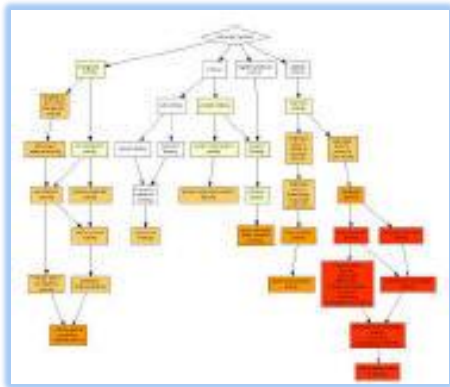LIFE SCIENCE
CORE FACILITIES

# Functional analysis of gene lists using Gene Ontology (GO)

Noa Wigoda

17.12.19

An Introduction to deep-sequencing analysis for biologists

# OUTLINE

- Single gene analysis / information

- Analysis of group of genes

- Gene ontology (GO)

- Enrichment analysis

  - Hypergeometric Test and Fisher exact test
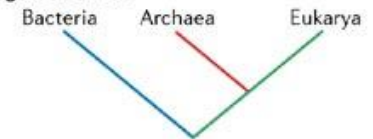
  - GO Independence Assumption
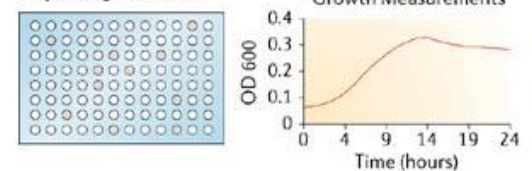


Genome sequence and annotation

Available literature
Pub Med

Phylogenetic data
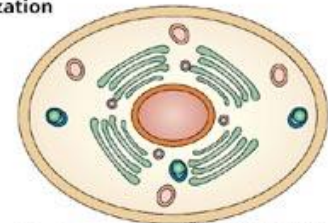Bacteria   Archaea   Eukarya

Physiological data
Growth Measurements

Databases
KEGG   EcoCyc

Localization
Signal sequences: PLLLLPISGSALP

Ask question which can be answered with a simple "Yes" or "No."

**20 Questions**

Is it part of a complex?

Is it a protein coding gene?

Is it a regulator – transcription factor?
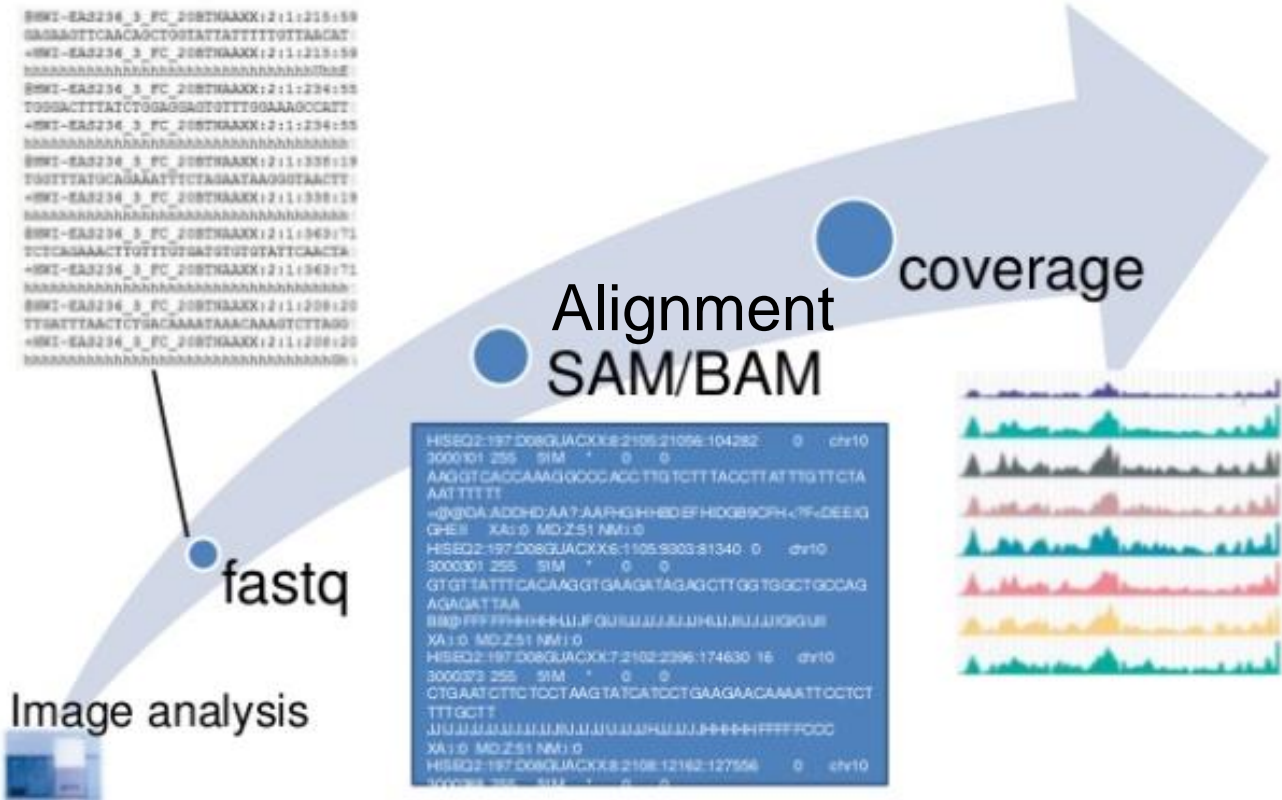
Is it in the nucleus?

Is it an enzyme?

Is it related to a disease?

All the answers are "attributes" or characteristics of the item (gene).

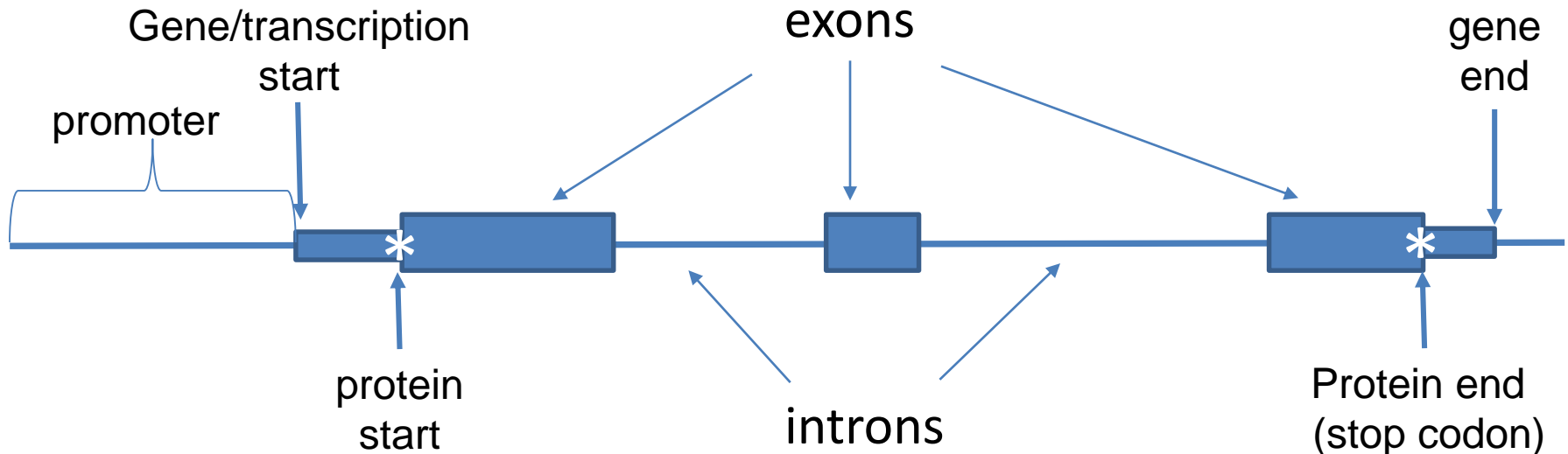# What have we done until now?



