



Practical Course on "Gene/Protein Functional Networks & Interactomes"

23 and 24 November 2015 - UCT, Cape Town, South Africa (thanks to Prof. Nicola Mulder)









DAY 1

Session 1 (9:30 - 12:30, 3h)

Bioinformatic tools for Functional Enrichment Analysis (FEA)

Session 2 (13:30 - 16:30, 3h)

Construction of gene functional networks

DAY2

Session 3 (9:30 - 12:30, 3h)

Protein interaction networks

Session 4 (13:30 - 16:30, 3h)

Construction and analysis of gene/protein networks

Dr. Javier De Las Rivas

Cancer Research Center (CiC-IBMCC, CSIC/USAL), Salamanca, Spain







Session 1 (9:30 - 12:30, 3h)
Bioinformatic tools for Functional Enrichment Analysis (FEA)
Session 2 (13:30 - 16:30, 3h)
Construction of gene functional networks

- Introduction to biological information and annotation spaces: GO, KEGG, Interpro
- Functional Enrichment Analysis (EA): from single to modular methods
 - Using EA tools to annotate gene lists:
 DAVID (single), GSEA (gene sets), GeneCodis (modular)
 - Sort out problems after EA: post-enrichment tool GeneTermLinker (postEA)
 - From co-annotation and enrichment to functional networks: networks construction using a R tool



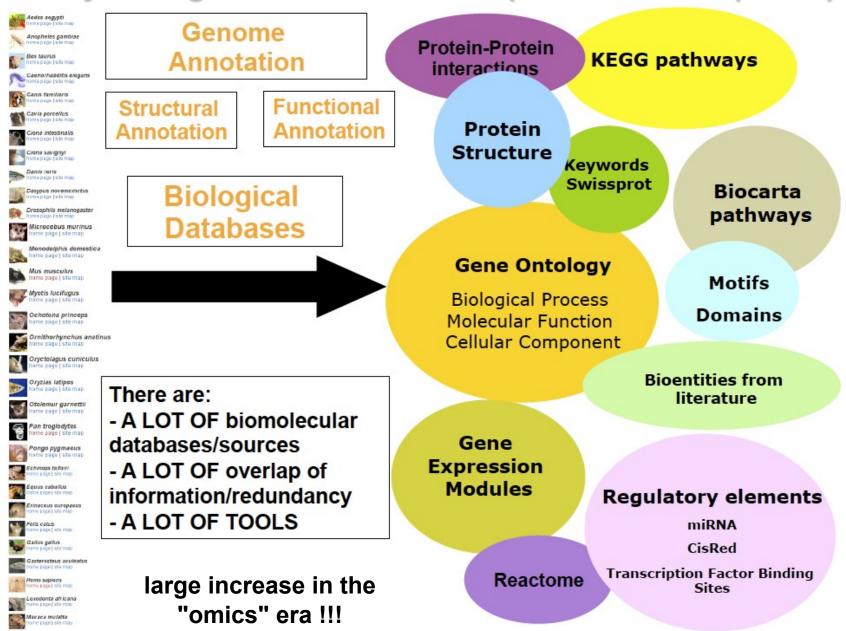




Session 1 (9:30 - 12:30, 3h)
Bioinformatic tools for Functional Enrichment Analysis (FEA)
Session 2 (13:30 - 16:30, 3h)
Construction of gene functional networks

- Introduction to biological information and annotation spaces:
 GO, KEGG, Interpro
- Functional Enrichment Analysis (EA): from single to modular methods
 - Using EA tools to annotate gene lists:
 DAVID (single), GSEA (gene sets), GeneCodis (modular)
 - Sort out problems after EA: post-enrichment tool GeneTermLinker (postEA)
 - From co-annotation and enrichment to functional networks: networks construction using a R tool

Many biological data resources (= annotation spaces)



3 orthogonal annotation spaces

Biological terms and definitions

Gene Ontology (GO) http://www.geneontology.org/

Biomolecular pathways and reactions

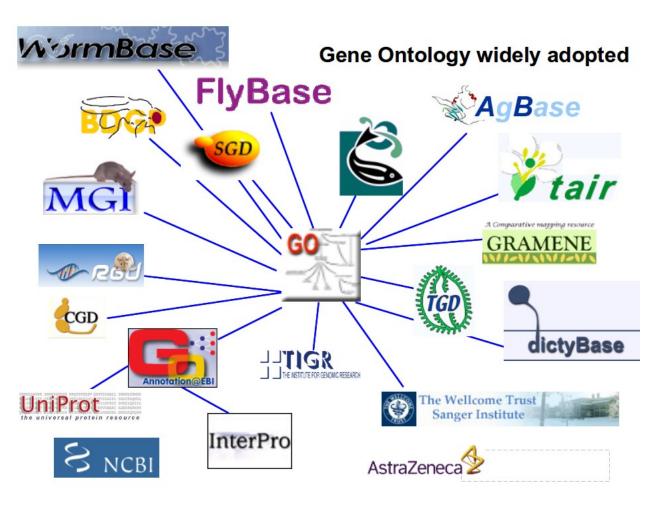
KEGG (pathways) http://www.genome.jp/kegg/

Biomolecular sequences and structure (domains, motifs)

InterPro (sequences & structure) http://www.ebi.ac.uk/interpro/



Gene Ontology (GO) http://www.geneontology.org/



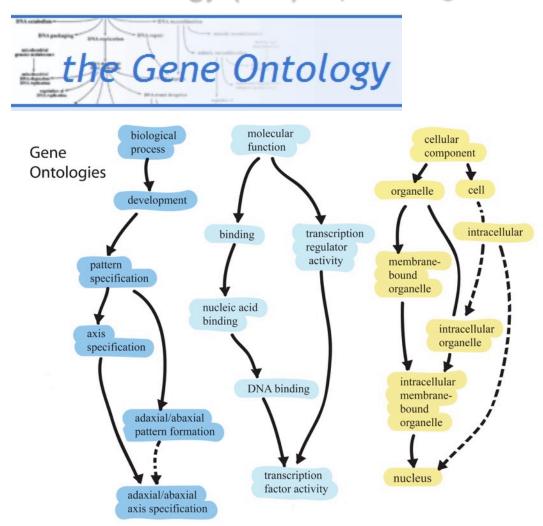
The Gene Ontology project is a major bioinformatics initiative with the aim of standardizing the representation of gene and gene product attributes across species and databases.

The project provides a controlled vocabulary of biological terms for describing gene product characteristics and gene product annotation data from GO Consortium members, as well as tools to access and process this data.

7



Gene Ontology (GO) http://www.geneontology.org/



Modified from Clark et al. (2005)

GO is organized in 3 independent hierarchies:

Biological Process Ontology (BP)

Biological and cellular processes where a given gene or gene product is involved.

Molecular Function Ontology (MF)

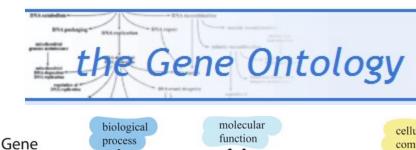
Molecular functions and activities that a given gene or gene product has.

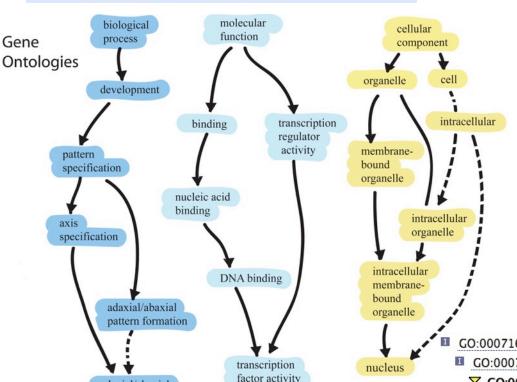
Cellular Component Ontology (CC)

Place or part of the cell where a given gene or gene product works most of its time.



Gene Ontology (GO) http://www.geneontology.org/





GO is organized in 3 independent hierarchies:

 Biological Process Ontology (BP)

Example, search for: "NOTCH" or "NOTCH signaling"

Find: GO:0007219 Notch signaling pathway

GO includes ancestors & children terms

GO:0007165 signal transduction [43440 gene products]

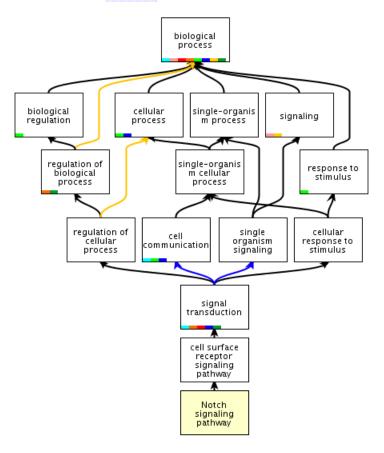
- GO:0007166 cell surface receptor signaling pathway [20166 gene products]
 - GO:0007219 Notch signaling pathway [870 gene products]
 - GO:0045746 negative regulation of Notch signaling pathway [154 gene products]

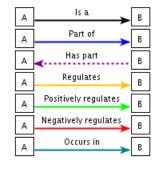
adaxial/abaxial axis specification

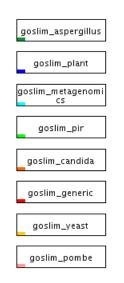


Gene Ontology (GO) http://www.geneontology.org/

View this term in QuickGO.







GO is organized in 3 independent hierarchies:

 Biological Process Ontology (BP)

Example, search for: "NOTCH" or "NOTCH signaling"

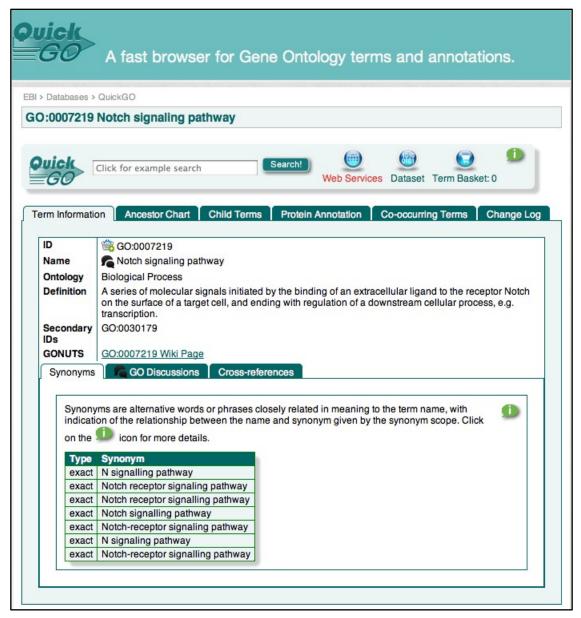
Find:

GO:0007219 Notch signaling pathway

GO includes ancestors & children terms

2187 gene productsassigned in all species (Oct.2014)1661 gene productsassigned in mammalia (Oct.2014)





http://www.ebi.ac.uk/QuickGO/

GO is organized in 3 independent hierarchies:

 Biological Process Ontology (BP)

Example, search for: "NOTCH" or "NOTCH signaling"

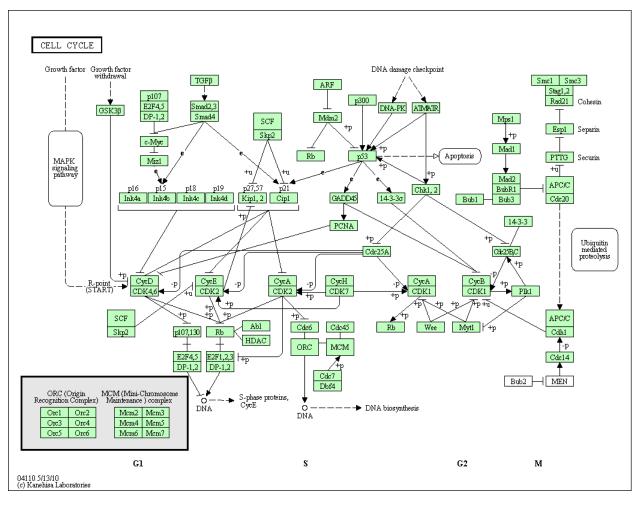
Find:

GO:0007219 Notch signaling pathway

GO includes ancestors & children terms

2187 gene productsassigned in all species (Oct.2014)1661 gene productsassigned in mammalia (Oct.2014)

KEGG (pathways) http://www.genome.jp/kegg/

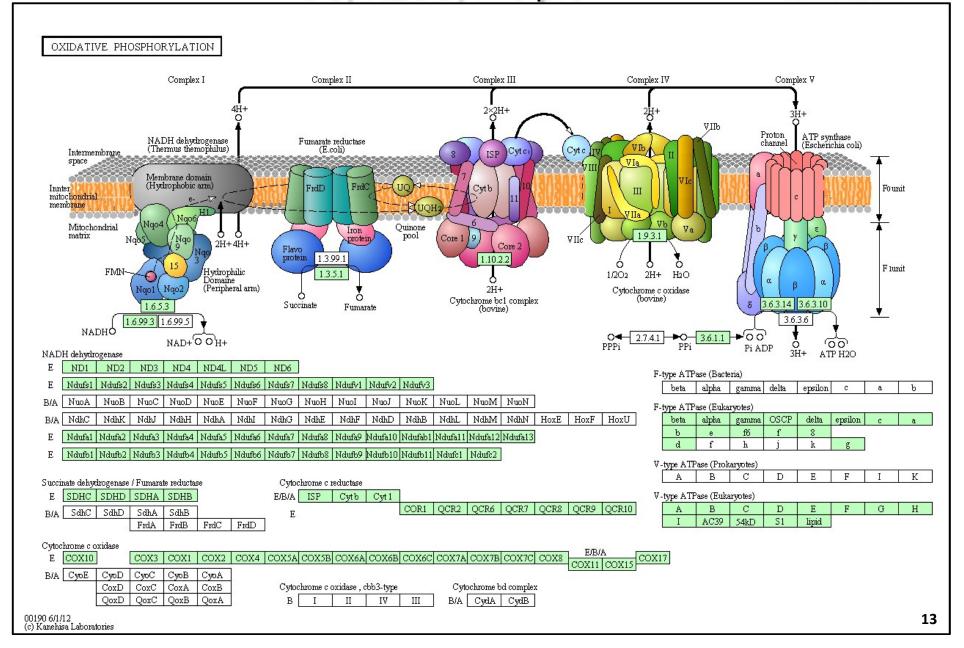




KEGG PATHWAY is a collection of manually drawn pathway maps representing our knowledge on the molecular interaction and reaction networks

- 1. Metabolism
- 2. Genetic Information Processing
- 3. Environmental Informat. Processing
- 4. Cellular Processes
- 5. Organismal Systems
- 6. Human Diseases
- 7. Drug Development



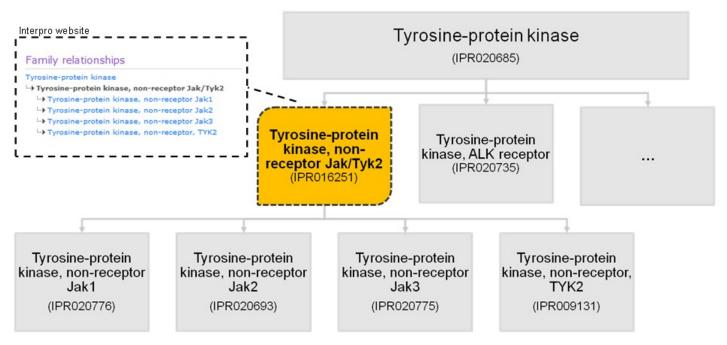




InterPro (sequences & structure) http://www.ebi.ac.uk/interpro/

InterPro provides **sequence/structure analysis** and classification of proteins into families using found protein **motifs**, protein **domains** and important **sites**. They combine protein signatures from a number of member databases into a single searchable resource to produce a integrated database and diagnostic tool.







InterPro (sequences & structure) http://www.ebi.ac.uk/interpro/

InterPro provides **sequence/structure analysis** and classification of proteins into families using found protein **motifs**, protein **domains** and important **sites**. They combine protein signatures from a number of member databases into a single searchable resource to produce a integrated database and diagnostic tool.



