

## Exploring the Semantics of Biomedical Ontologies



Francisco M Couto  
12 April 2011  
EBI External Seminar

## How to compare?

- Manually
  - Before computers
  - based on catalogues
- Problems
  - Digital Era
  - Information overload
  - UniProtKB/TreMBL (14 million sequences)



## New Entities

- Common step
  - **Compare** it with known entities
  - to transfer knowledge
- Common to any research area
  - Characterization of new proteins
  - Microarray Analysis
  - Information Retrieval in general



## Computational Comparison

- Using digital representations
- Structural Similarity
  - Based on the syntax of the representations
  - Protein sequence (BLAST)
- Semantic Similarity
  - Based on the implicit and explicit semantics (Ontologies)
  - Protein molecular function (GO Analysis)

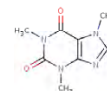


## Structural Similarity

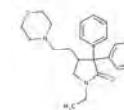
- Structure is “always” available
  - We normally start by identifying
    - how an entity look like
    - before finding what it does (meaning)
- Structure is less ambiguous
  - It’s easier to establish an agreement on describing
    - how an entity look like
    - than on what it does
- Performance and Scalability
  - String matching techniques (BLAST)



## Example: Similar Semantics



Caffeine (CHEBI:27732)



Doxapram (CHEBI:681849)

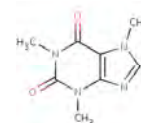
- Different structure but
- both central nervous system stimulants (CHEBI:35337)

## Semantic Similarity

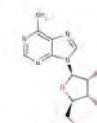
- Solve the problems of structural similarity
  - Two entities that look similar
    - **Do not imply**
    - They have the same meaning
  - And two entities with similar meaning
    - **Do not imply**
    - that they look similar
- When we want to find similar entities
  - Based on their meaning
  - **Not on** how they look like



## Example: Similar Structure



caffeine (CHEBI:27732)



adenosine (CHEBI:16335)

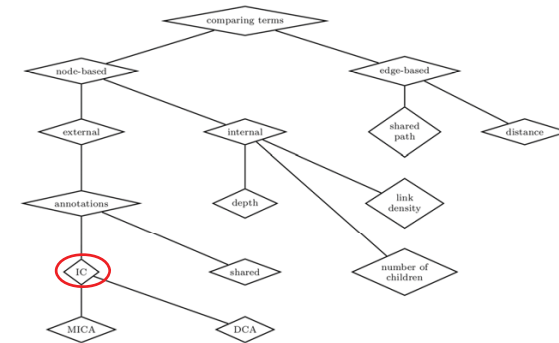
- Similar structure but
- different roles
  - adenosine is an anti-arrhythmia drug (CHEBI:38070)
  - nucleoside (CHEBI:18254)

## Semantic description

- Is not always available
- Semantic Web Initiative
  - inserting machine-readable metadata
- Ontologies
  - Serve as metadata vocabularies
  - to describe entities (annotations)
- Biomedical Ontologies
  - Large effort and involvement from the community
  - 101 ontologies at <http://www.obofoundry.org/>
  - 82 million annotations from UniProtKB-GOA



## Comparing Concepts

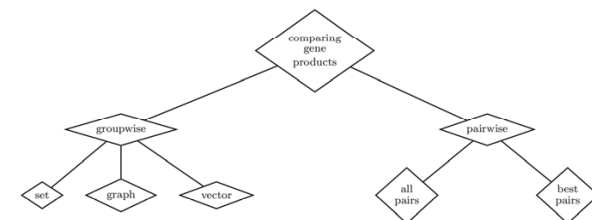


## Semantic similarity measures

- Input:
  - two ontology concepts
  - or two sets of terms annotating two entities
- Output:
  - a numerical value reflecting the closeness in meaning between them



## Comparing entities



## Information Content

- measures how specific and informative a concept is

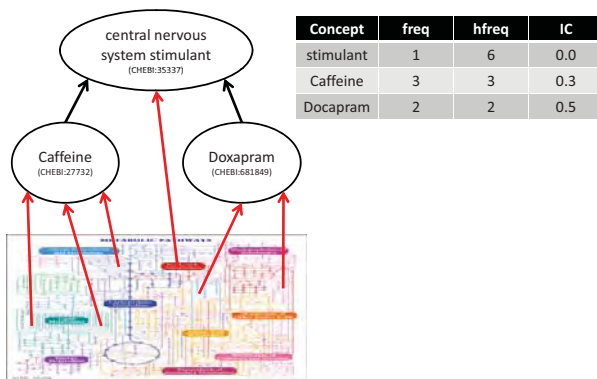
$$IC(c) = -\log\left(\frac{freq(c)}{maxFreq}\right)$$

