## BIMM 143

**Pathway Analysis and the Interpretation of Gene Lists**

Lecture 15

**Barry Grant**

UC San Diego

http://thegrantlab.org/bimm143

---



---

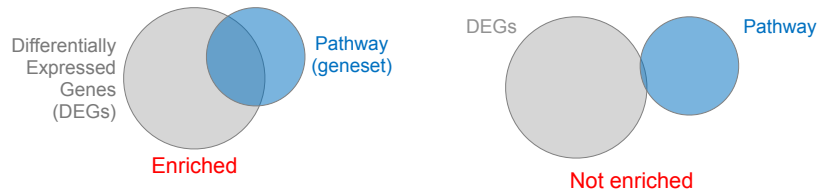My high-throughput experiment generated a long list of genes/proteins…

What do I do now? 🫤

---

Pathway analysis!
(a.k.a. geneset enrichment)

Use bioinformatics methods to help extract biological meaning from such lists…

## Pathway analysis (a.k.a. geneset enrichment)
### Principle



- Variations of the math: overlap, ranking, networks... ➢ *Not critical, different algorithms show similar performances*
- DEGs come from your experiment ➢ *Critical, needs to be as clean as possible*
- Pathway genes ("geneset") come from annotations ➢ *Important, but typically not a competitive advantage*

## Pathway analysis (a.k.a. geneset enrichment)
### Limitations

- **Post-transcriptional regulation** is neglected
- **Directionality** is hard to capture sensibly
  - e.g. IκBα/NF-κB
- **Tissue-specific** variations of pathways are not annotated
  - e.g. NF-κB regulates metabolism, not inflammation, in adipocytes
- **Size bias**: stats are influenced by the size of the pathway
- **Geneset annotation bias**: can only discover what is already known
- **Non-model organisms**: no high-quality genesets available
- Many pathways/receptors **converge** to few regulators
  - e.g. tens of innate immune receptors activate 4 TFs: NF-kB, AP-1, IRF3/7, NFAT

## Starting point for pathway analysis:
### Your gene list

- You have a list of genes/proteins of interest
- You have quantitative data for each gene/protein
  - Fold change
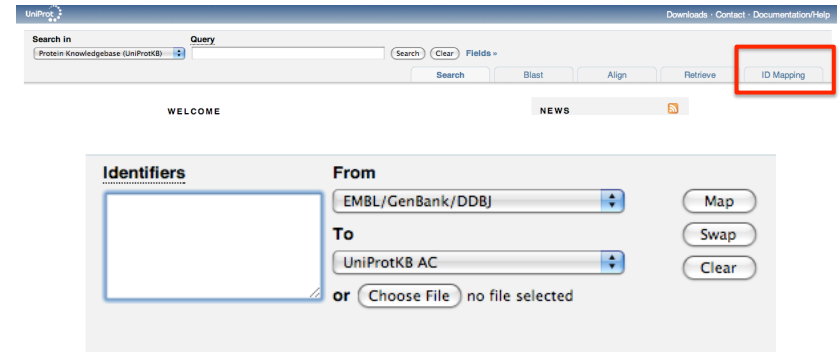  - p-value
  - Spectral counts
  - Presence/absence



## Translating between identifiers

- Many different identifiers exist for genes and proteins, e.g. UniProt, Entrez, etc.
- Sometimes you have to translate one set of ids into another
  - A program might only accept certain types of ids
  - You might have a list of genes with one type of id and info for genes with another type of id

## Translating between identifiers

- Many different identifiers exist for genes and proteins, e.g. UniProt, Entrez, etc.

- Sometimes you have to translate one set of ids into another

  - A program might only accept certain types of ids

  - You might have a list of genes with one type of id and info for genes with another type of id

- **Various web sites translate ids -> *best for small lists***

  - **UniProt < www.uniprot.org>; IDConverter < idconverter.bioinfo.cnio.es >**


## Translating between identifiers: UniProt < www.uniprot.org >




## Translating between identifiers

- Many different identifiers exist for genes and proteins, e.g. UniProt, Entrez, etc.