Search for NOTCH1 sequence: **NOTC1_HUMAN** (2555 aa)

Protein family membership In a Notch (IPR008297) In a Neurogenic locus Notch 1 (IPR022362) Domains and repeats In a Domain in Pomain in Repeat South (IPR008297) Domain in Pomain in Repeat South (IPR008297) Domain in Pomain in Repeat South (IPR008297) Domain in Pomain in Repeat







Hands-on: Practical Examples

Explore web sites: GO, KEGG, InterPro (& Pfam)

(NOTCH and NOTCH signaling)



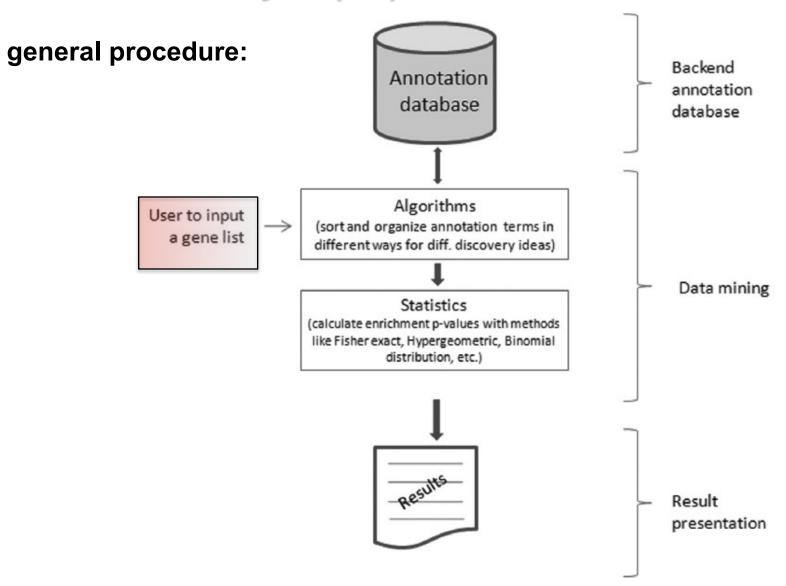




Session 1 (9:30 - 12:30, 3h)
Bioinformatic tools for Functional Enrichment Analysis (FEA)
Session 2 (13:30 - 16:30, 3h)
Construction of gene functional networks

- Introduction to biological information and annotation spaces: GO, KEGG, Interpro
- Functional Enrichment Analysis (EA): from single to modular methods
 - Using EA tools to annotate gene lists:
 DAVID (single), GSEA (gene sets), GeneCodis (modular)
 - Sort out problems after EA: post-enrichment tool GeneTermLinker (postEA)
 - From co-annotation and enrichment to functional networks: networks construction using a R tool







Enrichment tool name	Year release			Category	
FunSpec	2002	g:Profiler	2007	Hypergeometric	Class I
Onto-express	2002	ProbCD	2007	Yule's Q; Goodman-Kruskal's gamma; Cramer's T	Class I
EASE	2003	GOEAST	2008	Hypergeometric	Class I
FatiGO/FatiWise/FatiGO +	2003	GOHyperGAll	2008	Hypergeometric	Class I
FuncAssociate	2003	CatMap	2004	Permutations	Class II
GARBAN	2003	Godist	2004	Kolmogorov–Smirnov test	Class II
GeneMerge	2003	GO-Mapper	2004	Gaussian distribution; EQ-score	Class II
GoMiner	2003	iGA	2004	Permutations; hypergeometric; t-test; Z-score	Class II
MAPPFinder	2003	GSEA	2005	Kolmogorov–Smirnov-like statistic	Class II
CLENCH	2004	MEGO	2005	Z-score	Class II
GO::TermFinder	2004	PAGE	2005	Z-score	Class II
GOAL	2004	T-profiler	2005	t-Test	Class II
GOArray	2004	FuncCluster	2006	Fisher's exact	Class II
GOStat	2004	FatiScan	2007	Fisher's Exact	Class II
GoSurfer	2004	FINA	2007	Fisher's exact	Class II
OntologyTraverser	2004	GAzer	2007	Z-statistics; permutation	Class II
THEA	2004	GeneTrail	2007	Hypergeometric; Kolmogorov-Smirnov	Class II
BiNGO	2005	MetaGP	2007	Z-score	Class II
FACT	2005	Ontologizer	2004	Fisher's exact	Class III
gfinder	2005	POSOC	2004	POSET (a discrete math: finite partially ordered set)	Class III
Gobar	2005	topGO	2006	Fisher's exact	Class III
GOCluster	2005	GO-2D	2007	Hypergeometric; binomial	Class III
GOSSIP	2005	GENECODIS	2007	Hypergeometric; chi-square	Class III
L2L	2005	GOSim	2007	Resnik's similarity	Class III
WebGestalt	2005	PalS	2008	Percent	Class III
BayGO	2006	ProfCom	2008	Greedy heuristics	Class III
eGOn/GeneTools	2006	GOTM	2004	Hypergeometric	Class I,II
Gene Class Expression	2006	ermineJ	2005	Permutations; Wilcoxon rank-sum test	Class I,II
GOALIE	2006	DAVID	2003	Fisher's Exact (modified as EASE score)	Class I,III
and the second s	0.000	GOToolBox	2004	Hypergeometric; Fisher's exact; Binomial	Class I,III
		ADGO	2006	Z-statistic	Class II,III
Huang et al. (2009) NA	R	FunNet	2008	Unclear	Unclear

Functional	
Enrichment Analysis	(EA)

3 major types:

singular EA = SEA

selected gene list used to query different annotation spaces (one by one)

gene set EA = GSEA

not a list, but the entire genes ranked used to query different annotation spaces (one by one)

modular EA = MEA

selected gene list used to query multiple annotation spaces (at once)

Tool category Description

Class I: singular enrichment analysis (SEA)

Enrichment P-value is calculated on each term from the pre-selected interesting gene list. Then, enriched terms are listed in a simple linear text format. This strategy is the most traditional algorithm. It is still dominantly used by most of the enrichment analysis tools.

Class II: gene set enrichment analysis (GSEA) Entire genes (without pre-selection) and associated experimental values are considered in the enrichment analysis. The unique features of this strategy are: (i) No need to pre-select interesting genes, as opposed to Classes I and II; (ii) Experimental values integrated into *P*-value calculation.

Class III: modular enrichment analysis (MEA) This strategy inherits key spirit of SEA. However, the term—term/gene—gene relationships are considered into enrichment *P*-value calculation. The advantage of this strategy is that term—term/gene—gene relationship might contain unique biological meaning that is not held by a single term or gene. Such network/modular analysis is closer to the nature of biological data structure.

Huang et al. (2009) NAR





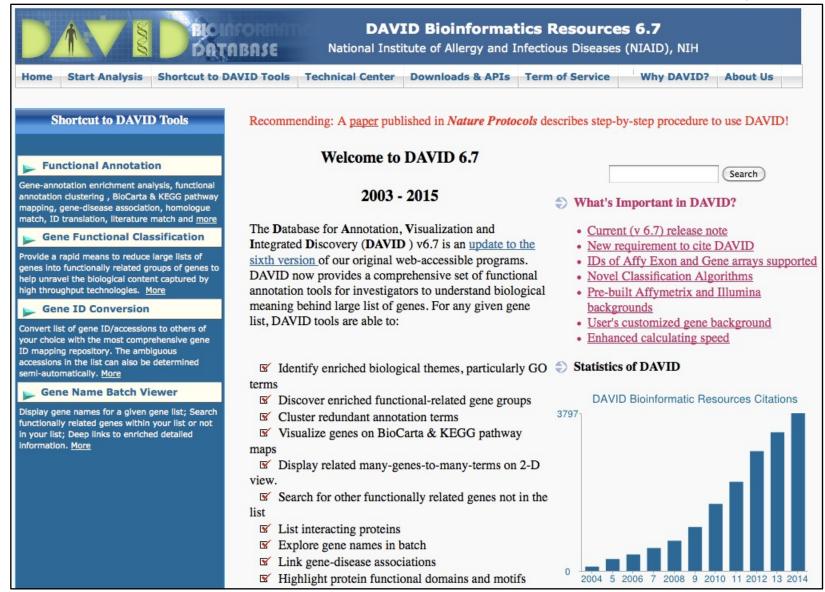


Session 1 (9:30 - 12:30, 3h)
Bioinformatic tools for Functional Enrichment Analysis (FEA)
Session 2 (13:30 - 16:30, 3h)
Construction of gene functional networks

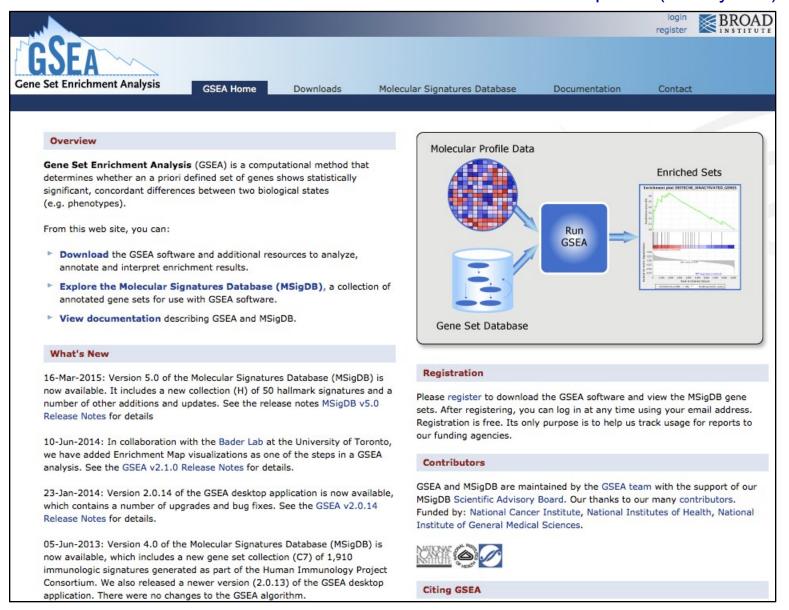
- Introduction to biological information and annotation spaces: GO, KEGG, Interpro
- Functional Enrichment Analysis (EA): from single to modular methods
 - Using EA tools to annotate gene lists:
 DAVID (single), GSEA (gene sets), GeneCodis (modular)
 - Sort out problems after EA: post-enrichment tool GeneTermLinker (postEA)
 - From co-annotation and enrichment to functional networks: networks construction using a R tool

Type 1: singular EA = SEA (DAVID)

selected gene list used to query different annotation spaces (one by one)



Type 2: gene set EA = GSEA (GSEA) entire genes ranked used to query different annotation spaces (one by one)



Type 3: modular EA = MEA (GeneCodis)

selected gene list used to query multiple annotation spaces (at once)

Analysis | Comparative Analysis | Web Services | Help | Release info | Other GeneCodis sites



Gene annotations co-ocurrence discovery

Modular and Singular Enrichment Analysis [?]

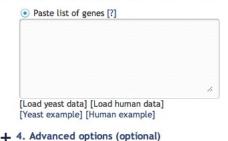
1. Select the organism [?]

Please select

2. Select the annotations [?] [Last update on Dec, 2011]

-Annotations-	
GO Biological Process	
GO Molecular Function	
GO Cellular Component	
GOSlim Process	
GOSlim Function	
GOSlim Component	
KEGG Pathways	
InterPro Motifs	
MicroRNA	
Omim Diseases	
Panther Pathways	
PharmGKB Drugs	
Pubmed	
Transcription Factors	

3. Paste your lists of genes [see allowed IDs]



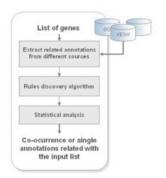
Or upload a file with the list of genes [?]

Choose File no file selected

gene annotations co-ocurrence discovery

What is GeneCodis

GeneCodis is a grid-based tool that integrates different sources of biological information to search for biological features (annotations) that frequently cooccur in a set of genes and rank them by statistical significance. It can be used to determine biological annotations or combinations of annotations that are significantly associated to a list of genes under study with respect to a reference list. GeneCodis executes the following workflow:



Why use GeneCodis?

Home | Analysis | Web Services | Help | Relase info | Other GeneCodis sites

Most of the currently available tools are designed to evaluate single annotations and they provide a list of annotations with their corresponding p-values without taking account the potential relationships among them. Finding relationships among annotations based on cooccurrence patterns can extend the understanding of the biological events associated to a given experimental

What is new in this version?

GeneCodis 2.0 includes a new and faster rule discovery algorithm based on well-known apriori algorithm. This new algorithm let reduce the minimum support needed to expand the range of results available. The whole system is running in a multi-grid environment that allow us to handle simultaneous queries from different users in less time. Moreover, GeneCodis 2.0 includes new funcionalities like searching for simple annotations related to the list of genes, visualizations, new annotations and more gene identifiers are supported.



This is the second version of Genecodis. If you want to use the old version, you can find it here

+ 5. Use your own annotations

3 major types:

singular EA = SEA (DAVID ncbi)
selected gene list used to query different
annotation spaces (one by one)

gene set EA = GSEA (GSEA Broad)
entire genes ranked used to query different
annotation spaces (one by one)

modular EA = MEA (GeneCodis CNB)
selected gene list used to query
multiple annotation spaces (at once)

Description Tool category Enrichment P-value is calculated Class I: singular on each term from the pre-selected enrichment interesting gene list. Then, analysis enriched terms are listed in a (SEA) simple linear text format. This strategy is the most traditional algorithm. It is still dominantly used by most of the enrichment analysis tools. Class II: Entire genes (without pre-selection) and associated experimental gene set enrichment values are considered in the analysis enrichment analysis. The unique (GSEA) features of this strategy are: (i) No need to pre-select interesting genes, as opposed to Classes I and II; (ii) Experimental values integrated into P-value calculation. Class III: This strategy inherits key spirit of modular SEA. However, the term-term/ enrichment gene-gene relationships are considered into enrichment P-value analysis calculation. The advantage of this (MEA) strategy is that term-term/genegene relationship might contain unique biological meaning that is not held by a single term or gene. Such network/modular analysis is closer to the nature of biological

data structure.

Huang et al. (2009) NAR

3 major types:

singular EA = SEA

selected gene list used to query different annotation spaces (one by one)

gene set EA = GSEA

entire genes ranked used to query different annotation spaces (one by one)

modular EA = MEA

selected gene list used to query multiple annotation spaces (at once)

statistical and theoretical basis of these methods



Functional Enrcihment Analysis (EA, Análisis de Enriquecimiento Funcional): is the identification of biological functions and processes that are over-represented in a given gene list or in a gene subset selected in a given study.

