# Pathway analysis: introduction and discussion

Francesca M. Buffa

Department of Oncology
University of Oxford

# Complex data structure
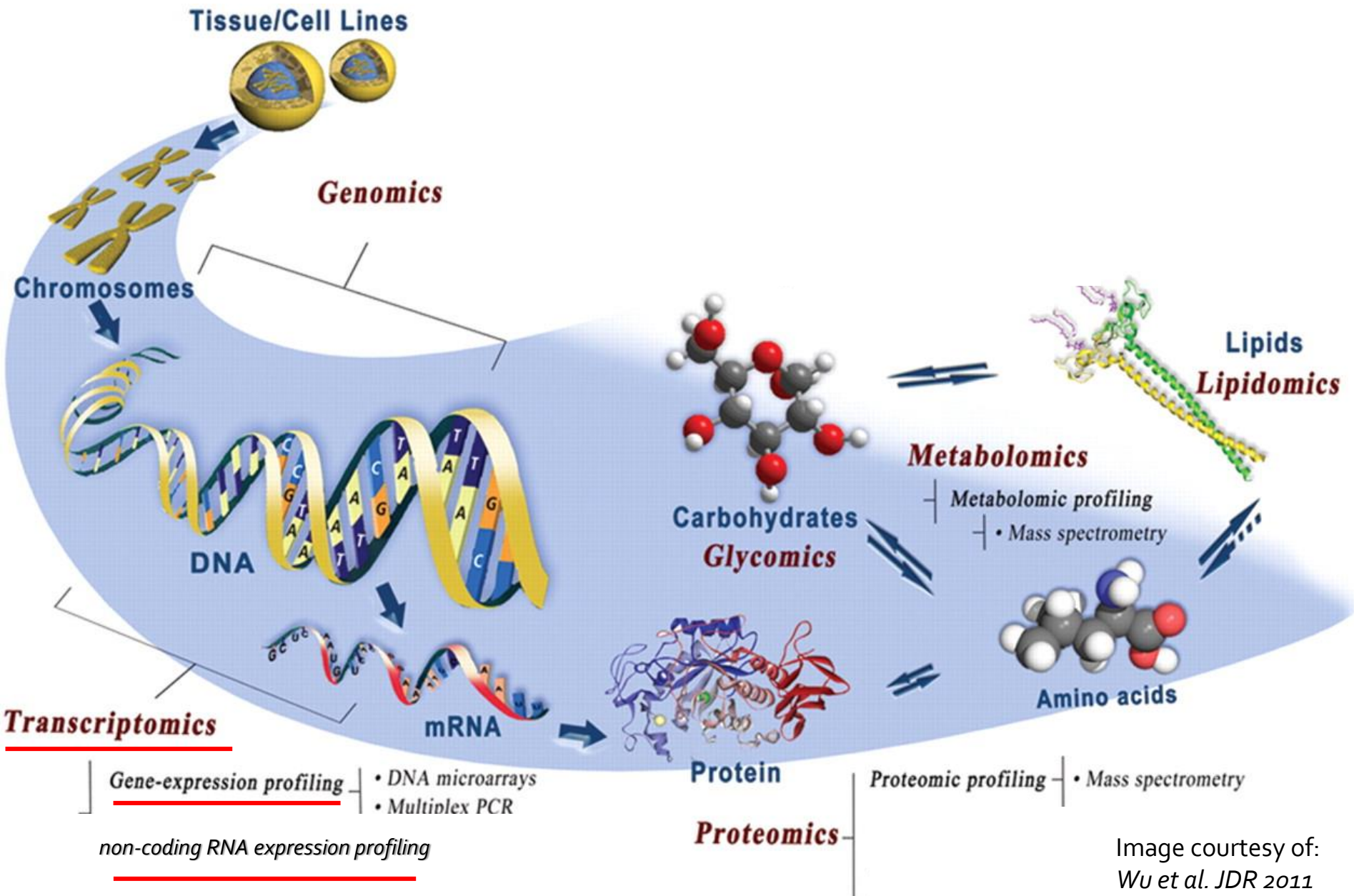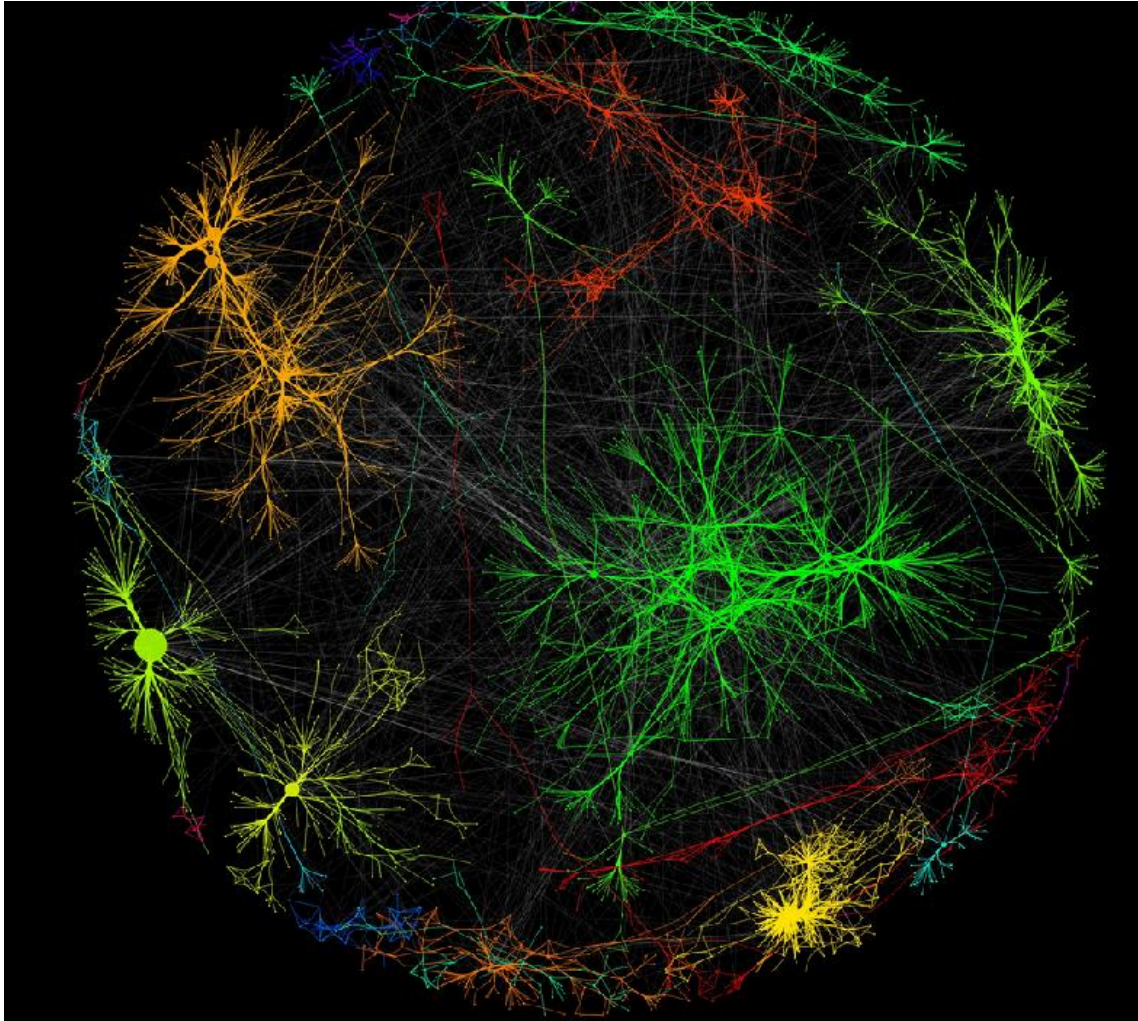


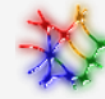Image courtesy of:
*Wu et al. JDR 2011*

# Cytoscape



### ClueGO
Creates and visualizes a functionally grouped network of

### CluePedia
CluePedia: A ClueGO plugin for pathway insights using integrated

### AgilentLiteratureSearch
Mines scientific literature to find publications related to search

### BiNGO
Calculates overrepresented GO terms in the network and display

### GeneMANIA
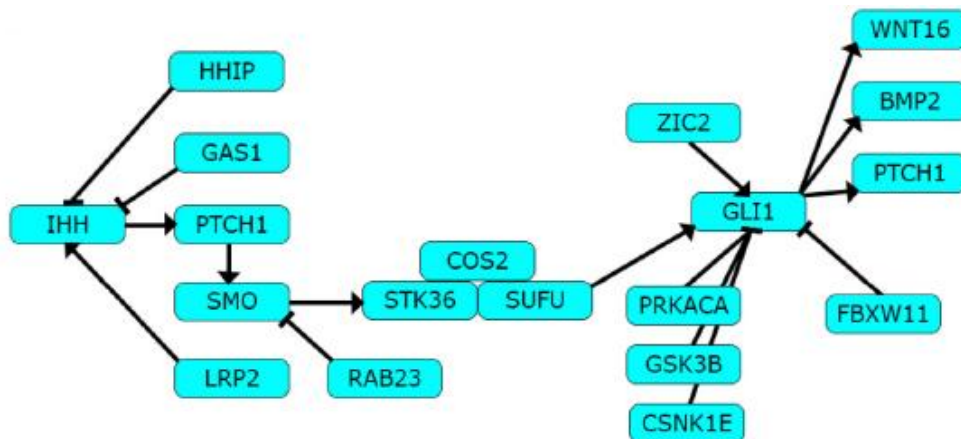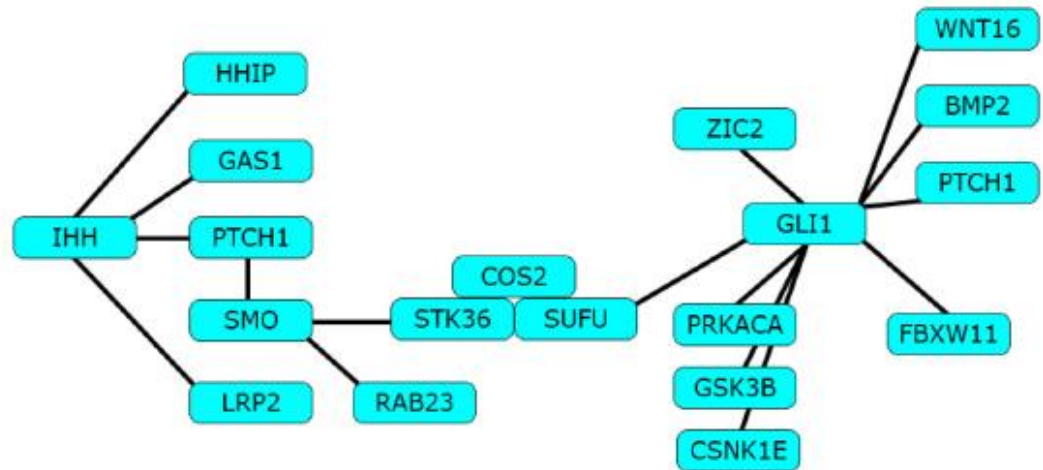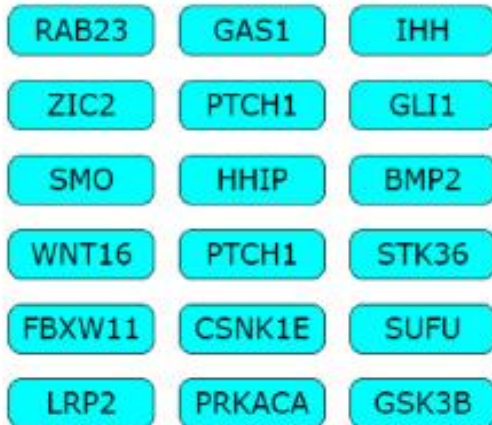Imports interaction networks from public databases from a list of

### clusterMaker2
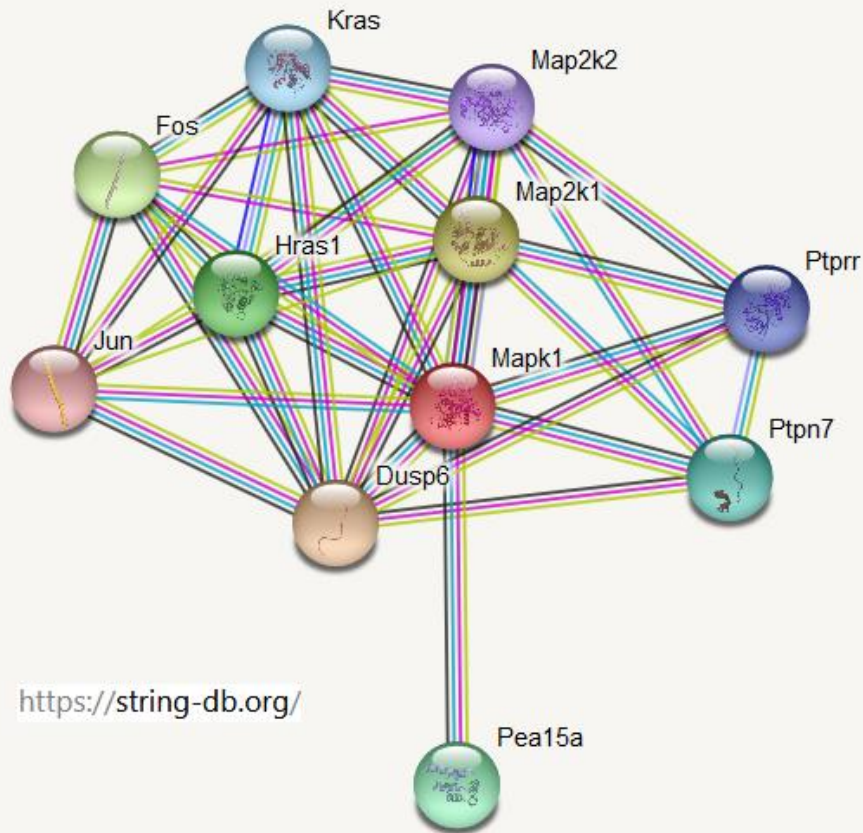Multi-algorithm clustering app for Cytoscape

# How do we represent a pathway

# STRING

# Kegg pathways

# Kegg pathways



**Nodes**
- Genes
- Group of genes
- Compounds
- Other networks

**Edges**
- Activation/Inhibition
- Expression
- Indirect
- Phosphorylation/Diphosphorilation
- Ubiquination
- Association/Dissocation

# Many repositories of biological pathways

# Many repositories of biological pathways

| | |
|---|---|
| 2471 Pathways | Many More! |
| >42 000 Pathways | WikiPathways |
| 505, 700 Pathways | Pathway Commons |
| 2132 (H. sapiens) | KEGG |
| | Reactome |

# Biological Pathway Exchange (BioPAX)



Before BioPAX

>100 DBs and tools
Tower of Babel

After BioPAX
Unifying language

# Pathway analysis



| affyID | gene.name | accession | unigene | FC | pfp |
|---|---|---|---|---|---|