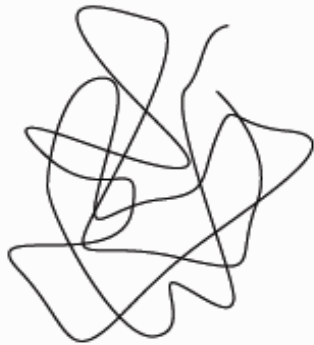Information flow of sequencing data

# What is a Gene ?

A gene is a region of DNA that encodes instructions for how the cell can make a gene product, which can be:

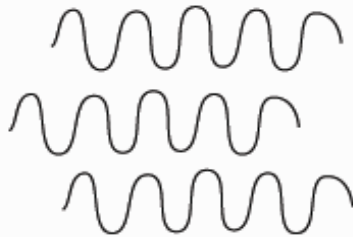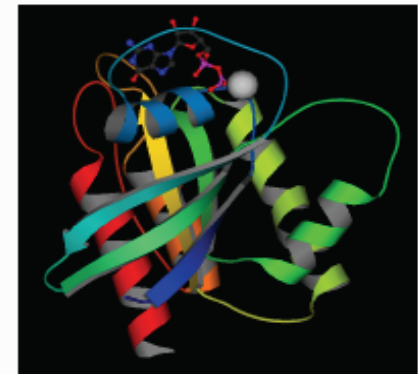- a protein
- a noncoding RNA.

# Data sources



| Genome | Transcripts | Protein |
| --- | --- | --- |

- There are several kinds of databases, looking at the genome, transcriptome or proteome level
- The mapping of the different names is not trivial

# Levels of annotation per gene

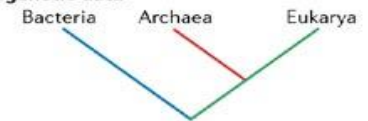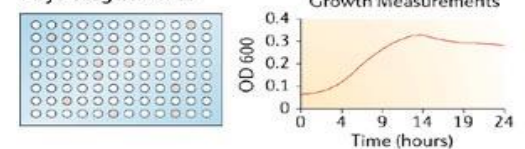| Level | Database |
|---|---|
| Sequence | GenBank<br>SwissProt (curated) |
| Metabolic pathways | Kegg<br>Transpath<br>MetaCyc |
| Literature | PubMed |
| Gene ontology (GO) | Biological process<br>Molecular function<br>Cell compartment |
| Integrated – Meta databases | GeneCards<br>Entrez Gene<br>OMIM<br>InterPro |



Genome sequence and annotation

Available literature
Pub Med

Phylogenetic data
Bacteria    Archaea    Eukarya

Physiological data
Growth Measurements
OD 600
0.4
0.3
0.2
0.1
0
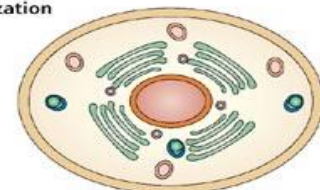0   4   9   14   19   24
Time (hours)

Databases
KEGG    EcoCyc

Localization

Signal sequences: *PLLLLPISGSALP*

# OUTLINE

- Single gene analysis / information

- **Analysis of group of genes**

- Gene ontology (GO)

- Enrichment analysis

  - Hypergeometric Test and Fisher exact test
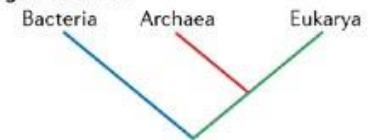
  - GO Independence Assumption



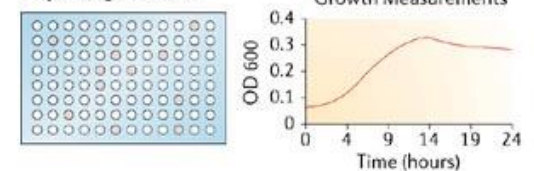Genome sequence and annotation
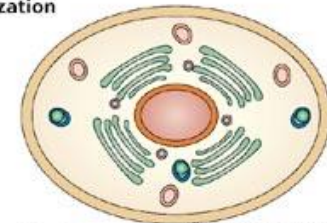
Available literature

Phylogenetic data
Bacteria  Archaea  Eukarya

Physiological data

Databases

Localization

Signal sequences: *PLLLLPISGSALP*
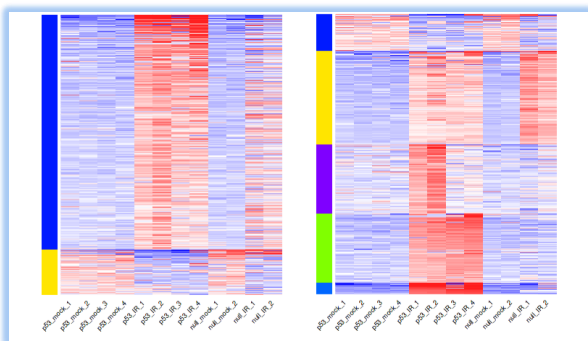
Copyright © 2006 Nature Publishing Group
**Nature Reviews | Genetics**

# What have we done until now?

A complex high-throughput experiment:
Deep Sequencing
Proteomics
Microarrays
…

## What did we get?

## Lists of genes



**Clusters** of differential genes

Up regulated     Down regulated

# Functional Genomics:
# Find the biological meaning

- Take a list of "interesting" genes and find their biological meaning

- Requires a reference set of "biological knowledge"

-  Linking between genes and biological function:

  ❖ Gene ontology: GO

  ❖ Pathways databases

# The problem

- Vast amounts of biological data

- Different names/terms for the same concepts

  For example: the same function can be called

  translation or protein synthesis.

- Cross-species comparison is difficult

# Part of the solution

# Gene Ontology

# What is Ontology?

1700s

Ontology (from the Greek…) is the philosophical study of the nature of being, existence or reality in general, as well as of the basic categories of being and their relations.

Traditionally listed as a part of the major branch of philosophy known as metaphysics, ontology deals with questions concerning what entities exist or can be said to exist, and how such entities can be grouped, related within a hierarchy, and subdivided according to similarities and differences.

- The Gene Ontology (GO) project is a collaborative effort to address the need for consistent descriptions of gene products in different databases.

- Gene ontology is an annotation system

- The project provides the controlled and consistent vocabulary of terms and gene product annotations, i.e. terms occur only once, and there is a dictionary of allowed words

http://geneontology.org/

# Why use GO?

- The goal of the GeneOntology (GO) project is to provide a uniform way to describe the functions of gene products from organisms across all kingdoms of life and thereby enable analysis of genomic data.

- bio-ontologies such as GO make domain knowledge available to both humans and computers.

- GO provides the ability to group gene products to some high level term.

There are three structured, controlled vocabularies (ontologies) that use terms to describe gene products in a species-independent manner:

- Biological processes
  - A recognized series of events, must have more than one distinct steps
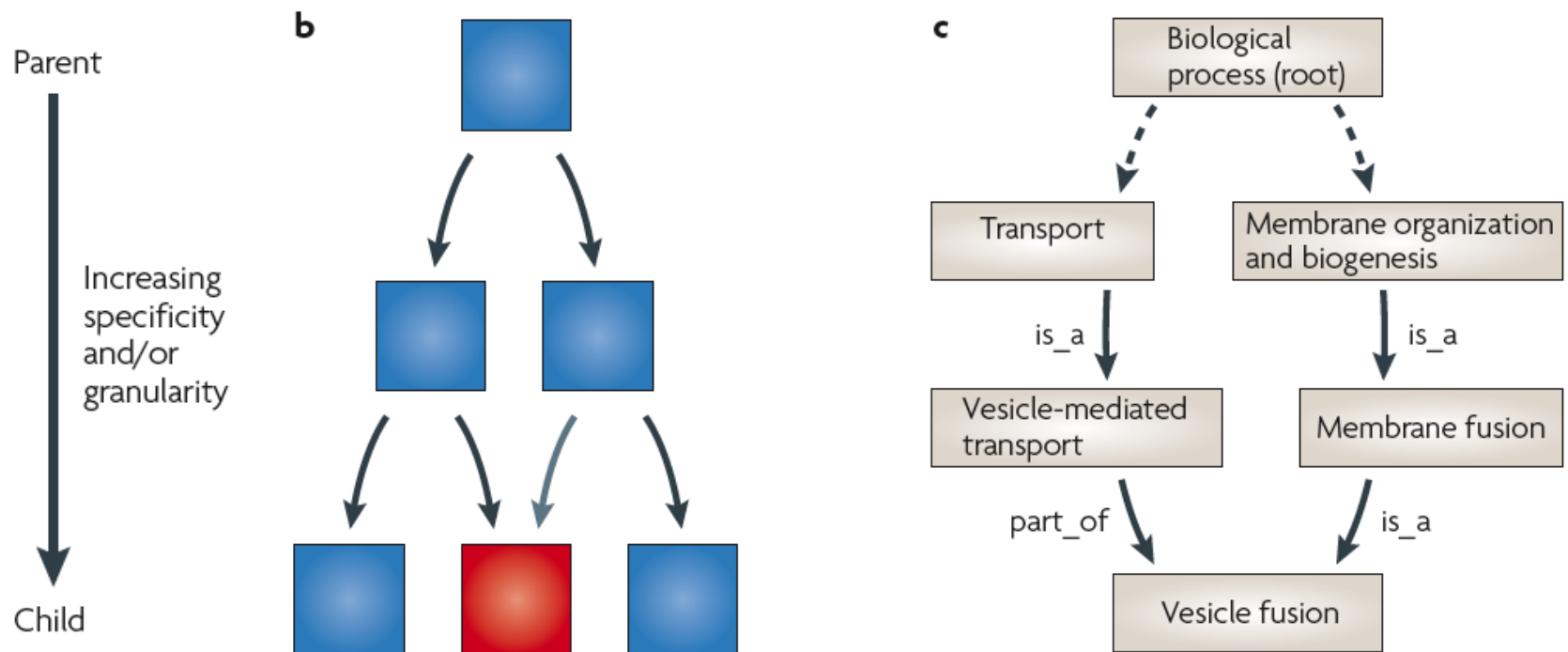  - Examples: cell division, pyrimidine metabolic process