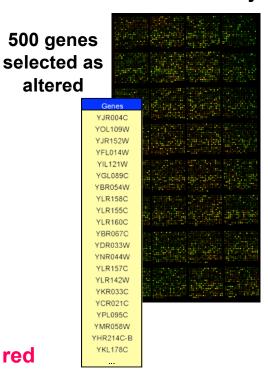


Find all the **genes** that are **annotated** to a given biological term and evaluate how **relevant** is such functional annotation

¿Which is the **probability** that such annotation is by chance (i.e. random)? The **p-value** allows to calculate the **statistical probability** of a given selection and allows to fix confidence thresholds (for example, set the test to allow less than 5% expected selection by chance)

If in 10,000 genes total => 100 genes are red
In 500 genes selected in one study => 5 genes are expected red
so any ratio found ≤ 1% (0.01) is normal and equal to by chance

10,000 genes total analysed in a 'omic' study





AATF	cell adhesion, ribosome biogenesis, apoptosis,
BAG1	neuron differentiation, apoptosis, anti-apoptosis,
CLIC4	cytoplasm, actin cytoskeleton, membrane
ATRIP	nucleus, DNA repair, protein binding,
NFKB1	signal trasduction, apoptosis, response to
MDM1	nucleus
OFD1	centrosome, oxidoreductase activity,
MAD1L1	mitotic metaphase, cytosol, spindle,
PDCD1	apoptosis, humoral immune response, protein
SNCG	axon, cell soma, adult locomotory behavior

QUERY gene list studied: 4 genes of 10 (in the gene list) annotated to 'apoptosis' REFERENCE gene list (i.e. whole proteome): 452 genes of 29,905 (universe) annotated to 'apoptosis'

Statistical test for significant enrichment in the gene list studied

'apoptosis' is an ENRICHED TERM with p-value = 1.12x10-5

3 major types:

singular EA = SEA

selected gene list used to query different annotation spaces (one by one)

gene set EA = GSEA

entire genes ranked used to query different annotation spaces (one by one)

modular EA = MEA

selected gene list used to query multiple annotation spaces (at once)

These methods alocate groups of **genes** to **terms** in a **contingency table** (with **YES** or **NO** categories) and then they apply the:

Fisher Exact test or Hypergeometric test

that calculate the probability using the **contingency table**of frequency data cross-classified according the two
categorical variables of assignment or not:
YES = if the **gene** in the list is assigned to the **term**NO = if the **gene** in the list is not-assigned to the **term**

Real	Event observed					
Event	Yes	No	Marginal total			
Yes	Hit	False alarm	Real Yes			
No	Miss	Correct non-event	Real No			
Marginal total	Obs Yes	Obs No	Sum total			
Real Event observed						
Event	Yes	No	Marginal total			
Arrograms 2	a	b	a + b			
Yes						
Yes No	С	d	c + d			

3 major types:

singular EA = SEA

selected gene list used to query different annotation spaces (one by one)

gene set EA = GSEA

entire genes ranked used to query different annotation spaces (one by one)

modular EA = MEA

selected gene list used to query multiple annotation spaces (at once)

These methods alocate groups of **genes** to **terms** in a **contingency table** (with **YES** or **NO** categories) and then they apply the:

Fisher Exact test or Hypergeometric test

that calculate the probability using the **contingency table**of frequency data cross-classified according the two
categorical variables of assignment or not:
YES = if the **gene** in the list is assigned to the **term**NO = if the **gene** in the list is not-assigned to the **term**

Gene	Gene Selected in Query List YES / NO					
Annotated White	Yes	No	Marginal total			
Yes	Hit	False alarm	W Yes			
No	Miss	Correct non-event	W No			
Marginal total	Obs Yes	Obs No	Sum total			
	1					
Gene Annotated	Gene Selected in Query List YES / NO					
White	Yes	No	Marginal total			
Yes	a	b	a + b			
No	С	d	c + d			
Marginal total	a+c	b + d	a+b+c+d=n			

Hypergeometric test

The classical application of the hypergeometric distribution is **sampling without replacement**. Think of an urn with two types of marbles, black ones and white ones. Define drawing a white marble as a success and drawing a black marble as a failure (analogous to the binomial distribution). If the variable N describes the number of **all marbles in the urn** (see contingency table below) and K describes the number of **white marbles**, then N - K corresponds to the number of **black marbles**. In this example X is the random variable whose outcome is K, the number of white marbles actually drawn in the experiment. This situation is illustrated by the following contingency table:

	drawn	not drawn	total
white marbles	k	K-k	K
black marbles	n – k	N + k - n - K	N – K
total	n	N – n	N

Now, assume (for example) that there are 5 white and 45 black marbles in the urn. Standing next to the urn, you close your eyes and draw 10 marbles without replacement. What is the probability that exactly 4 of the 10 are white? Note that although we are looking at success/failure, the data are not accurately modeled by the binomial distribution, because the probability of success on each trial is not the same, as the size of the remaining population changes as we remove each marble.

This problem is summarized by the following contingency table:

	drawn	not drawn	total		
white marbles	k = 4	K - k = 1	K = 5		
black marbles	n-k=6	N+k-n-K=39	N-K=45		
total	n = 10	N - n = 40	N = 50		

The probability of drawing exactly k white marbles can be calculated by the formula

$$P(X=k) = f(k; N, K, n) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}.$$

Hence, in this example calculate

$$P(X=4) = f(4;50,5,10) = \frac{\binom{5}{4}\binom{45}{6}}{\binom{50}{10}} = \frac{5 \cdot 8145060}{10272278170} = 0.003964583...$$

Universe

5 red and 45 black marbles in the urn of 50

&

Selection

4 red and 6 black of 10 select marbles







Hands-on: Practical Examples

Enrichment statistical tests
(Hypergeometric test and Fisher's Exact test)

run R Protocol: Protocol_ 1_ EnrichmentTests.R

3 major types:

singular EA = SEA (DAVID ncbi)
selected gene list used to query different
annotation spaces (one by one)

gene set EA = GSEA (GSEA Broad)
entire genes ranked used to query different
annotation spaces (one by one)

modular EA = MEA (GeneCodis CNB)
selected gene list used to query
multiple annotation spaces (at once)

Description Tool category Enrichment P-value is calculated Class I: singular on each term from the pre-selected enrichment interesting gene list. Then, analysis enriched terms are listed in a (SEA) simple linear text format. This strategy is the most traditional algorithm. It is still dominantly used by most of the enrichment analysis tools. Class II: Entire genes (without pre-selection) and associated experimental gene set enrichment values are considered in the analysis enrichment analysis. The unique (GSEA) features of this strategy are: (i) No need to pre-select interesting genes, as opposed to Classes I and II; (ii) Experimental values integrated into P-value calculation. Class III: This strategy inherits key spirit of modular SEA. However, the term-term/ enrichment gene-gene relationships are considered into enrichment P-value analysis calculation. The advantage of this (MEA) strategy is that term-term/genegene relationship might contain unique biological meaning that is not held by a single term or gene. Such network/modular analysis is closer to the nature of biological data structure.

Huang et al. (2009) NAR

3 major types:

singular EA = SEA

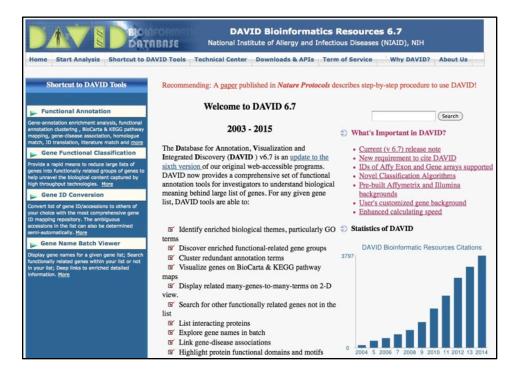
selected gene list used to query different annotation spaces (one by one)

gene set EA = GSEA

entire genes ranked used to query different annotation spaces (one by one)

modular EA = MEA

selected gene list used to query multiple annotation spaces (at once)



Type 1: singular EA = SEA (DAVID) selected gene list used to query different annotation spaces (one by one)

Type 1: singular EA = SEA (DAVID) selected gene list used to query different annotation spaces (one by one)

94 ch	art records						E	Down	load File
Sublis	t <u>Category</u> \$	<u>Te</u>	<u>ırm</u> :	RT	Genes	Count	<u>%</u> 🗧	P-Value	Benjamini;
	GOTERM_BP_FAT	defense response		RT		22	22.7	6.2E-11	5.1E-8
	GOTERM_BP_FAT	immune response		RT		20	20.6	2.1E-8	8.5E-6
	GOTERM_CC_FAT	extracellular region		<u>RT</u>		30	30.9	3.2E-6	5.6E-4
	GOTERM_CC_FAT	anchored to membrane		RT		10	10.3	7.3E-6	6.5E-4
	GOTERM_BP_FAT	response to bacterium		RT	_	9	9.3	2.1E-5	5.9E-3
	GOTERM_CC_FAT	extracellular region part		RT		18	18.6	5.3E-5	3.1E-3
	GOTERM_BP_FAT	defense response to bacterium		RT	_	7	7.2	5.9E-5	1.2E-2
	GOTERM_MF_FAT	carbohydrate binding		RT		10	10.3	7.1E-5	1.6E-2
	GOTERM_BP_FAT	response to wounding		RT		13	13.4	7.8E-5	1.3E-2
	GOTERM_CC_FAT	secretory granule		RT		8	8.2	1.1E-4	4.9E-3
	GOTERM_BP_FAT	inflammatory response		RT	_	10	10.3	1.5E-4	2.0E-2
	GOTERM_CC_FAT	extracellular space		RT		14	14.4	2.3E-4	8.3E-3
	GOTERM_CC_FAT	cytoplasmic vesicle		RT		13	13.4	4.8E-4	1.4E-2
	GOTERM_CC_FAT	cytoplasmic membrane-bounded vesicle		RT		12	12.4	4.8E-4	1.2E-2
	GOTERM_CC_FAT	membrane-bounded vesicle		RT		12	12.4	6.3E-4	1.4E-2
	annotation space	terms				genes		p	-values







Hands-on: Practical Examples

Protein_ SETs_ 2015.xls (106g hs, 175g hs, yeast 11pathways, yeast 59g5pc)

run DAVID-FAC

3 major types:

singular EA = SEA

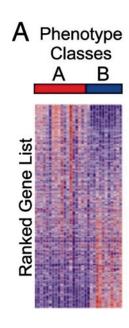
selected gene list used to query different annotation spaces (one by one)

gene set EA = GSEA

entire genes ranked used to query different annotation spaces (one by one)

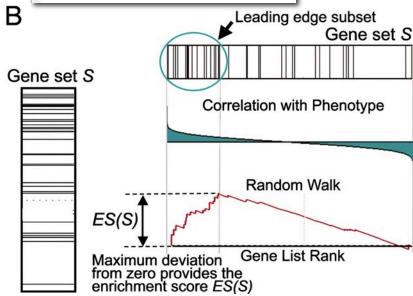
modular EA = MEA

selected gene list used to query multiple annotation spaces (at once)



Subramanian et al. (2005) PNAS







Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles

Aravind Subramaniana, Pablo Tamayoa, Vamsi K. Mootha, Sayan Mukherjeed, Benjamin L. Eberta, Michael A. Gillette*، , Amanda Paulovich^a, Scott L. Pomeroy^h, Todd R. Golub^{a, e}, Eric S. Lander^{a,c,i,k}, and Jill P. Mesirov^{a,k}

*Broad institute of Massachusetts Institute of Technology and Harvard. 320 Charles Street, Cambridge, MA 02141; 'Oppartment of Systems Biology, Alpert 536, Harvard Medical School, 200 Longwood Avenue, Booton, MA 02466', "Institute for Genome Sciences and Policy, Center for Interdisciplinary Engineering, Medicine, and Applied Sciences, Duke University, 10' 15 Science Drive, Durham, NC 2709, "Oppartment of Medical Oncore of Medical Oncore institute, 48 Binney Street, Booton, MA 02113; 'Division of Pulmonary and Critical Care Medicine, Masschusetts General Hospital, 55 Flust Street, Booton, MA 02113; 'Division of Pulmonary and Critical Care Medicine, Masschusetts General Hospital, 55 Flust Street, Booton, Care, Properties of Center, 1105 Gardieve Avenue North, 2023, P.O. Sox 19203, Seattle, WA 88190-1014, "Oppartment of Neurology, Enders 200, Children's Hospital, Hervard Medical School, 300 Longwood Avenue, Boston, MA 02115; 'Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02142; and Whitehead Institute for Biomedical Research, Massachusetts Institute of Technology, Cambridge, MA 02142.

Contributed by Eric S. Lander, August 2, 2005

Although genomewide RNA expression analysis has become a routine tool in biomedical research, extracting biological insight from such information remains a major challenge. Here, we describe a powerful analytical method called Gene Set Enrichment Analysis (GSEA) for interpreting gene expression data. The method derives its power by focusing on gene sets, that is, groups of genes that share common biological function, chromosomal location, or regulation. We demonstrate how GSEA yields insights into several cancer-related data sets, including leukemia and lung cancer. Notably, where single-gene analysis finds little similarity between two independent studies of patient survival in lung cancer, GSEA reveals many biological pathways in common. The GSEA method is embodied in a freely available software package, together with an initial database of 1,325 biologically defined gene sets.

enomewide expression analysis with DNA microarrays has become a mainstay of genomics research (1, 2). The challenge no longer lies in obtaining gene expression profiles, but rather in interpreting the results to gain insights into biological mechanisms. In a typical experiment, mRNA expression profiles are generated

for thousands of genes from a collection of samples belonging to one of two classes, for example, tumors that are sensitive resistant to a drug. The genes can be ordered in a ranked list L, according to their differential expression between the classes. The challenge is to extract meaning from this list.

A common approach involves focusing on a handful of genes at the top and bottom of L (i.e., those showing the largest difference) to discern telltale biological clues. This approach has a few major

(i) After correcting for multiple hypotheses testing, no individual gene may meet the threshold for statistical significance, because the relevant biological differences are modest relative to the noise inherent to the microarray technology.

(ii) Alternatively, one may be left with a long list of statistically significant genes without any unifying biological theme. Interpretation can be daunting and ad hoc, being dependent on a biologist's area of expertise.

(iii) Single-gene analysis may miss important effects on pathways. Cellular processes often affect sets of genes acting in concert. An increase of 20% in all genes encoding members of a metabolic pathway may dramatically alter the flux through the pathway and may be more important than a 20-fold increase in a single gene.

(iv) When different groups study the same biological system, the list of statistically significant genes from the two studies may show distressingly little overlap (3).

To overcome these analytical challenges, we recently developed a method called Gene Set Enrichment Analysis (GSEA) that 0 2005 by The National Academy of Sciences of the USA

www.pnas.org/cgi/doi/10.1073/pnas.0506580102

evaluates microarray data at the level of gene sets. The gene sets are defined based on prior biological knowledge, e.g., published information about biochemical pathways or coexpression in previous experiments. The goal of GSEA is to determine whether members of a gene set S tend to occur toward the top (or bottom) of the list L, in which case the gene set is correlated with the phenotypic class distinction.

We used a preliminary version of GSEA to analyze data from muscle biopsies from diabetics vs. healthy controls (4). The method revealed that genes involved in oxidative phosphorylation show reduced expression in diabetics, although the average decrease per gene is only 20%. The results from this study have been independently validated by other microarray studies (5) and by in vivo functional studies (6).

Given this success, we have developed GSEA into a robust technique for analyzing molecular profiling data. We studied its characteristics and performance and substantially revised and generalized the original method for broader applicability.

In this paper, we provide a full mathematical description of the GSEA methodology and illustrate its utility by applying it to several diverse biological problems. We have also created a software package, called GSEA-P and an initial inventory of gene sets (Molecular Signature Database, MSigDB), both of which are freely

Overview of GSEA. GSEA considers experiments with genomewide expression profiles from samples belonging to two classes, labeled 1 or 2. Genes are ranked based on the correlation between their expression and the class distinction by using any suitable metric

Given an a priori defined set of genes S (e.g., genes encoding products in a metabolic pathway, located in the same cytogenetic band, or sharing the same GO category), the goal of GSEA is to determine whether the members of S are randomly distributed throughout L or primarily found at the top or bottom. We expect

Freely available online through the PNAS open access option

Abbreviations: ALL acute lymphoid leukemia; AML, acute myeloid leukemia; ES, enrich-ment sone; FDR, faibe discovery rate; GSEA, Gene Set Enrichment Analysis; MAPK, mitogen-activated protein kinase; MSigDB, Molecular Signature Database; NES, normalized enrichment score.

See Commentary on page 15278.

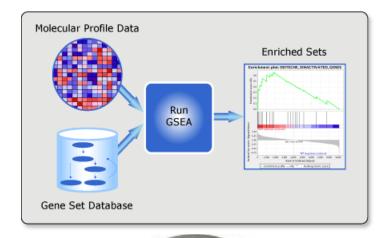
⁶A.S. and P.T. contributed equally to this work.