| 1419476_at | Adamdec1 | NM_021475 | Mm.36742 | 27.31 | 0 |
| 1448162_at | Vcam1 | BB250384 | Mm.76649 | 26.58 | 2.44 |
| 1419128_at | Itgax | NM_021334 | Mm.22378 | 17.06 | 3.83 |
| 1415989_at | Vcam1 | BB250384 | Mm.76649 | 13.57 | 0 |
| 1418776_at | 5830443L24Rik | NM_029509 | Mm.42526 1 | 11.76 | 0 |

# Many ways to approach pathway analysis



Functional Pathway Analysis

Input

Over-Representation Analysis (ORA)
Differential Expression Analysis → Differentially Expressed (DE) Genes → Number of DE and Reference Genes in Each Pathway

Functional Class Scoring (FCS)
Gene-level Statistics → Gene-set (Pathway) Statistics

Pathway Topology (PT)
DE Genes or Gene-level Statistics → Pathway Impact Factor
Pathway Topology
• Number of Reactions
• Position of Gene
• Type of Reaction

Pathway Database

Assess Pathway Significance

# Over-representation analysis

We have mapped our significantly differentially expressed genes to pathways. So we can start to interpret our results.

How likely is it that if we consider a random set of genes we will observe these pathways?

# Over-representation analysis

Are there any pathways that have a larger than expected subset of our selected genes in their annotation list?

Gene universe

# Over-representation analysis

Are there any pathways that have a larger than expected subset of our selected genes in their annotation list?
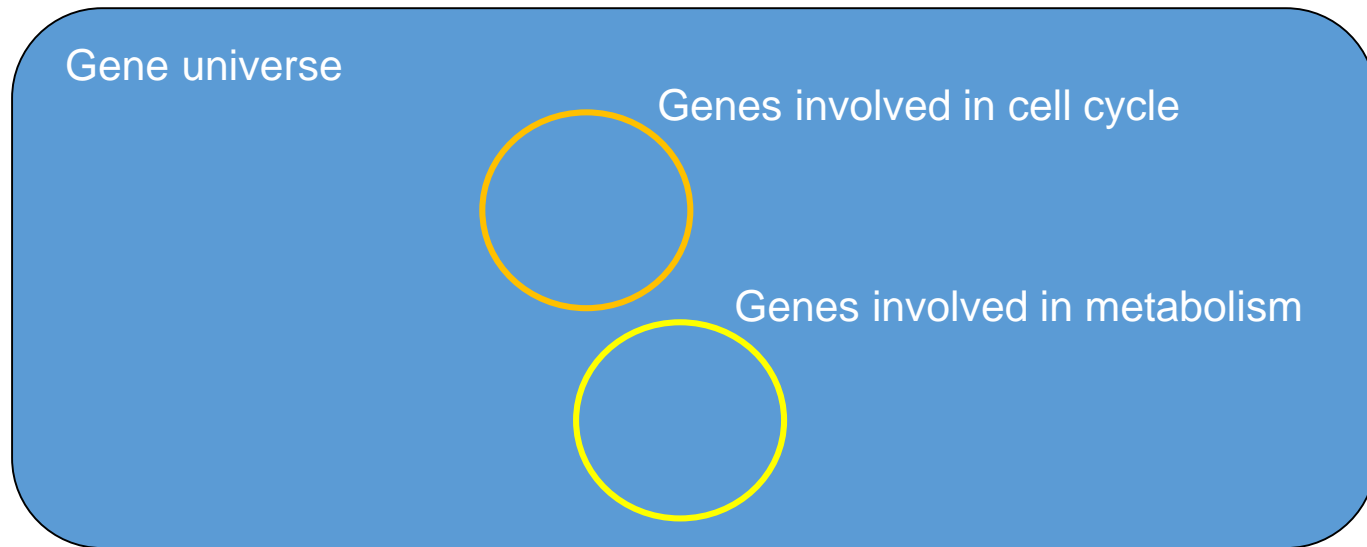
# Over-representation analysis

Are there any pathways that have a larger than expected subset of our selected genes in their annotation list?

# Over-representation analysis

Are there any pathways that have a larger than expected subset of our selected genes in their annotation list?



