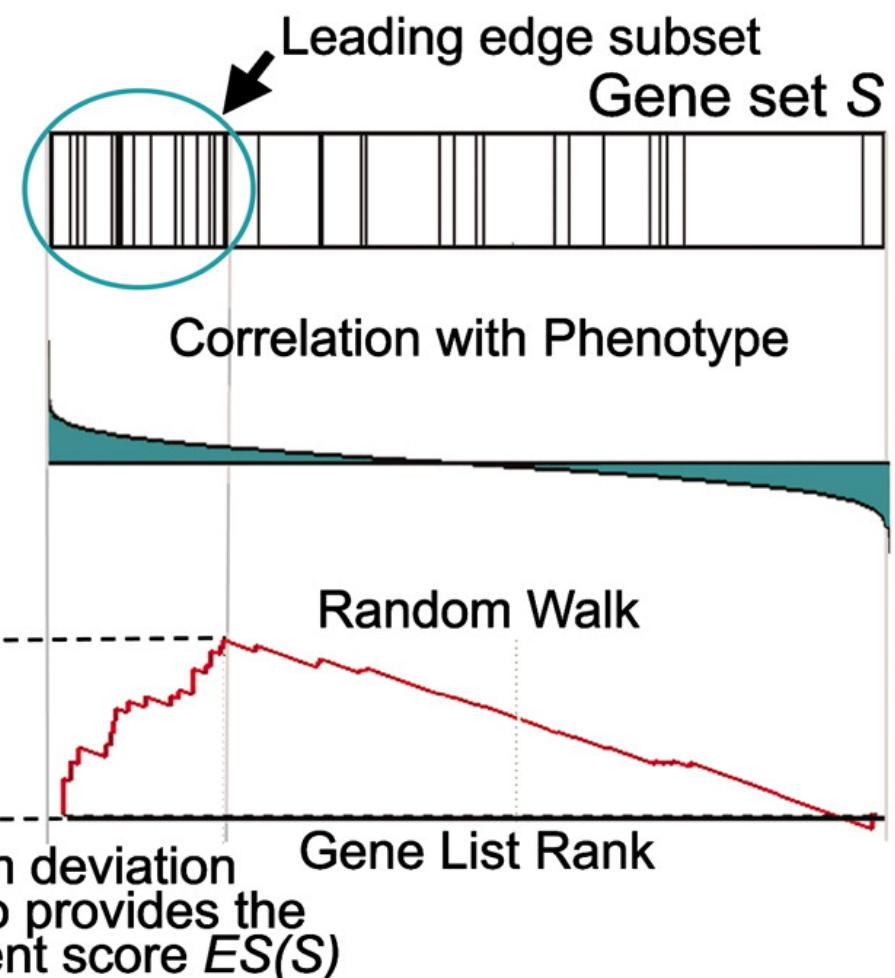
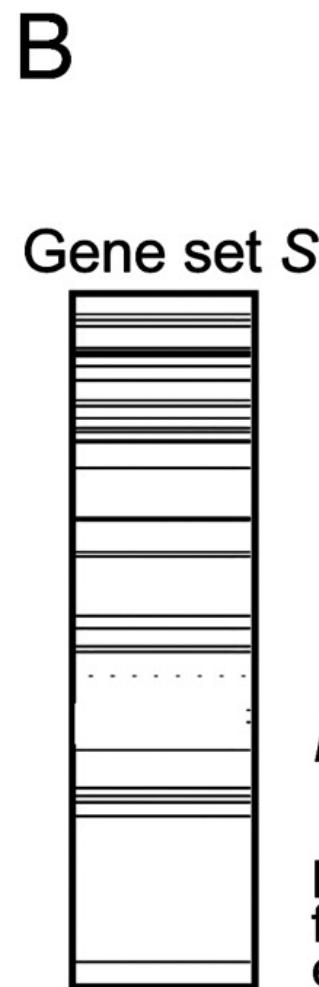
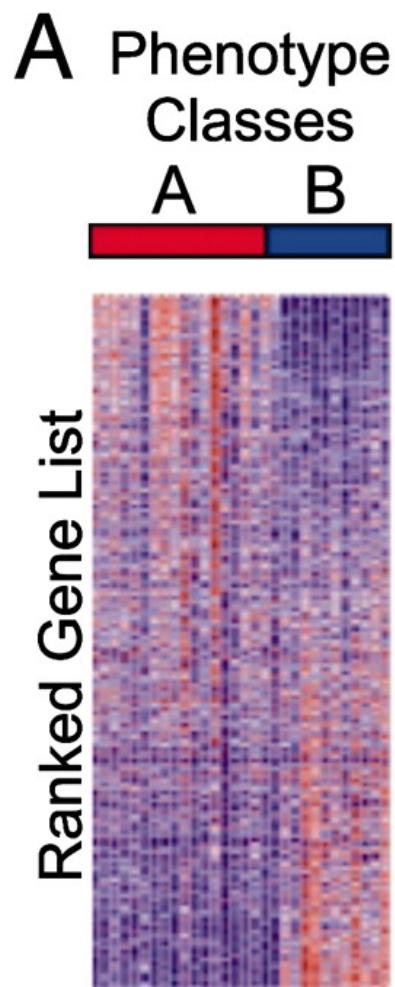
- Approach
  - Try to capture weaker **coordinated changes** that can cumulatively be informative
  - Compute **gene-level statistics**: significance of deviation from background distribution
    - t-test, Z-score, Q-statistic, ANOVA, signal-to-noise, correlation with phenotype
  - Aggregate gene-level statistics into a single **pathway-level statistic**
    - Univariate or multivariate (account for dependencies b/w genes)
    - E.g. KS test, Wilcoxon rank sum (Mann-Whitney), ANCOVA
  - Assess pathway-level **significance**
    - **Competitive** null hypothesis (permute pathway labels)  
vs. **self-contained** null hypothesis (permute class labels)

# Functional Class Scoring (FCS)

- Advantages
  - Do not require arbitrary threshold for significance of groups
  - Use all available molecular measurements for pathway analysis
  - Detect coordinated changes in pathway components
  - Account for dependence between genes in pathway
- Limitations
  - Analyzes each pathway independently
  - May use changes to rank genes in a pathway only

*(Exception: some methods use sum or mean of gene-level scores to compute pathway-level score)*

# GSEA: Gene Set Enrichment Analysis



# Controlling errors

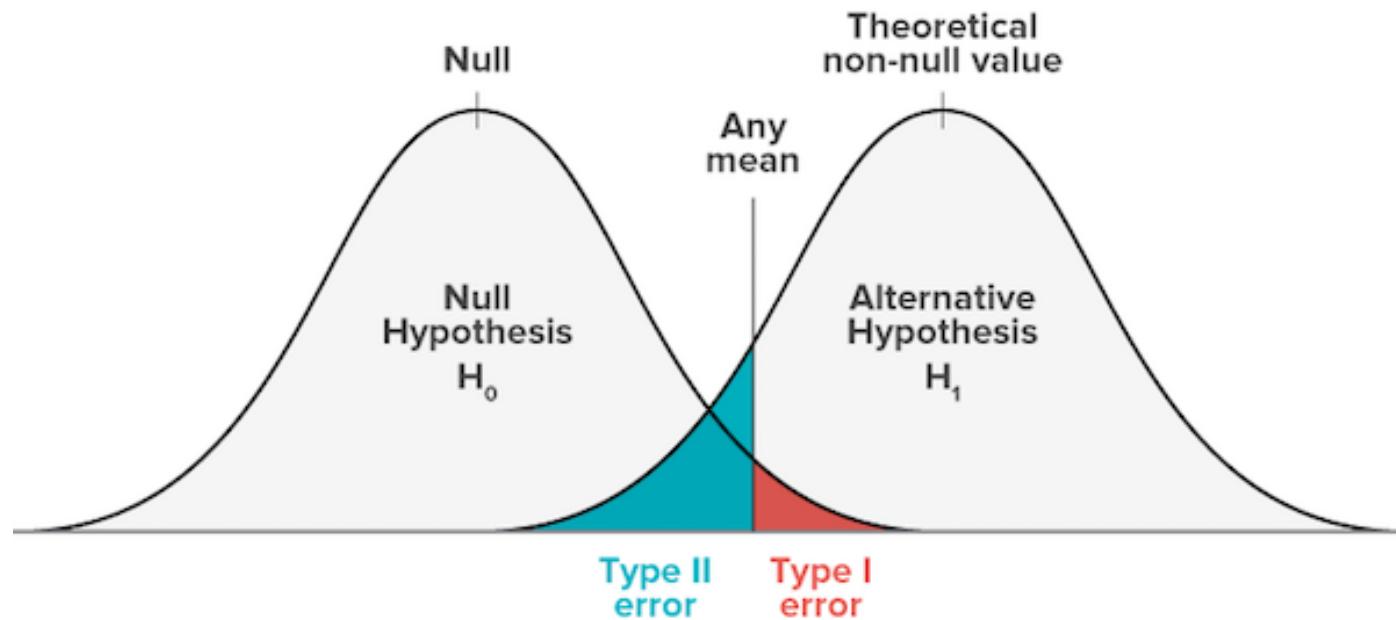
| Table of error types                     |                | Null hypothesis ( $H_0$ ) is         |                                      |
|--|----------------|--------------------------------------|--------------------------------------|
|  |                | True                                 | False                                |
| Decision About Null Hypothesis ( $H_0$ ) | Reject         | Type I error<br>(False Positive)     | Correct inference<br>(True Positive) |
|  | Fail to reject | Correct inference<br>(True Negative) | Type II error<br>(False Negative)    |