^kTo whom correspondence may be addressed. E-mail: lander@broad.mit.edu or

PNAS | October 25, 2005 | vol. 102 | no. 43 | 15545-15550

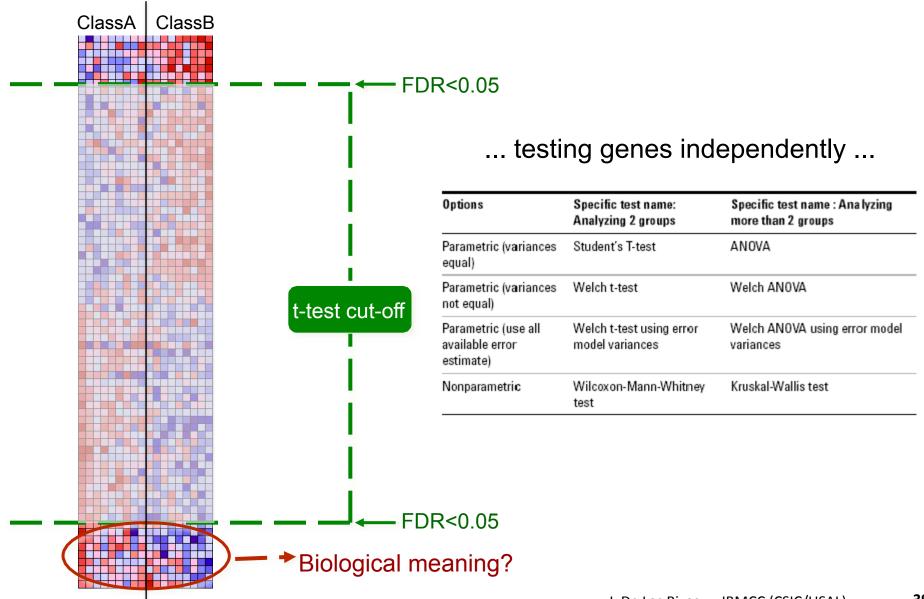
GSEA

MIT - Broad Institute

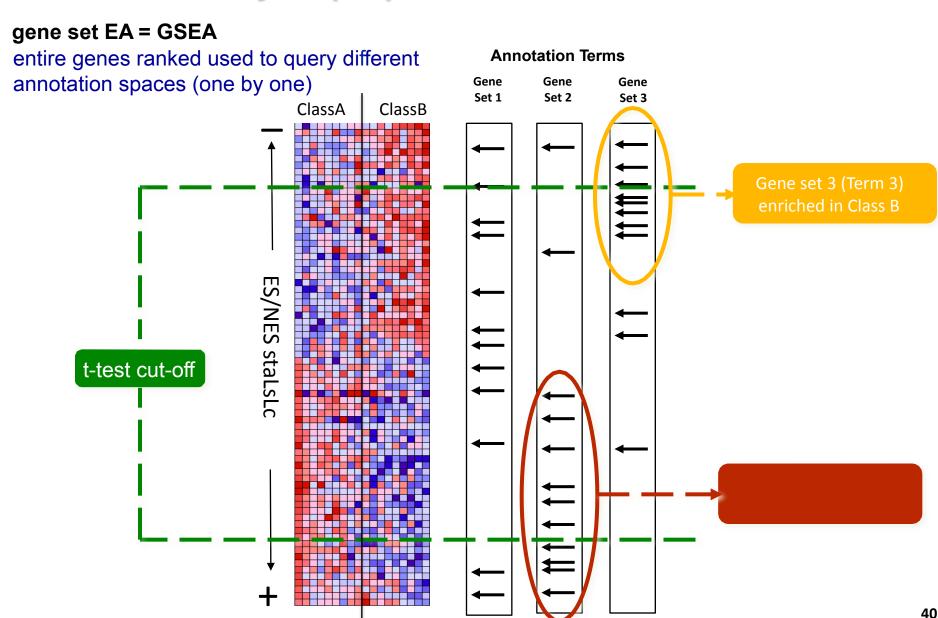




How GSEA works?



How GSEA works?



How GSEA works?

gene set EA = GSEA

entire genes ranked used to query different annotation spaces (one by one)

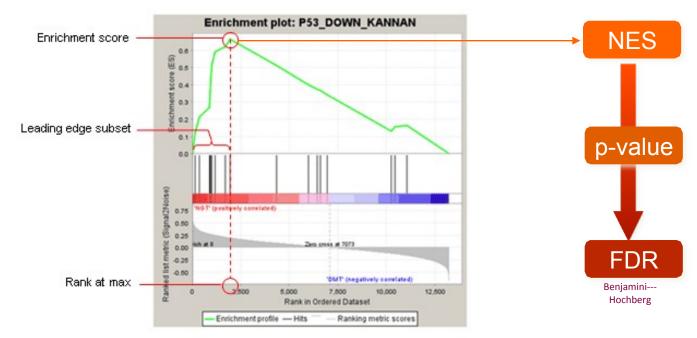
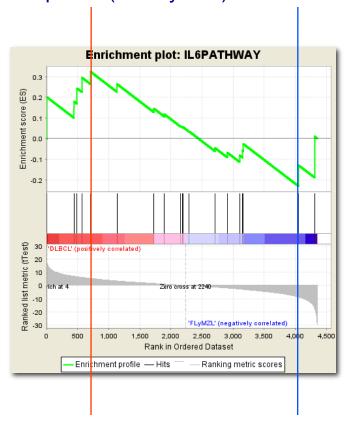


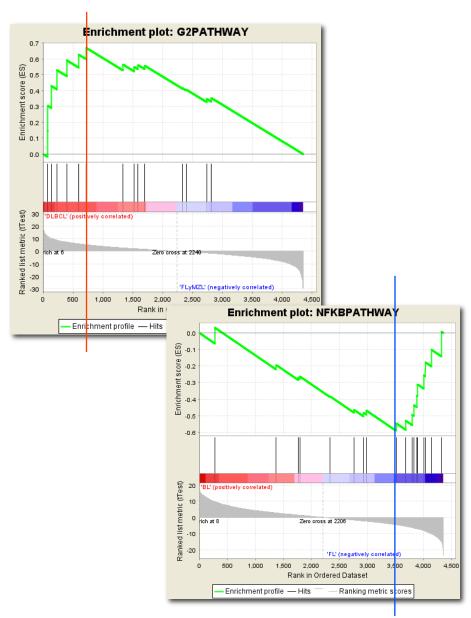
Fig 1: Enrichment plot: P53_DOWN_KANNAN
Profile of the Running ES Score & Positions of GeneSet Members on the Rank Ordered List

How GSEA works?

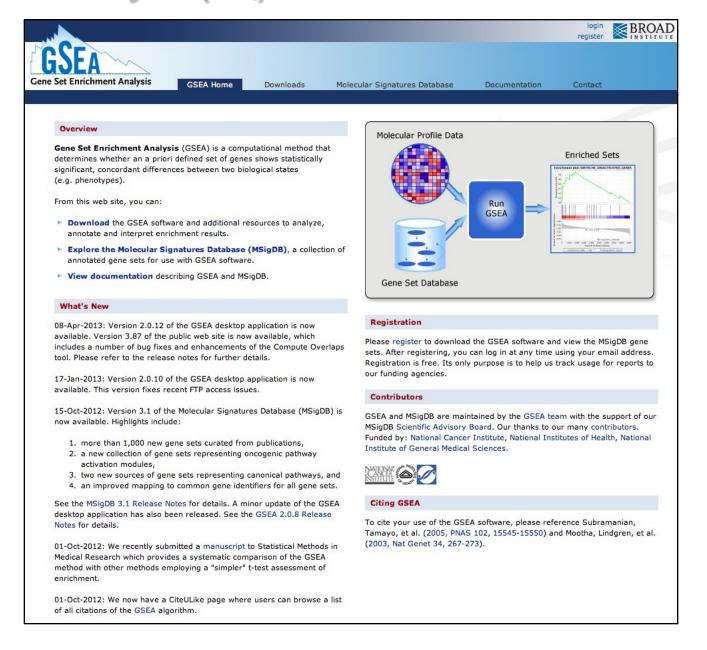
gene set EA = GSEA

entire genes ranked used to query different annotation spaces (one by one)



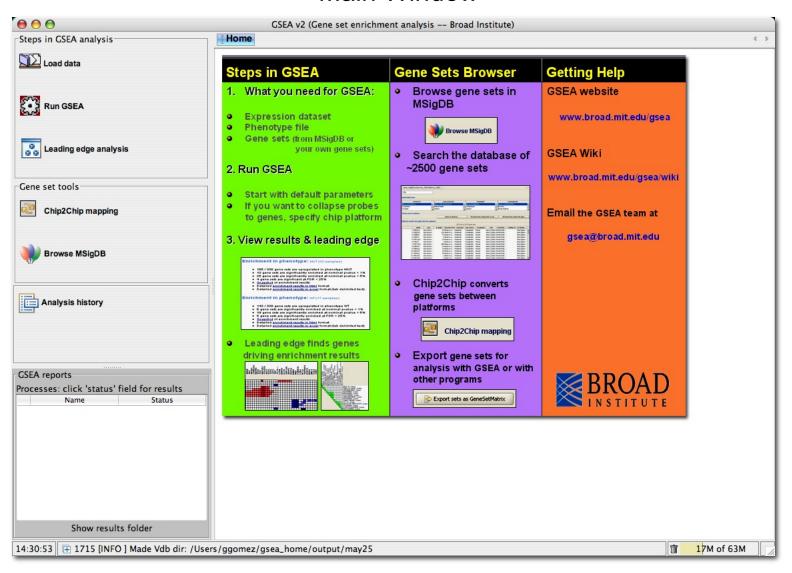


http://www.broad.mit.edu/gsea/



GSEA software

Main Window



GSEA software

Downloads

	de and the Molecular Signatures Database (MSigDB) are freely available to individuals in both acade see the GSEA/MSigDB license for more details.	emia and industry for
Software		
	A software. All options implement exactly the same algorithm. Usage recommendations and installa entations of GSEA require Java 1.6 or higher. If your computer has Java 1.5 and cannot upgrade to	
javaGSEA Desktop Application	 Easy-to-use graphical user interface Runs on any desktop computer (Windows, Mac OS X, Linux etc.) that supports Java1.6+ Produces richly annotated reports of enrichment results Integrated gene sets browser to view gene set annotations, search for gene sets and map gene sets between platforms The GSEA team suggests always starting GSEA by using these Launch buttons, or by clicking the icon that the application installs on your desktop, in order to ensure optimal memory allocation 	Launch with 512Mb memory £ Launch Launch with 1Gb memory £ Launch
javaGSEA Java Jar file	► Command line usage ► Runs on any platform that supports Java1.6+ ► We recommend using the 'Launch' buttons above instead of this mode for most users	download gsea2-2.07.jar
GSEA Java Source Code Java source files	 100% Java implementation of GSEA Incorporate GSEA into your own data analysis pipeline Programmatically call the open source GSEA java API 	download gsea2_distrib-2.04.zip
R-GSEA R Script	 ► Usage from within the R programming environment ► Easily inspect, learn and tweak the algorithm ► Incorporate GSEA into your own data analysis pipeline ► Programmatically call the open source GSEA R API ► Click here to learn more about the R-GSEA script 	download GSEA-P-R.1.0.zip
GenePattern GSEA Module	 Use GSEA from within GenePattern Use GSEA in concert with a large suite of other analytics found in GenePattern (a powerful and flexible analysis platform developed at the Broad Institute) 	GenePattern site













Session 1 (9:30 - 12:30, 3h)
Bioinformatic tools for Functional Enrichment Analysis (FEA)
Session 2 (13:30 - 16:30, 3h)
Construction of gene functional networks

- Introduction to biological information and annotation spaces: GO, KEGG, Interpro
- Functional Enrichment Analysis (EA): from single to modular methods
 - Using EA tools to annotate gene lists:
 DAVID (single), GSEA (gene sets), GeneCodis (modular)
 - Sort out problems after EA: post-enrichment tool GeneTermLinker (postEA)
 - From co-annotation and enrichment to functional networks: networks construction using a R tool

Functional	
Enrichment Analysis (EA)

3 major types: (examples)

singular EA = SEA

selected gene list used to query different annotation spaces (one by one)

gene set EA = GSEA

entire genes ranked used to query different annotation spaces (one by one)

modular EA = MEA

selected gene list used to query multiple annotation spaces (at once)

Tool category Description

Class I: singular enrichment analysis (SEA) Enrichment P-value is calculated on each term from the pre-selected interesting gene list. Then, enriched terms are listed in a simple linear text format. This strategy is the most traditional algorithm. It is still dominantly used by most of the enrichment analysis tools.

Class II: gene set enrichment analysis (GSEA) Entire genes (without pre-selection) and associated experimental values are considered in the enrichment analysis. The unique features of this strategy are: (i) No need to pre-select interesting genes, as opposed to Classes I and II; (ii) Experimental values integrated into *P*-value calculation.

Class III: modular enrichment analysis (MEA) This strategy inherits key spirit of SEA. However, the term—term/gene—gene relationships are considered into enrichment *P*-value calculation. The advantage of this strategy is that term—term/gene—gene relationship might contain unique biological meaning that is not held by a single term or gene. Such network/modular analysis is closer to the nature of biological data structure.

Huang et al. (2009) NAR

3 major types:

singular EA = SEA (DAVID ncbi)
selected gene list used to query different
annotation spaces (one by one)

gene set EA = GSEA (GSEA Broad)
entire genes ranked used to query different
annotation spaces (one by one)

modular EA = MEA (GeneCodis CNB)
selected gene list used to query
multiple annotation spaces (at once)

Description Tool category Enrichment P-value is calculated Class I: singular on each term from the pre-selected enrichment interesting gene list. Then, analysis enriched terms are listed in a (SEA) simple linear text format. This strategy is the most traditional algorithm. It is still dominantly used by most of the enrichment analysis tools. Class II: Entire genes (without pre-selection) and associated experimental gene set enrichment values are considered in the analysis enrichment analysis. The unique (GSEA) features of this strategy are: (i) No need to pre-select interesting genes, as opposed to Classes I and II; (ii) Experimental values integrated into P-value calculation. Class III: This strategy inherits key spirit of modular SEA. However, the term-term/ enrichment gene-gene relationships are considered into enrichment P-value analysis calculation. The advantage of this (MEA) strategy is that term-term/genegene relationship might contain unique biological meaning that is not held by a single term or gene. Such network/modular analysis is closer to the nature of biological

data structure.