- Cellular components
  - Where a gene product is located (an anatomical structure)
  - Examples: nucleus, proteasome

- Molecular functions
  - describes activities, such as catalytic or binding activities

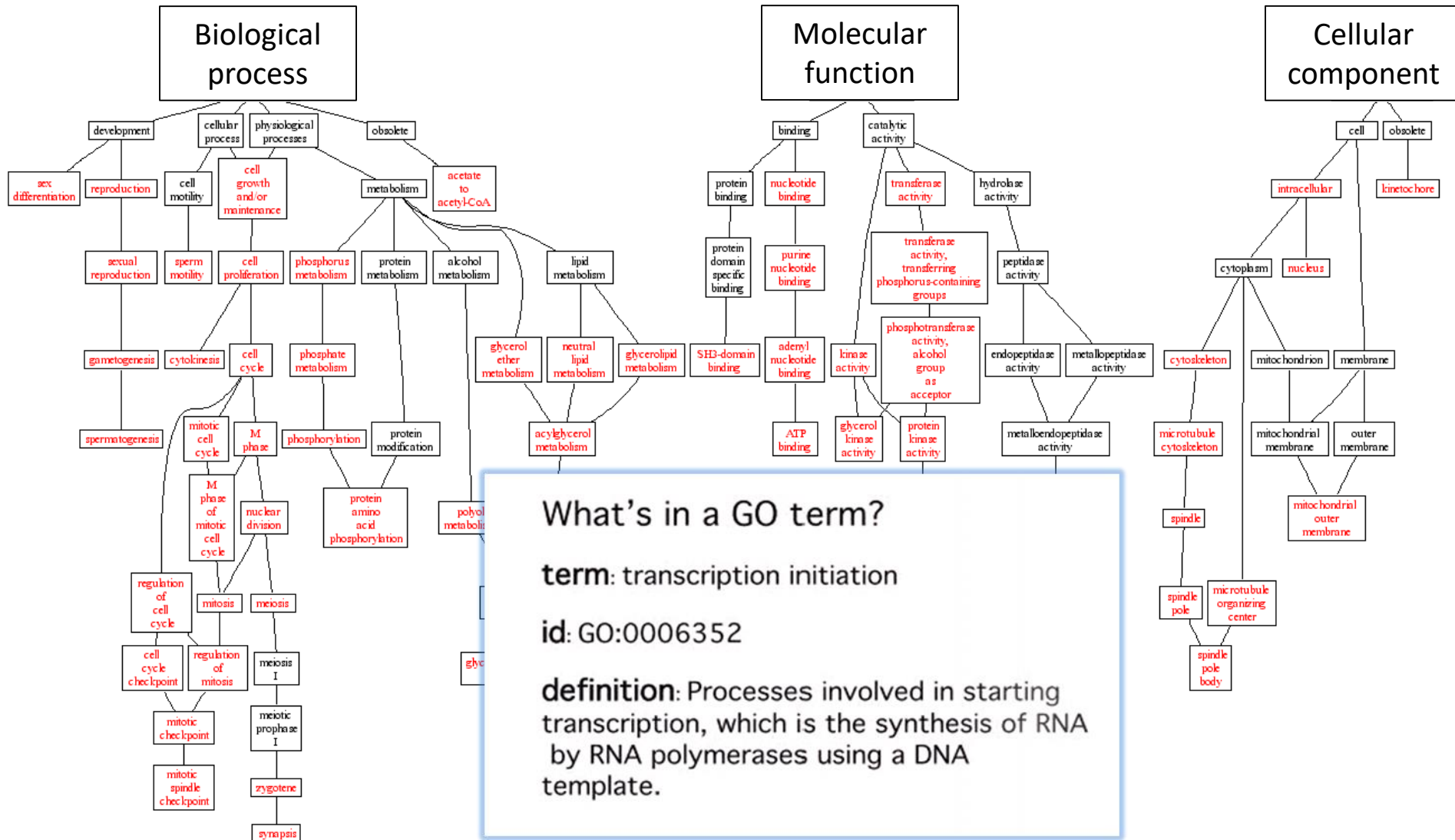# Gene ontology is represented as a directed acyclic graph (DAG)



Taken from: Nature Reviews Genetics 9:509-515 (2008)
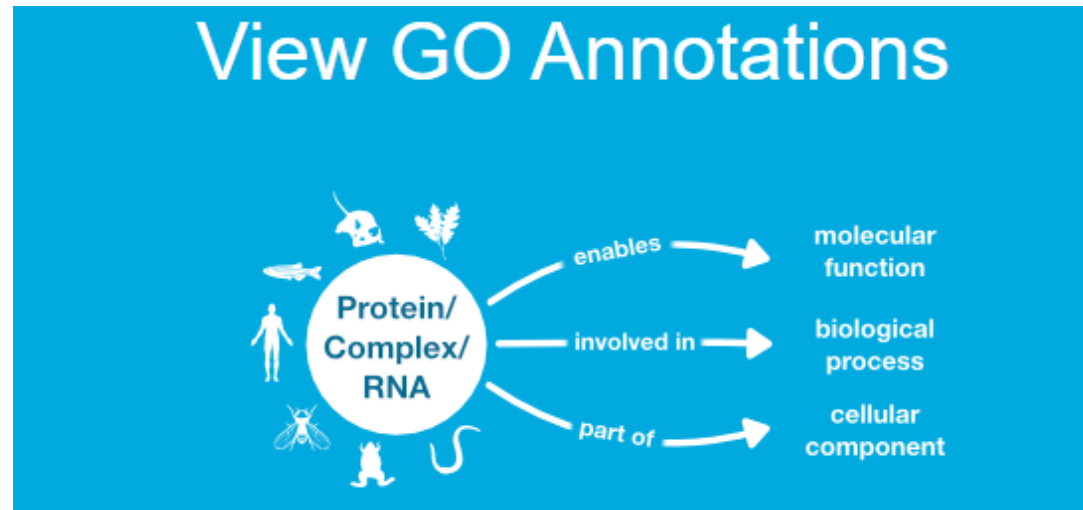
# Directed Acyclic Graph (DAG)

- A child can have more than one parent

  - parents are closer to the root and are more general

  - children are further from the root and more specific

- There are no cycles - there is a root

- It is a directed graph

- You can skip levels in the graph

# Example

# GO annotation

- A GO annotation is a statement about the function of a particular gene.

- GO annotations are associations made between gene products or protein complexes and the GO terms that describe them.