# Controlling errors

|  |   | predicted condition   |  |   |   |
|--|---|---|--|---|---|
| total population   |   | prediction positive   | prediction negative  | Prevalence = $\frac{\sum \text{condition positive}}{\sum \text{total population}}$  |   |
| true condition   | condition positive  | True Positive (TP)  | False Negative (FN)<br>(type II error)                               | True Positive Rate (TPR), Sensitivity, Recall,<br>Probability of Detection<br>$= \frac{\sum \text{TP}}{\sum \text{condition positive}}$ | False Negative Rate (FNR),<br>Miss Rate<br>$= \frac{\sum \text{FN}}{\sum \text{condition positive}}$        |
|  | condition negative  | False Positive (FP)<br>(Type I error)   | True Negative (TN)   | False Positive Rate (FPR), Fall-out,<br>Probability of False Alarm<br>$= \frac{\sum \text{FP}}{\sum \text{condition negative}}$         | True Negative Rate (TNR),<br>Specificity (SPC)<br>$= \frac{\sum \text{TN}}{\sum \text{condition negative}}$ |
| Accuracy<br>$= \frac{\sum \text{TP} + \sum \text{TN}}{\sum \text{total population}}$ | Positive Predictive Value (PPV),<br>Precision<br>$= \frac{\sum \text{TP}}{\sum \text{prediction positive}}$ | False Omission Rate (FOR)<br>$= \frac{\sum \text{FN}}{\sum \text{prediction negative}}$       | Positive Likelihood Ratio (LR+)<br>$= \frac{\text{TPR}}{\text{FPR}}$ | Diagnostic Odds Ratio (DOR)<br>$= \frac{\text{LR}^+}{\text{LR}^-}$  |   |
|  | False Discovery Rate (FDR)<br>$= \frac{\sum \text{FP}}{\sum \text{prediction positive}}$                    | Negative Predictive Value (NPV)<br>$= \frac{\sum \text{TN}}{\sum \text{prediction negative}}$ | Negative Likelihood Ratio (LR-)<br>$= \frac{\text{FNR}}{\text{TNR}}$ |   |   |

⇒ *Precision-recall (AUC)*  
 ⇒ *Sensitivity-specificity*  
 ⇒ *Log odds ratio*

# Effect size, error rates, power



|               |          | Reality   |   |
|---------------|----------|---|---|
|               |          | Positive  | Negative  |
| Study Finding | Positive | True Positive<br>(Power)<br>( $1-\beta$ )             | False Positive<br><b>Type I Error</b><br>( $\alpha$ ) |
|               | Negative | False Negative<br><b>Type II Error</b><br>( $\beta$ ) | True Negative   |

# GSEA: Gene Set Enrichment Analysis

- **Step 1: Compute enrichment score (ES)**
  - The degree to which a set  $S$  is overrepresented at the extremes (top or bottom) of the entire ranked list  $L$
  - Procedure: Walk down  $L$ , increase (decrease) running-sum statistic if gene is (not) in  $S$ 
    - The magnitude of the increment depends on the correlation of the gene with the phenotype
  - ES = maximum deviation from 0
- **Step 2: Estimate significance of ES (P value)**
  - Permutation test: generate null distribution by permuting class (phenotype) labels
    - Permuting class labels preserves gene-gene correlation structure
  - Empirical P-value is observed ES relative to null distribution
- **Step 3: Calculate false discovery rate (FDR)**
  - Normalize ES for each gene set to account for size of set (NES)
  - Compare tails of observed and null distributions

# FWER vs FDR

**Null hypothesis testing methods adjusted for multiple comparisons (multiple hypothesis testing)**

- **Family-wise error rate (e.g. Bonferroni correction):**
  - Control the probability of *at least one* Type I error
  - Most conservative
- **False Discovery rate (Benjamini Hochberg):**
  - Control the expected proportion of rejected null hypotheses that are false (incorrect rejections)
  - Less conservative, more Type I errors
  - Greater power (probability that the test correctly rejects the null hypothesis when the alternative hypothesis is true)

# Pathway Topology (PT)

- Approach
  - Use annotation databases that catalog **interactions**
    - Activation, inhibition, physical, metabolic
    - Cellular compartment
  - Same steps as FCS (functional class scoring), BUT use **pathway topology** to compute gene-level statistics
    - Similarity between genes based on **proximity** in network (metabolic networks)
    - E.g. correlation, covariance weighted by graph distance
  - IF: impact factor (signaling pathways)
    - Model signaling pathways as graphs (nodes, edges)
    - Gene-level statistic is PF: perturbation factor
      - Expression AND linear function of PF for all genes in pathway
    - IF is sum of all PFs for pathway

# Pathway Topology (PT)

- Limitations
  - Networks are not necessarily the same in each cell
  - Condition-specific changes not known or annotated
  - Do not model dynamic states
  - Do not model interdependence between pathways

⇒ *All methods will improve as resolution, accuracy, and completeness of annotations increase*



# PathCORE-T: identifying and visualizing globally co-occurring pathways in large transcriptomic compendia

Kathleen M. Chen<sup>1</sup> , Jie Tan<sup>2</sup> , Gregory P. Way<sup>1</sup> , Georgia Doing<sup>3</sup> , Deborah A. Hogan<sup>3</sup> and Casey S. Greene<sup>1\*</sup>

\* Correspondence: [csgreene@upenn.edu](mailto:csgreene@upenn.edu)

<sup>1</sup>Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine, University of Pennsylvania, 3400 Civic Center Blvd., Philadelphia, PA 19104, USA  
Full list of author information is available at the end of the article

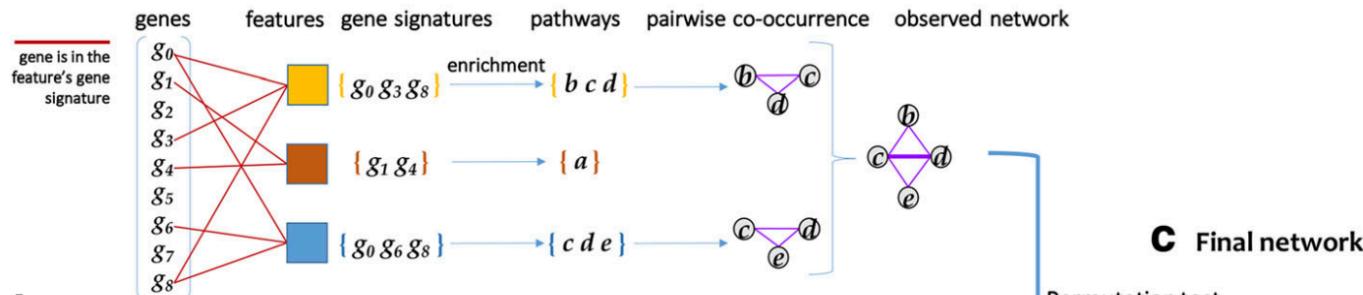
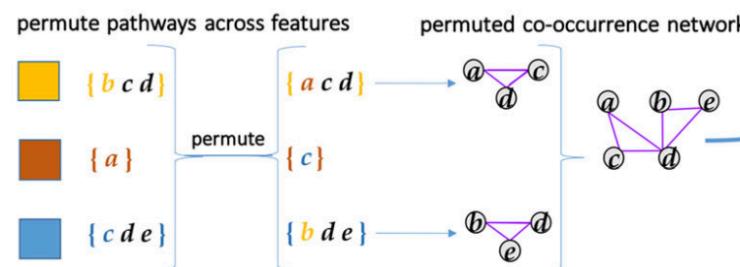
## Abstract

**Background:** Investigators often interpret genome-wide data by analyzing the expression levels of genes within pathways. While this within-pathway analysis is routine, the products of any one pathway can affect the activity of other pathways. Past efforts to identify relationships between biological processes have evaluated overlap in knowledge bases or evaluated changes that occur after specific treatments. Individual experiments can highlight condition-specific pathway-pathway relationships; however, constructing a complete network of such relationships across many conditions requires analyzing results from many studies.