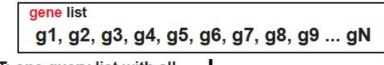
Huang et al. (2009) NAR



problems:

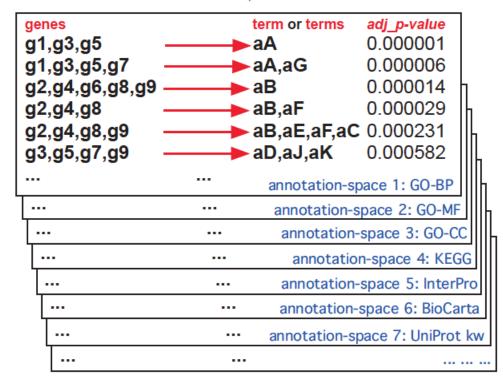
- 1.there are many annotation spaces that **overlap**
- 2.this functional overlapping produces **repetitive** results
- 3.for a **single** query **gene list** SEA and GSEA tools produced many **result lists** that are highly **redundant**

General approach of the Enrichment Analysis (EA) tools: SEA, MEA, GSEA



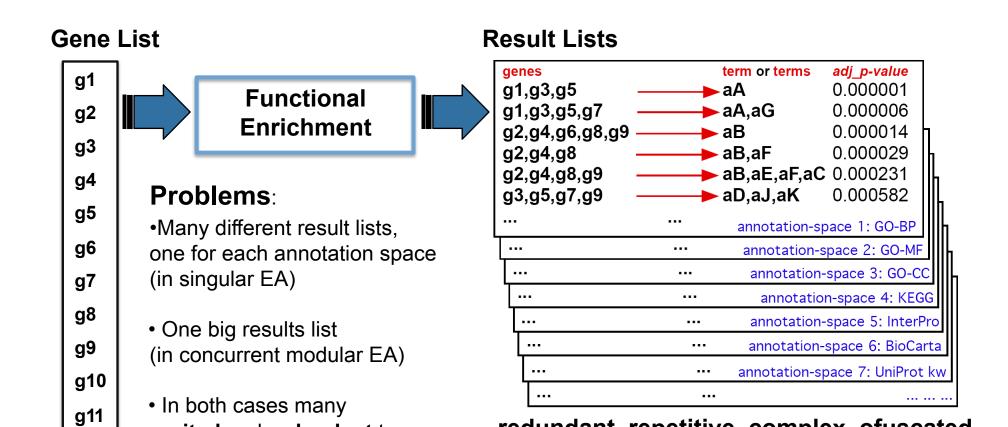
INPUT: one query list with all significant genes selected (one-dimension, simple)





repited and redundant terms



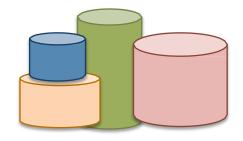


redundant, repetitive, complex, ofuscated

results



Examples of: **redundant**, **repetitive** ... inherent problems to the annotation spaces:



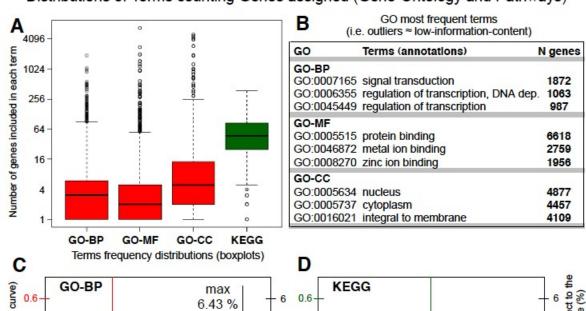
- 3 Equivalent terms, with the same biological meaning
 GO:0007049 cell cycle === GO:0022402 cell cycle process
- ③ Redundancy of terms, repeated in different annotation spaces
 GO:0007049: cell cycle === KEGG hsa04110: cell cycle
- ③ Bias due to highly frequent promiscuous terms that are unspecific
 GO:0005515: regulation of biological process includes
 ≈ 45% of all annotated human genes in GO-BP

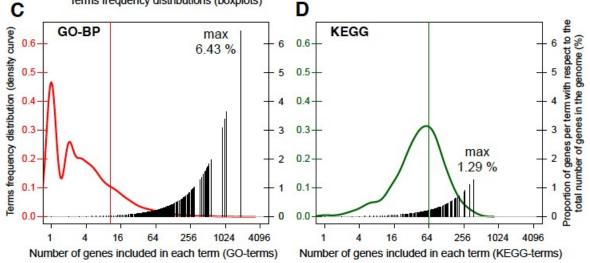


problems:

- 1.there are many annotation spaces that **overlap**
- 2.this functional overlapping produces **repetitive** results
- 3.for a **single** query **gene list** SEA and GSEA tools produced many **result lists** that are highly **redundant**

Terms (biological annotations) Distributions of Terms counting Genes assigned (Gene-Ontology and Pathways)



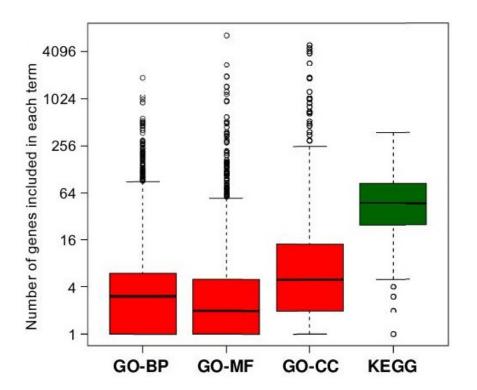


showing most frequent terms in GO and KEGG



Analysis of the distribution of **genes per term** in different annotation spaces

Distributions of number genes per term in each annotation space



GO terms with a very high number of genes annotated that have **low-information-content** and easily become **over-represented**

GO	Terms (annotations)	Genes (N)
GO-BP		
GO:0007165 GO:0006355 GO:0045449	signal transduction regulation of transcription, DNA dep. regulation of transcription	1872 1063 987
GO-MF		
GO:0005515	protein binding	6618
GO:0046872	metal ion binding	2759
GO:0008270	zinc ion binding	1956
GO-CC		
GO:0005634	nucleus	4877
GO:0005737	cytoplasm	4457
GO:0016021	integral to membrane	4109



Distribution of genes per term (GeneOntology terms and Pathway terms)

problems:

Most Frequent Terms (highly used ≈ low-information-content)

GO-BP GO:0006355 GO:0007165 GO:0006350	N genes 1869 1700 1516	Terms regulat. cellular transcription signal transduction cellular transcription
GO-MF GO:0005515 GO:0046872 GO:0008270	5306 2370 2251	protein binding metal ion binding zinc ion binding
GO-CC GO:0005634 GO:0005737 GO:0016021	4930 4264 4013	nucleus cytoplasm integral to membrane

Comparative gene-size of the annotation-terms that include at least 3 genes (GO-BP versus KEGG)

00000	760	495	422		4	02	376	5	35	7	- 1	348		31	2
1700			267	22	4	198	191	ı	186	1	82	1	77	1	68
signal transduction	570	472	258	21	9	167	143	14	1	136	13	14	12	5 1	124
			250			162	123	12	2 1	119	11	4	112	1	111
10000	506	450	239	21	4	159	110	94	93	9	2	89	8		86
cellular			227	21	2	153	105	85	76	75	73 54		72	71 51	69
transcription	503	434	237	21	5	151	99	82		61 59					
			232	20	5		98	80	64	57	48				
874	497	424				147	97	78	63	56	47 46				
organismal development			227	20	2	146	96	77		55					

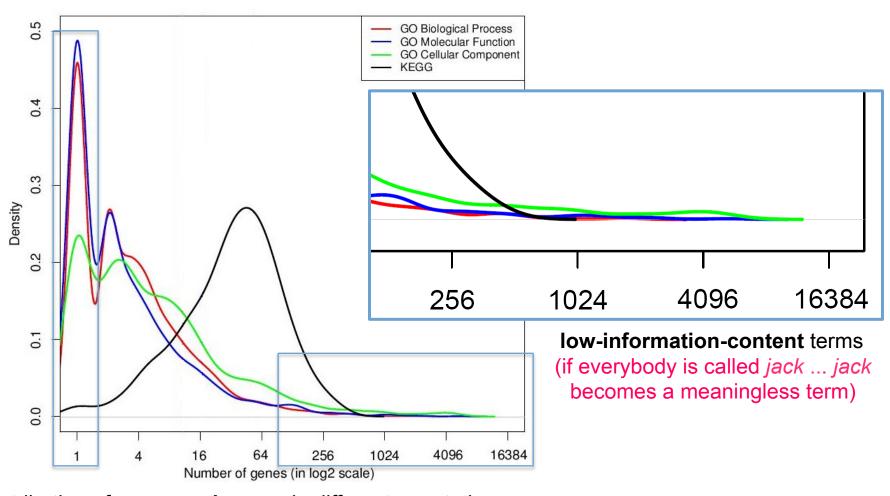
GO-BP (human) aprox.1800 terms terms with n° genes > [Med+2.5*SD] Highly used terms (red) ≈ low-information-content 33 terms that include 36.8 % of genes

363	246	161	130	129	1	124	12	3	114	4	10	08
olfactory transduction		150	98	87	74	73		72	7	70	6	8
320	207	150	96	85	67	66		65	6	4	6	53
pathways in cancer	195	145	90	84	62	55	54		52	51		50
	193		93	83	61	49	44	43	42		11	39
MAPK	176	133	92	83	58	48	37	33	29	25	24	27
singnaling				82	57	47	36 35	32	22	18		14
252	175	131	89	81	56	45	34	30		17 16		

KEGG (human) aprox. 200 terms terms with no genes > [Med+2.5*SD]
 Highly used terms (red) ≈ low-information-content 9 terms that include 20.8 % of genes



Analysis of the distribution of **genes per term** in different annotation spaces



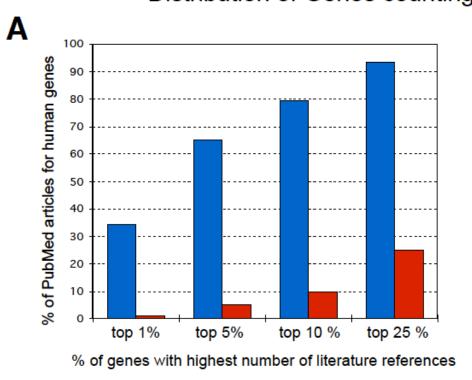
Distribution of terms and genes in different annotation spaces



problems: there is also a clear knowledge bias in the scientific literature

Genes (biological entities): Distribution of Genes counting number of articles (PubMed)

В



	Genes	N articles (PubMed)
1	TP53	26891
2	TNF (TNF-a)	21616
3	ÌNS	20298
4	NFKB1	16441
5	FOS (c-FOS)	16164
6	MYC (c-MYC)	12839
7	CD4	12427
8	TGFB1	12207
9	IFNG	12185
10	IL1B	11767
11	EGF	10493
12	BCL2	10123
13	CALCA	10006
14	IL6	9321
15	IL2	9186
16	VEGFA	8651
17	ESR1	8564
18	CAT	8274
19	FOSB	7227
20	IL4	7013

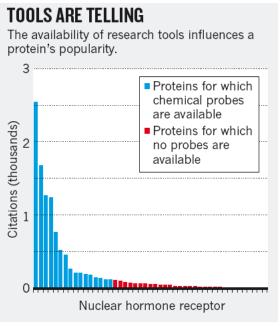


problems: there is also a clear knowledge bias in the scientific literature



Too many roads not taken

Most protein research focuses on those known before the human genome was mapped. Work on the slew discovered since, urge Aled M. Edwards and his colleagues.

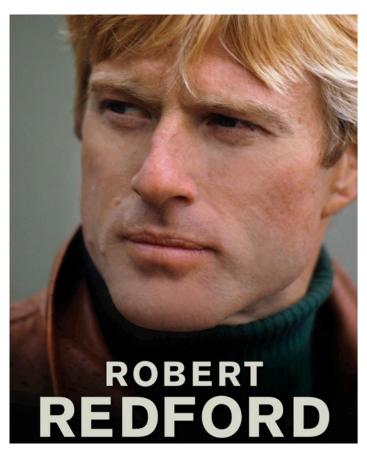


Edwards et al. (2011) SCIENCE



problems: knowledge bias in the scientific literature === everybody wants
to publish about celebrity proteins ('robert redford' et al.)

Genes (biological entities):
Distribution of Genes counting number of articles (PubMed)



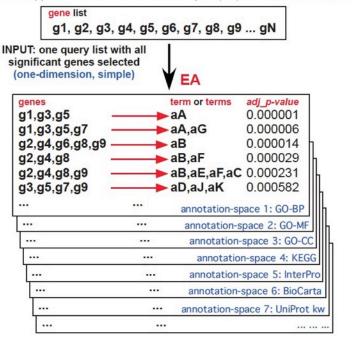
	Genes	N articles (PubMed)
1	TP53	26891
2	TNF (TNF-a)	21616
3	INS	20298
4	NFKB1	16441
5	FOS (c-FOS)	16164
6	MYC (c-MYC)	12839
7	CD4	12427
8	TGFB1	12207
9	IFNG	12185
10	IL1B	11767
11	EGF	10493
12	BCL2	10123
13	CALCA	10006
14	IL6	9321
15	IL2	9186
16	VEGFA	8651
17	ESR1	8564
18	CAT	8274
19	FOSB	7227
20	IL4	7013

solution:

design a new system of multiple EA that avoids redundant terms and provides a **unified result** based in functional convergence of **genes** & **terms** =

GeneTerm Linker tool

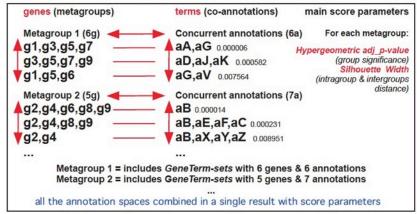
General approach of the Enrichment Analysis (EA) tools: SEA, MEA, GSEA



OUTPUT: many lists of genes/annotations highly redundant (enriched gene-sets) (multiple-dimensions, complex)



Approach beyond EA: reciprocal linkage of genes & terms to find functional association



OUTPUT: one list of metagroups built using significant coherent gene-term linkage and removing redundant non-informative sets (one-dimension, simple)



solution:

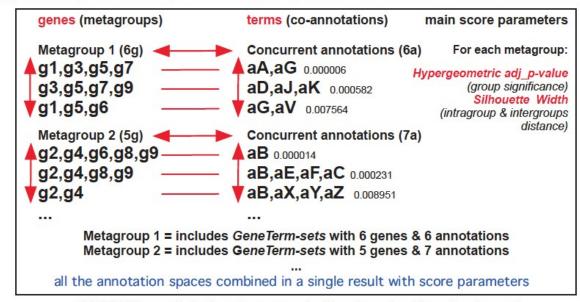
design a new system of multiple EA that avoids redundant terms and provides a unified result based in functional convergence of genes & terms =

GeneTerm Linker tool

OUTPUT: many lists of genes/annotations highly redundant (enriched gene-sets) (multiple-dimensions, complex)



Approach beyond EA: reciprocal linkage of genes & terms to find functional association



OUTPUT: one list of metagroups built using significant coherent gene-term linkage and removing redundant non-informative sets (one-dimension, simple)



solution: GeneTerm Linker tool

	Term Linker ssociation by non-redundant reciprocal linkage
1. Input a list of genes of interest:	3. Organism:
	Homo sapiens 💠
	4. Annotation Spaces:
	 ☐ GO Biological Process ☐ GO Molecular Function
	GO Cellular Component
[Human example] [Yeast example]	☐ KEGG Pathways ☐ InterPro Motifs And Domains
[numan example] [reast example]	
2. Input a list of genes of reference (optional):	5. Minimum Support: 4 💠
	6. Email address (optional):
	(Submit analysis) (Reset)



singular EA = SEA :: DAVID http://david.abcc.ncifcrf.gov/
selected gene list used to query different
annotation spaces (one by one)

gene set EA = GSEA :: GSEA http://www.broadinstitute.org/gsea/ entire genes ranked used to query different annotation spaces (one by one)

modular EA = MEA :: GeneCodis http://genecodis.dacya.ucm.es/ selected gene list used to query multiple annotation spaces (at once)

Beyond EA:: GeneTerm Linker selected gene list used to query multiple annotation spaces (at once) avoiding repetition and redundancy and linking genes and terms

http://gtlinker.dacya.ucm.es/







Hands-on: Practical Examples

Protein_ SETs_ 2014.xls (106g hs, 175g hs, yeast 11pathways, yeast 59g5pc)

run GeneTerm Linker







Session 1 (9:30 - 12:30, 3h)
Bioinformatic tools for Functional Enrichment Analysis (FEA)
Session 2 (13:30 - 16:30, 3h)
Construction of gene functional networks

- Introduction to biological information and annotation spaces: GO, KEGG, Interpro
- Functional Enrichment Analysis (EA): from single to modular methods
 - Using EA tools to annotate gene lists:
 DAVID (single), GSEA (gene sets), GeneCodis (modular)
 - Sort out problems after EA: post-enrichment tool GeneTermLinker (postEA)
 - From co-annotation and enrichment to functional networks: networks construction using a R tool

GeneTerm Linker a post enrichment tool



OPEN @ ACCESS Freely available online

Fontanillo et al. (2011) PLoS ONE



Functional Analysis beyond Enrichment: Non-Redundant Reciprocal Linkage of Genes and Biological Terms

Celia Fontanillo¹⁹, Ruben Nogales-Cadenas²⁹, Alberto Pascual-Montano², Javier De Las Rivas¹*

1 Cancer Research Center (CiC-IBMCC, CSIC/USAL), Campus Miguel de Unamuno, Salamanca, Spain, 2 National Center of Biotechnology (CNB, CSIC), Campus de Cantoblanco UAM, Madrid, Spain

Abstract

Functional analysis of large sets of genes and proteins is becoming more and more necessary with the increase of experimental biomolecular data at *omic*-scale. Enrichment analysis is by far the most popular available methodology to derive functional implications of sets of cooperating genes. The problem with these techniques relies in the redundancy of resulting information, that in most cases generate lots of trivial results with high risk to mask the reality of key biological events. We present and describe a computational method, called *GeneTerm Linker*, that filters and links enriched output data identifying sets of associated genes and terms, producing metagroups of coherent biological significance. The method uses fuzzy reciprocal linkage between genes and terms to unravel their functional convergence and associations. The algorithm is tested with a small set of well known interacting proteins from yeast and with a large collection of reference sets from three heterogeneous resources: multiprotein complexes (CORUM), cellular pathways (SGD) and human diseases (OMIM). Statistical *Precision*, *Recall* and balanced *F-score* are calculated showing robust results, even when different levels of random noise are included in the test sets. Although we could not find an equivalent method, we present a comparative analysis with a widely used method that combines enrichment and functional annotation clustering. A web application to use the method here proposed is provided at http://gtlinker.cnb.csic.es.