  - attributed to a source
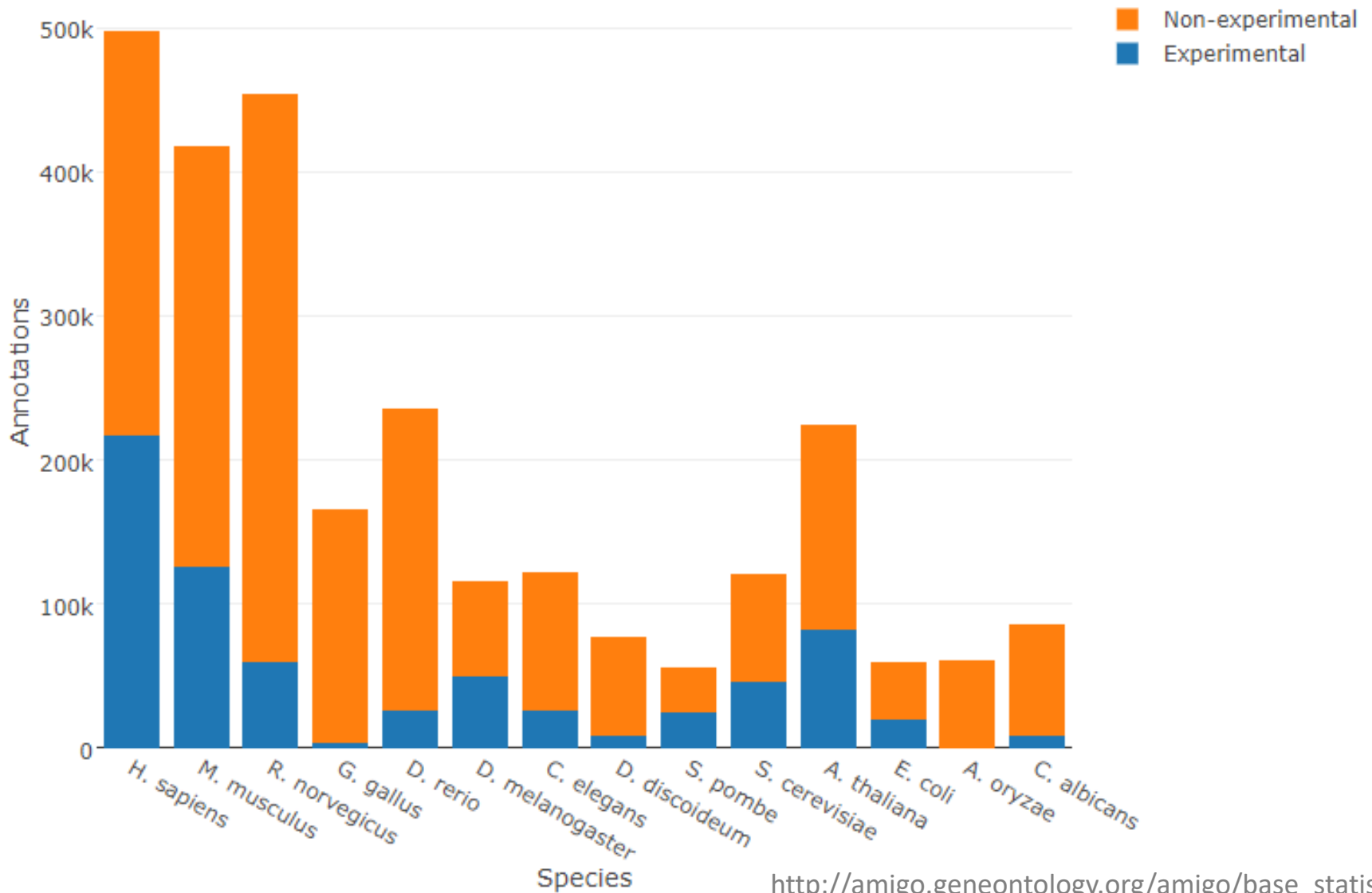  - indicate the evidence upon which it is based.

# Evidence codes
## not all annotations are created equal

| | | | |
|---|---|---|---|
| HTP | EXP | Inferred from Experiment | **BLAST** |
| | | Inferred from High Throughput Experiment | |

| | | |
|---|---|---|
| HDA | IDA | Inferred from Direct Assay |
| | IPI | Inferred from Physical Interaction |
| HMP | IMP | Inferred from Mutant Phenotype |
| HGI | IGI | Inferred from Genetic Interaction |
| HEP | IEP | Inferred from Expression Pattern |

| | |
|---|---|
| ISS | Inferred from Sequence/Structural Similarity |
| TAS | Traceable Author Statement |
| NAS | Non-traceable Author Statement |
| IC | Inferred by Curator |
| ND | No Data available |

| | |
|---|---|
| IEA | Inferred from electronic annotation |

http://geneontology.org/docs/guide-go-evidence-codes/

# Type of annotation per species



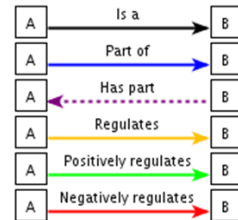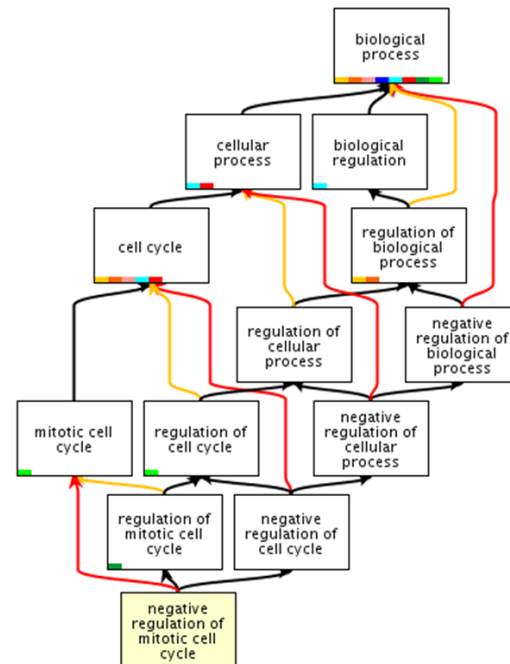http://amigo.geneontology.org/amigo/base_statistics

# Ontology Relations

Defines the relationships (the arrows) between the ontology terms.

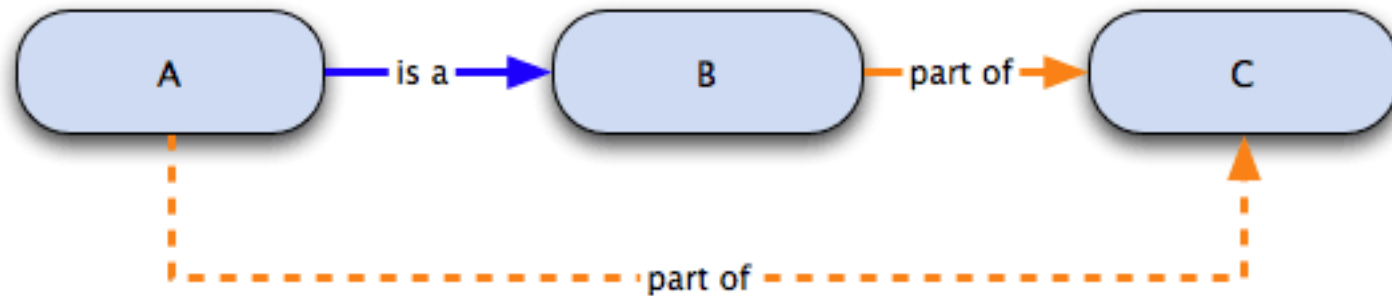There are three types of relationships:
- ❖ is_a
- ❖ part_of
- ❖ regulates:
  - • positively regulates
  - • negatively regulates

# Ontology Relations

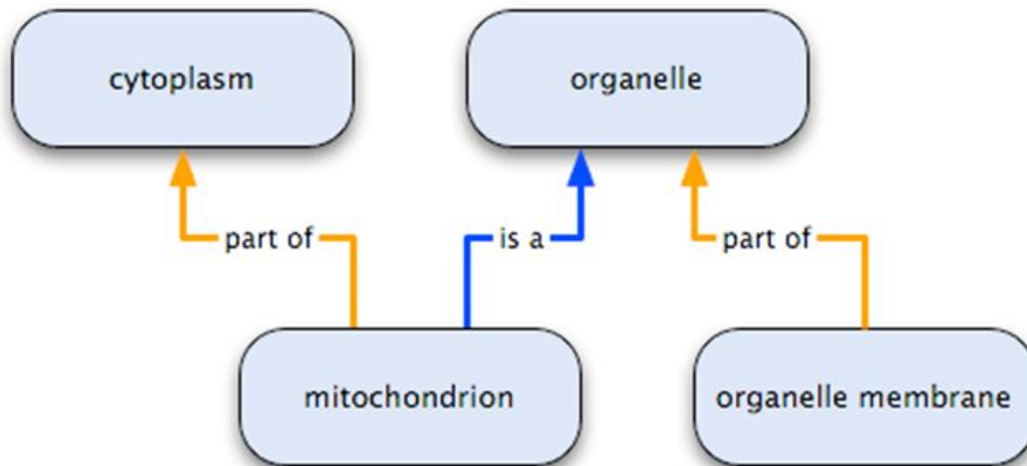- is_a is a simple class-subclass relationship

   Example: nuclear chromosome is_a chromosome.