**Results:** We developed PathCORE-T framework by implementing existing methods to identify pathway-pathway transcriptional relationships evident across a broad data compendium. PathCORE-T is applied to the output of feature construction algorithms; it identifies pairs of pathways observed in features more than expected by chance as *functionally co-occurring*. We demonstrate PathCORE-T by analyzing an existing eADAGE model of a microbial compendium and building and analyzing NMF features from the TCGA dataset of 33 cancer types. The PathCORE-T framework includes a demonstration web interface, with source code, that users can launch to (1) visualize the network and (2) review the expression levels of associated genes in the original data. PathCORE-T creates and displays the network of globally co-occurring pathways based on features observed in a machine learning analysis of gene expression data.

**Conclusions:** The PathCORE-T framework identifies transcriptionally co-occurring pathways from the results of unsupervised analysis of gene expression data and visualizes the relationships between pathways as a network. PathCORE-T recapitulated previously described pathway-pathway relationships and suggested experimentally testable additional hypotheses that remain to be explored.

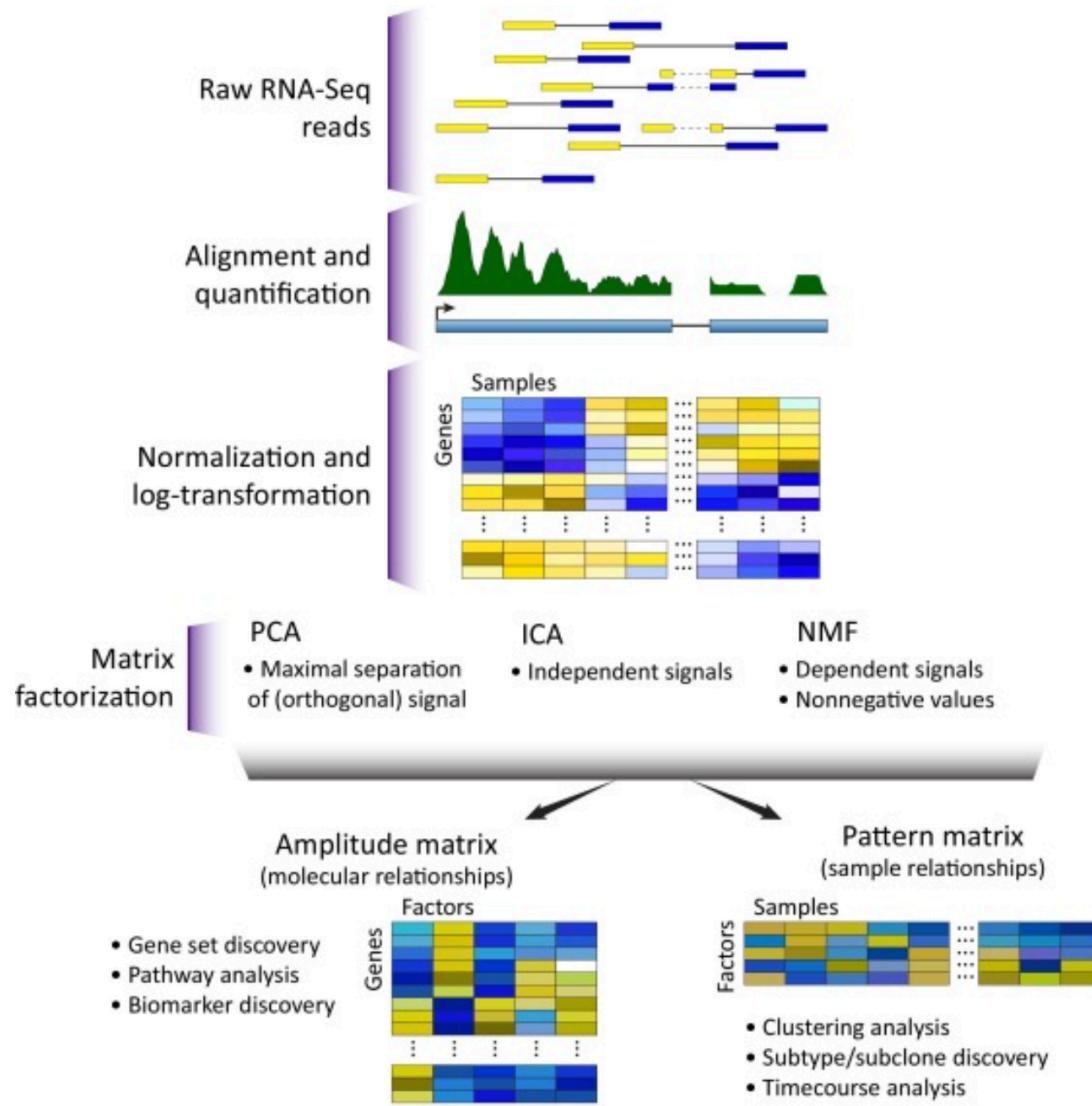
**Keywords:** Gene expression, Unsupervised feature construction, Crosstalk, Pathway interactions

**a** Constructing the observed co-occurrence network**b** Constructing one permuted co-occurrence network**C** Final network

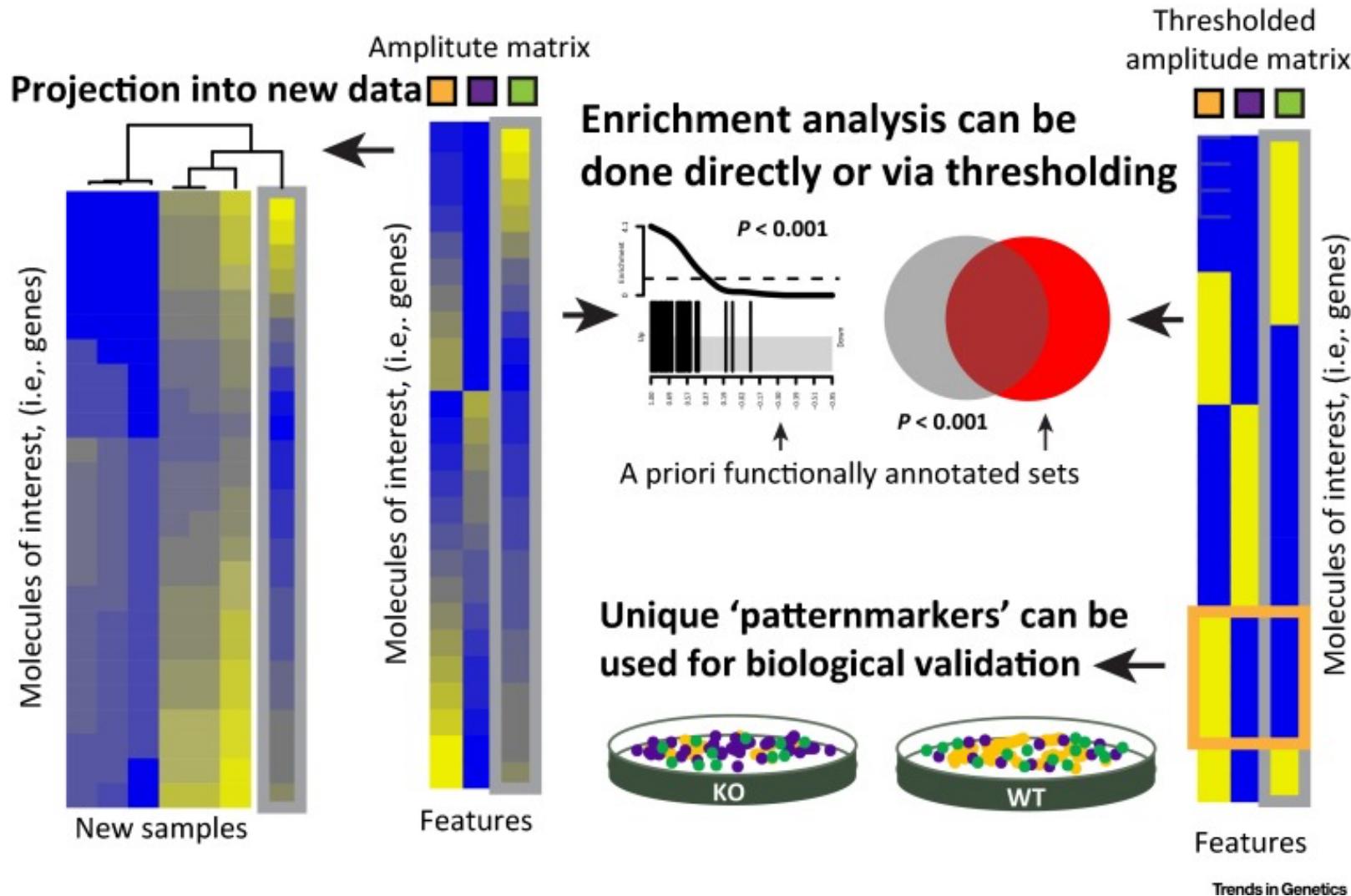
Permutation test  
evaluate the significance of each observed edge