- Inversely proportional to frequency in a given corpus
- The frequency is also propagated to its ancestors
  - IC proportional to the depth of a concept
- Extrinsic IC
  - number of entities mapped to each concept
- Intrinsic IC
  - number of children

## Similarity based on IC

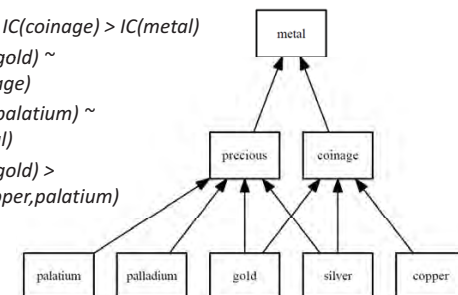
- Similarity proportional to the
  - IC of the most informative common ancestor (MICA)
    - Shared information between two concepts
    - Resnik
  - Weighted Jaccard index of two sets of concepts
    - Shared information between two entities
    - SimGIC

## Example



## Example

- $IC(copper) > IC(coinage) > IC(metal)$
- $Sim(copper, gold) \sim IC(coinage)$
- $Sim(copper, palatium) \sim IC(metal)$
- $Sim(copper, gold) > Sim(copper, palatium)$



## Applications

- Chemical Classification
- Chemical entity recognition and mapping
- Ontology matching and extension
- Enzyme family coherency assessment
- Epidemic and VPH information retrieval

## Classification Problems

- BBB:
  - Do compounds cross the blood brain barrier?
- P-gp:
  - Are compounds substrates to the P-glycoprotein?
- estrogen:
  - Are compounds ligands to an estrogen receptor?



## CHEMICAL CLASSIFICATION

João D Ferreira, Francisco Couto 2010: Semantic Similarity for Automatic Classification of Chemical Compounds. PLoS Computational Biology 9(6), e1000937

## Structural Similarity

- Fingerprints and Machine Learning

Testing set	Classification system	Accuracy	Reference
BBB	Artificial Neural Networks	75.7%	(Doniger et al., 2002)
	Random Forest	80.9%	(Svetnik et al., 2003)
	Support Vector Machines	81.5%	(Doniger et al., 2002)
P-gp	Four-point Pharmacophore	62.7%	(Penzotti et al., 2002)
	Support Vector Machines	79.4%	(Xue et al., 2004)
	Random Forest	80.6%	(Svetnik et al., 2003)
estrogen	Decision Forest	~80%	(Tong et al., 2003)
	Random Forest	82.8%	(Svetnik et al., 2003)

## Hybrid Approach

- Based on Semantic Similarity
  - CHEBI and KEGG pathways

$$\text{sim}_{\text{hybrid}} = \alpha \cdot \text{sim}_{\text{structural}} + (1 - \alpha) \cdot \text{sim}_{\text{semantic}}$$

Set	Chym Parameters	Accuracy
BBB <sub>p</sub>	FP3, simGIC, all, $\alpha = 0.28$	90.9%
P-gp <sub>p</sub>	FP4, simUL, all, $\alpha = 0.66$	87.7%
estrogen <sub>p</sub>	FP4, simGIC, role, $\alpha = 0.42$	84.2%

## New Predictions

Compound		Set	Coefficient
ID	Name		
1015	orthanilic acid	BBB <sub>p</sub>	0.503
2654	aminoglutethimide	BBB <sub>p</sub>	0.489
2089	O-methylserotonin	BBB <sub>p</sub>	0.477
3638	chloroquine	BBB <sub>p</sub>	0.475
2430	aconitine	P-gp <sub>p</sub>	0.474
1883	4-hydroxystyrene	estrogen <sub>p</sub>	0.577
5078	flavonol*	estrogen <sub>p</sub>	0.577
5262	galangin	estrogen <sub>p</sub>	0.577

\* This compound is a false positive.