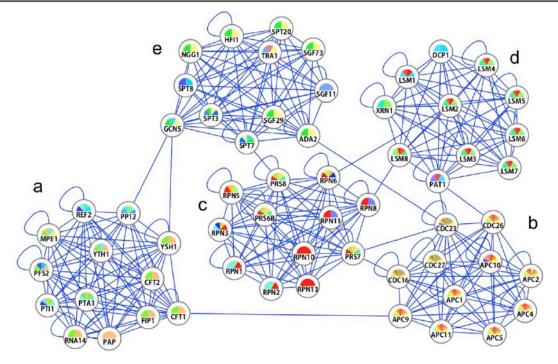
Citation: Fontanillo C, Nogales-Cadenas R, Pascual-Montano A, De Las Rivas J (2011) Functional Analysis beyond Enrichment: Non-Redundant Reciprocal Linkage of Genes and Biological Terms. PLoS ONE 6(9): e24289. doi:10.1371/journal.pone.0024289



Protein Complexes (yeast)	Number of proteins
a. mRNA cleavage and polyadenylation specificity factor complex CFT1,CFT2,FIP1,GLC7,MPE1,PAP1,PFS2,PTA1,PTI1,REF2,RNA14,YSH1,YTH1	13
b. anaphase-promoting complex APC1,APC2,APC4,APC5,APC9,APC11,CDC16,CDC23,CDC26,CDC27,DOC1	11
C. proteasome, 19/22S regulator complex RPN1,RPN2,RPN3,RPN5,RPN6,RPN8,RPN10,RPN11,RPN13,RPT1,RPT3,RPT6	12
d. U6 snRNP complex DCP1,KEM1,LSM1,LSM2,LSM3,LSM4,LSM5,LSM6,LSM7,LSM8,PAT1	11
e. SAGA complex ADA2,GCN5,HFI1,NGG1,SGF11,SGF29,SGF73,SPT3,SPT7,SPT8,SPT20,TRA1	12

Testing the method with a set of **59** yeast nuclear **proteins** working in **5** known **complexes**

Bader et al (2003)



Reconstruction of the complexes using experimental protein interaction data (from APID and APID2NET)

Prieto et al (2006) Hernández-Toro et al (2007)

Cic

Validation example

Query list of 59 yeast proteins

GeneTerm Linker results: 5 metagroups

Gene Metagroups found (include all the related groups)	GENES found	p value	CoAnnotaGons (GOBP, GOMF, GOCC, KEGG, InterPro)
Metagroup 1; SilhouePe: 0.529 GLC7, REF2, YTH1, FIP1, PAP1, PFS2, CFT1, RNA14, PTI1, PTA1, MPE1, CFT2, YSH1	13 (59)	2.26E 26	GO:0005847:mRNA cleavage and polyadenylaLon specificity factor complex (CC), GO:0006378:mRNA polyadenylaLon (BP), GO:0006379:mRNA cleavage (BP), GO:0031123:RNA 3'end processing (BP), GO:0006353:transcripLon terminaLon (BP),
Metagroup 2; SilhouePe: 0.823 CDC23, APC5, CDC16, APC2, APC1, DOC1, APC9, APC11, APC4, CDC27, CDC26, GLC7, TRA1	11 (59)	4.85E 24	04111:Cell cycle – yeast, GO:0004842:ubiquiLnprotein ligase acLvity (MF),GO:0000070:mitoLc sister chromaLd segregaLon (BP),GO: 0016567:protein ubiquiLnaLon (BP),GO:0000022:mitoLc spindle elongaLon (BP),GO:0005680:anaphasepromoLng complex (CC),
Metagroup 3; SilhouePe: 0.609 RPN13, RPN8, RPN1, RPT1, RPN3, RPN10, RPN11, RPN2, RPN5, RPT3, RPT6, RPN6, APC2, DOC1	14 (59)	8.11E 15	GO:0000502:proteasome complex (CC), 03050:Proteasome, GO: 0034515:proteasome storage granule (CC), GO:0008541:proteasome regulatory parLcle, lid subcomplex (CC), GO:0004175:endopepLdase acLvity (MF), GO:0030234:enzyme regulator acLvity (MF)
Metagroup 4; SilhouePe: 0.750 LSM5, DCP1, PAP1, LSM3, LSM8, LSM6, LSM1, LSM4, LSM7, LSM2, PTA1, KEM1, PAT1, YSH1	14 (59)	1.31E 14	03018:RNA degradaLon, GO:0008033:tRNA processing (BP), GO: 0030529:ribonucleoprotein complex (CC), GO:0046540:U4/U6 x U5 trisnRNP complex (CC), GO:0000398:nuclear mRNA splicing, via spliceosome (BP), IPR001163:LikeSm ribonucleoprotein (LSM), GO:0005688:U6 snRNP (CC),
Metagroup 5; SilhouePe: 0.570 NGG1, HFI1, TRA1, SPT20, SGF29, SPT7, SGF73, SPT8, SGF11, GCN5, SPT3, ADA2, RPN6, CFT2	14 (59)	5.26E 12	GO:0000124:SAGA complex (CC), GO:0016568:chromaln modificalon (BP), GO:0016573:histone acetylalon (BP), GO:0046695:SLIK (SAGAlike) complex (CC), GO:0003702:RNA polymerase II transcriplon factor aclivity (MF), GO: 0004402:histone acetyltransferase aclivity (MF),



a. mRNA cleavage and polyadenilation specificity factor complex: **13 genes**

Gene Metagroups found (include all the related groups)	GENES found InterPr	pvalu o)	e CoAnnotaGons (GOBP, GOMF, GOCC, KEGG,
Metagroup 1; SilhouePe: 0.529 GLC7, REF2, YTH1, FIP1, PAP1, PFS2, CFT1, RNA14, PTI1, PTA1, MPE1, CFT2, YSH1	13 (59) (BP),	2.26E26	GO:0005847:mRNA cleavage and polyadenylaGon specificity factor complex (CC), GO:0006378:mRNA; polyadenylaGon (BP), GO:0006379:mRNA cleavage GO:0031123:RNA 3'end processing (BP), GO: 0006353:transcripLon terminaLon
Metagroup 2; SilhouePe: 0.823 CDC23, APC5, CDC16, APC2, APC1, DOC1, APC9, APC11, APC4, CDC27, eD626g&L67, TRA1	11 (59)		(PAP);Cell cycle – yeast, GO:0004842:ubiquiLnprotein ligase acLvity (MF),GO:0000070:mitoLc sister chromaLd segregaLon (BP),GG0:0000022:mitoLc spindle (BP),GO:0005680:anaphasepromoLng complex (CC),
Metagroup 3; SilhouePe: 0.609 RPN13, RPN8, RPN1, RPT1, RPN3, RPN10, RPN11, RPN2, RPN5, RPT3, RPT6, RPN6, APC2, DOC1	14 (59)	8.11E	GO:0000502:proteasome complex (CC), 03050:Proteasome, GO: 0034515:proteasome storage granule (CC), GO:0008541:proteasome regulatory parLcle, lid subcomplex (CC), GO:0004175:endopepLdase acLvity (MF), GO:0030234:enzyme regulator acLvity (MF)
Metagroup 4; SilhouePe: 0.750 LSM5, DCP1, PAP1, LSM3, LSM8, LSM6, LSM1, LSM4, LSM7, LSM2, PTA1, KEM1, PAT1, YSH1	14 (59)	1.31E	03018:RNA degradaLon, GO:0008033:tRNA processing (BP), GO: 0030529:ribonucleoprotein complex (CC), GO:0046540:U4/U6 x U5 trisnRNP complex (CC), GO:0000398:nuclear mRNA splicing, via spliceosome (BP), IPR001163:LikeSm ribonucleoprotein (LSM), GO:0005688:U6 snRNP (CC),
Metagroup 5; SilhouePe: 0.570 NGG1, HFI1, TRA1, SPT20, SGF29, SPT7, SGF73, SPT8, SGF11, GCN5, SPT3, ADA2, RPN6, CFT2	14 (59)	5.26E	GO:0000124:SAGA complex (CC), GO:0016568:chromaln modificalon (BP), GO:0016573:histone acetylalon (BP), GO:0046695:SLIK (SAGAlike) complex (CC), GO:0003702:RNA polymerase II transcriplon factor acluity (MF), GO: 0004402:histone acetyltransferase acluity (MF),



b. Anaphase-promoting complex: **11 genes**

Gene Metagroups found (include all the related groups)	GENES found interPr	pvalu o)	e CoAnnotaGons (GOBP, GOMF, GOCC, KEGG,
Metagroup 1; SilhouePe: 0.529 GLC7, REF2, YTH1, FIP1, PAP1, PFS2, CFT1, RNA14, PTI1, PTA1, MPE1, CFT2, YSH1	13 (59) 0006353	2.26E26	GO:0005847:mRNA cleavage and polyadenylaLon specificity factor complex (CC), GO:0006378:mRNA polyadenylaLon (BP), GO:0006379:mRNA cleavage (BP), GO:00031123:RNA 3'end processing (BP), GO:
Metagroup 2; SilhouePe: 0.823 CDC23, APC5, CDC16, APC2, APC1, DOC1, APC9, APC11, APC4, CDC27, CDC26, GLC7, TRA1	11 (59) elong	4.85E24	04111:Cell cycle – yeast, GO:0004842:ubiquiLn protein ligase acLvity (MF), GO:0000070:mitoGc sistgregh@ma(BB), GO:0000022:mitoLc spindle (BP),GO:0005680:anaphasepromoGng complex (CC),
Metagroup 3; SilhouePe: 0.609 RPN13, RPN8, RPN1, RPT1, RPN3, RPN10, RPN11, RPN2, RPN5, RPT3, RPT6, RPN6, APC2, DOC1	14 (59)	8.11E	GO:0000502:proteasome complex (CC), 03050:Proteasome, GO: 0034515:proteasome storage granule (CC), GO:0008541:proteasome regulatory parLcle, lid subcomplex (CC), GO:0004175:endopepLdase acLvity (MF), GO:0030234:enzyme regulator acLvity (MF)
Metagroup 4; SilhouePe: 0.750 LSM5, DCP1, PAP1, LSM3, LSM8, LSM6, LSM1, LSM4, LSM7, LSM2, PTA1, KEM1, PAT1, YSH1	14 (59)	1.31E	03018:RNA degradaLon, GO:0008033:tRNA processing (BP), GO: 0030529:ribonucleoprotein complex (CC), GO:0046540:U4/U6 x U5 tri snRNP complex (CC), GO:0000398:nuclear mRNA splicing, via spliceosome (BP), IPR001163:LikeSm ribonucleoprotein (LSM), GO:0005688:U6 snRNP (CC),
Metagroup 5; SilhouePe: 0.570 NGG1, HFI1, TRA1, SPT20, SGF29, SPT7, SGF73, SPT8, SGF11, GCN5, SPT3, ADA2, RPN6, CFT2	14 (59)	5.26E	GO:0000124:SAGA complex (CC), GO:0016568:chromaln modificalon (BP), GO:0016573:histone acetylalon (BP), GO:0046695:SLIK (SAGAlike) complex (CC), GO:0003702:RNA polymerase II transcriplon factor aclvity (MF), GO: 0004402:histone acetyltransferase aclvity (MF),



12

Gene Metagroups

	found (include all the related groups)	GENES found interpr	ı pvaiu	e CoAnnotaGons (GOBP, GOMF, GOCC, KEGG,
	Metagroup 1; SilhouePe: 0.529 GLC7, REF2, YTH1, FIP1, PAP1,			GO:0005847:mRNA cleavage and polyadenylaLon specificity factor complex (CC), GO:0006378:mRNA polyadenylaLon (BP), GO:0006379:mRNA cleavage
	PFS2, CFT1, RNA14, PTI1, PTA1, MPE1, CFT2, YSH1	13 (59) 0006353	2.26E26	(BP), GO:0031123:RNA 3'end processing (BP), GO:
	Metagroup 2; SilhouePe: 0.823 CDC23, APC5, CDC16, APC2, APC1, DOC1, APC9, APC11, APC4, CDC27, €D626g€LL€₹, TRA1	11 (59)		04111:Cell cycle – yeast, GO:0004842:ubiquiLnprotein ligase acLvity (MF),GO:0000070:mitoLc sister chromaLd segregaLon (BP),GO:0000022:mitoLc spindle (BP),GO:0005680:anaphasepromoLng complex (CC),
	Metagroup 3; SilhouePe: 0.609 RPN13, RPN8, RPN1, RPT1, RPN3, RPN10, RPN11, RPN2, RPN5, RPT3, RPT6, RPN6, APC2, DOC1	14 (59)	8.11E	GO:0000502:proteasome complex (CC), 03050:Proteasome, GO:0034515:proteasome storage granule (CC), GO:0008541:proteasome regulatory
	,	15		parLcle, lid subcomplex (CC), GO: 90041 t75: (nntl o);e pJ_0 l 2 0s @ 30234:enzyme regulator
	Metagroup 4; SilhouePe: 0.750 LSM5, DCP1, PAP1, LSM3, LSM8, LSM6, LSM1, LSM4, LSM7, LSM2, PTA1, KEM1, PAT1, YSH1	14 (59)	1.31E	銀紅印料 (全角目)Lon, GO:0008033:tRNA processing (BP), GO: 0030529:ribonucleoprotein complex (CC), GO:0046540:U4/U6 x U5 trisnRNP complex (CC), GO:0000398:nuclear mRNA splicing, via spliceosome (BP), IPR001163:LikeSm ribonucleoprotein (LSM), GO:0005688:U6 snRNP (CC),
	Metagroup 5; SilhouePe: 0.570 NGG1, HFI1, TRA1, SPT20, SGF29, SPT7, SGF73, SPT8, SGF11, GCN5, SPT3, ADA2, RPN6, CFT2	14 (59)	5.26E	GO:000124:SAGA complex (CC), GO:0016568:chromaln modificalon (BP), GO:0016573:histone acetylalon (BP), GO:0046695:SLIK (SAGAlike) complex (CC), GO:0003702:RNA polymerase II transcriplon factor acluity (MF), GO: 0004402:histone acetyltransferase acluity (MF),