**Fig. 2** The approach implemented in PathCORE-T to construct a pathway co-occurrence network from an expression compendium. **a** A user-selected feature extraction method is applied to expression data. Such methods assign each gene a weight, according to some distribution, that represents the gene's contribution to the feature. The set of genes that are considered highly representative of a feature's function is referred to as a feature's gene signature. The gene signature is user-defined and should be based on the weight distribution produced by the unsupervised method of choice. In the event that the weight distribution contains both positive and negative values, a user can specify criteria for both a positive and negative gene signature. A test of pathway enrichment is applied to identify corresponding sets of pathways from the gene signature(s) in a feature. We consider pathways significantly overrepresented in the same feature to co-occur. Pairwise co-occurrence relationships are used to build a network, where each edge is weighted by the number of features containing both pathways. **b**  $N$  permuted networks are generated to assess the statistical significance of a co-occurrence relation in the graph. Here, we show the construction of one such permuted network. Two invariants are maintained during a permutation: (1) pathway side-specificity (if applicable, e.g. positive and negative gene signatures) and (2) the number of distinct pathways in a feature's gene signature. **c** For each edge observed in the co-occurrence network, we compare its weight against the weight distribution generated from  $N$  (default: 10,000) permutations of the network to determine each edge's  $p$ -value. After correcting the  $p$ -value by the number of edges observed in the graph using the Benjamini—Hochberg procedure, only an edge with an adjusted  $p$ -value below alpha (default: 0.05) is kept in the final co-occurrence network.

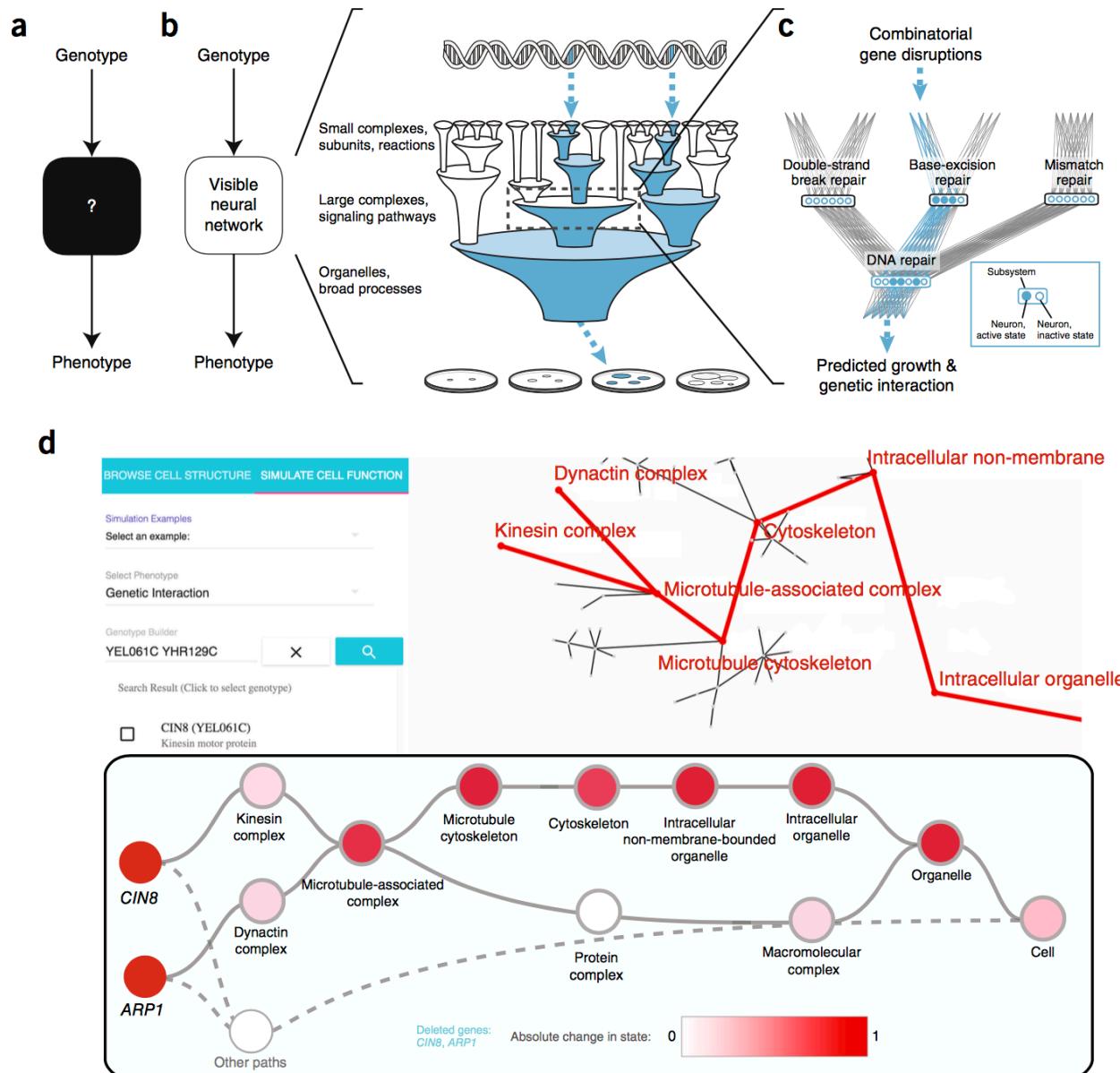
# Enrichment analysis: feature extraction by matrix factorization



# Enrichment analysis: feature extraction by matrix factorization



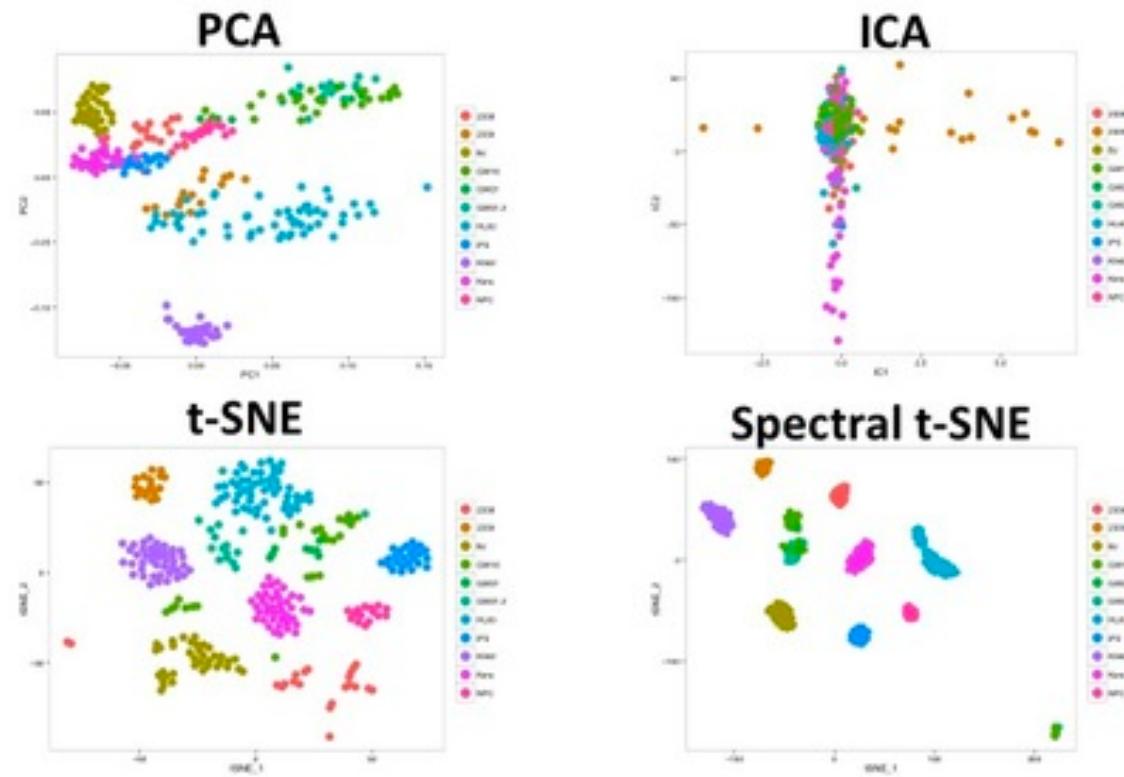
# Dcell: Visual neural network for prediction of phenotypes and GIs