Two-way table:

|  | Selected | Universe |
|---|---|---|
| In Pathway | ? | ? |
| Not In Pathway | ? | ? |
| Total | ? | ? |

# Over-representation analysis

Are there any pathways that have a larger than expected subset of our selected genes in their annotation list?
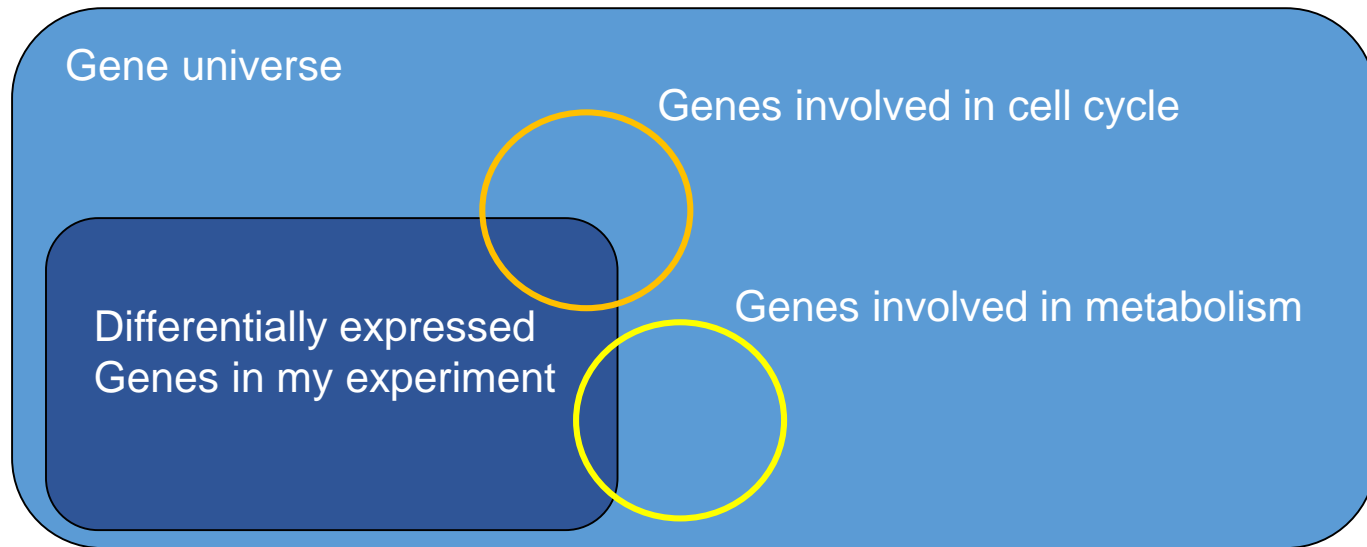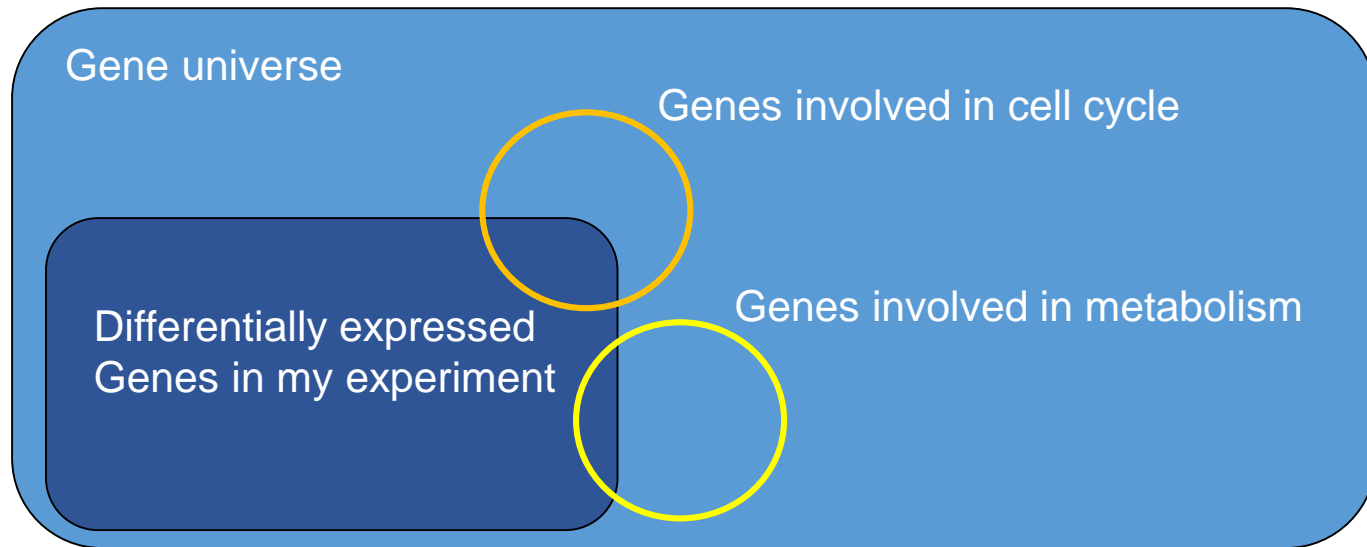
Gene universe

Genes involved in cell cycle

Genes involved in metabolism

Differentially expressed Genes in my experiment

Two-way table:

|  | Selected | Universe |
|---|---|---|
| In Pathway | 22 | 7500 |
| Not In Pathway | 28 | 22500 |
| Total | 50 | 30000 |

# Over-representation analysis

Are there any pathways that have a larger than expected subset of our selected genes in their annotation list?



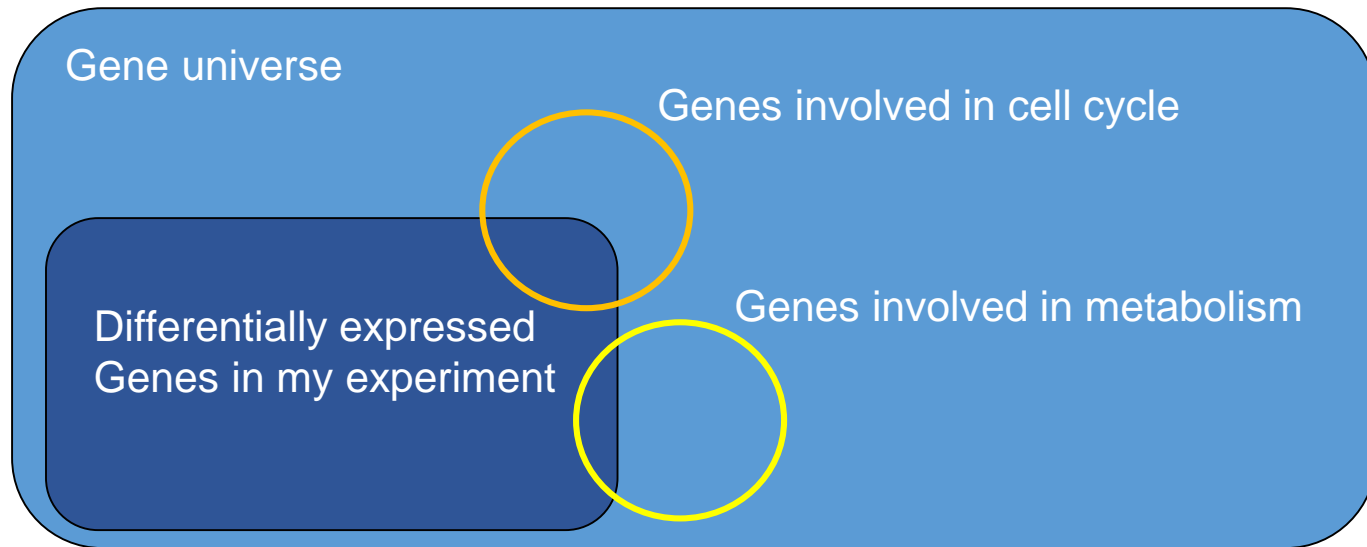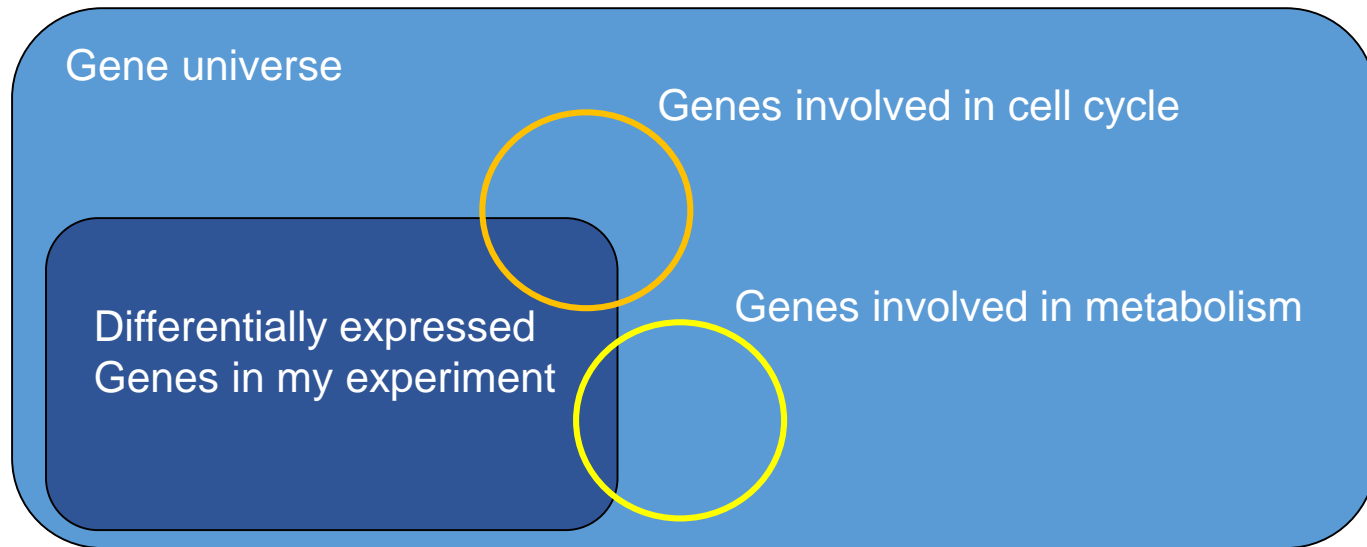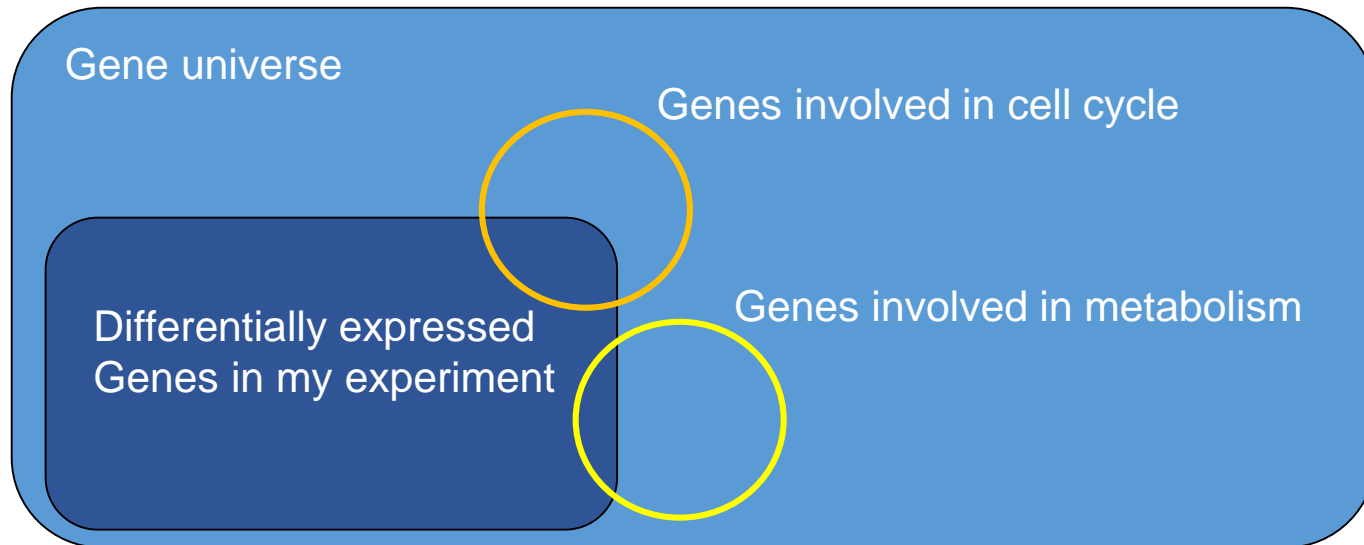Two-way table:

