Protein function prediction

(focusing on methods not based on the detection of sequence similarity)

Manuel J. Gómez

Bioinformatics Lab Centro de Astrobiología, CSIC-INTA

ONTOLOGIES IN MOLECULAR BIOLOGY

Controlled and structured vocabularies, constructed with TWO purposes: proposing standard collections of terms, and, organizing the knowledge of a given field around its language: the relations between the terms are supposed to reflect the biological reality.

• Enzyme Commission Nomenclature.

EC 1	Oxidoreductases.
EC 1. 1	Acting on the CH-OH group of donors.
EC 1. 1. 1	With NAD(+) or NADP(+) as acceptor.
EC 1. 1. 2	With a cytochrome as acceptor.

- MeSH (Medical Subject Headings) terms: NLM controlled vocabulary.
- Gene Ontology: developed and maintained by a consortium (GO Consortium) of laboratories and institutions involved in molecular biology database management.

ONTOLOGIES IN MOLECULAR BIOLOGY



GENE ONTOLOGY

Biological terms have been groupd in three ontologies:

- Molecular functions
- Cellular processes
- Cellular components

Within each ontology, the terms are related hierarchically.

The relation between parent and child terms can be of two types:

- Part of
- Instance of

Most molecular biology databases have joined this initiative, and have included annotations following this standard.



Gene Ontology

Mycoplasma pneumoniae proteome. Distribution of GO anotations.

GO Classification for M. pneumoniae							
Term	Proteins						
<u>GO:0003674</u>	431	62.7%					
<u>GO:0003676</u>	nucleic acid binding	119	17.3%				
<u>GO:0030528</u>	transcription regulator activity	5	0.7%				
<u>GO:0003754</u>	chaperone activity	10	1.4%				
<u>GO:0003824</u>	catalytic activity	245	35.6%				
<u>GO:0015070</u>	toxin activity	1	0.1%				
<u>GO:0005194</u>	cell adhesion molecule activity	1	0.1%				
<u>GO:0005198</u>	structural molecule activity	56	8.1%				
<u>GO:0005215</u>	transporter activity	64	9.3%				
<u>GO:0005488</u>	binding	211	30.7%				
<u>GO:0005554</u>	$molecular_function$ unknown	65	9.4%				
<u>GO:0008150</u>	biological_process	343	49.9%				
<u>GO:0008152</u>	metabolism	280	40.7%				
<u>GO:0006810</u>	transport	65	9.4%				
<u>GO:0006950</u>	response to stress	16	2.3%				
<u>GO:0007049</u>	cell cycle	22	3.2%				
<u>GO:0007154</u>	cell communication	16	2.3%				
<u>GO:0007275</u>	development	5	0.7%				
<u>GO:0007582</u>	physiological process	335	48.7%				
<u>GO:0005575</u>	cellular_component	277	40.3%				
<u>GO:0005576</u>	extracellular	1	0.1%				
<u>GO:0005623</u>	cell	274	39.8%				
<u>GO:0005941</u>	unlocalized	5	0.7%				

FUNCTION PREDICTION PROTOCOL: SEQ SIM

Based on the transfer of function from proteins that are similar at sequence level





GeneQuiz (EBI) (Andrade et al 1999)

PEDANT (Andrade et al 1999)

MAGPIE (Gaasterland et al, 1996)

ENSEMBL (Clamp et al, 2003)

RiceGAAS (Sakata el al, 2002)



DISTRIBUTION OF DIFFERENT LEVELS OF ANNOTATION IN MICROBIAL GENOMES





Manuel J. Gómez

FUNCTION PREDICTION PROTOCOL: SEQ SIM + STRUCT SIM

Based on sequence similarity and structural analyses, mainly to identify remote homologs



Known / Predicted structure

DETECTION OF SIMILARITY OF PROTEIN 1D FEATURES

- Several protein features that can be calculated or predicted from the primary sequence are used frequently to get hints about protein function.
- These may include:
 - Amino acid composition and isoelectric point.
 - Presence of signal peptides and trans-membrane segments.
 - Secondary structure elements.
 - Post translational modification sites.
- The challenge is to devise methods to find correlations between these properties and protein functions.