A really nice introduction to neural  
networks in 4 short videos

# Enrichment Analysis: Seurat

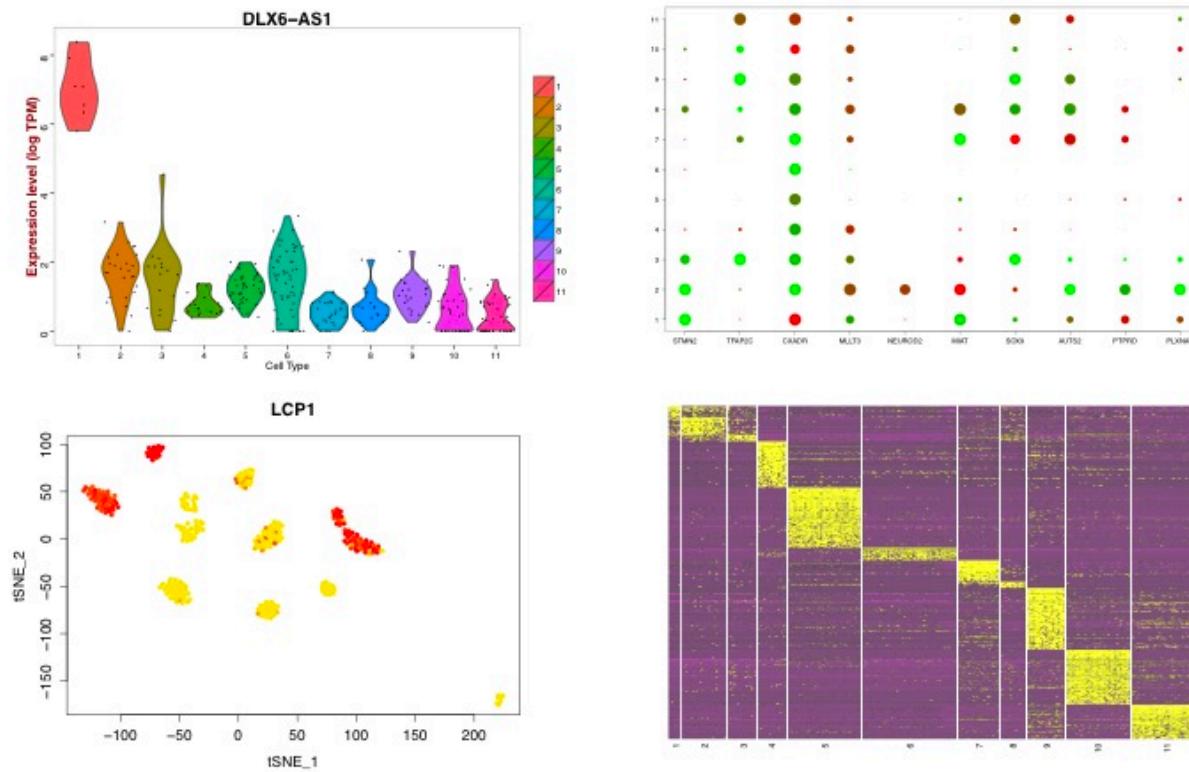
## Part 1 – Cluster single cells



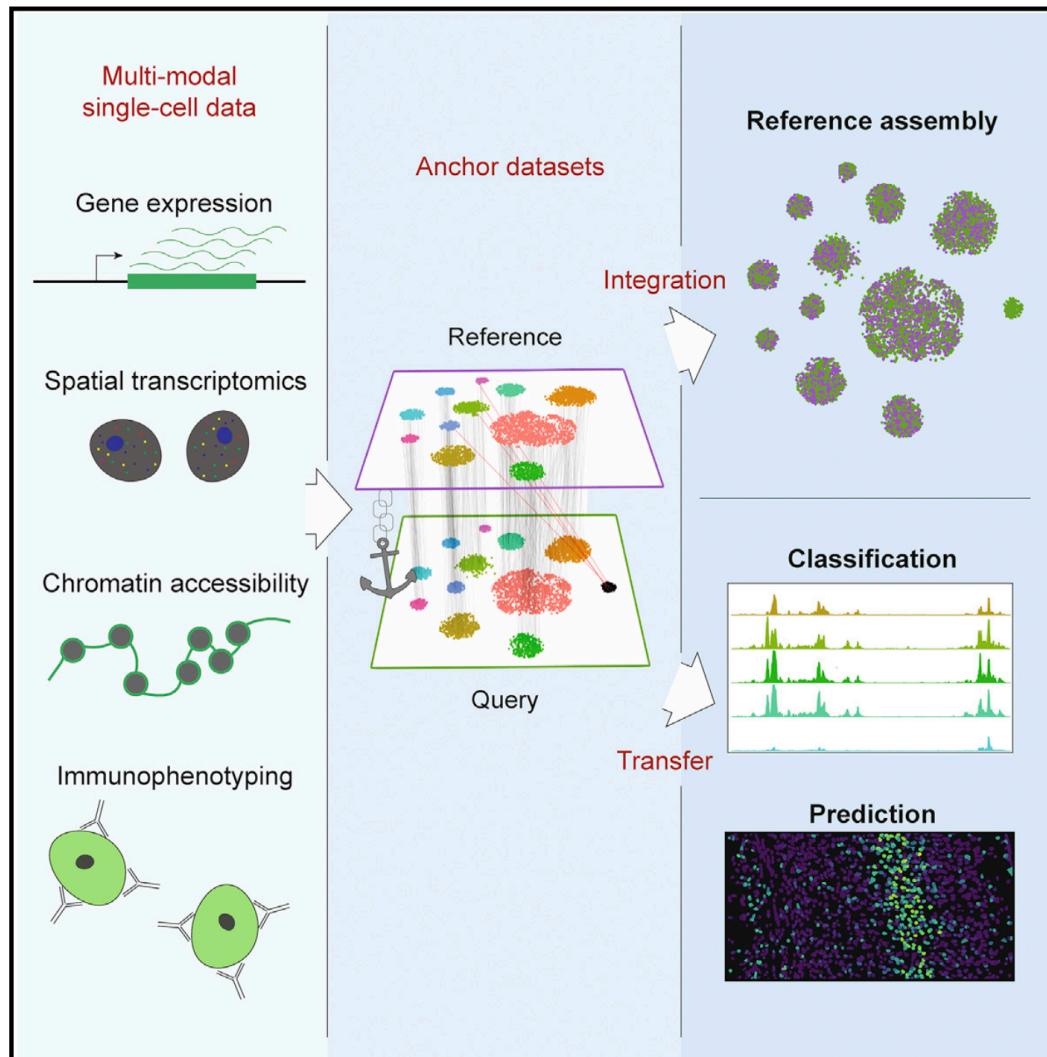
⇒Now use UMAP

# Enrichment Analysis: Seurat

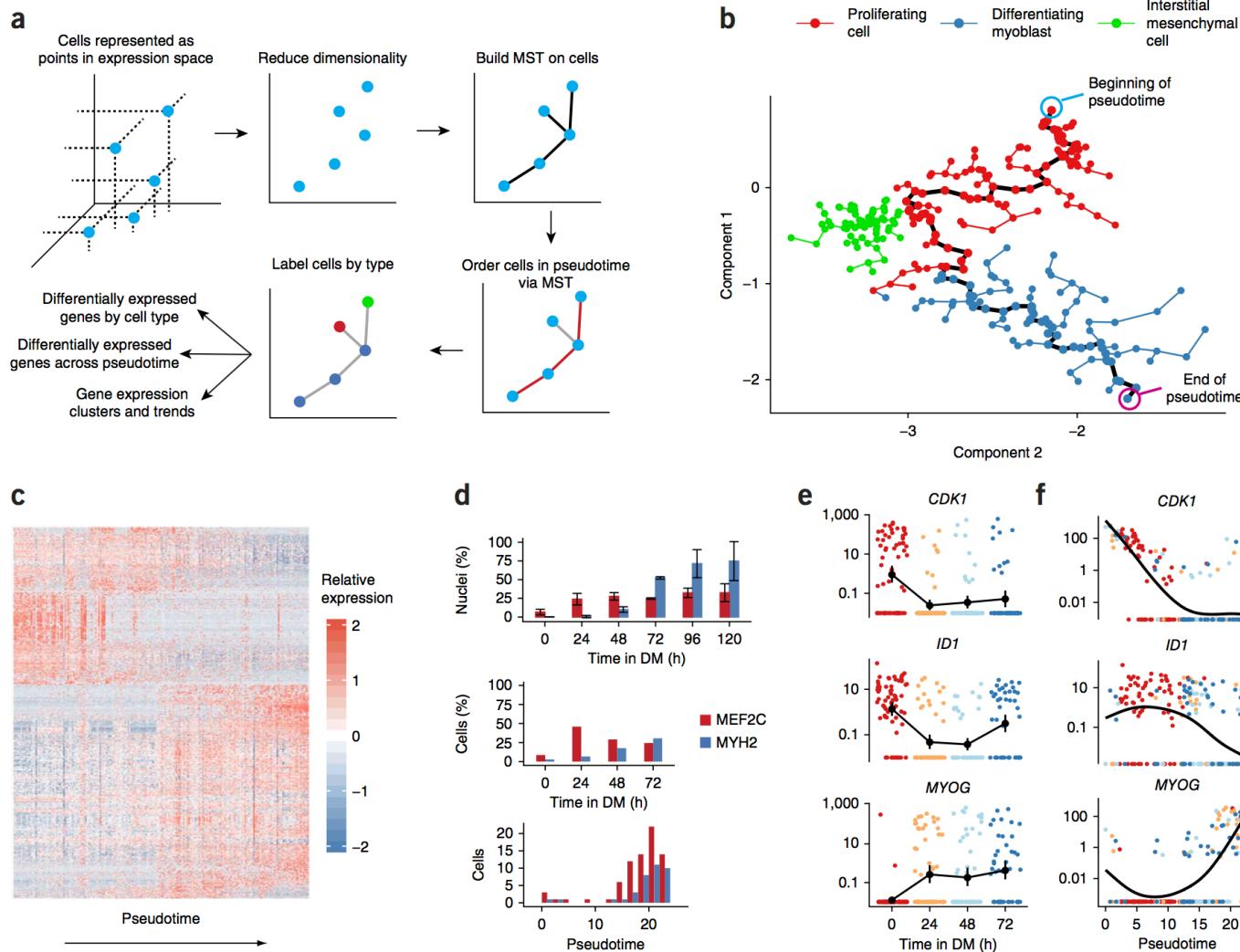
## Part 2 – Discover and visualize markers