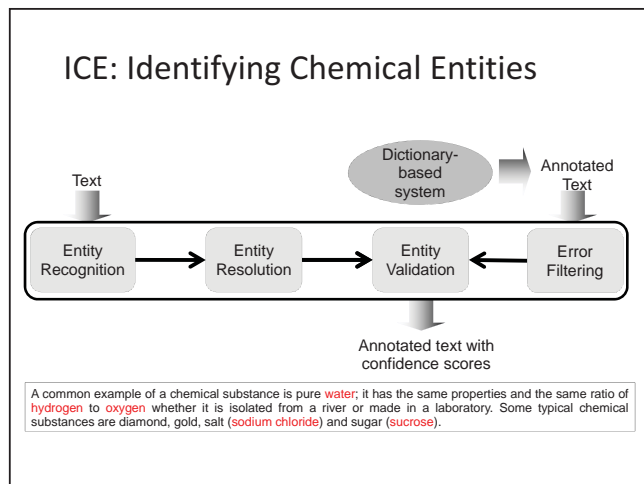
## Structural vs. Semantic

Alpha	BBB <sub>p</sub>	P-gp <sub>p</sub>	estrogen <sub>p</sub>
0.0	81.0%	74.1%	73.3%
0.1	86.9%	74.1%	74.3%
0.2	88.9%	79.0%	74.3%
0.3	<b>90.2%</b>	76.5%	79.2%
0.4	88.2%	81.5%	<b>84.2%</b>
0.5	85.0%	84.0%	83.2%
0.6	83.0%	85.2%	78.2%
0.7	83.0%	<b>86.4%</b>	81.2%
0.8	81.0%	82.7%	76.2%
0.9	77.1%	84.0%	71.3%
1.0	71.9%	85.2%	79.2%



## CHEMICAL ENTITY RECOGNITION AND MAPPING

Tiago Grego, Piotr Pezik, Francisco Couso, Dietrich Rebholz-Schuhmann, Identification of Chemical Entities in Patent Documents, 3rd International Workshop on Practical Applications of Computational Biology and Bioinformatics (IWPACB'09) 2009



### Example

- “Despite a lack of data regarding their efficacy, both **caffeine** and **doxapram** have been recommended for treatment of hypercapnia in equine neonates with central nervous system damage.” (PMID: 18371030)
- The fact that caffeine and doxapram are semantically similar
  - is an evidence for being correctly identified

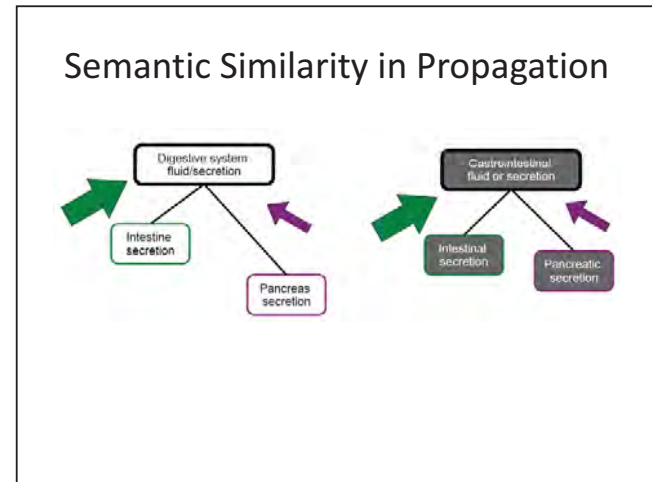
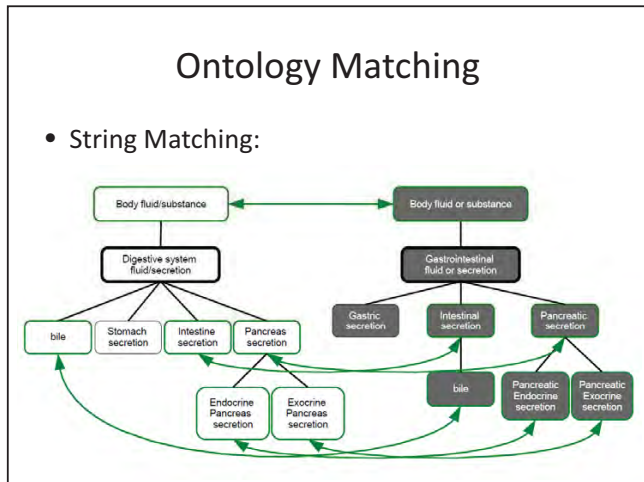
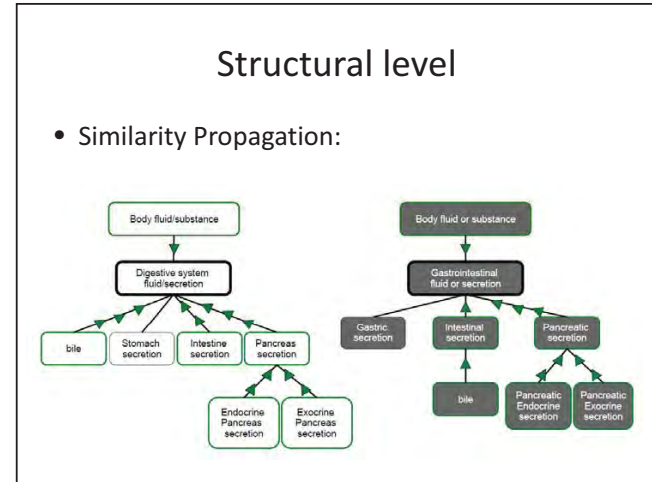
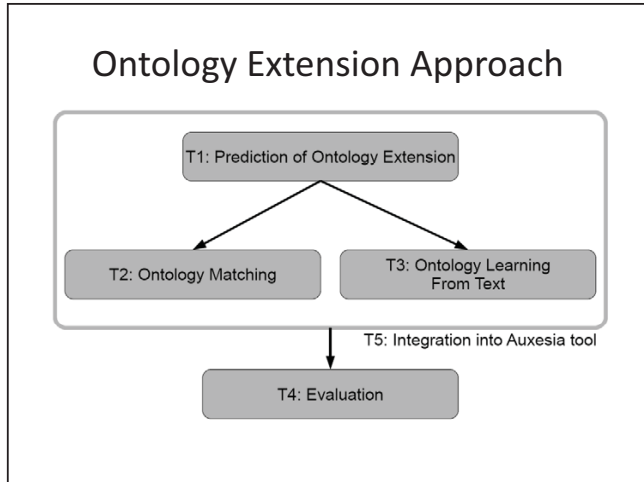
### Entity Validation