c. proteasome, regulator complex: 12 genes



Gene Metagroups found (include all the related groups)	GENES found interPr	pvalu o)	e CoAnnotaGons (GOBP, GOMF, GOCC, KEGG,
Metagroup 1; SilhouePe: 0.529 GLC7, REF2, YTH1, FIP1, PAP1, PFS2, CFT1, RNA14, PTI1, PTA1, MPE1, CFT2, YSH1	13 (59) 0006353	2.26E26	GO:0005847:mRNA cleavage and polyadenylaLon specificity factor complex (CC), GO:0006378:mRNA polyadenylaLon (BP), GO:0006379:mRNA cleavage (BP), GO:0031123:RNA 3'end processing (BP), GO:
Metagroup 2; SilhouePe: 0.823 CDC23, APC5, CDC16, APC2, APC1, DOC1, APC9, APC11, APC4, CDC27, eD626g起L67, TRA1	11 (59)		04111:Cell cycle – yeast, GO:0004842:ubiquiLnprotein ligase acLvity (MF),GO:0000070:mitoLc sister chromaLd segregaLon (BP),GSG7:protein ubiquiLnaLon(BP),GO:0000022:mitoLcspindle (BP),GO:0005680:anaphasepromoLng complex (CC),
Metagroup 3; SilhouePe: 0.609 RPN13, RPN8, RPN1, RPT1, RPN3, RPN10, RPN11, RPN2, RPN5, RPT3, RPT6, RPN6, APC2, DOC1	14 (59)	8.11E	GO:0000502:proteasome complex (CC), 03050:Proteasome, GO: 0034515:proteasome storage granule (CC), GO:0008541:proteasome regulatory parLcle, lid subcomplex (CC), GO:0004175:endopepLdase acLvity (MF), GO:0030234:enzyme regulator acLvity (MF)
Metagroup 4; SilhouePe: 0.750 LSM5, DCP1, PAP1, LSM3, LSM8, LSM6, LSM1, LSM4, LSM7, LSM2, PTA1, KEM1, PAT1, YSH1	14 (59)	1.31E	03018:RNA degradaLon, GO:0008033:tRNA processing (BP), GO:0030529:ribonucleoprotein complex (CC), GO: 0046540:U4/U6 x U5 trisnRNP complex (CC), GO: 0000398:nuclear mRNA splicing, via spliceosome (BP), IPR001163:LikeSm ribonucleoprotein (LSM), GO: 0005688:U6 snRNP (CC),
Metagroup 5; SilhouePe: 0.570 NGG1, HFI1, TRA1, SPT20, SGF29, SPT7, SGF73, SPT8, SGF11, GCN5, SPT3, ADA2, RPN6, CFT2	14 (59)	5.26E	GO:0000124:SAGA complex (CC), GO:0016568:chromaln modificalon (BP), GO:0016573:histone acetylalon (BP), GO:0046695:SLIK (SAGAlike) complex (CC), GO:0003702:RNA polymerase II transcriplon factor aclvity (MF), GO: 0004402:histone acetyltransferase aclvity (MF),

d. U6 snRNP complex: **11 genes**



Gene Metagroups found (include all the related groups)	GENES found interPr	pvalu o)	e CoAnnotaGons (GOBP, GOMF, GOCC, KEGG,
Metagroup 1; SilhouePe: 0.52 GLC7, REF2, YTH1, FIP1, PAP1, PFS2, CFT1, RNA14, PTI1, PTA1, MPE1, CFT2, YSH1	, 13 (59)	2.26E26	GO:0005847:mRNA cleavage and polyadenylaLon specificity factor complex (CC), GO:0006378:mRNA polyadenylaLon (BP), GO:0006379:mRNA cleavage (BP), GO:0031123:RNA 3'end processing (BP), GO: terminaLon (BP),
Metagroup 2; SilhouePe: 0.82 CDC23, APC5, CDC16, APC2, A DOC1, APC9, APC11, APC4, CI eb626g可见可, TRA1	APC1,		04111:Cell cycle — yeast, GO:0004842:ubiquiLnprotein ligase acLvity (MF),GO:0000070:mitoLc sister chromaLd segregaLon (BP),GO:0000022:mitoLc spindle (BP),GO:00005680:anaphasepromoLng complex (CC),
Metagroup 3; SilhouePe: 0.60 RPN13, RPN8, RPN1, RPT1, RP RPN10, RPN11, RPN2, RPN5, RPT6, RPN6, APC2, DOC1	PN3.	8.11E	GO:0000502:proteasome complex (CC), 03050:Proteasome, GO: 0034515:proteasome storage granule (CC), GO:0008541:proteasome regulatory parLcle, lid subcomplex (CC), GO:0004175:endopepLdase acLvity (MF), GO:0030234:enzyme regulator acLvity (MF)
Metagroup 4; SilhouePe: 0.75 LSM5, DCP1, PAP1, LSM3, LSM LSM6, LSM1, LSM4, LSM7, LS PTA1, KEM1, PAT1, YSH1	И8, 143	1.31E	03018:RNA degradaLon, GO:0008033:tRNA processing (BP), GO: 0030529:ribonucleoprotein complex (CC), GO:0046540:U4/U6 x U5 trisnRNP complex (CC), GO:0000398:nuclear mRNA splicing, via spliceosome (BP), IPR001163:LikeSm ribonucleoprotein (LSM), GO:0005688:U6 snRNP (CC),
Metagroup 5; SilhouePe: 0.57 NGG1, HFI1, TRA1, SPT20, SG SPT7, SGF73, SPT8, SGF11, GG SPT3, ADA2, RPN6, CFT2	GF29, CN5,	5.26E	GO:0000124:SAGA complex (CC), GO:0016568:chromaLn modificaLon (BP), GO:0016573:histone acetylaLon (BP), GO:0046695:SLIK (SAGAlike) complex (CC), GO: 0003702:RNA polymerase II transcripLon factor acLvity (MF), GO:0004402:histone acetyltransferase acLvity (MF),

e. SAGA complex: 12 genes

Validation with other gene/protein sets from multiple biological sources



Use 3 sets of genes/proteins grouped by different biological criteria:

Comprehensive Resource of Mammalian protein complexes
 (CORUM) genes in 10 mammal multi-protein complexes



③ Saccaromyces Genome Database (SGD) genes in 10 known yeast pathways



Online Mendelian Inheritance in Man (OMIM) genes in 10 human diseases
Online Mendelian Inheritance in Man



Validation with other gene/protein sets from multiple biological sources



Use 3 sets of genes/proteins grouped by different biological criteria:

- © Comprehensive Resource of Mammalian protein complexes (CORUM) genes in 10 mammal multi-protein complexes
- ③ Saccaromyces Genome Database (SGD) genes in 10 known yeast pathways
- Online Mendelian Inheritance in Man (OMIM) genes in 10 human diseases

We evaluate the signal considering Proise tolerance (20% noise added): $\frac{Precission}{TP + FP} = \frac{1}{TP + FP}$

③ Recall



	COMPLEXES (from CORUM db, human)	GENES in Ref.	GENES Tested	GENES Found	Common GENES	Precision (%)	Recall (%)	F-score (%)	adjusted p.value	TERMS Found	TERMS Found (only first shown)
1c	C complex spliceosome	80	96	68	68	100.00	85.00	91.89	5.25E-138	6	GO:0005681:spliceosomal
2c	Mediator (transcriptional coactivator) complex	32	39	28	28	100.00	87.50	93.33	3.96E-064	10	GO:0016592:mediator complex
3c	Proteasome (20S/26S)	22	27	22	22	100.00	100.00	100.00	5.34E-063	12	03050:Proteasome
4c	RNA polymerase II (RNAPII)	26	32	24	24	100.00	92.31	96.00	1.21E-059	12	GO:0006350:transcription
5c	F1F0-ATP synthase (EC 3.6.3.14), mitochondrial	16	20	14	14	100.00	87.50	93.33	2.63E-045	12	GO:0005753:mitochondrial ATPase
6c	DAB complex, transcription preinitiation complex	16	20	16	16	100.00	100.00	100.00	3.20E-042	6	03022:Basal transcription factors
7c	Exosome	11	14	11	11	100.00	100.00	100.00	2.47E-040	4	GO:0006364:rRNA processing
8c	eIF3 complex, eukaryotic initiation of translation factor-3	13	16	11	11	100.00	84.62	91.67	8.98E-038	4	GO:0005852:eukar. translation initiation factor .
9с	Nup 107-160 nuclear pore subcomplex	9	11	9	9	100.00	100.00	100.00	1.34E-033	6	GO:0005635:nuclear envelope
10c	CENP-A NAC-CAD kinetochore complex	13	16	13	13	100.00	100.00	100.00	2.32E-028	2	GO:0005694:chromosome
	with 20% noise (% of random genes included)		+20% noise	aver	age values=	100.00	93.69	96.62			
	DISEASES (from OMIM db, human)	GENES in Ref.	GENES Tested	GENES Found	Common	Precision (%)	Recall (%)	F-score (%)	adjusted p.value	TERMS Found	TERMS Found (only first shown)
1d	Retinitis pigmentosa	51	62	39	38	97.44	74.51	84.44	1.15E-066	1	GO:0007601:visual perception (BP)
2d	Deafness (autosomal dominant and recesive)	84	101	31	31	100.00	36.90	53.91	4.11E-053	6	GO:0007605:sensory perception of sound
3d	Cardiomyopathy (dilated, familial and hypertrophic)	44	53	19	19	100.00	43.18	60.32	1.48E-031	5	GO:0008307:structural constituent of muscle
4d	Epidermolysis bullosa (dystrophica, simplex, junctional)	11	14	11	11	100.00	100.00	100.00	2.77E-024	2	GO:0031581:hemidesmosome assembly
	Congenital disorder of glycosylation (type I and II)	23	28	11	11	100.00	47.83	64.71	1.18E-023	1	00510:N-Glycan biosynthesis
	Muscular dystrophy (congenital, limb-girdle, rigid spine)	25	30	11	11	100.00	44.00	61.11	4.45E-019	4	GO:0007517:muscle organ development
7d	Glycogen storage disease	19	23	9	9	100.00	47.37	64.29	2.53E-018	3	00500:Starch and sucrose metabolism
	Leigh syndrome	8	10	7	7	100.00	87.50	93.33	1.16E-017	10	GO:0006120:mitochondrial electron transport
	Acute Leukemia (lymphoblastic ALLs & myeloid AMLs)	37	45	9	9	100.00	24.32	39.13	7.71E-016	2	05221:Acute myeloid leukemia
10d	Diabetes mellitus (type 1 or 2, gestational, neonatal)	13	16	5	5	100.00	38.46	55.56	1.01E-010	2	GO:0005975:carbohydrate metabolic process
	with 20% noise (% of random genes included)		+20% noise	aven	age values=	99.74	54.41	67.68			· ·
	PATHWAYS (from SGD db, yeast)	GENES in Ref.	GENES Tested	GENES Found	Common GENES	Precision (%)	Recall (%)	F-score (%)	adjusted p.value	TERMS Found	TERMS Found (only first shown)
	gluconeogenesis	22	27	22	22	100.00	100.00	100.00	7.05E-047	7	GO:0006094:gluconeogenesis
	TCA cycle, aerobic respiration	22	27	23	22	95.65	100.00	97.78	2.36E-037	10	00020:Citrate cycle (TCA cycle)
3p	sphingolipid metabolism	19	23	16	16	100.00	84.21	91.43	7.64E-033	4	00600:Sphingolipid metabolism
4p	de novo biosynthesis of purine nucleotides	35	42	21	21	100.00	60.00	75.00	2.02E-030	4	00230:Purine metabolism
5p	lipid-linked oligosaccharide biosynthesis	12	15	12	12	100.00	100.00	100.00	1.38E-029	5	GO:0005783:endoplasmic reticulum
6р	ergosterol biosynthesis	11	14	11	11	100.00	100.00	100.00	1.02E-028	7	GO:0006696:ergosterol biosynthetic process
7p	superpathway of glucose fermentation	14	17	13	13	100.00	92.86	96.30	9.18E-027	5	00010:Glycolysis / Gluconeogenesis
8р	fatty acid biosynthesis, initial steps	17	21	13	12	92.31	70.59	80.00	9.93E-027	6	GO:0006631:fatty acid metabolic process
	inositol phosphate biosynthesis	19	23	11	11	100.00	57.89	73.33	1.02E-025	3	00562:Inositol phosphate metabolism
9p									0.000 040		00.0000700
	folate biosynthesis	18	22	10	9	90.00	50.00	64.29	6.98E-019	6	GO:0006730:one-carbon metabolic process





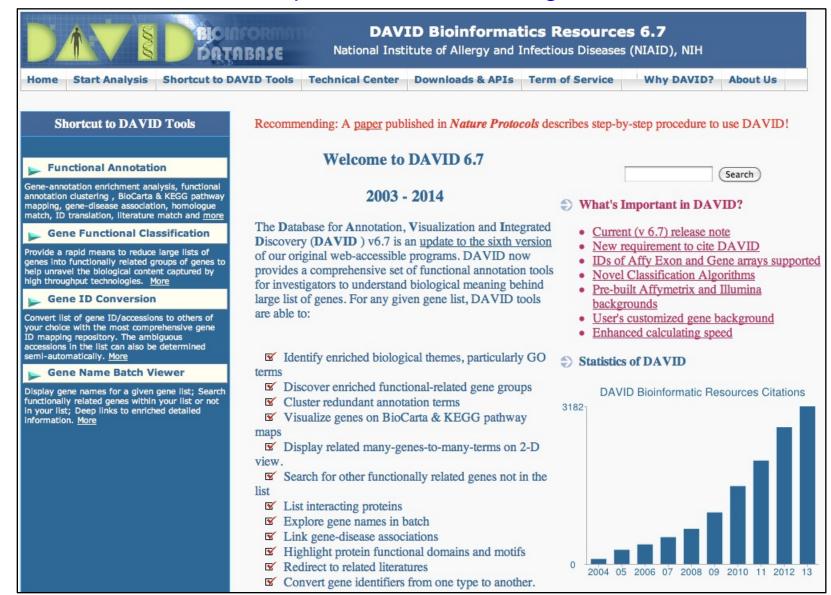


Session 1 (9:30 - 12:30, 3h)
Bioinformatic tools for Functional Enrichment Analysis (FEA)
Session 2 (13:30 - 16:30, 3h)
Construction of gene functional networks

- Introduction to biological information and annotation spaces: GO, KEGG, Interpro
- Functional Enrichment Analysis (EA): from single to modular methods
 - Using EA tools to annotate gene lists:
 DAVID (single), GSEA (gene sets), GeneCodis (modular)
 - Sort out problems after EA: post-enrichment tool
 GeneTermLinker (postEA) vs DAVID FAC (single+clustering)
 - From co-annotation and enrichment to functional networks: networks construction using a R tool



http://david.abcc.ncifcrf.gov/







http://david.abcc.ncifcrf.gov/

Huang et al (2007) Genome Biol.

Open Access

Software

The DAVID Gene Functional Classification Tool: a novel biological module-centric algorithm to functionally analyze large gene lists Da Wei Huang**, Brad T Sherman**, Qina Tan*, Jack R Collins†, W Gregory Alvord*, Jean Roayaei*, Robert Stephens†, Michael W Baseler§, H Clifford Lane¶ and Richard A Lempicki*

Addresses: *Laboratory of Immunopathogenesis and Bioinformatics, Clinical Services Program, SAIC-Frederick, Inc., National Cancer Institute

at Frederick, Frederick, Frederick, Frederick, M Frederick, MD 21702, U USA. *Laboratory of Im 20892, USA.

PROTOCOL

Huang et al (2009) Nature Protoc.

Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources

Da Wei Huang^{1,2}, Brad T Sherman^{1,2} & Richard A Lempicki¹

¹Laboratory of Immunopathogenesis and Bioinformatics, Clinical Services Program, SAIC-Frederick Inc., National Cancer Institute at Frederick, Frederick, Maryland 21702, USA. ²These authors contributed equally to this work. Correspondence should be addressed to R.A.L. (rlempicki@mail.nih.gov) or D.W.H. (huangdawei@mail.nih.gov)

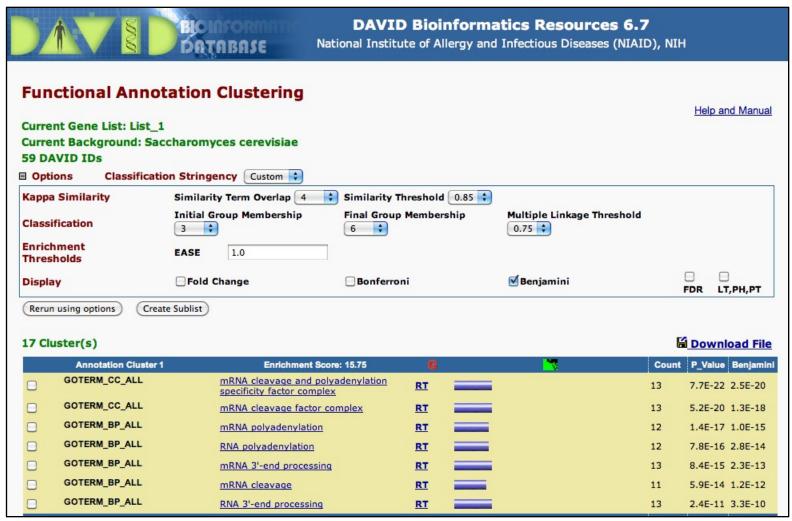
Published online 18 December 2008; doi:10.1038/nprot.2008.211

DAVID bioinformatics resources consists of an integrated biological knowledgebase and analytic tools aimed at systematically extracting biological meaning from large gene/protein lists. This protocol explains how to use DAVID, a high-throughput and integrated data-mining environment, to analyze gene lists derived from high-throughput genomic experiments. The procedure first requires uploading a gene list containing any number of common gene identifiers followed by analysis using one or more text and pathway-mining tools such as gene functional classification, functional annotation chart or clustering and functional annotation table. By following this protocol, investigators are able to gain an in-depth understanding of the biological themes in lists of genes that are enriched in genome-scale studies.



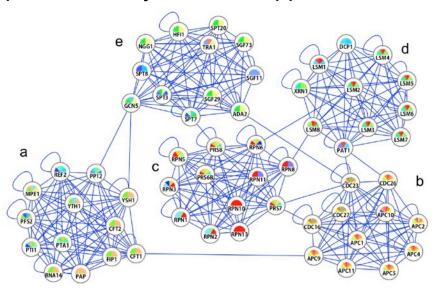
DAVID bioinformatics suite provides a **post-enrichment** tool:

DAVID Functional Annotation Clustering (DAVID - FAC), after single enrichment analysis it clusters the terms based in the genes that they share.





Comparative analysis of two approaches that provide groups of genes and terms



Testing a set of **59** yeast nuclear **proteins** working in **5** known **complexes**

	GeneTerm Linker	DAVID FAC (by default)	DAVID FAC (tuned to find 5 groups)
Total groups reference	5	5	5
Total groups found	5	15	5
Accuracy	0.952	0.311	0.884
Jaccard coefficient	0.769	0.213	0.562



	GeneTerm Linker	DAVID FAC (by default)	DAVID FAC (tuned to find 5 groups)
Total groups reference	5	5	5
Total groups found	5	15	5
Accuracy	0.952	0.311	0.884
Jaccard coefficient	0.769	0.213	0.562

Accuracy: percentage of genes grouped correctly.

$$Accuracy = \frac{(TP + TN)}{TP + TN + FP + FN}$$

Jaccard coefficient: (measures similarity between sample sets) proportion of genes that belong to the same complex with respect to the total number of genes that are known to belong such complex or to any other.

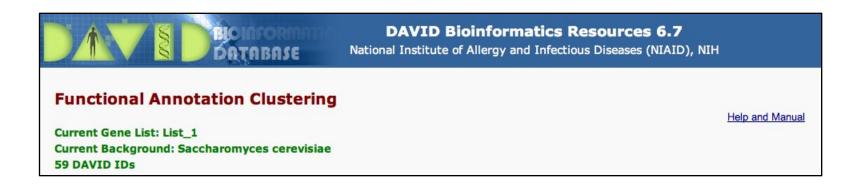
$$Coefficiente\ de\ Jaccard = \frac{TP}{TP + FP + FN}$$

96



DAVID bioinformatics suite provides a **post-enrichment** tool:

DAVID Functional Annotation Clustering (DAVID - FAC), after single enrichment analysis it clusters the terms based in the genes that they share.



DAVID Functional Annotation Clustering (DAVID - FAC)

- 1.It groups/clusters genes or terms in an independent way
- 2.It does not considers the enrichment *p-value* of each gene-term set to rank them during the clustering
- 3. It does not eliminates redundancies







Hands-on: Practical Examples

Protein_ SETs_ 2014.xls (106g hs, 175g hs, yeast 11pathways, yeast 59g5pc)

run DAVID-FAC & GeneTerm Linker







Session 1 (9:30 - 12:30, 3h)
Bioinformatic tools for Functional Enrichment Analysis (FEA)
Session 2 (13:30 - 16:30, 3h)
Construction of gene functional networks

- Introduction to biological information and annotation spaces: GO, KEGG, Interpro
- Functional Enrichment Analysis (EA): from single to modular methods
 - Using EA tools to annotate gene lists:
 DAVID (single), GSEA (gene sets), GeneCodis (modular)
 - Sort out problems after EA: post-enrichment tool GeneTermLinker (postEA)
 - From co-annotation and enrichment to functional networks:
 networks construction using a R tool

GeneTerm Linker web tool



Home | Help | Webservice

. Input a list of genes of interest:	3. Organism:
	Homo sapiens
	4. Annotation Spaces:
luman example] [Yeast example]	 GO Biological Process GO Molecular Function GO Cellular Component KEGG Pathways InterPro Motifs And Domains
. Input a list of genes of reference (optional):	5. Minimum Support: 4 💠
	6. Email address (optional):
	Submit analysis Reset

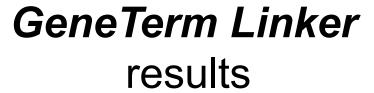
[Medline] [Online version]

GeneTerm Linker results



FuncLonal analysis results for a **100 genes** signature of **Acute LinfoblasGc Leukemia** (AML)

	Genes \$	#list \$	#ref_list \$	adjusted ¢ pValue	Silhouette Width	Terms
Metagroup 3 Show details			104(34208)	1.04382e-05	0.8589	GO:0010628 positive regulation of gene expression (BP)
Metagroup 2 Show details	BMP2 CMTM8 CXCL2 CXCL3 GNAI1 PF4 PPBP TNFSF4	8(94)	310(34208)	GO:0006935 chemotaxis (BP) GO:0008009 chemokine activity (MF) 04060 Cytokine-cytokine receptor interaction 04062 Chemokine signaling pathway IPR001811 Chemokine interleukin-8-like domain IPR002473 CXC chemokine, interleukin 8 IPR001089 CXC chemokine		GO:0008009 chemokine activity (MF) 04060 Cytokine-cytokine receptor interaction 04062 Chemokine signaling pathway IPR001811 Chemokine interleukin-8-like domain IPR002473 CXC chemokine, interleukin 8
Metagroup 4 Show details	GNAI1 LRRK2 RAB31 RAB32 RASD1	5(94)	165(34208)	9.53831e-05	0.7143	GO:0003924 GTPase activity (MF) IPR005225 Small GTP-binding protein domain IPR001806 Ras GTPase IPR003579 Ras small GTPase, Rab type
Metagroup 5 Show details	COL5A1 CTGF ECM1 LGALS3 PXDN	5(94)	272(34208)	0.000944947	0.7112	GO:0031012 extracellular matrix (CC) GO:0005578 proteinaceous extracellular matrix (CC)
Metagroup 1 Show details	LRRK2 MME PSD3 SHANK3 STXBP5 SYT1	6(94)	77(34208)	7.41975e-08	0.4167	GO:0008021 synaptic vesicle (CC) GO:0045202 synapse (CC)





FuncLonal analysis results for a 100 genes signature of Acute LinfoblasGc Leukemia

(AML)

(AIVIL)				
Genes	#list	#ref_list	adjusted pValue	Terms
CXCL2 CXCL3 PF4 PPBP	4(94)	9(34208)	8.59786e-07	GO:0006955 immune response (BP) GO:0006935 chemotaxis (BP) GO:0008009 chemokine activity (MF) GO:0005615 extracellular space (CC) GO:0005576 extracellular region (CC) 04060 Cytokine-cytokine receptor interaction 04062 Chemokine signaling pathway IPR001811 Chemokine interleukin-8-like domain IPR002473 CXC chemokine, interleukin 8 IPR001089 CXC chemokine
CMTM8 CXCL2 CXCL3 PF4 PPBP	5(94)	65(34208)	1.8799e-05	GO:0006935 chemotaxis (BP) GO:0005615 extracellular space (CC)
CXCL2 CXCL3 PF4 PPBP TNFSF4	5(94)	76(34208)	3.19339e-05	GO:0006955 immune response (BP) GO:0005615 extracellular space (CC) 04060 Cytokine-cytokine receptor interaction
BMP2 CXCL2 CXCL3 PF4 PPBP TNFSF4	6(94)	146(34208)	4.20382e-05	GO:0005615 extracellular space (CC) 04060 Cytokine-cytokine receptor interaction
CXCL2 CXCL3 GNAI1 PF4 PPBP	5(94)	186(34208)	0.000567261	04062 Chemokine signaling pathway

GeneTerm Linker results



Functional analysis results for a 100 genes signature of Acute LinfoblasGc Leukemia

(AME) tatistics of the analysis

Total number of genes sent to the analysis: 100

Number of genes recognized by the enrichment tool: 94

Total number of **genes** not-included in the enrichment analysis: 13 (show/hide)

Genes not present in the metagroups, only annotated to generic terms: 56 (show/hide)

Total number of **genes** included in the metagroups: 31 (show/hide)

Total number of significant **GeneTerm-sets** provided by the enrichment analysis: 98

Download the initial Enrichment Analysis (text format)

Number of filtered GeneTerm-sets that include generic terms in the enrichment: 80

Download the GeneTerm-sets filtered as generic (text format)

Total number of **GeneTerm-sets** of the enrichment used: 18

NOTE: Maximum number of **GeneTerm-sets** analysed in the web = 1000

Number of redundant GeneTerm-sets of the enrichment filtered: 2

Final number of **GeneTerm-sets** of the enrichment included in the metagroups: 16





The results of the gene---term **metagroups** generated by *GeneTerm Linker* can be used in a further analysis to build **funcGonal networks**

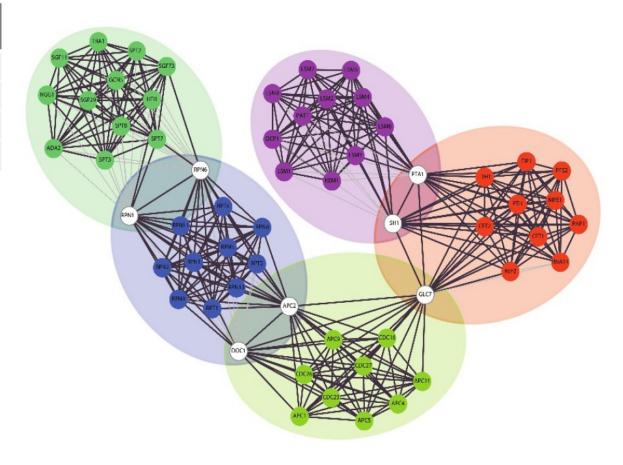
	gts 1	gts 2	gts 3	gts 4
g1	0	0	1	1
g2	0	0	1	0
g3	1	1	0	1
g4	1	0	1	0

Common gene term sets adjacency matrix (g :: gts)



	Mg 1	Mg 2	Mg 3	Mg 4
g1	0	1	0	1
g2	1	0	1	0
g3	0	1	0	1
g4	1	0	1	0

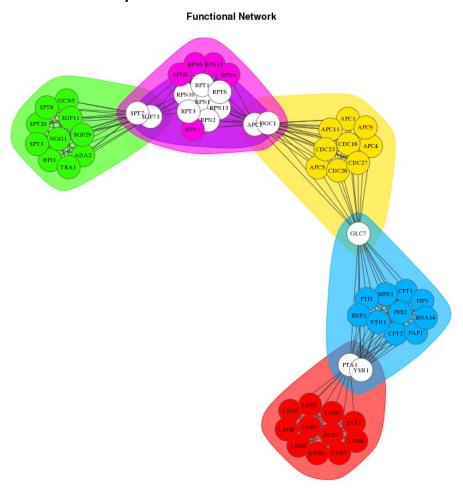
Common **metagroups** adjacency matrix (g :: Mg)







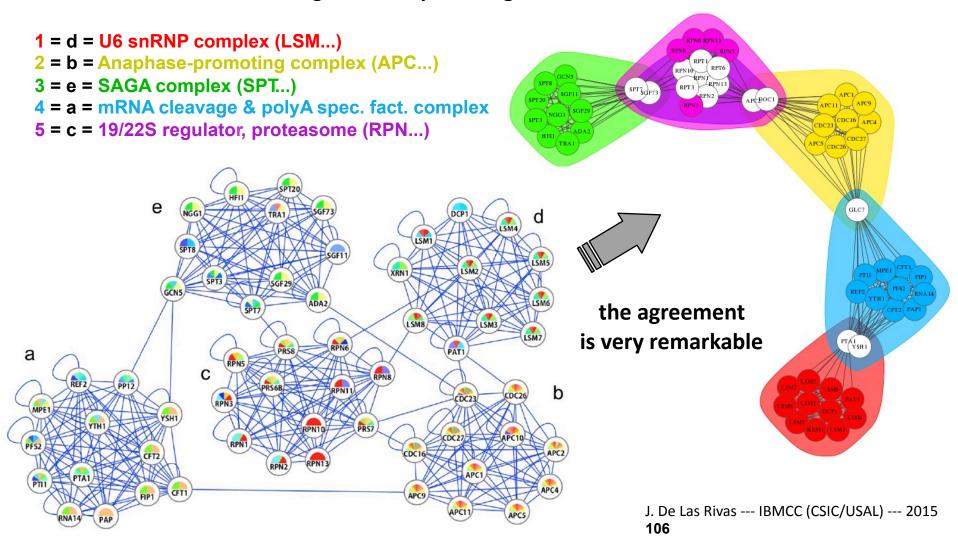
The results of the gene---term **metagroups** generated by *GeneTerm Linker* can be used in a further analysis to build **funcGonal networks**



Functional Network derived from GeneTerm Linker



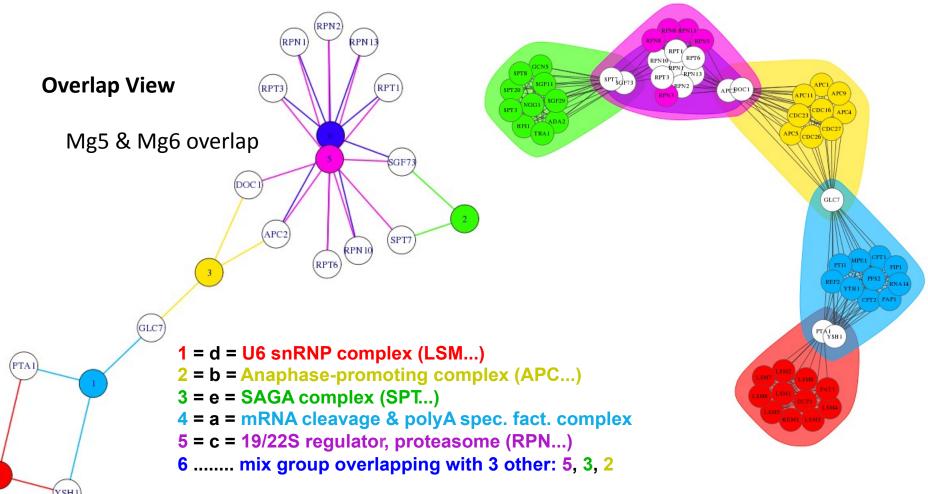
Build *Func&ondNetworks*:: from an experiment---based **protein interacGon network** to a knowledge---based **protein/gene funcGonal network**



Functional Network derived from GeneTerm Linker



Build *Func&ondNetworks*:: from an experiment---based **protein interacGon network** to a knowledge---based **protein/gene funcGonal network**



Functional Network derived from GeneTerm Linker



Build Func&ondNetworks :: from an experiment---based protein interacGon network to a knowledge---based **protein/gene funcGonal network**