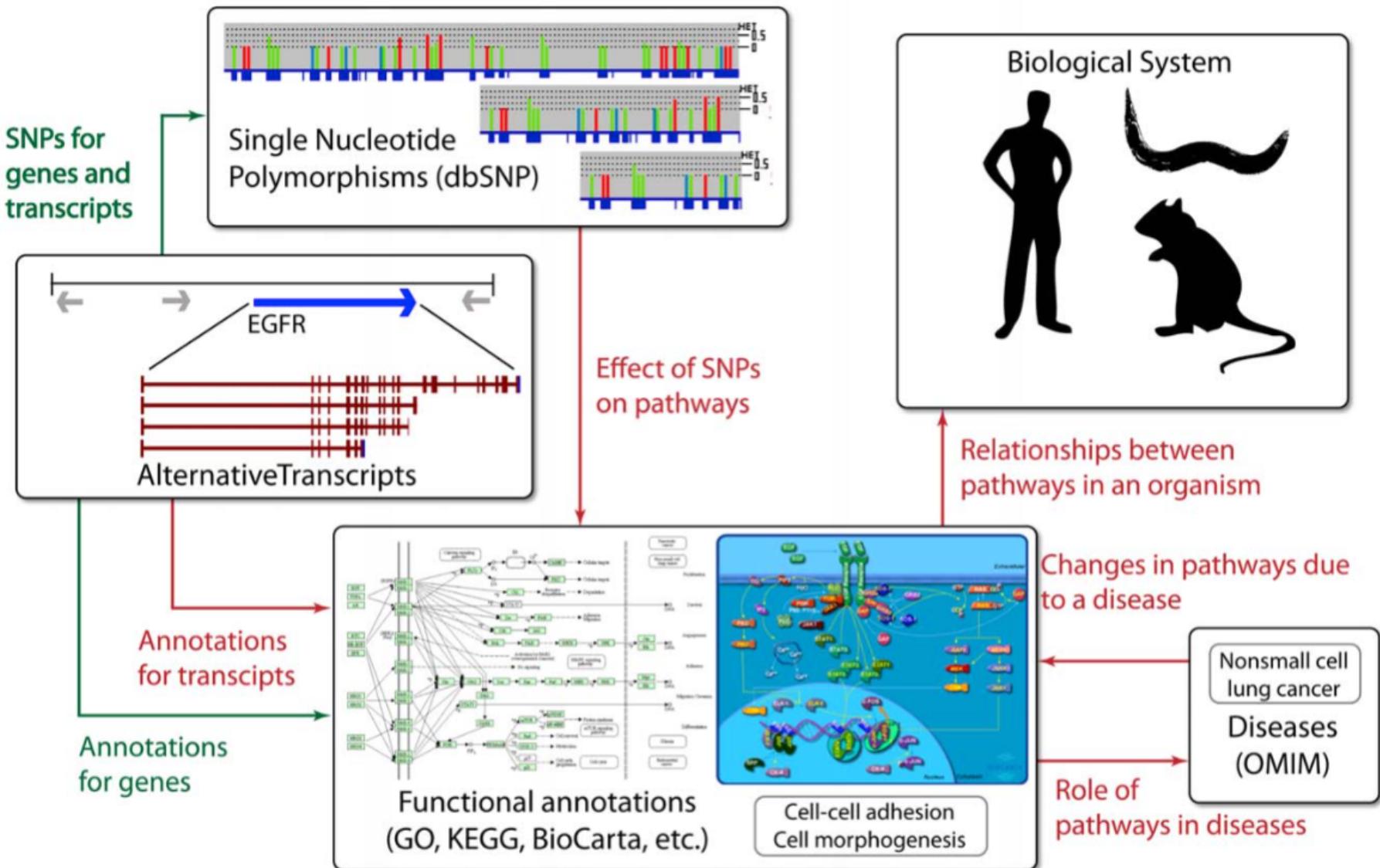
# Enrichment Analysis: Seurat

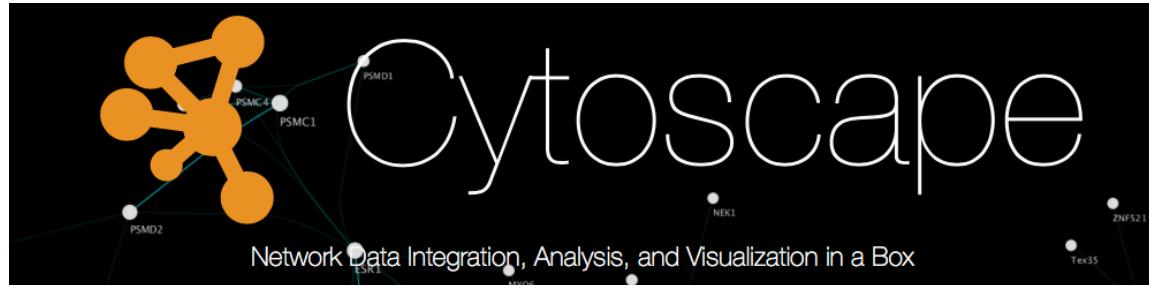


# Enrichment Analysis: Pseudotime



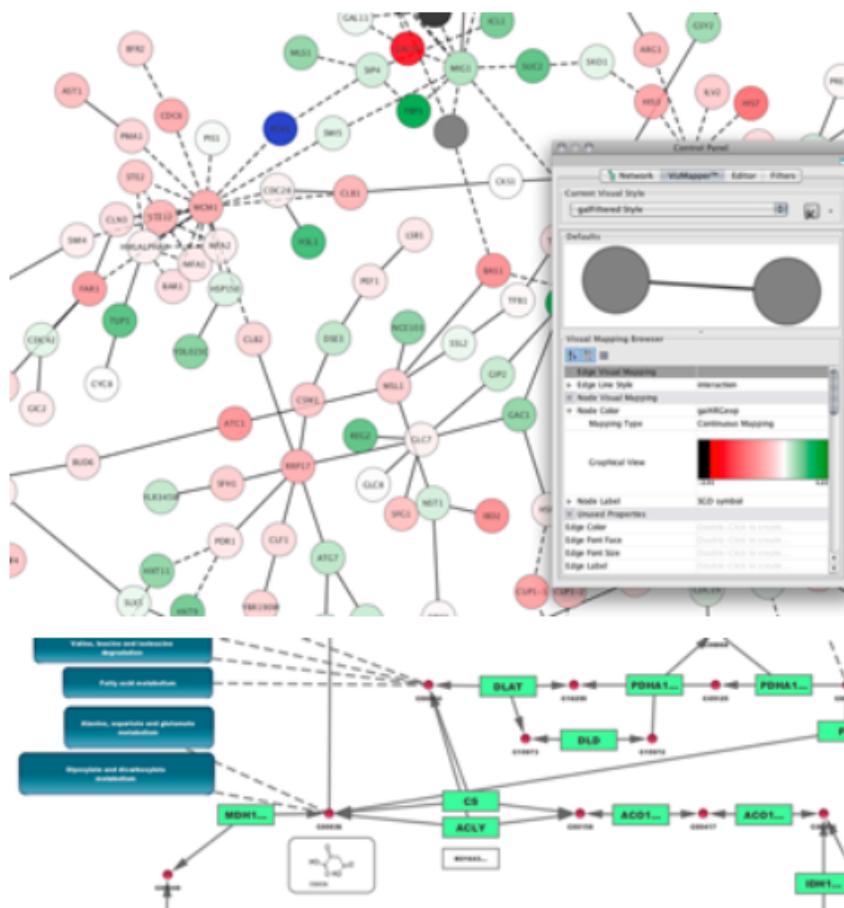
# Outstanding challenges





Cytoscape supports many use cases in molecular and systems biology, genomics, and proteomics:

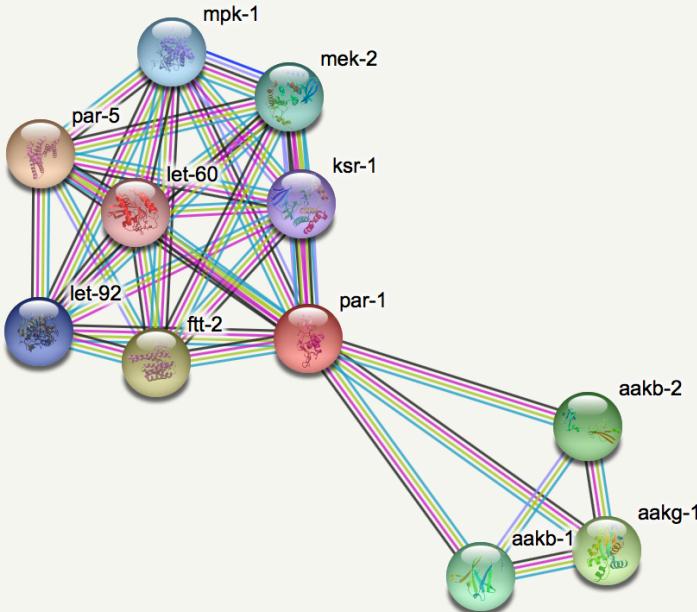
- Load molecular and genetic interaction data sets in many standards formats
  - Project and integrate global datasets and functional annotations
  - Establish powerful visual mappings across these data
  - Perform advanced analysis and modeling using [Cytoscape Apps](#)
  - Visualize and analyze human-curated pathway datasets such as [WikiPathways](#), [Reactome](#), and [KEGG](#).



# STRING database

Version: 10.5      [LOGIN](#) | [REGISTER](#)

 **STRING**      [Search](#)      [Download](#)      [Help](#)      [My Data](#)



The diagram illustrates a complex network of protein interactions. Nodes represent proteins, and edges represent interactions. Nodes are color-coded by cluster: purple (top), brown (left), red (center), green (bottom right), and blue (bottom left). Node labels include: mpk-1, mek-2, ksr-1, par-1, let-60, let-92, ftt-2, par-5, aakb-1, aakb-2, and aakg-1. Edges are colored and weighted, showing the strength and type of interactions between these proteins.

[Viewers](#) [Legend](#) [Settings](#) [Analysis](#) [Exports](#) [Clusters](#) [More](#) [Less](#)

## Nodes:

Network nodes represent proteins

*splice isoforms or post-translational modifications are collapsed, i.e. each node represents all the proteins produced by a single, protein-coding gene locus.*

### Node Color



colored nodes:  
query proteins and first shell of interactors



white nodes:  
second shell of interactors

### Node Content



empty nodes:  
proteins of unknown 3D structure



filled nodes:  
some 3D structure is known or predicted

## Edges:

Edges represent protein-protein associations

*associations are meant to be specific and meaningful, i.e. proteins jointly contribute to a shared function; this does not necessarily mean they are physically binding each other.*

### Known Interactions



from curated databases



experimentally determined

### Predicted Interactions



gene neighborhood