A dotted line means an inferred relationship, e.g. one that has not been expressly stated

# Ontology Relations

- part_of represent part-whole relationships;
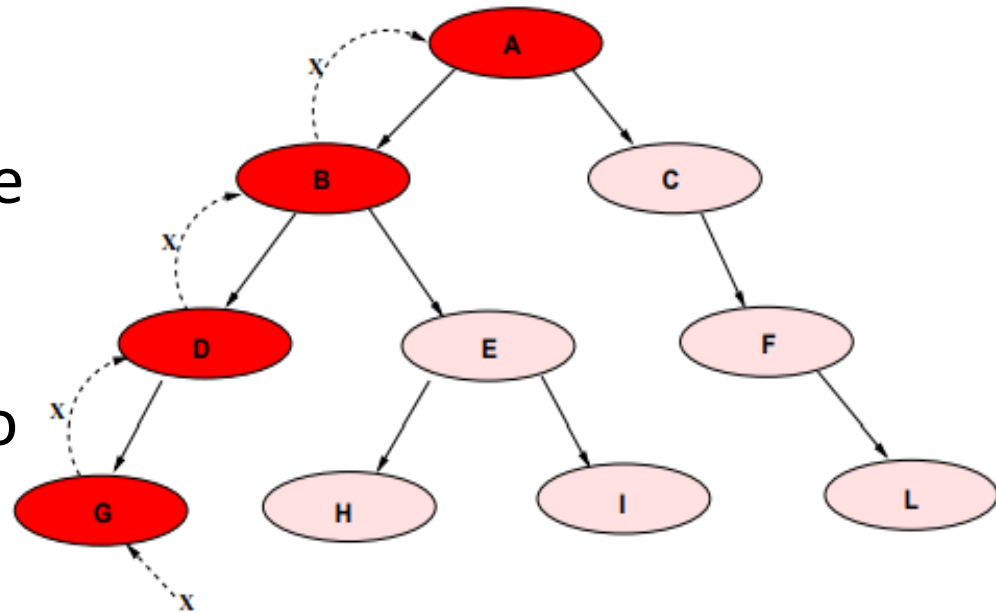  C part_of D means that whenever C is present,
  it is always a part of D.



Example: nucleus part_of cell; nuclei are always part of a cell, but not all cells have nuclei.

mitochondrion has two parents: it *is an* organelle and it is *part of* the cytoplasm;
organelle has two children: mitochondrion *is an* organelle, and organelle membrane is *part of* organelle

# Ontology Structure

Every GO term obeys "the true path rule":

- If a child term describes the gene product, then all its ancestors (parent) terms must also apply to that gene product.

- If a gene is not annotated to a term, it cannot be annotated to its offsprings.

# AmiGO

## a web application to query, browse and visualize ontologies



http://amigo.geneontology.org

# Available GO Information

Current ontology statistics: as of Dec, 2019:

**44674 terms**, 100.0% defined

- 29,380 Biological process terms

- 11,113 Molecular function terms

- 4,181 Cellular component terms

- 2711 obsolete terms (not included in figures above)

20 Questions

Which attribute is not a GO term?

Is it part of a complex?

Is it a protein coding gene?

Is it a regulator – transcription factor?

Is it in the nucleus?

Is it an enzyme?

Is it related to a disease?

All the answers are "attributes" or characteristics of the item (gene).

**GENE**ONTOLOGY
Unifying Biology

There are three structured, controlled vocabularies (ontologies) that use terms to describe gene products in a species-independent manner:

- Biological processes
  - must have more than one distinct steps
  - Examples: signal transduction,

- Cellular components
  - an anatomical structure
  - Examples: nucleus, proteasome

- Molecular functions
  - describes activities, such as catalytic or binding activities

Is it part of a complex?

Is it a protein coding gene?

Is it a regulator – transcription factor?

Is it in the nucleus?

Is it an enzyme?

Is it related to a disease?

*Not a GO term...*

# What is not GO?

- Gene products: e.g. cytochrome c is not in the ontologies, but attributes of cytochrome c, such as oxidoreductase activity, are

- Processes, functions or components that are unique to mutants or diseases: e.g. oncogenesis

- Attributes of sequence such as intron/exon parameters

- Protein domains or structural features

- Protein-protein interactions

- Environment, evolution and expression

- It is not **complete,** it is done "by hand" by curators

- A pathway

# GO Pitfalls

- Not complete

- Computational annotations

- NOT qualifier

- Identifier flagged as 'obsolete', some tools do not update their databases

# OUTLINE

- Single gene analysis / information

- Analysis of group of genes

- Gene ontology (GO)

- Enrichment analysis

  - Hypergeometric Test  and Fisher exact test
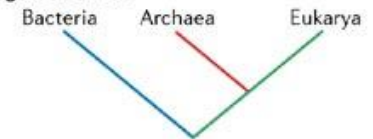
  - GO Independence Assumption



Genome sequence and annotation

Available literature
Pub Med

Phylogenetic data
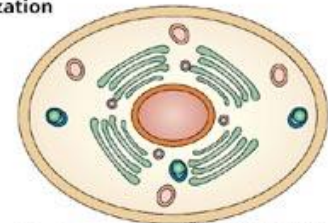Bacteria    Archaea    Eukarya

Physiological data
Growth Measurements
OD 600
Time (hours)

Databases
KEGG    EcoCyc

Localization

Signal sequences: *PLLLLPISGSALP*

Copyright © 2006 Nature Publishing Group
**Nature Reviews | Genetics**

# Two-class Design



**Expression Matrix**

Class-1    Class-2

**Selection by Threshold**

UP

DOWN

**Genes Ranked by Differential Statistic**

UP

DOWN

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

(modified) **bio**informatics.ca
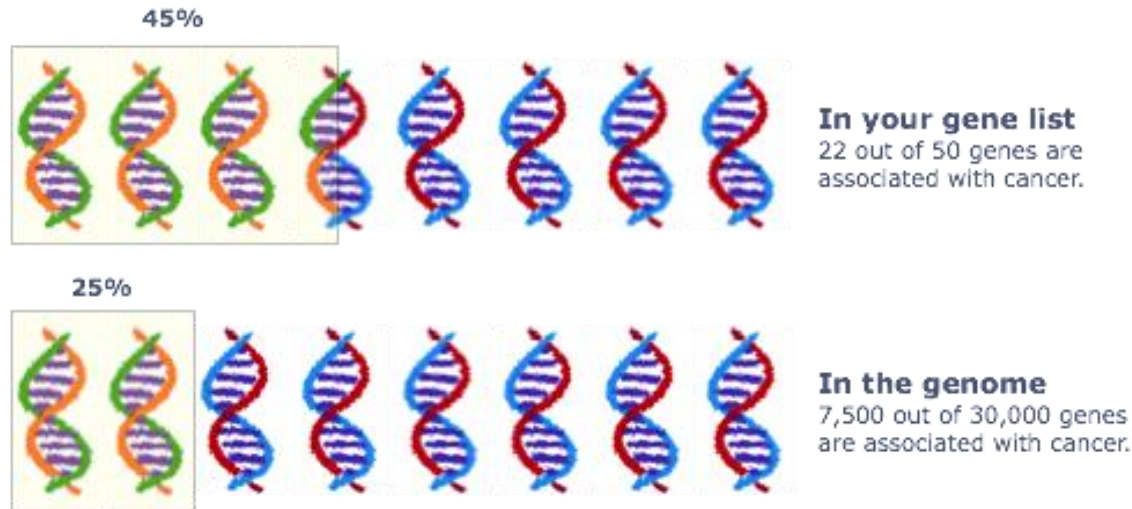
# What is functional enrichment?

- It is a measure of how much a group of gene products is found in our data set

- It requires some type of background measure, as a basis for comparison

- What we look at is how many we have (observed) as opposed to how many we would expect to see at random, given our background.

# Background

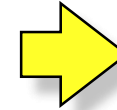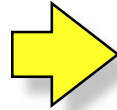The choice of an appropriate background is critical to get meaningful results



45%

In your gene list
22 out of 50 genes are associated with cancer.