- Sometimes you have to translate one set of ids into another

  - A program might only accept certain types of ids

  - You might have a list of genes with one type of id and info for genes with another type of id

- Various web sites translate ids -> *best for small lists*

  - UniProt < www.uniprot.org>; IDConverter < idconverter.bioinfo.cnio.es >

- **VLOOKUP in Excel - *good if you are an excel whizz - I am not!***

  - **Download flat file from Entrez, Uniprot, etc; Open in Excel; Find columns that correspond to the 2 IDs you want to convert between; Sort by ID; Use vlookup to translate your list**


## Translating between identifiers: Excel VLOOKUP

VLOOKUP(lookup_value, table_array, col_index_num)
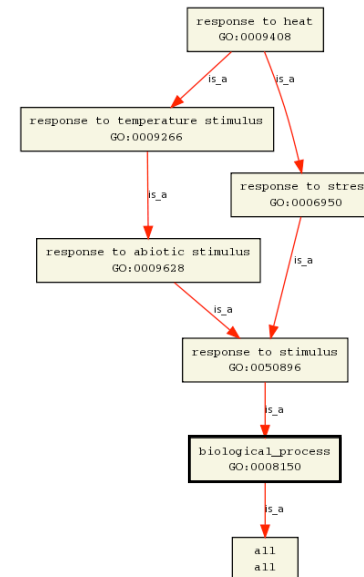
## Translating between identifiers

- Many different identifiers exist for genes and proteins, e.g. UniProt, Entrez, etc.

- Sometimes you have to translate one set of ids into another

  - A program might only accept certain types of ids

  - You might have a list of genes with one type of id and info for genes with another type of id

- Various web sites translate ids -> *best for small lists*

  - UniProt < www.uniprot.org >; IDConverter < idconverter.bioinfo.cnio.es >

- VLOOKUP in Excel -> *good if you are an excel whizz - I am not!*

  - Download flat file from Entrez, Uniprot, etc; Open in Excel; Find columns that correspond to the 2 ids you want to convert between; Use vlookup to translate your list

- Use the **merge()** or **mapIDs()** functions in **R** - <u>fast</u>, *versatile* & <u>reproducible!</u>

  - Also **clusterProfiler::bitr()** function and many others… [Link to clusterProfiler vignette]
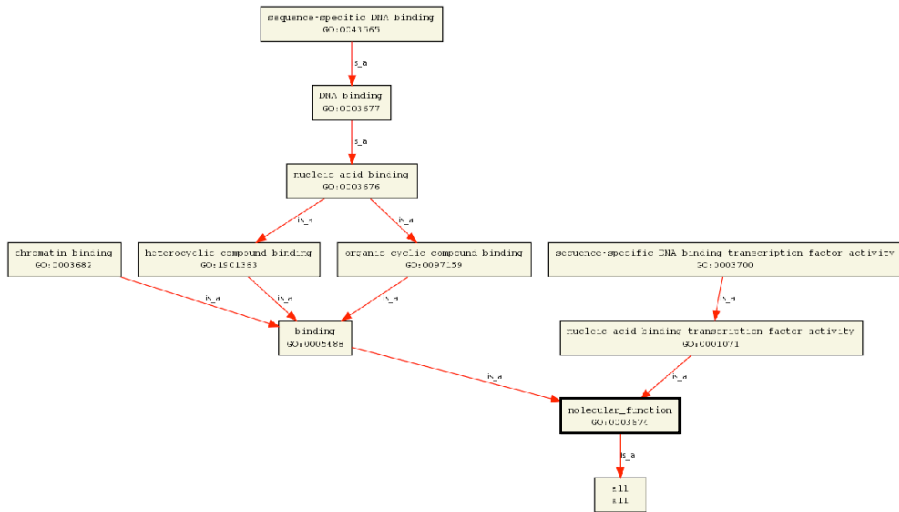
## What functional set databases do you want?

- Commonly used

  - **Gene Ontology (GO)**

  - **KEGG Pathways** (mostly metabolic)

  - **GeneGO MetaBase**

  - **Ingenuity Pathway Analysis (IPA)** INGENUITY SYSTEMS

  - **MSigDB** (gene sets based on chromosomal position, cis-regulatory motifs, GO terms, etc)

- Many others...

