# BIO214 Lecture 9

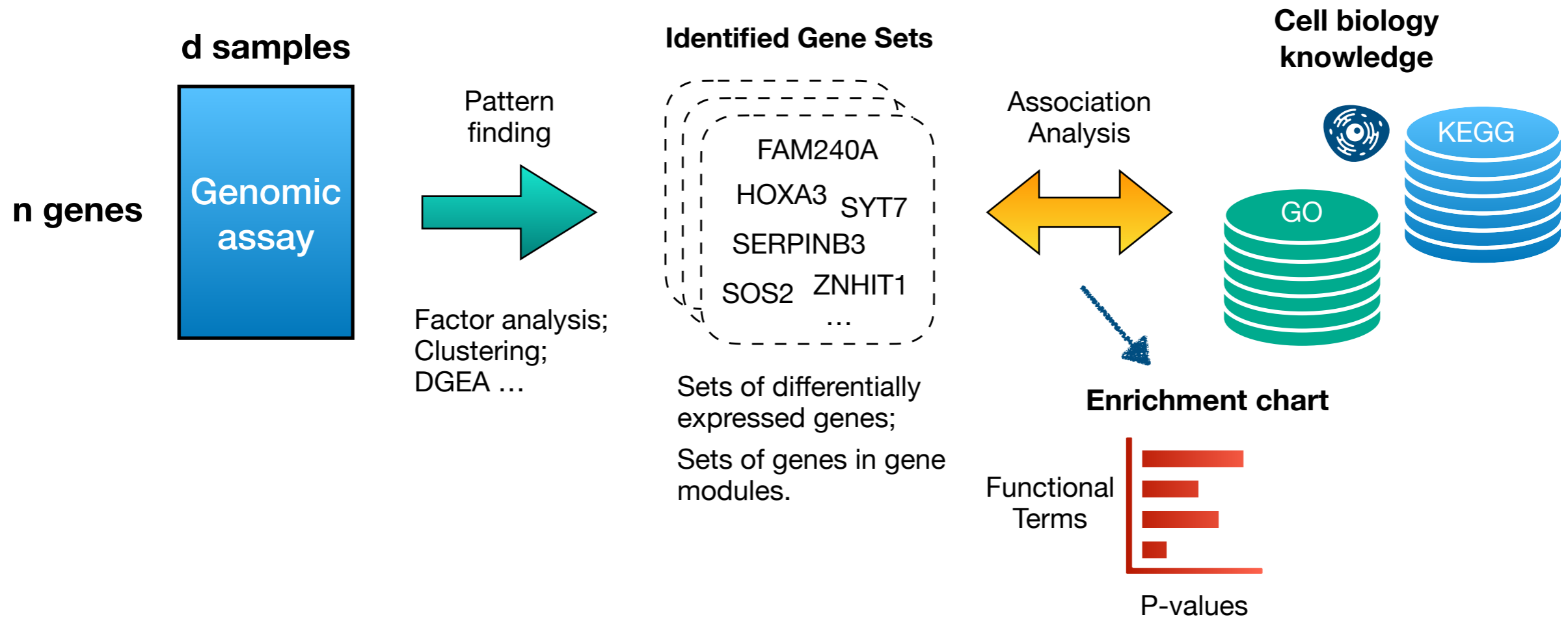## Bioinformatics-II

*Genomic knowledge representation*

**Zhen Wei; 2023-Feb-14**

# Outline

- Gene set annotation

- Range based annotation

- Network basics

- Co-expression network

# Gene set annotation

# Why we need functional annotations?



- Annotations are stored knowledge from previous biological experiments.

- Functional annotations are essential for the interpretation of gene sets obtained from the upstream analysis.

- Gene set enrichment is calculated via the statistical association between gene functions and gene sets.