25%

In the genome
7,500 out of 30,000 genes are associated with cancer.
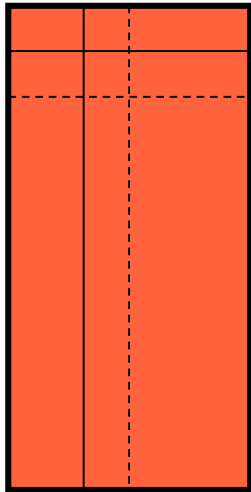
Fold of enrichment = 45% / 25% = 1.8

You should use all the genes detected by the method used in your experiment,

not all the genes in the genom, if possible.

# Enrichment test
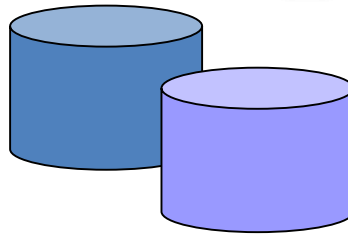
**RNA-seq experiment
(gene expression table)**

**Enrichment Table**

| Spindle | 0.00001 |
| Apoptosis | 0.00025 |

ENRICHMENT
TEST

**Gene-set
Databases
(GO)**

**bio**informatics.ca

# Enrichment test

**RNA-seq experiment
(gene expression table)**

**Gene list
(e.g UP)**

**Gene-set
Databases**

**Background genes
(array genes not significant)**

# Enrichment test

# Enrichment test

# Enrichment test

**Significant genes (e.g UP)**

**Overlap between gene list and gene-set**

**Random samples of array genes**

*Is this overlap larger than expected by **random sampling** the array genes?*

**Background genes (array genes not significant)**
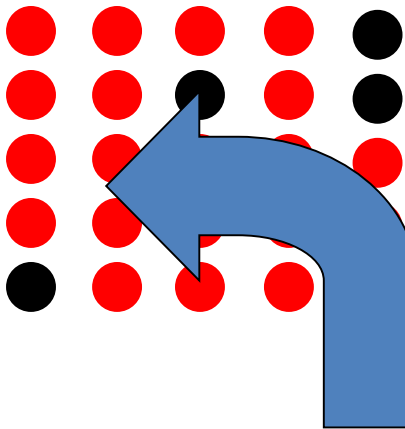
bioinformatics.ca

# Enrichment analysis

- Given:

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

  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?*

- Details:

  – Where do the gene lists come from?

  – How to assess "surprisingly" (statistics)

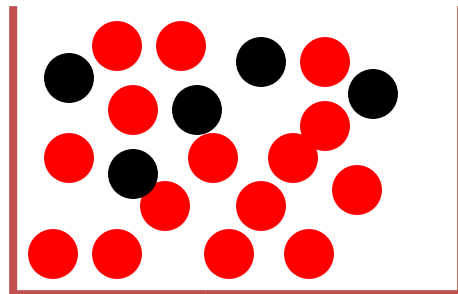  – How to correct for repeating the tests

# Randomization test

Random draws



... 7,834 draws later ...

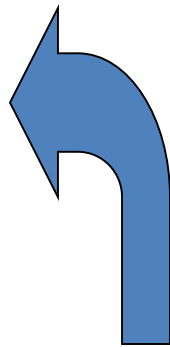*Expect a random draw with observed enrichment once every 1 / P-value draws*

Background population:
500 black genes
4500 red genes

# Fisher's exact test

## a.k.a., the hypergeometric test

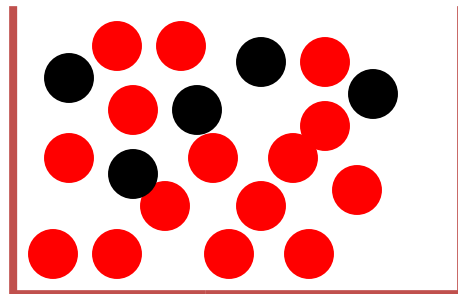Gene list

● RRP6
● MRD1
● RRP7
● RRP43
● RRP42

**Null hypothesis:** List is a random sample from population

**Alternative hypothesis:** More black genes than expected

Background population:
500 black genes
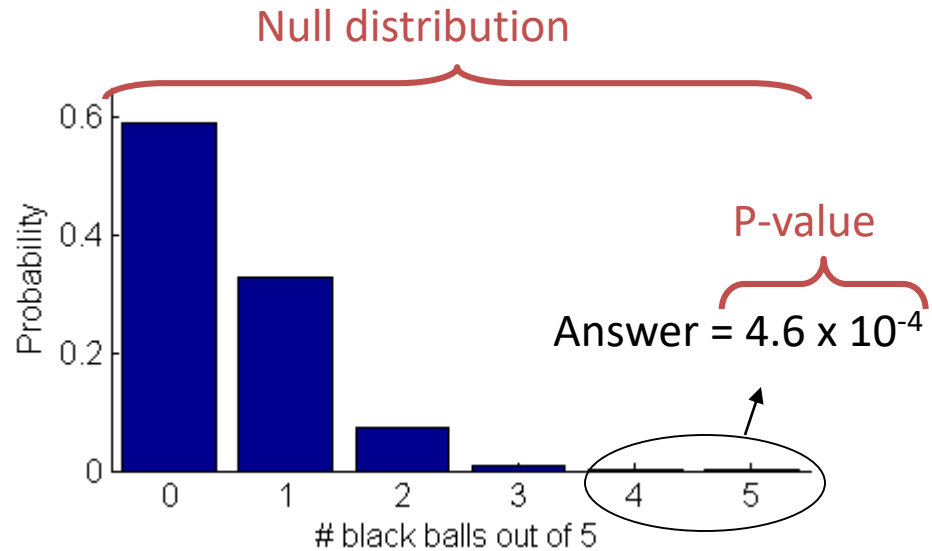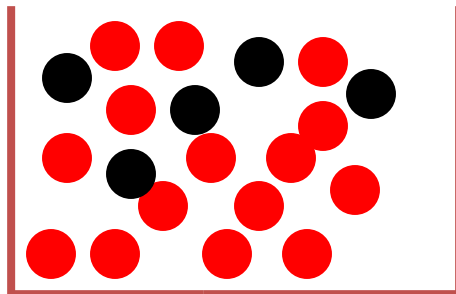4500 red genes

# Fisher's exact test

## a.k.a., the hypergeometric test

Gene list

Null distribution

- ● RRP6
- ● MRD1
- ● RRP7
- ● RRP43
- ● RRP42

P-value

Answer = $4.6 \times 10^{-4}$

Background population:
500 black genes
4500 red genes

# Problems working with large data sets

- The more comparisons we make, the more there is a chance that we will get random hits

- We need to correct for multiple tests, using statistical methods such as Bonferroni, FDR (Benjamini)