| | Selected | Universe |
|---|---|---|
| In Pathway | 22 | 7500 |
| Not In Pathway | 28 | 22500 |
| Total | 50 | 30000 |

Fold enrichment = (22/50) / (7500/30000) = 45% / 25% = 1.8

# The Hypergeometric test



"The probability of drawing up to *x* of a possible *K* items in *N* drawings without replacement from a group of *M* objects"

X = the number of differentially expressed genes belonging to the pathway

K = the number of genes belonging to the pathway

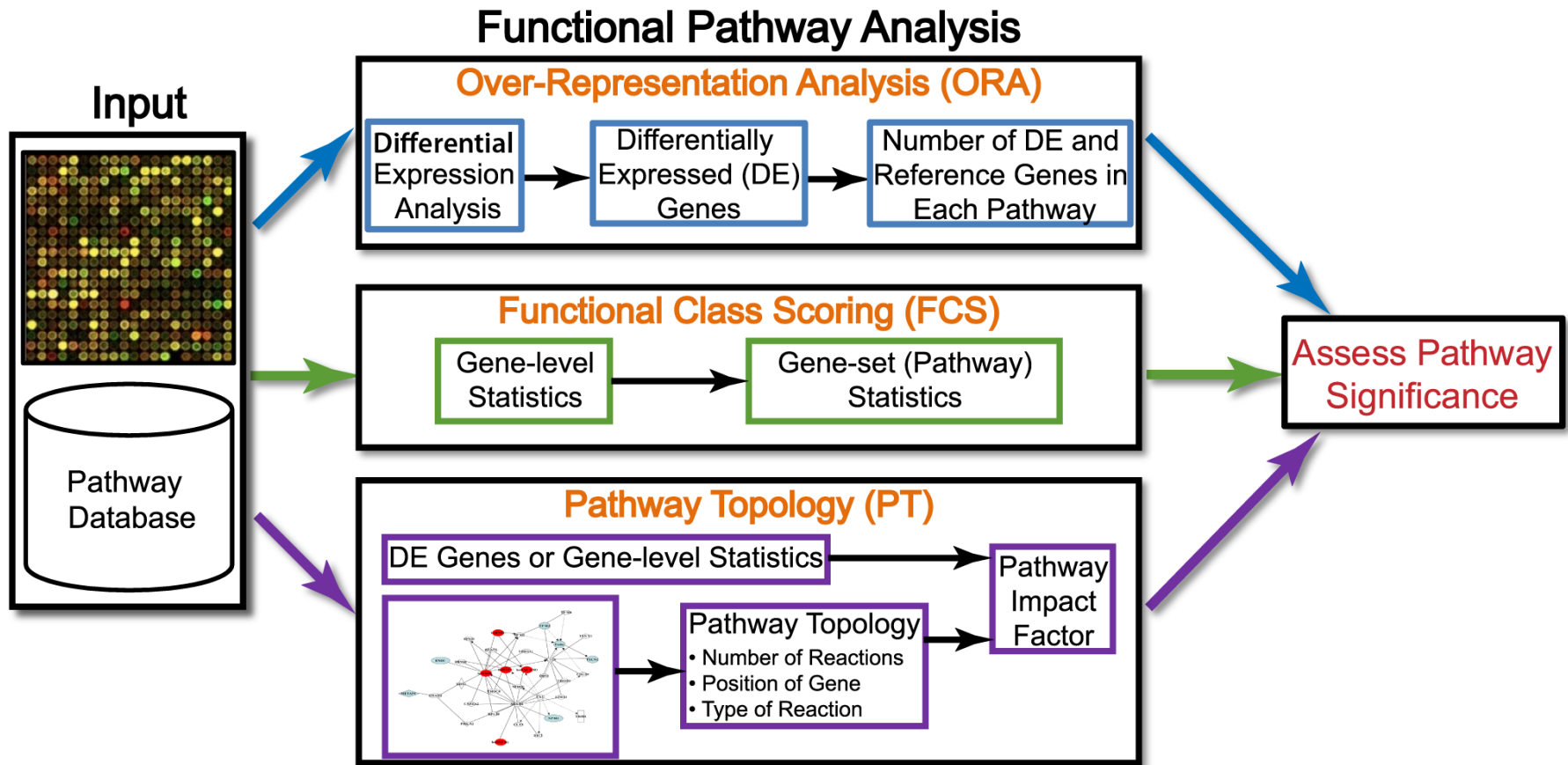N = the differentially expressed genes (or selected genes)

M = the universe

$$p = F(x \mid M, K, N) = \sum_{i=0}^{x} \frac{\binom{K}{i}\binom{M-K}{N-i}}{\binom{M}{N}}$$

# Limitations

- Not always clear how to define the universe

- The over-representation analysis is independent of the changes measured. All genes are treated equally.