  - Enzyme Classification, Pfam families

  - Open Biomedical Ontologies (OBO, www.obofoundry.org)

## GO database < www.geneontology.org >

- **What function does HSF1 perform?**

  - *response to heat; sequence-specific DNA binding; transcription; etc*

- **Ontology** => a structured and controlled vocabulary that allows us to annotate gene products consistently, interpret the relationships among annotations, and can easily be *handled by a computer*

- GO database consists of 3 ontologies that describe gene products in terms of their associated **biological processes**, **cellular components** and **molecular functions**
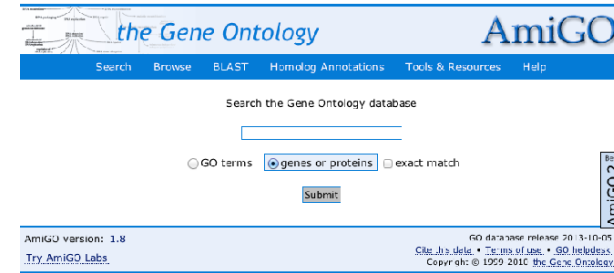


- Terms are nodes

- Relationships are edges

- Parent terms are more general

- Terms can have multiple parents

## GO Annotations

- GO is not a database of genes/proteins or sequences

- Gene products get annotated with GO terms by organism specific databases, such as Flybase, Wormbase, MGI, ZFIN, UniProt, etc
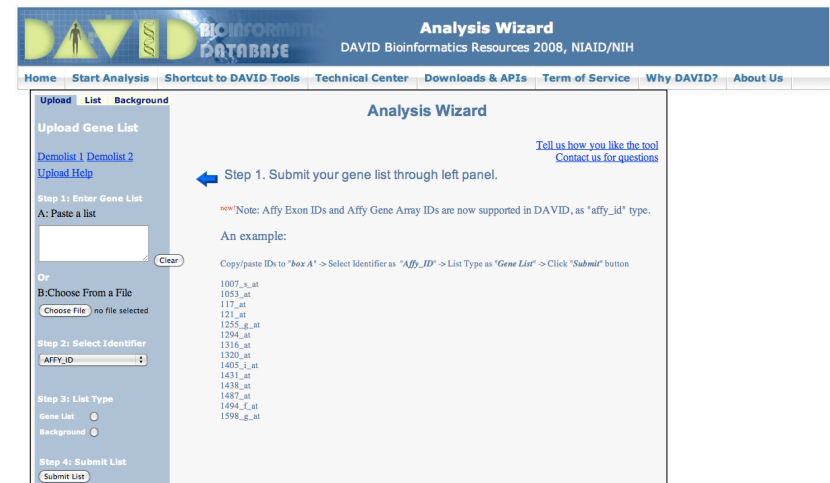
- Annotations are available through AmiGO < amigo.geneontology.org >



## GO evidence codes



| Evidence code | Evidence code description | Source of evidence | Manually checked | Current number of annotations* |
|---|---|---|---|---|
| IDA | Inferred from direct assay | Experimental | Yes | 71,050 |
| IEP | Inferred from expression pattern | Experimental | Yes | 4,598 |
| IGI | Inferred from genetic interaction | Experimental | Yes | 8,311 |
| IMP | Inferred from mutant phenotype | Experimental | Yes | 61,549 |
| IPI | Inferred from physical interaction | Experimental | Yes | 17,043 |
| ISS | Inferred from sequence or structural similarity | Computational | Yes | 196,643 |
| RCA | Inferred from reviewed computational analysis | Computational | Yes | 103,792 |
| IGC | Inferred from genomic context | Computational | Yes | 4 |
| IEA | Inferred from electronic annotation | Computational | No | 15,687,382 |
| IC | Inferred by curator | Indirectly derived from experimental or computational evidence made by a curator | Yes | 5,167 |
| TAS | Traceable author statement | Indirectly derived from experimental or computational evidence made by the author of the published article | Yes | 44,564 |
| NAS | Non-traceable author statement | No 'source of evidence' statement given | Yes | 25,656 |
| ND | No biological data available | No information available | Yes | 132,192 |
| NR | Not recorded | Unknown | Yes | 1,185 |

*October 2007 release

**Use and misuse of the gene ontology annotations**
Seung Yon Rhee, Valerie Wood, Kara Dolinski & Sorin Draghici
*Nature Reviews Genetics* **9**, 509-515 (2008)

## DAVID at NIAID < david.abcc.ncifcrf.gov >

# DAVID

- Notice that you can pick a *Background* (Universe)

**Analysis Wizard**

Tell us how you like the tool
Contact us for questions

Step 1. Successfully submitted gene list
Current Gene List: Uploaded List_2
Current Background: HOMO SAPIENS

Step 2. Analyze above gene list with one of DAVID tools

Which DAVID tools to use?