# **Gene Ontology**: use general biological knowledge to annotate genes

"An ontology is a formal representation of a body of knowledge within a given domain."
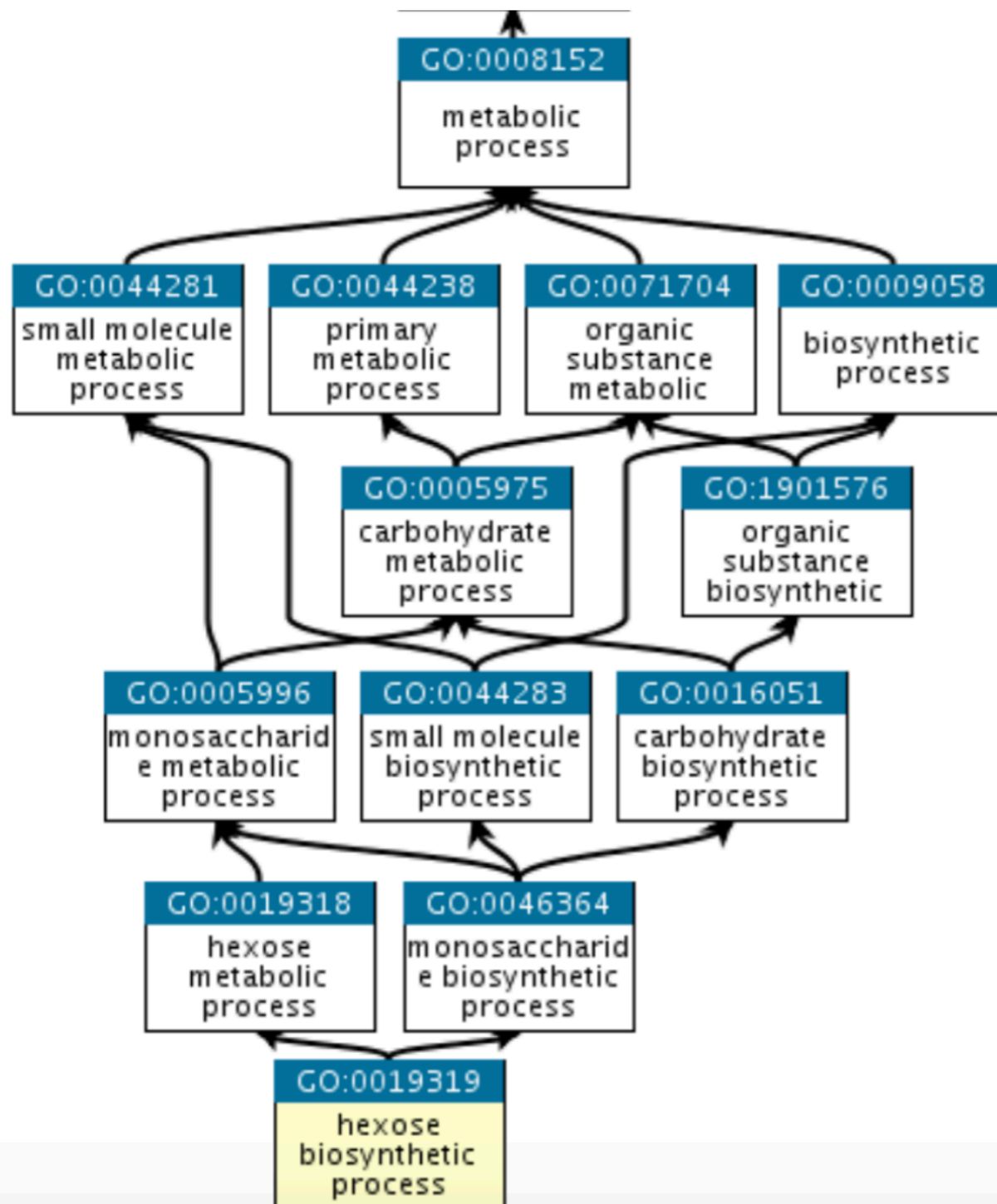
—- From Gene Ontology website

The Gene Ontology (GO) describes our knowledge of the biological domain with respect to three aspects:

- **Molecular Function**: Molecular-level activities performed by gene products.

- **Cellular Component**: The locations relative to cellular structures in which a gene product performs a function.

- **Biological Process**: The larger processes, or 'biological programs' accomplished by multiple molecular activities.