- Only the most significant genes are used which causes information loss

- Genes are assumed to be independent and the correlation structure is ignored

- Pathways are assumed to be independent

# Many ways to approach pathway analysis

Khatri P, Sirota M, Butte AJ (2012) Ten Years of Pathway Analysis: Current Approaches and Outstanding Challenges. PLoS Comput Biol 8(2): e1002375. doi:10.1371/journal.pcbi.1002375
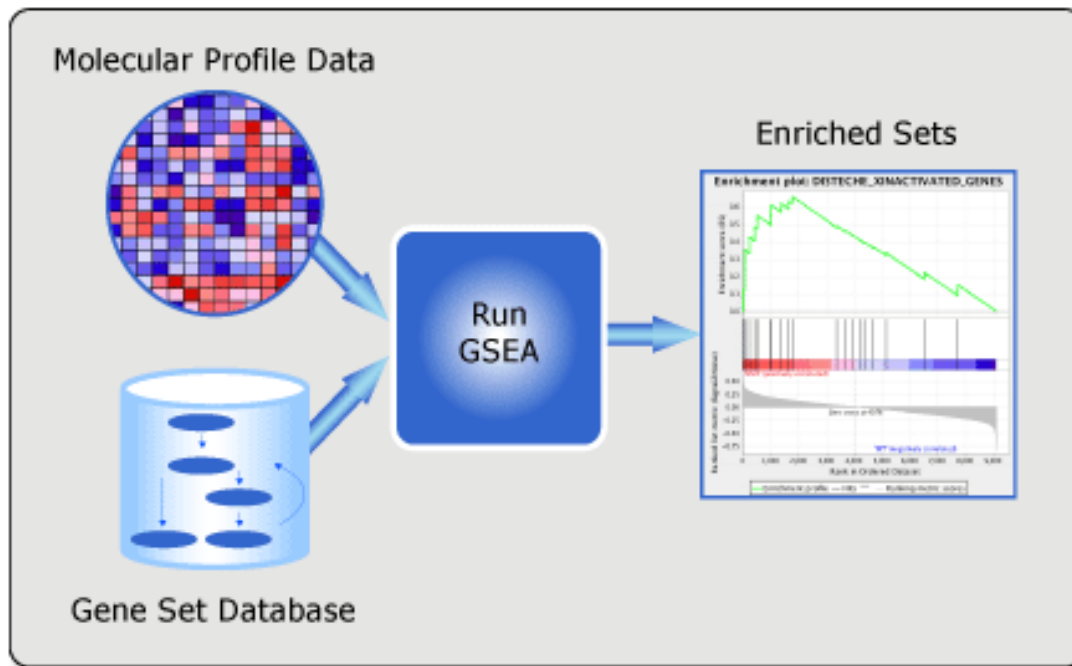http://www.ploscompbiol.org/article/info:doi/10.1371/journal.pcbi.1002375
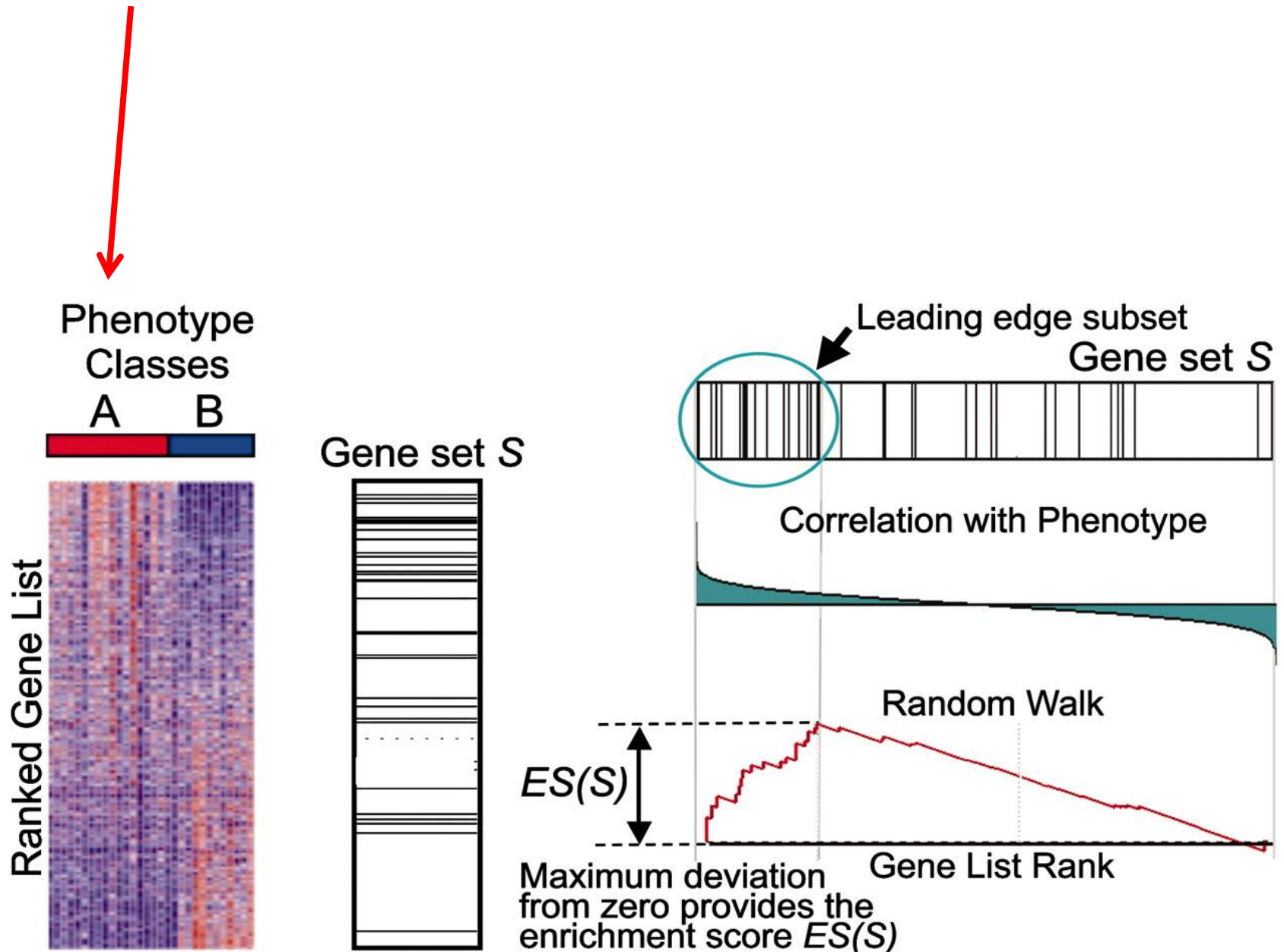
# Gene Set Enrichment Analysis

Tests whether an a priori defined set of genes shows statistically significant, concordant differences between two biological states (e.g. phenotypes)
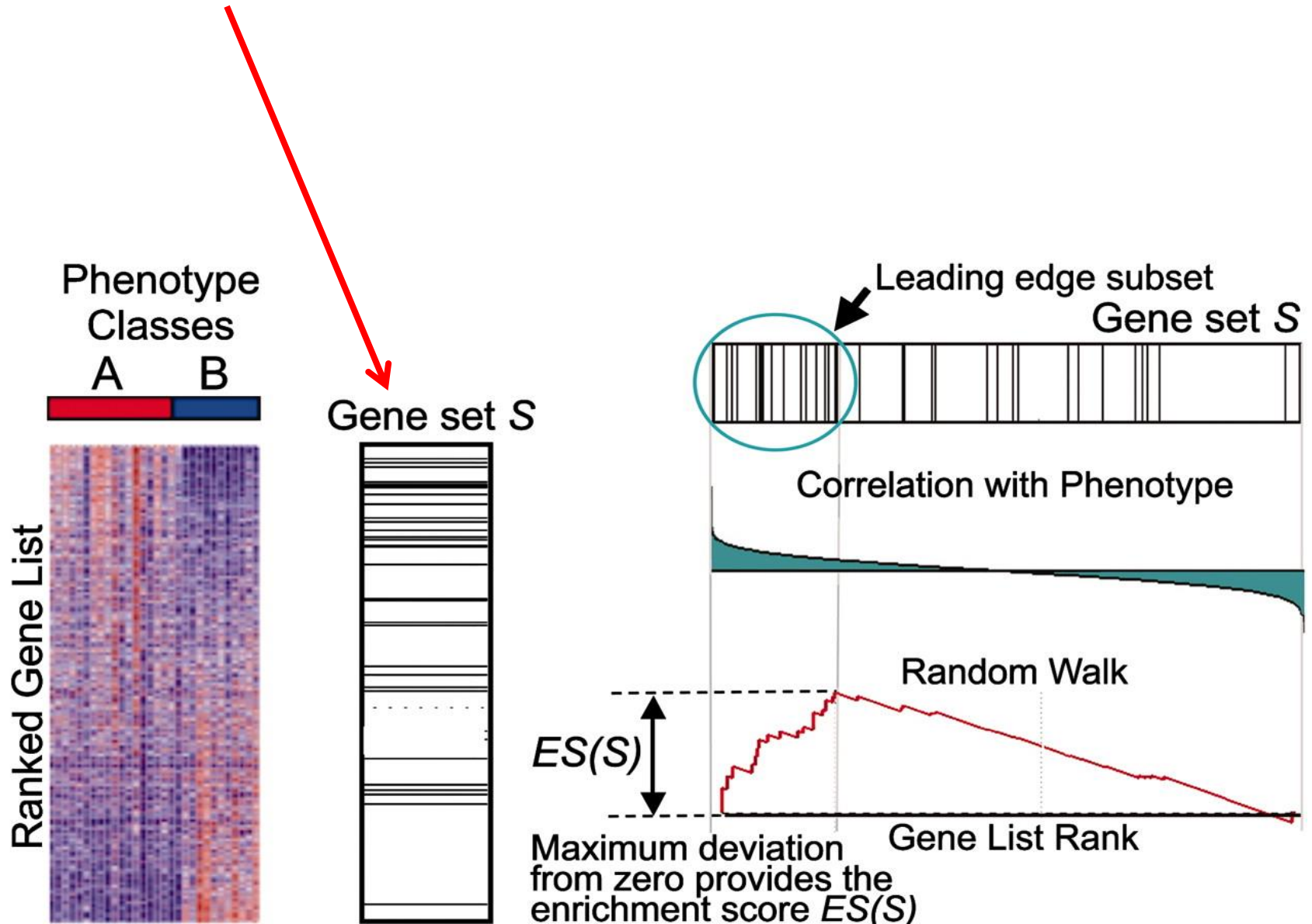
Hypothesis:
Pathway regulation can be detected either by looking at large changes in individual genes or by looking at coordinated changes in sets of functionally related genes.
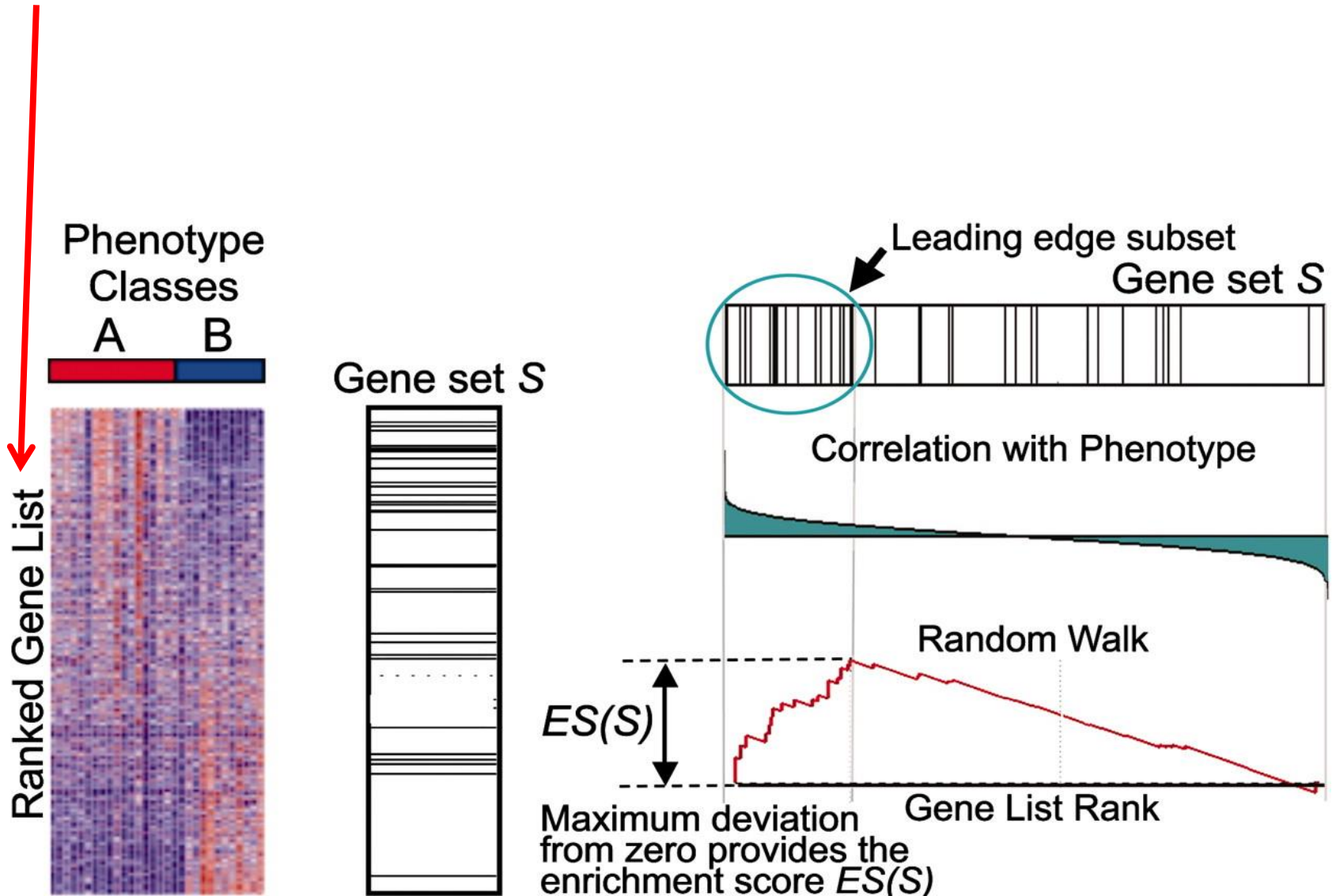