Around 5500 human protein sequences, classified according to: •EC number •Euclid

•Gene Ontology

Prediction or calculation of features that can be derived directly from the primary sequence

 \rightarrow

Training of a Neural Network in the association between features and function

Jensen et al. JMB (2002)

Abbriviation	Encoding	Description
ec	single value	Extinction coefficient predicted by <u>ExPASy ProtParam</u>
gravy	single value	Hydrophobicity predicted by <u>ExPASy ProtParam</u>
meg	single value	Number of negatively charged residues counted by <u>ExPASy ProtParam</u>
npos	single value	Number of positively charged residues counted by <u>ExPASy ProtParam</u>
nglyc	potential in 5 bins	N–glycosylation sites predicted by <u>NetNGlyc</u>
oglyc	potential-threshold in 10 bins	GalNAc O – glycosylations predicted by <u>NetOGlyc</u>
pest	fraction in 10 bins	PEST rich regions identified by <u>PESTfind</u>
phosST	potential in 10 bins	Serine and threonine phosporylations predicted by <u>NetPhos</u>
phosY	potential in 10 bins	Tyrosine phosporylations predicted by <u>NetPhos</u>
psipred 💋	helix, sheet, coil in 5 bins	Predicted secondary structure from <u>PSI – Pred</u>
psort	20 probabilities	Subcellular location predtions by <u>PSORT</u>
signalp	meanS, maxY, log(cleavage pos)	Signal peptide predictions made by <u>SignalP</u>
tmhmm	inside, outside, membrane in 5 bins	Transmembrane helix predictions made by <u>TMHMM</u>

DETECTION OF SIMILARITY OF PROTEIN 1D FEATURES



DETECTION OF SIMILARITY OF PROTEIN 1D FEATURES

(a) 1.0 Amino acid biosynthesis Biosynthesis of cofactors Cell envelope 0.8 Cellular processes Central intermediary metabolism Energy metabolism Fatty acid metabolism False positive rate 0.6 Purines and pyrimidines **Regulatory functions** Replication and transcription Translation Transport and binding 0.4 0.2 0.0 0.8 1.0 0.2 0.4 0.6 0.0Sensitivity

Sensitivity vs specificity in the prediction of functional roles, according to Euclid

Manuel J. Gómez

METHODS BASED ON ANALYSES OF PROTEIN STRUCTURE

Several methods have been proposed or used:

- Similarity, at a general structural level, to proteins of known function (detection of remote homologs by fold recognition).
- Similarity at local structural level (common structure of the active site, in enzymes that are not related)
- Prediction of binding sites (cavities) followed by prediction of ligands with docking algorithms.
- Prediction of interactions with other proteins by analysis of protein surfaces to identify potential interaction sites, followed prediction of interaction partners by docking methods.

The utility of the last three types of methods to predict function is, for the moment, anecdotical.

Eisenstein Curr Op Biotec (2000)

DISTRIBUTION OF DIFFERENT LEVELS OF ANNOTATION IN MICROBIAL GENOMES





Manuel J. Gómez

FUNCTION PREDICTION PROTOCOL: SEQ & STRUCT SIM + INT

Based on sequence similarity, structural analyses and information about interacting partners.



- In general the prediction of function based on the detection of functional interactions follows the concept presented often as "guilty by association".
- In the absence of any detectable indication of homology:
 - If protein A is a regulator of, or it is regulated by, protein B;
 - or protein A interacts physically with protein B;
 - or protein A and protein B are co-expressed;

then, we can assume that the function of both proteins is also related.

FUNCTION PREDICTION METHODS



FUNCTIONAL RELATION OR INTERACTION

We can consider four sources of information about functional interactions:

- Functional Genomics
- Proteomics
- Text mining (information extraction)
- Comparative sequence genomics

Functional Genomics



Clustering of genes according to their pattern of expression: UPGMA versus SOTA

J. Dopazo

Bioinformatics Lab. CAB, CSIC-INTA





Large scale detection of protein interactions by the Yeast 2 hybrid system





C. elegans Walhout et al. (2000) Science 287:116-122

S. cerevisiae Uetz et al. (2000) Nature 403:623-631 Ito et al. (2000) PNAS 97:1143-1147

METHODS BASED ON THE DETECTION **Proteomics: OF FUNCTIONAL INTERACTIONS** TAP, Tandem 55 **Affinity Purification** а BP)-TEV site-(Protein Spacer-PCR product 56 Homologous Gene recombination 57 targeting Chromosome Gene Gene 109 Fusion NH₂ Protein 60 - Spacer-CBP)-TEV site-(Protein En montene Ortopiden Milochon Nucleus Vacuolar Membra 000 212 STO1 ADES MAS Mapl10 3cale Play ma> eres ECM16

CTR9

Cellzome, **Nature 2002**

Text mining (information extraction)

The current list of rules includes .protein [word]* [verb] [word]* protein .[verb] of [word]* protein [word]* [by,to] [word]* protein .[noun] of [word]* protein [word]* [by,with] [word]* protein .[noun] between [word]* protein [word]* and [word]* protein .protein [word]* protein [word]* [complex/es, dimer, heterodimer] .complex formed between [word]* protein [word]* and [word]* protein .complex/es of [word]* protein [word]* and [word]* protein .protein [word]* forms a complex with [word]* protein

Sub-rules are used for incorporating particular cases

1.1.	protein [word]* [verb] [word]* but not [word]* protein
1.2.	protein [word]* cannot [word]* [verb] [word]* protein
1.3.	protein [word]* does not [word]* [verb] [word]* protein
1.4.	protein [word]* did not [word]* [verb] [word]* protein
1.5.	protein [word]* was not [word]* [verb] [word]* protein
1.6.	protein [word]* not [word]* [verb] [word]* by [word]* protein
1.7.	protein [word]* not required for [word]* [verb] [word]* protein
1.8.	protein [word]* failed to [word]* [verb] [word]* protein

Rules to identify information about interactions in scientific texts.