Functional Annotation Tool
- Functional Annotation Clustering
- Functional Annotation Chart
- Functional Annotation Table

Gene Functional Classification Tool
Gene ID Conversion Tool
Gene Name Batch Viewer

**Gene List Manager**

Select to limit annotations by one or more species  Help

- Use All Species -
HOMO SAPIENS(4402)
SYNTHETIC CONSTRUCT(5)

Select

**List Manager**  Help

Uploaded List_2

Select List to:
Use  Rename
Remove  Combine

Show Gene List new!

Upload  List  Background

---

# DAVID

- *Functional Annotation Tool*

**Annotation Summary Results**

Help and Tool Manual

Current Gene List: Uploaded List_3   2320 DAVID IDs
Current Background: HOMO SAPIENS   Check Defaults ☑  Clear All

- Main Accessions (0 selected)
- Other Accessions (0 selected)
- Gene Ontology (4 selected)
- Protein Domains (3 selected)
- Pathways (3 selected)
- General Annotations (0 selected)
- Functional Categories (3 selected)
- Protein Interactions (0 selected)
- Literature (0 selected)
- Disease (1 selected)
- Tissue Expression

**Combined View for Selected Annotation**

Functional Annotation Clustering new!  ←
Functional Annotation Chart
Functional Annotation Table

---

# DAVID

- Specify functional sets

**Annotation Summary Results**

Help and Tool Manual

Current Gene List: Uploaded List_3   2320 DAVID IDs
Current Background: HOMO SAPIENS   Check Defaults ☐  Clear All

- Main Accessions (0 selected)
- Other Accessions (0 selected)
- Gene Ontology (1 selected)

| | | | | |
|---|---|---|---|---|
| ☐ GOTERM_BP_1 | 71% | 1669 | Chart | |
| ☐ GOTERM_BP_2 | 71% | 1652 | Chart | |
| ☐ GOTERM_BP_3 | 69% | 1609 | Chart | |
| ☐ GOTERM_BP_4 | 65% | 1519 | Chart | |
| ☑ GOTERM_BP_5 | 61% | 1432 | Chart | |
| ☐ GOTERM_BP_ALL | 71% | 1669 | Chart | |
| ☐ GOTERM_CC_1 | 75% | 1754 | Chart | |
| ☐ GOTERM_CC_2 | 75% | 1741 | Chart | |
| ☐ GOTERM_CC_3 | 75% | 1741 | Chart | |
| ☐ GOTERM_CC_4 | 70% | 1634 | Chart | |
| ☐ GOTERM_CC_5 | 69% | 1605 | Chart | |
| ☐ GOTERM_CC_ALL | 75% | 1754 | Chart | |
| ☐ GOTERM_MF_1 | 75% | 1745 | Chart | |
| ☐ GOTERM_MF_2 | 73% | 1716 | Chart | |
| ☐ GOTERM_MF_3 | 65% | 1523 | Chart | |

---

# DAVID

- Let's look at the *Functional Annotation Chart*

**Annotation Summary Results**

Help and Tool Manual

Current Gene List: Uploaded List_3   2320 DAVID IDs
Current Background: HOMO SAPIENS   Check Defaults ☑  Clear All

- Main Accessions (0 selected)
- Other Accessions (0 selected)
- Gene Ontology (4 selected)
- Protein Domains (3 selected)
- Pathways (3 selected)
- General Annotations (0 selected)
- Functional Categories (3 selected)
- Protein Interactions (0 selected)
- Literature (0 selected)
- Disease (1 selected)
- Tissue Expression

**Combined View for Selected Annotation**

Functional Annotation Clustering new!  ←
Functional Annotation Chart
Functional Annotation Table

## DAVID

- *Functional Annotation Chart*



**Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources**
Da Wei Huang, Brad T Sherman & Richard A Lempicki
*Nature Protocols* **4**, 44 - 57 (2009)