http://www.broadinstitute.org/gsea/

1) Define phenotype classes (e.g. Cells in Hypoxia or Normoxia)



Phenotype Classes
A  B

Ranked Gene List

Gene set S

Leading edge subset
Gene set S

Correlation with Phenotype

Random Walk

ES(S)

Gene List Rank
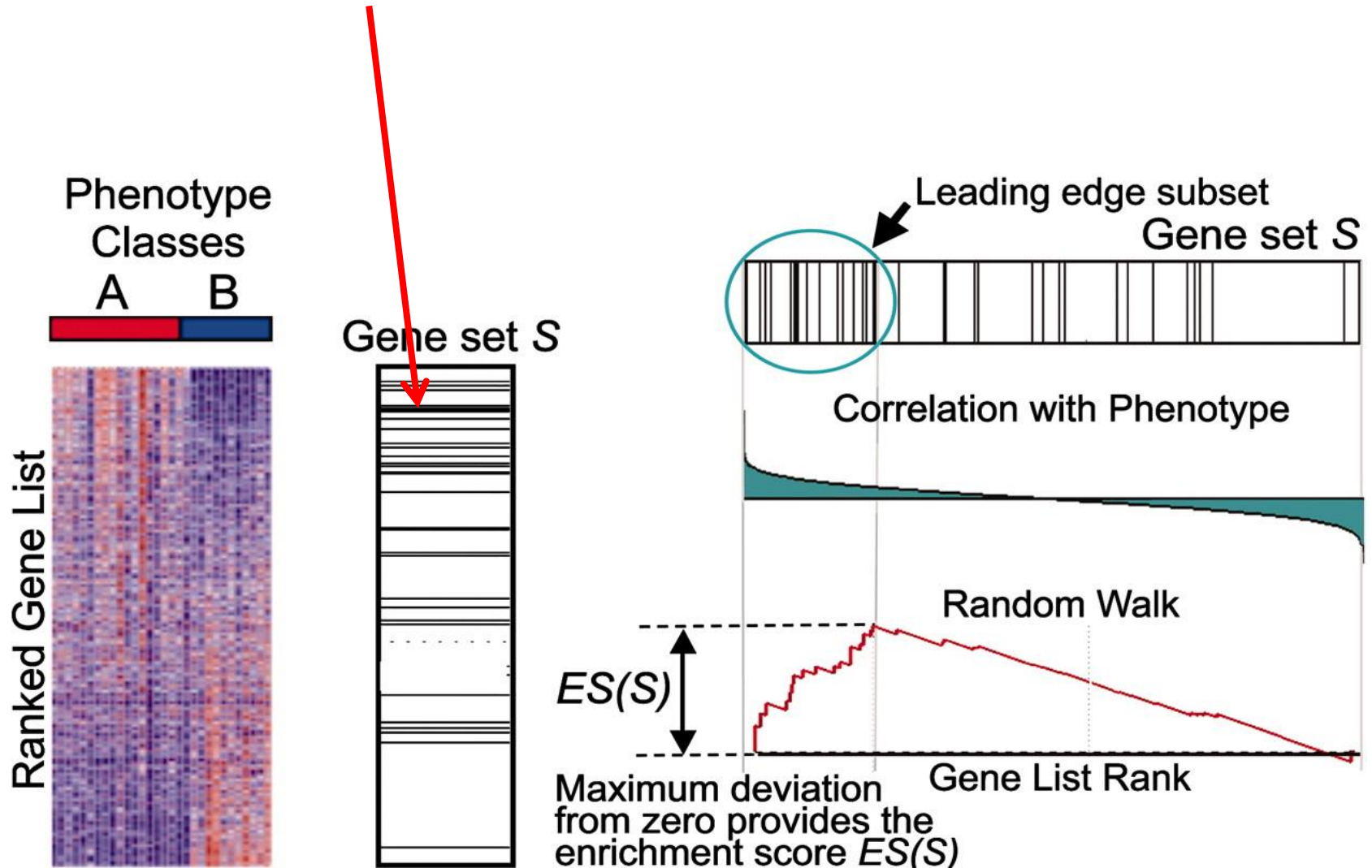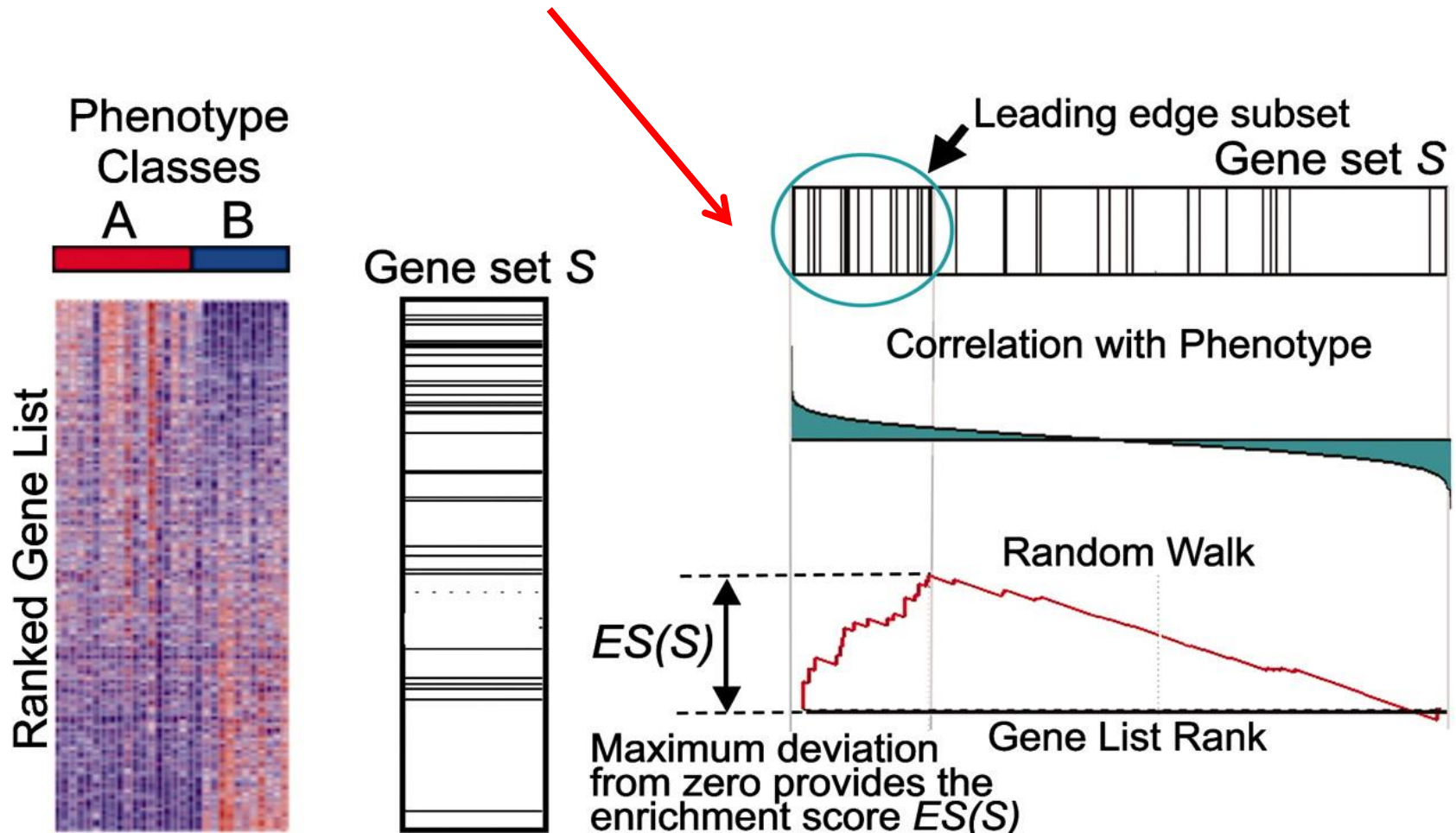
Maximum deviation from zero provides the enrichment score ES(S)