- Statistical significance doesn't necessarily mean biological significance

# Beyond Fisher's Exact Test

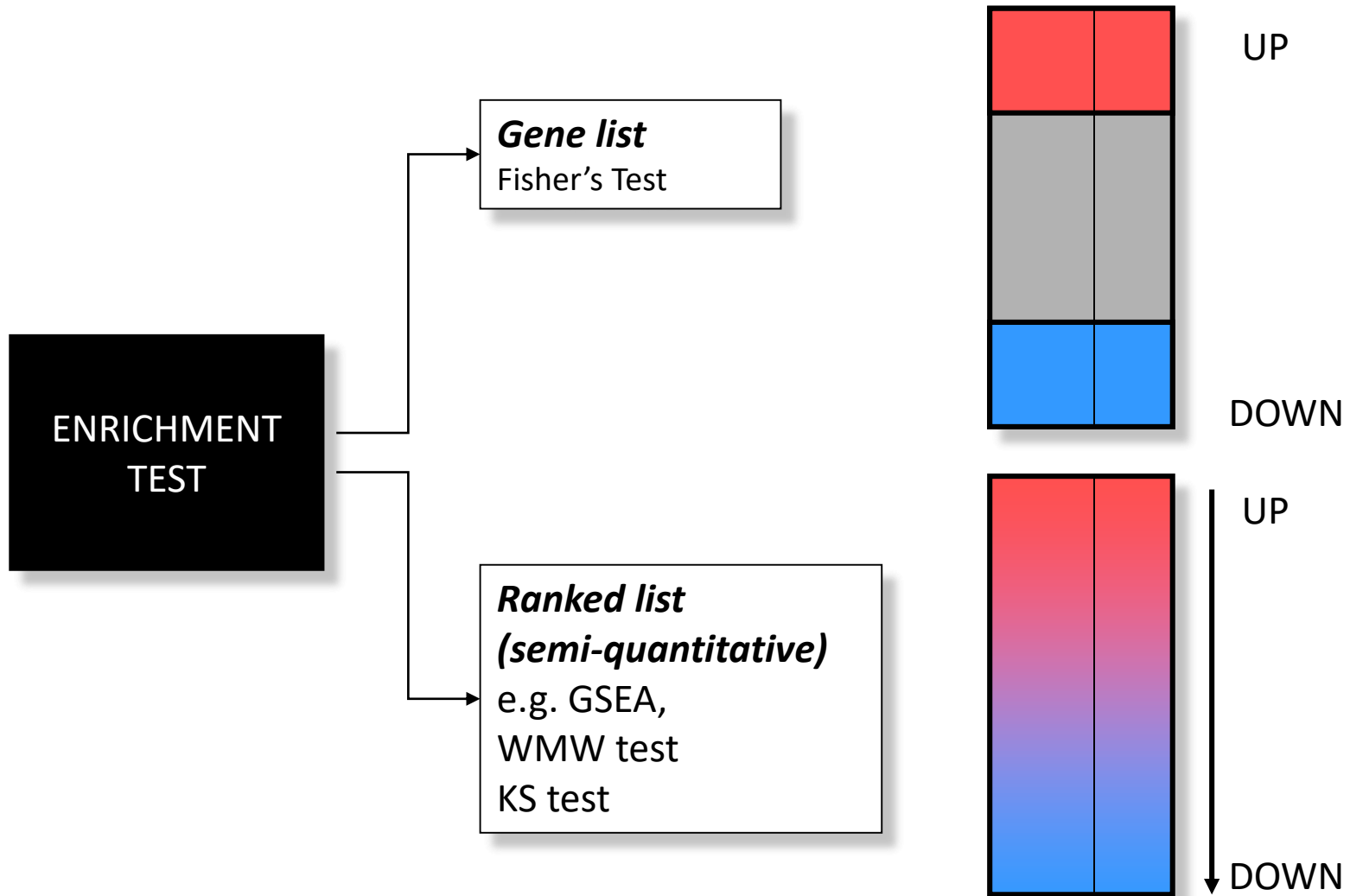Possible problems with Fisher's Exact Test:

- No "natural" value for the threshold
- Different results at different threshold settings
- Possible loss of statistical power due to thresholding
    - No resolution between significant signals with different strengths
    - Weak signals neglected

Solution: enrichment tests based on ranked lists

# Beyond Fisher's Exact Test



**Gene list**
Fisher's Test

**Ranked list (semi-quantitative)**
e.g. GSEA,
WMW test
KS test

ENRICHMENT TEST

UP

DOWN

UP

DOWN

WMW - the Wilcoxon-Mann-Whitney test    KS - the Kolmogorov–Smirnov test

# OUTLINE

- Single gene analysis / information

- Analysis of group of genes

- Gene ontology (GO)

- Enrichment analysis

  - Hypergeometric Test  and Fisher exact test
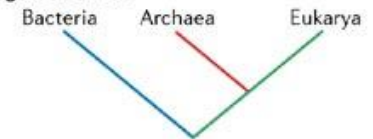
  - GO Independence Assumption
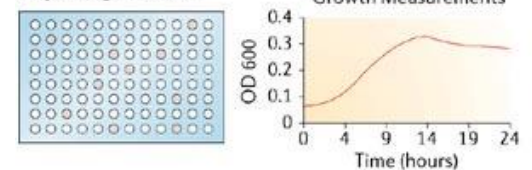