gene fusions



gene co-occurrence

### Others



textmining



co-expression



protein homology

## Your Input:

par-1 Serine/threonine-protein kinase par-1 ; Required for cytoplasmic partitioning and asymmetric cell division in early embryogenesis (PubMed-7758115). Phosphorylates and restricts the asymmetry effector mex-5 (and possibly also mex-6) to the anterior cytoplasm of the zygote (PubMed-18842813). Regulates mes-1 expression during early embryogenesis (PubMed-11003841). Critical role in postembryonic vulval morphogenesis (PubMed-12490197). Involved in the establishment of neuronal polarity (PubMed-20023164) (1216 aa)

| Neighborhood | Gene Fusion | Cocurrence | Coexpressions | Experiments | Databases | Textmining | [Homology] | Score |
|--------------|-------------|------------|---------------|-------------|-----------|------------|------------|-------|
|--------------|-------------|------------|---------------|-------------|-----------|------------|------------|-------|

## Predicted Functional Partners:

- |  |       |       |
|--|-------|-------|
| par-5 14-3-3-like protein 1 (248 aa)   | ● ● ● | 0.981 |
| ftt-2 14-3-3-like protein 2 ; Required for extension of life-span by sir-2.1 (PubMed-16777605). Promotes nuclear export of yap-1 ...   | ● ● ● | 0.973 |
| aakg-1 Protein AAKG-1 (582 aa)   | ● ● ● | 0.967 |
| aakb-2 AMP-Activated Kinase Beta subunit family member (aakb-2) (274 aa)   | ● ● ● | 0.955 |
| aakb-1 AMP-Activated Kinase Beta subunit family member (aakb-1) (269 aa)   | ● ● ● | 0.951 |
| mek-2 Dual specificity mitogen-activated protein kinase kinase mek-2 ; Functions in the let-60 Ras signaling pathway; acts down... .   | ● ● ● | 0.931 |
| mpk-1 Mitogen-activated protein kinase mpk-1 ; Function in let-60 Ras signaling pathway; acts downstream of lin-45 raf kinase, b...    | ● ● ● | 0.924 |
| let-92 LETHAL family member (let-92) (318 aa)  | ● ● ● | 0.922 |
| ksr-1 Kinase Suppressor of activated Ras family member (ksr-1) ; Serine/threonine-protein kinase which positively regulates Ra...      | ● ● ● | 0.919 |
| let-60 Ras protein let-60 ; The level of let-60 controls the switch between vulval and hypodermal cell fates during C.elegans vulva... | ● ● ● | 0.915 |

## Your Current Organism:

Caenorhabditis elegans

NCBI taxonomy Id: [6239](#)