1) Define phenotype classes (e.g. Cells in Hypoxia or Normoxia)
2) Define a gene set of interest (e.g. Genes in the HIF1a pathway)



Phenotype Classes
A   B

Ranked Gene List

Gene set $S$

Leading edge subset
Gene set $S$

Correlation with Phenotype

Random Walk

$ES(S)$

Maximum deviation from zero provides the enrichment score $ES(S)$

Gene List Rank

1) Define phenotype classes (e.g. Cells in Hypoxia or Normoxia)
2) Define a gene set of interest (e.g. Genes in the HIF1a pathway)
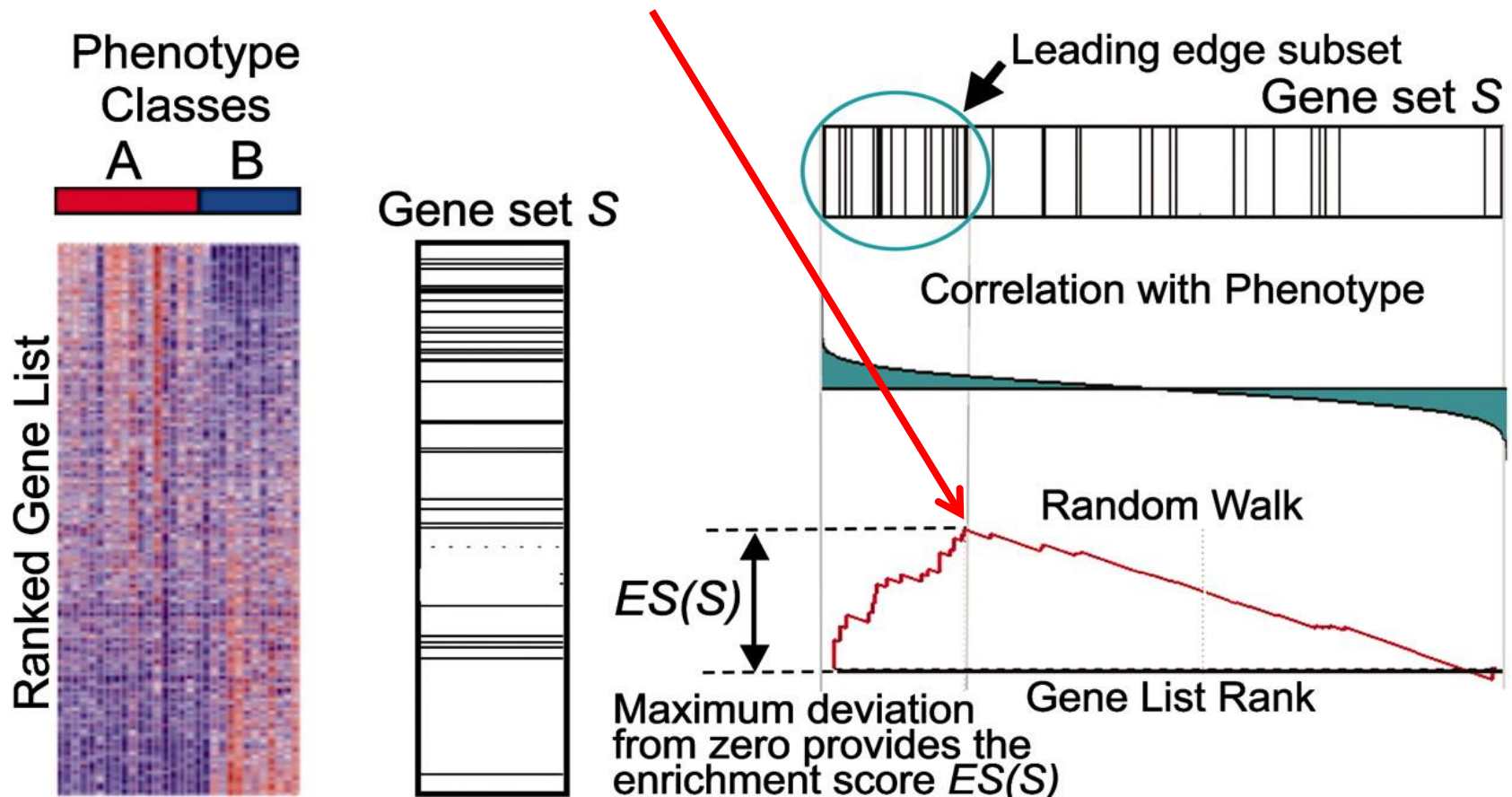3) Sort genes based on their differential expression between classes

1) Define phenotype classes (e.g. Cells in Hypoxia or Normoxia)
2) Define a gene set of interest (e.g. Genes in the HIF1a pathway)
3) Sort genes based on their differential expression between classes
4) Tag genes from the set S within the sorted list

1) Define phenotype classes (e.g. Cells in Hypoxia or Normoxia)
2) Define a gene set of interest (e.g. Genes in the HIF1a pathway)
3) Sort genes based on their differential expression between classes
4) Tag genes from the set S within the sorted list
5) Walk down the list, for each gene: if gene is in S running-sum statistic up, if not down. (The magnitude of the increment depends on FC)

1) Define phenotype classes (e.g. Cells in Hypoxia or Normoxia)
2) Define a gene set of interest (e.g. Genes in the HIF1a pathway)
3) Sort genes based on their differential expression between classes
4) Tag genes from the set S within the sorted list
5) Walk down the list, for each gene: if gene is in S running-sum statistic up, if not down. (The magnitude of the increment depends on FC)
6) ES is the maximum deviation from zero of this random walk

# Improvement on over-representation

- No need to define an arbitrary threshold for selecting significant genes

- The molecular measurements of the actual changes are not ignored but used in order to detect coordinated changes in the expression of genes in the same pathway.

- Coordinate changes are considered: the dependence between genes in a pathway is accounted for

# Limitations

▪ Pathway are considered independent. However a gene can function in more than one pathway, meaning that pathways can cross and overlap.

▪ Most methods use ranks instead of the actual changes (exception exist: gene set analysis http://statweb.stanford.edu/~tibs/GSA/ but only available as R function at the moment)

▪ The nature of the functional link between genes , the strength of the evidence for this link, the role of the genes in the pathway are not considered, only the list of genes in a pathway is used

# Many ways to approach pathway analysis

# MinePATH



Available at [minepath.org](minepath.org) (Koumakis et al., 2012)

# Kegg pathways



MAPK SIGNALING PATHWAY

# Kegg pathways



**Nodes**
- Genes
- Group of genes
- Compounds
- Other networks

**Edges**
- Activation/Inhibition
- Expression
- Indirect
- Phosphorylation/Diphosphorilation
- Ubiquination
- Association/Dissocation

# System overview

# System overview

Discretise the gene expression data

**0: down-regulated genes**
**1: up-regulated genes**



**Data Collection** — **Pre-process** — **Annotation** — **Combine sources** — **Analyze Data** — **Visualize**

MicroArray

Discretize

g1=hsa1
g2=hsa2

Pathway(s)

Deco

|  | Sample 1 | Sample 2 | Sample 3 | ... |
|---|---|---|---|---|
| Gene A | 7.8 | 2.5 | 5.9 | ... |
| Gene B | 7.8 | 6.4 | 7.0 | ... |
| Gene C | 6.3 | 8.2 | 6.4 | ... |
| ... | ... | ... | ... | ... |

Discretised data

|  | Sample 1 | Sample 2 | Sample 3 | ... |
|---|---|---|---|---|
| Gene A | 1 | 0 | 1 | ... |
| Gene B | 1 | 1 | 1 | ... |
| Gene C | 0 | 1 | 0 | ... |
| ... | ... | ... | ... | ... |

# System overview

Decompose gene networks

**Data Collection**

MicroArray



**Pre-process**

Discretize



Pathway(s)



Decompose



Clone
multi-probes

# System overview

Combine discretized gene expression data and decomposed sub-paths

# Combine discretized gene expression data and decomposed sub-paths

**Activation**

A → B

**Inhibition**

A —| B

|  | Sample 1 | Sample 2 | Sample 3 | … |
|---|---|---|---|---|
| Gene A | 1 | 0 | 1 | … |
| Gene B | 1 | 1 | 1 | … |
| Gene C | 0 | 1 | 0 | … |
| … | … | … | … | … |