For example, the gene product "cytochrome c" can be described by the **molecular function** *oxidoreductase activity*, the **biological process** *oxidative phosphorylation*, and the **cellular component** *mitochondrial matrix*.
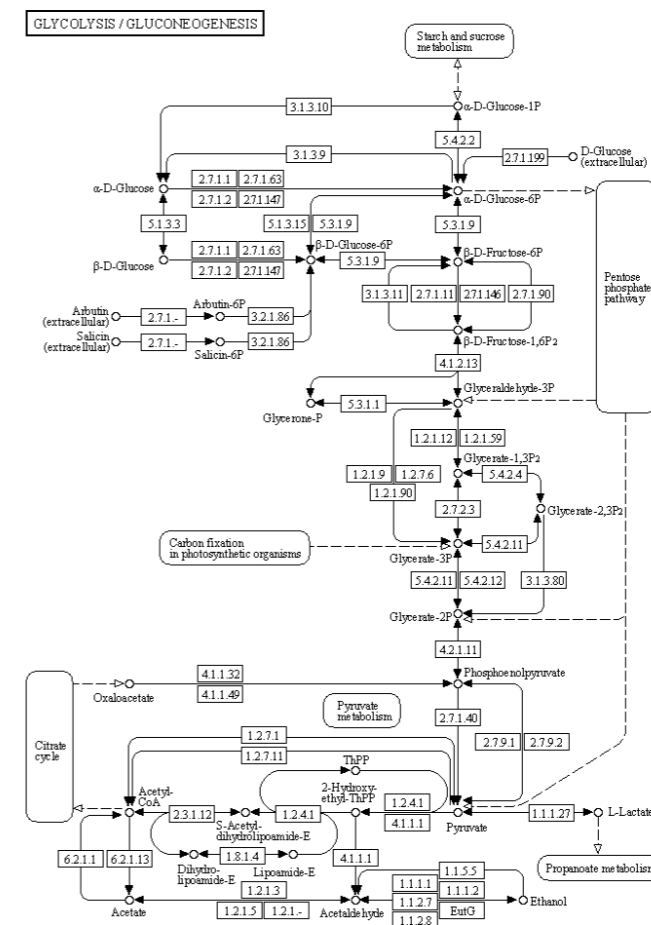
# GO Graph



- The Gene Ontology (GO) is represented as a graph with terms as nodes and relationships between terms as edges.

- GO is hierarchical, with more specific child terms and more general parent terms.

- Terms can have multiple parent terms.
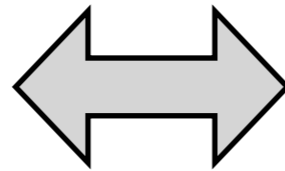
# KEGG: gene annotation via signaling pathway



- KEGG is a database for understanding biological systems.

- KEGG pathway maps are molecular interaction/reaction networks represented in terms of KEGG Orthology groups.

- These maps can help generalize experimental evidence from one organism to others based on genomic information.

The KEGG pathway map for glycolysis / gluconeogenesis

# Calculating statistical association between annotations

| | User Genes | Genome |
|---|---|---|
| In Pathway | 10 | 30 |
| Not Pathway | 4 | 60 |

⟷

Can be viewed as

| | Balls drawed | Balls in urn |
|---|---|---|
| Black balls | 10 | 30 |
| White balls | 4 | 60 |

Hypergeometric pmf

$$P(X) = \frac{\binom{k}{x}\binom{N-k}{n-x}}{\binom{N}{n}}$$

- **Fisher's exact test** is often used to calculate p-value of association between gene sets and functional terms.

- The p-value is calculated by the hypergeometric distribution:

  1. Enumerate all possible 2 by 2 tables that are as or more associated than the observed given fixed margins (column and row sums).

  2. Use hypergeometric distribution to calculate the probabilities of each table, sum them up and you will get the p-value.

# Range based annotations

# Range based annotations

## Transcript annotation

**Information stored in transcript annotations (GTF/GFF files)**



**The genome coordinate of a locus:**
**(Not necessarily a gene)**

Chr9; [1653482−1654482]; + strand

**Genomic Feature Extraction**

**Overlapping with 5'UTR, CDS, 3'UTR, intron, exon or not?**

**Properties of the overlapped features, e.g. length, GC content, state of evolutionary conservation.**

- Gene & Transcript annotations from GTF/GFF files are often used to annotate range based genomic experiments (e.g. peaks from CHIP-Seq).

# Other range based annotations

**Epigenetic markers**



- **ENCODE** stands for ENCyclopedia Of DNA Elements.
- It's a database that collects high-quality data about epigenetic markers, expressed transcripts, and epitranscriptomic markers.
- ENCODE uses strict and well-documented data processing pipelines to ensure data quality.
- Researchers can use the epigenetic markers from ENCODE to annotate their own experiments.