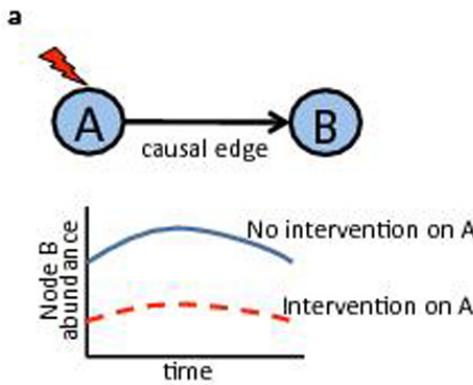
Other names: *C. elegans*, *Caenorhabditis elegans*, *Rhabditis elegans*, *nematode*

# Gene regulatory network (GRN) inference

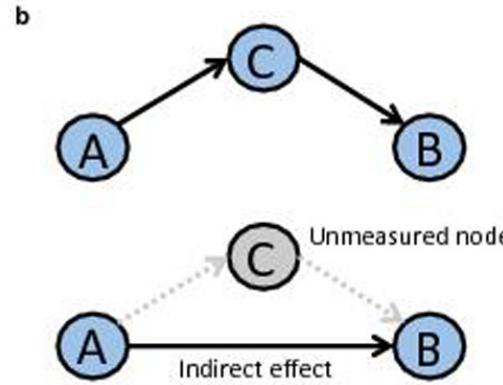
- Integrate many types of molecular and perturbation studies to generate predictive models
- Gene expression, physical interactions, chromatin accessibility, histone marks, TF-gene interactions...
- Combine network priors and measured features from many large-scale datasets
- Mathematical models (large linear systems)

# DREAM challenge: inferring causal networks

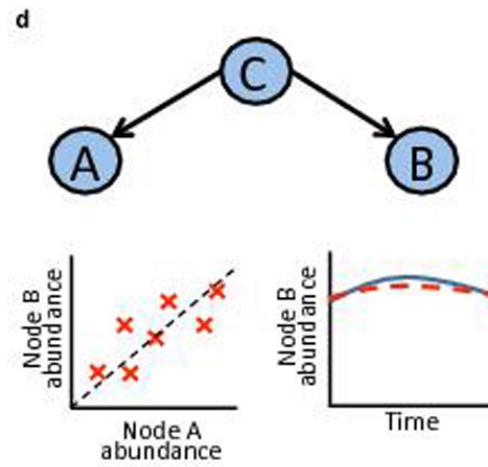
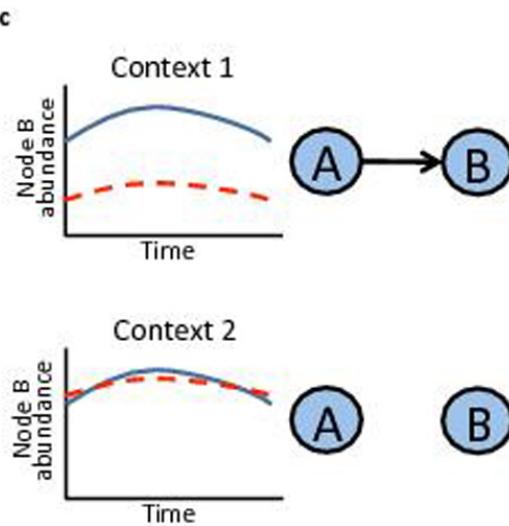
Direct, causal  
(A influences B)



Indirect, causal  
(A influences B)

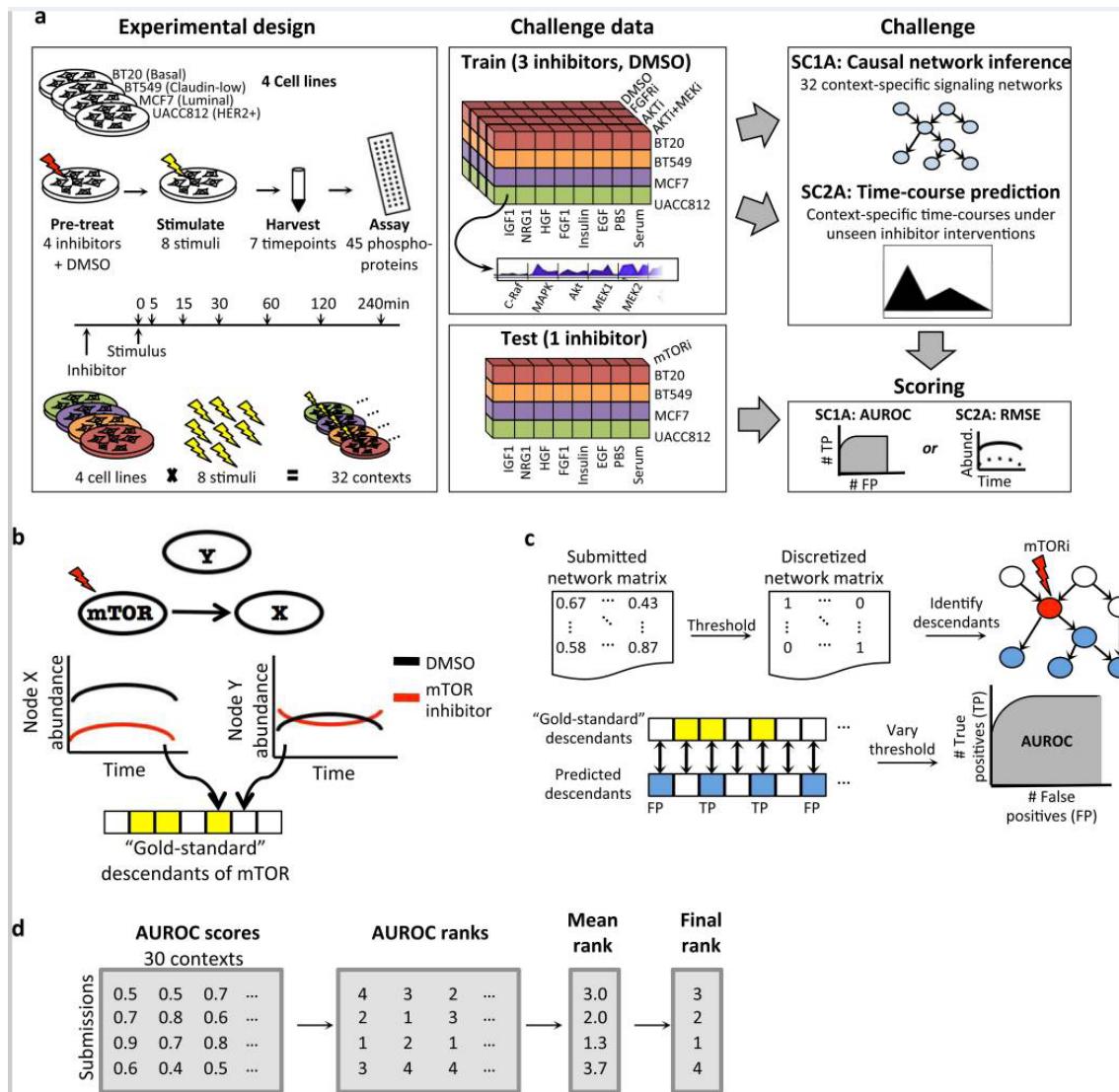


Context-dependent interaction



Indirect,  
correlated but  
not causal  
(A does NOT  
influence B)

# DREAM challenge: inferring causal networks



## Gene regulatory network inference software tools |...

[+ Submit tool](#)

Elucidating gene regulatory network (GRN) from large scale experimental data remains a central challenge in systems biology. The advent of high-throughput data generation technologies has allowed researchers to fit theoretical models to experimental data on gene-expression profiles.

**CMGRN / Constructing Multilevel Gene Regulatory Networks**

(2) 0 discussions

An integrative web server to unravel hierarchical interactive networks at different regulatory levels. The developed method used the Bayesian network modeling to infer causal interrelationships among...

**MIPRIP / Mixed Integer linear Programming based Regulatory Interaction Pre**

(0) 0 discussions

Predicts regulators of a gene of interest from gene expression profiles of the samples under study and known regulator binding information (from e.g. ChIP-seq/ChIP-chip databases). MIPRIP is...

**AMIGO2**

(0) 0 discussions

A MATLAB environment to solve mathematical optimization problems which in dynamic modeling and control of biological systems. AMIGO2 is organized in four main modules: the pre-processor, the...

**SpidermiR**

(0) 0 discussions

Offers an easy access to both Gene Regulatory Networks (GRNs) and miRNAs to the end user. SpidermiR integrates co-expression, physical interaction, co-localization, genetic influence, pathways, and...

---

RELATED WEBSITES

---



Gene regulatory networks  
Wikipedia

---

INFORMATION

---



Share the latest tools and make your expertise visible

[CONTRIBUTE](#)