# Combine discretized gene expression data and decomposed sub-paths

**Activation**

| | | B | |
|---|---|---|---|
| | | ON | OFF |
| A | ON | ✓ | ✗ |
| | OFF | ✗ | ✗ |

A ➡ B

*Logic AND*

**Inhibition**
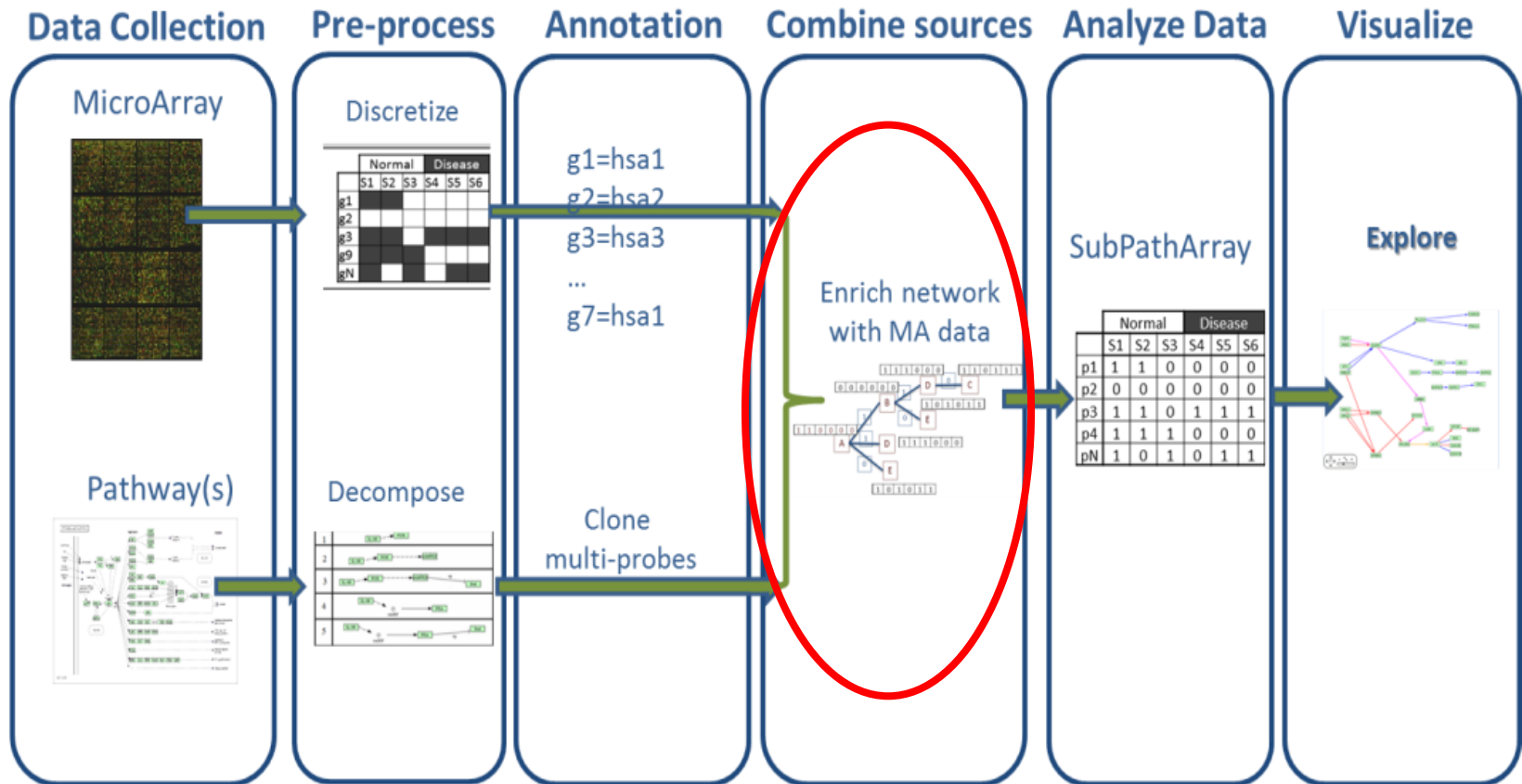
| | | B | |
|---|---|---|---|
| | | ON | OFF |
| A | ON | ✗ | ✓ |
| | OFF | ✓ | ✗ |

A ⊣ B

*Logic XOR*

| | Sample 1 | Sample 2 | Sample 3 | ... |
|---|---|---|---|---|
| Gene A | 1 | 0 | 1 | ... |
| Gene B | 1 | 1 | 1 | ... |
| Gene C | 0 | 1 | 0 | ... |
| ... | ... | ... | ... | ... |

Interactions as logical operators

# System overview

Evaluate sub-paths

# System overview

Visualize the results

# Visualizing the best pathways
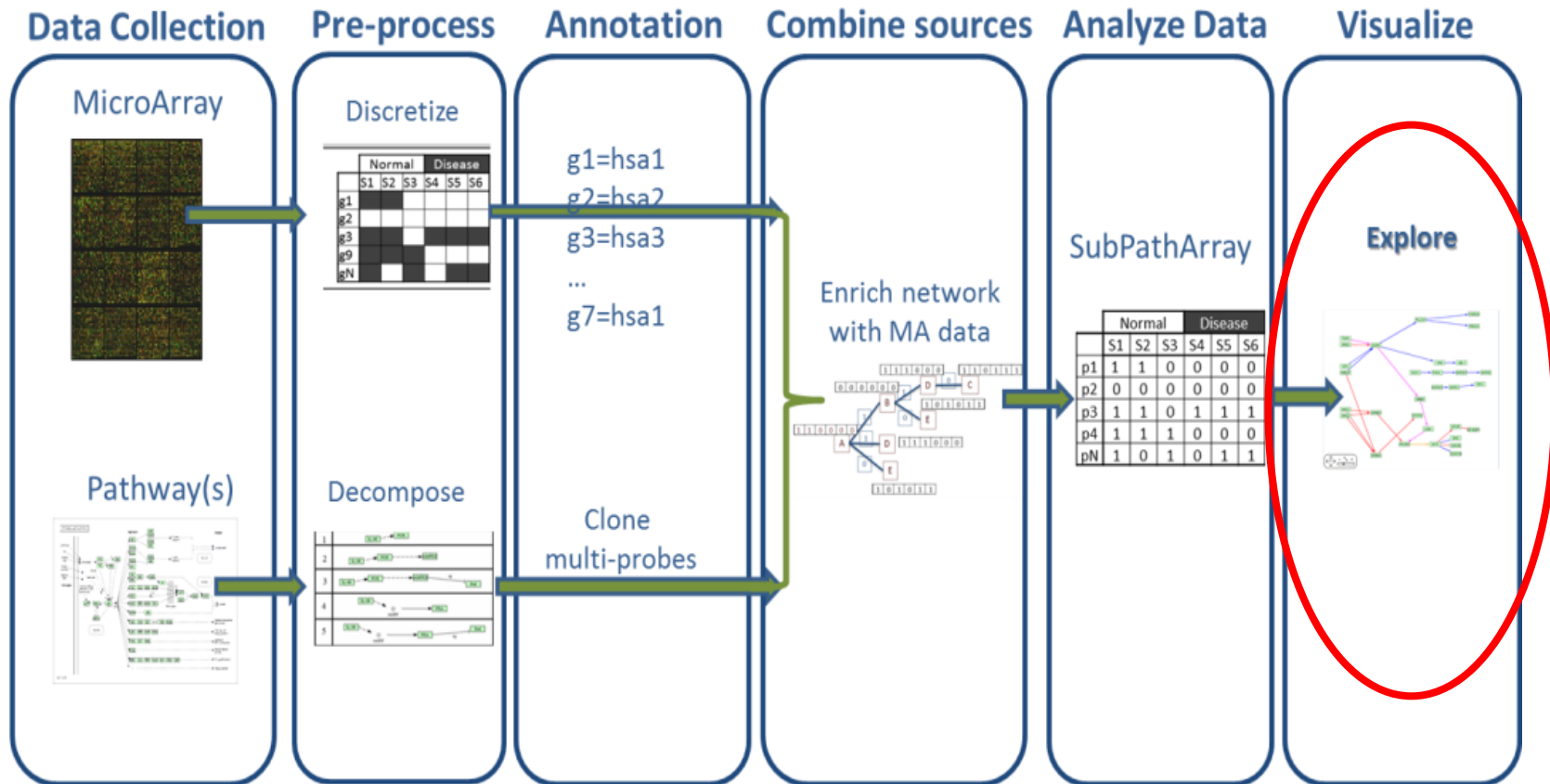
- '**Red**' is used to encode sub-path relations that are active for phenotype 1 (Class 1)

- '**Blue**' for relations that are active for phenotype 2 (Class 2)

- '**Magenta**' for relations holding for both phenotypes

- '**Orange**' for relations that are "always-active"

- "**Yellow**" for the association/disassociation relations

- '**Grey**' for inactive relations.

# There are still many gaps



SNPs for genes and transcripts

Single Nucleotide Polymorphisms (dbSNP)

EGFR

AlternativeTranscripts

Effect of SNPs on pathways

Annotations for transcripts

Annotations for genes

Functional annotations (GO, KEGG, BioCarta, etc.)

Cell-cell adhesion
Cell morphogenesis

Biological System

Relationships between pathways in an organism

Changes in pathways due to a disease

Nonsmall cell lung cancer

Diseases (OMIM)

Role of pathways in diseases