# Introduction to biological networks

# Correlational v.s. causal gene networks

Undirected      Directed      Cyclic      Acyclic

| Types | Description | Example |
|---|---|---|
| Correlational graphs (undirected graph) | Represent the positive / negative correlation between genes. The significantly correlated genes are linked by an undirected edge. | PPI network, gene co-expression network |
| Cause-effect graph (directed graph) | Describe the relationship of causality between genes, such as a gene is changed upon the action of another gene. The direction of the arrowed edge represents cause and effect. | Cell signaling network, epigenetic regulatory network |

# Representation: adjacency matrix

Adjacency matrix: $\mathbf{A}$

A  B

C

D

$$\begin{array}{c} \\ A \\ B \\ C \\ D \end{array} \begin{array}{cccc} A & B & C & D \\ \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix} \end{array}$$

From (node i)

To (node j)

- **Nodes:** A, B, C, D
- **Edges**: A <-> B, B <-> C, C <-> D, D <-> A
- The whole network structure can be specified by an $n \times n$ **adjacency matrix**, where $n$ is the number of nodes.
- $\mathbf{A}_{ij} = 1$ if nodes $i$ and $j$ are connected from $i$ to $j$.
- Can be softly weighted by probabilities, i.e. $\mathbf{A}_{ij} \in [0,1]$

# Degree of node



Degree of node C = 3

- **Neighbors** are pairs of nodes connected by an edge.

- **Degree** ($k$) of a node counts the number of edges connecting its neighbors to it.

- Degrees for an undirected graph can be calculated by the row sums of the adjacency matrix.

# Random v.s. Scale-free Network

- The distribution of degrees over a graph reveals essential network properties.

- In **random network**, edges are added to node pairs with equal probabilities.

- The degree distribution for random network is Poisson distribution.

- In **scale-free network**, the probability of adding a new edge from node $i$ to a new node increases as the degree of node $i$ increases.

- The degree distribution for scale free network is power distribution.



A Random network

Aa

B Scale-free network

Ba

Ab

Most vertices have degree close to the 'average degree.'

P(k)

k

Bb

P(k)

Hubs = high degree vertices.

1    10    100    1,000

k

$P(k) \sim k^{-\lambda}; \quad \lambda < 3 \Rightarrow$

scale free

# Properties of scale-free network

- Average steps between a random pair of nodes in a graph of size $n$:

  - For a red[random network], the average path length is $\sim \log(n)$

  - For a red[scale-free network], the average path length is $\sim \log(\log(n))$

There by, information transfer is more efficient on a scale-free network.

- When "attacks" are made by removing nodes from the graph:

  - If the failures happened randomly, the scale-free network is more likely to survive than the random network.

  - If the failures are targeted toward the **hub nodes** (the nodes with highest degree), then the scale-free network is more vulnerable than the random network.

# Essential proteins are hub-nodes

- A protein is essential if its knock-down is lethal.
- In yeast PPI network, the proteins with higher degree (more direct interactions with other proteins) are more likely to be essential proteins.
- 2240 edges are formed among 1870 nodes (proteins) in yeast PPI network.
- 93% of proteins have degrees < 3, among them, 21% are essential to yeast survival.
- 0.7 % of proteins have > 15 degree, and 62% of those are essential.
- The overall correlation coefficient between lethality and connectivity is 0.76.

# Co-expression network

# How to construct gene network from gene expression levels?

# Workflow of co-expression network analysis



- Pairwise correlation used to construct network
- Clustering identifies modules
- Differential co-expression analysis identifies regulatory genes
- Guilt-by-association approach identifies potential disease genes

# Limitation of Pearson correlation



- Pearson correlation cannot capture non-linear interaction (last row).

# Performances of different network inference methods

Marbach, Daniel, et al. "Wisdom of crowds for robust gene network inference." *Nature methods* 9.8 (2012): 796-804.

# GINIE3: a high performing network inference algorithm



Top 1 performance in DREAM4 competition.

To create a gene regulatory network in GINIE3:

- For each gene, train Random Forest predictors ($f_j$) with its expression levels as output and other genes' levels as input.

- For each predictors, rank all input genes by feature importance.

- Combine the rankings of all predictors to get the edge scores for network's regulatory links.

Huynh-Thu, Vân Anh, et al. "Inferring regulatory networks from expression data using tree-based methods." PloS one 5.9 (2010): e12776.