Genome sequence and annotation

Available literature

PubMed
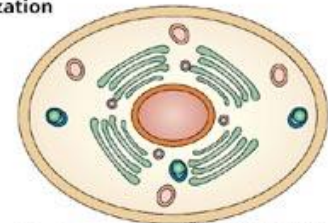
Phylogenetic data

Physiological data

Databases

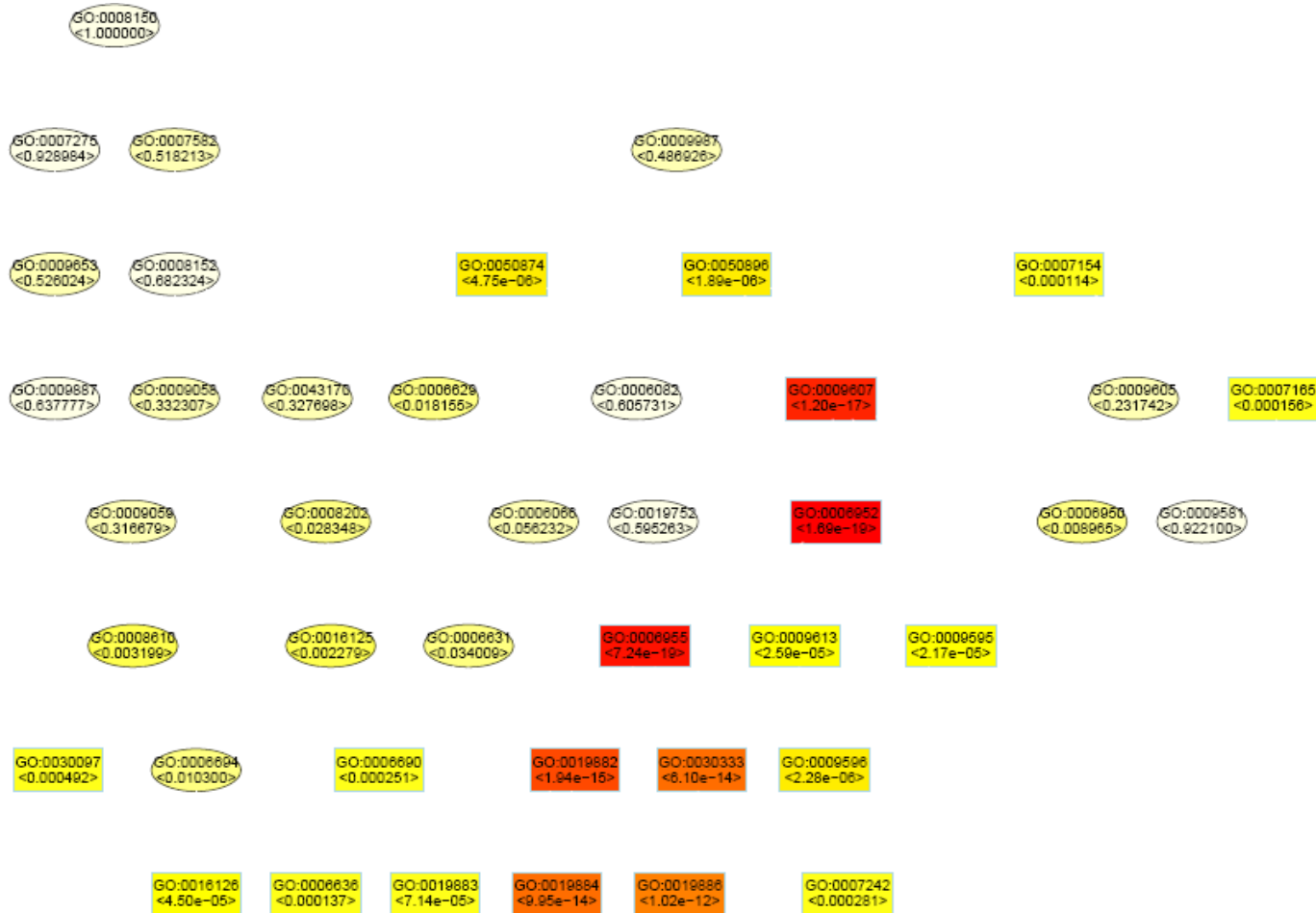KEGG    EcoCyc

Localization

Signal sequences: PLLLLPISGSALP
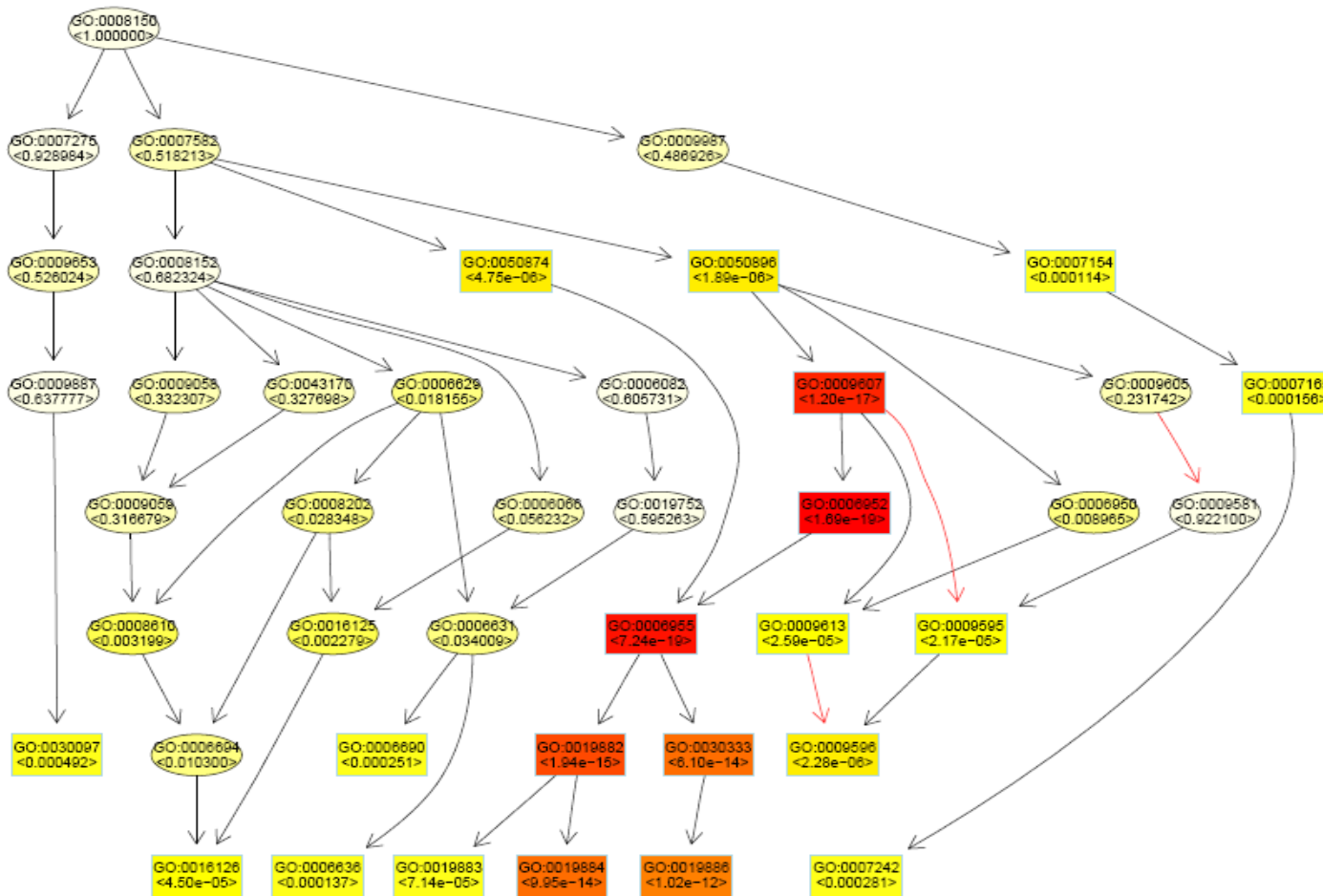
# Term-for-term

- The most common type of analysis

- Each term is considered independently of its neighbors in the GO tree

- Compares observed to expected and calculates significance

# GO Independence Assumption



Note: The coloring of the nodes represent the *relative* significance of the GO terms: dark red is the most significant, light yellow is the least significant from the graph

# GO Independence Assumption



Note: The coloring of the nodes represent the *relative* significance of the GO terms: dark red is the most significant, light yellow is the least significant from the graph

# Algorithms review

➤ **classic algorithm**

- Calculate significance of each GO term independently.
- Adjust pvalues for multiple testing (Bonferroni, FDR, etc.).
- Kolmogorov-Smirnov test can easily be used in this case
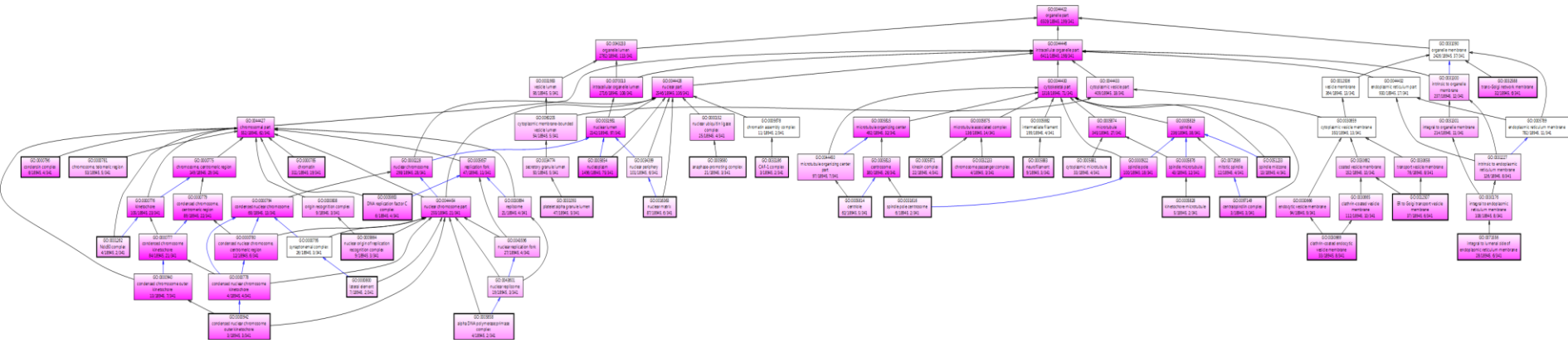
➤ **elim algorithm**

- Nodes are processed bottom-up in the GO graph.
- It iteratively removes the genes annotated to significant GO terms from more general GO terms.
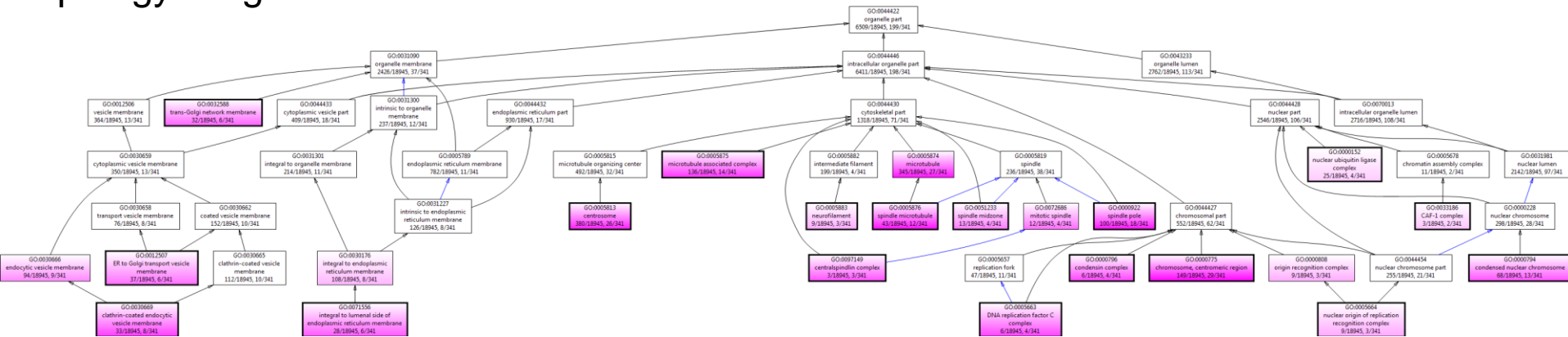- Intuitive and simple to interpret.

➤ **weight algorithm**

- The genes obtain weights that denote the gene relevance in the significant nodes.
- To decide if a GO term $u$ better represents the interesting genes, the enrichment score of node $u$ is compared with the scores of its children.
- Children with a better score than $u$ better represent the interesting genes; their significance is increased
- Children with a lower score than $u$ have their significance reduced.

Alexa A, Rahnenführer J, Lengauer T. Bioinformatics. 2006 Jul 1;22(13):1600-7

# Same input data – different results....

Term for term



Topology weighted

# *Thanks to:*

Dr. Esti Feldmeser & Dr. Shifra Ben Dor for

interchanging and improving slides

for
your attention
Questions?