## GSEA < www.broadinstitute.org/gsea >

- Download GSEA desktop application



- Excellent tutorial, user's guide and example datasets to work through
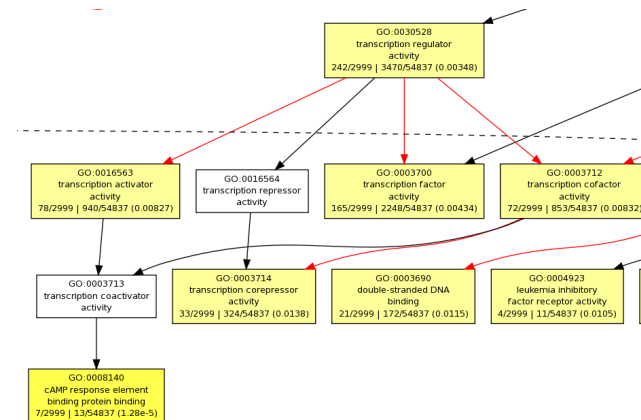
Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles
Aravind Subramanian, Pablo Tamayo, Vamsi K. Mootha, Sayan Mukherjee, Benjamin L. Ebert, Michael A. Gillette, ...
*PNAS* **102**, 15545-15550 (2005)

## Overlapping functional sets

- Many functional sets overlap, in particular those from databases that are hierarchical in nature (e.g. GO)

- Hierarchy enables:

  - Annotation flexibility (e.g. allow different degrees of annotation completeness based on what is known)

  - Computational methods to "understand" function relationships (e.g. ATPase function is a subset of enzyme function)

- Unfortunately, this also makes functional profiling trickier

## GOEast < omicslab.genetics.ac.cn/GOEAST >

- Graphical view of enriched GO terms and their relationships

## GO SLIMs

- Cut-down versions of the GO ontologies containing a subset of the terms in the whole GO

- GO FAT (DAVID):

  - filters out very broad GO terms based on a measured specificity of each term

## DAVID Functional Annotation Clustering

- Based on shared genes between functional sets



## Want more?

- **GeneGO** < portal.genego.com >

  - MD/PhD curated annotations, great for certain domains (eg, Cystic Fibrosis)

  - Nice network analysis tools

  - Email us for access

- **Oncomine** < www.oncomine.org >

  - Extensive cancer related expression datasets

  - Nice concept analysis tools

  - Research edition is free for academics, Premium edition $$$

- **Lots of other Bioconductor packages in this area!**

# Hands-on time!

https://bioboot.github.io/bimm143_W18/lectures/#15

Do it Yourself!

## Advice:
## Figure out "**What do I want to do with my list?**"

- Organize/summarize data for presentation or manuscript
    - DAVID: GO_FAT -> Functional Annotation Clustering -> Pick threshold

- Infer biological processes from the list
    - DAVID: Functional Annotation Chart -> explore functional databases and see which make sense
    - GSEA: Select MSigDB sets of interest -> e.g., immunologic signatures
    - Use domain specific database it at all possible!

- Find "missing" genes/proteins not detected by experiment
    - ConceptGen: Gene-gene enrichment

---

## Pathway analysis (a.k.a. geneset enrichment)
## **Principle**



- Variations of the math: overlap, ranking, networks... ➢ *Not critical, different algorithms show similar performances*
- DEGs come from your experiment ➢ *Critical, needs to be as clean as possible*
- Pathway genes ("geneset") come from annotations ➢ *Important, but typically not a competitive advantage*

---

## Pathway analysis (a.k.a. geneset enrichment)
## **Limitations**

- **Post-transcriptional regulation** is neglected

- **Directionality** is hard to capture sensibly
    - e.g. IκBα/NF-κB

- **Tissue-specific** variations of pathways are not annotated
    - e.g. NF-κB regulates metabolism, not inflammation, in adipocytes

- **Size bias**: stats are influenced by the size of the pathway

- **Geneset annotation bias**: can only discover what is already known

- **Non-model organisms**: no high-quality genesets available

- Many pathways/receptors **converge** to few regulators
    - e.g. tens of innate immune receptors activate 4 TFs: NF-kB, AP-1, IRF3/7, NFAT