E. coli INTERACTIONS DATABASE PROTOTYPE

Interaction networks identified by extracting information from the literature.



by C. Blaschke

Comparative Sequence Genomics

- Methods based on the analysis of evolutionary information stored in sequences.
- There are FIVE methods, so far
 - Conservation of gene context or neighbourhood
 - Detection of gene fusions
 - Similarity of gene trees
 - Identification of correlated mutations
 - Similarity of phylogenetic profiles

Conservation of gene context or neighborhood



Manuel J. Gómez

Detection of gene fussions, to infer physical interactions



Manuel J. Gómez



Detection of correlated mutations



Similarity of phylogenetic profiles



NON ORTHOLOGOUS GENE DISPLACEMENT

Anti-correlation of phylogenetic profiles



Morett et al, 2003

DATABASES OF PROTEIN INTERACTIONS

Predicted interactions, by comparative genomics (integration of several methods):

EcID

Predictome

STRING

Prolinks

Experimentally detected protein interactions:

DIP BIND MIPS

PIMrider

Experimentally detected and predicted interactions: IntAct (EBI).

Netscape: The E Coli Predicted Protein Interactions Database. File Edit View Go Bookmarks Options Directory Window Help Back Forward Home Edit Reload Print Find Company Print Print Print Print Company Print P									
An application of the "in silico two-l entire E Coli genome. Predicted Interaction partners for <u>G1786676 DP3X_ECOLI</u> <u>P06710</u> DNA POLYMERASE III SUBUNITS GAMMA AND TAU (EC 2.7.7.7) [GENE: DNAX AND DNAZ]									
[Home] [Info] [Contact] [Search Hel]	Z-Score	Database Links	Protein	Run ECID for it					
	5.763	G1790860 CREB_ECOLI P08368	TRANSCRIPTIONAL REGULATORY PROTEIN CREB [GENE: CREB]	>>					
Enter the E. Coli protein you want to search i	3.687	G1789556 RS15_ECOLI P02371	30S RIBOSOMAL PROTEIN S15 [GENE: RPSO OR SECC]	>>					
Search for polymerase as	3.625	G1790408 BIRA_ECOLI P06709	BIRA BIFUNCTIONAL PROTEIN (BIOTIN OPERON REPRESSOR) (BIOTIN[ACETYL- COA-CARBOXYLASE] SYNTHETASE) (EC 6.3.4.15) (BIOTINPROTEIN LIGASE) [GENE: BIRA OR BIOR OR DHBB]	>>					
Last update: xxx/xxx/xxx.	3.440	G1789701 RL6_ECOLI P02390	50S RIBOSOMAL PROTEIN L6 [GENE: RPLF]	>>					
7-0	3.184	G1789691 RS4_ECOLI P02354	30S RIBOSOMAL PROTEIN S4 [GENE: RPSD OR RAMA]	>>					
	3.119	G2367268 THDF_ECOLI P25522	THIOPHENE AND FURAN OXIDATION PROTEIN THDF [GENE: THDF OR TRME]	>>					
	2.900	<u>G2367200</u> TRUB_ECOLI P09171	TRNA PSEUDOURIDINE SYNTHASE B (EC 4.2.1.70) (TRNA PSEUDOURIDINE 55 SYNTHASE) (PSI55 SYNTHASE) (PSEUDOURIDYLATE SYNTHASE) (URACIL HYDROLYASE) (P35 PROTEIN) [GENE: TRUB OR P35]	>>					
	2.608	G1787357 MFD_ECOLI P30958	TRANSCRIPTION-REPAIR COUPLING FACTOR (TRCF) [GENE: MFD]	>>					

by F. Pazos, 1999

Bioinformatics Lab. CAB, CSIC-INTA

STRING: a database of predicted functional associations between proteins.

Von Mering et al. (2003) Nucleic Acids Research, 31 (1):258.

Integrates information about:

- Phylogenetic profiles
- Conservation of gene context
- Existence of gene fusions

With that information, functional associations between COGs are established. The functional associations may be interpreted as physical or regulatory interactions or participation in the same metabolic pathway.

The network of interactions is huge and may include all COGs, depending of the stringency used to calculate it.

In any case, that network can be broken down in modules that may represent individual metabolic pathways, regulatory circuits or protein complexes.



Manuel J. Gómez

This presentation contains material from:

Christian Blaschke, BioAlma Florencio Pazos, CNB Alfonso Valencia, CNB Manuel J Gómez, CAB