- The scope of a scientific publication is typically narrow
  - Specific protein, specific metabolic pathway, specific disease, etc
- Thus it is expected for the entities present in a document to be related
  - Measure semantic similarity between the entities



### ONTOLOGY MATCHING AND EXTENSION

Cátia Pesquita, Cosmin Stroe, Isabel F. Cruz, Francisco Couto, BLOOMS on AgreementMaker: results for OAEI 2010 - ISWC Workshop on Ontology Matching 2010





## OAEI 2011 anatomy results

- Adult Mouse Anatomy (2744 classes)
- NCI Thesaurus (3304 classes) describing the human anatomy.

System	precision	recall	f-measure
AgreementMaker (F-measure)	90.3%	85.3%	87.7%
AgreementMaker (precision)	96.2%	75.1%	84.3%
BLOOMS-XLDB	96.7%	72.5%	82.9%

## Goals

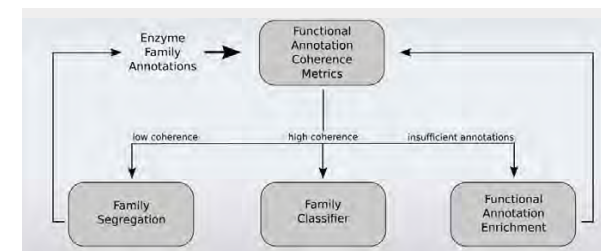
- Measure and Improve
  - functional annotation coherence of protein families.
- Enrich
  - under-annotated families with functional annotations.
- Classify
  - novel sequences into richly annotated protein families.



## ENZYME FAMILY COHERENCY ASSESSMENT

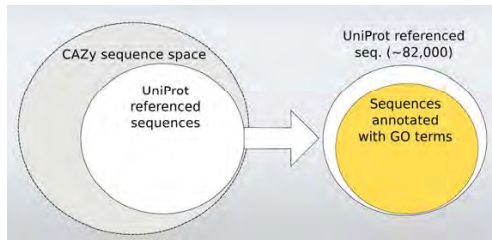
Hugo Bastos, Tiago Grego, Francisco Couto, Pedro M. Coutinho, Enzyme family coherence assessment: validation and prediction. JB'2009 - Challenges in Bioinformatics p. 26-30, Lisbon, Portugal, November, 2009.c

## Approach



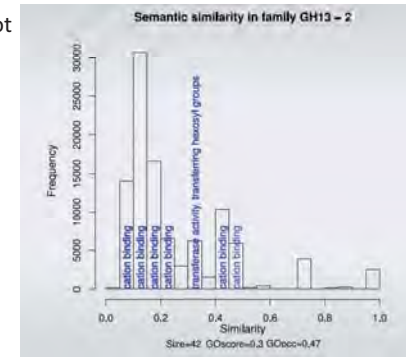
### CAZy: case-study

- Database specialized in catalytic modules of enzymes that degrade, modify or create glycosidic bonds and modules associated to carbohydrate adhesion

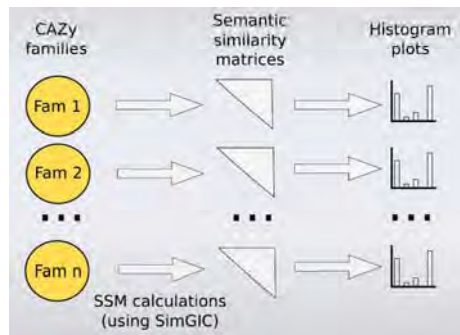


### Plots

- Histograms plot frequency of protein pairs within SSM ranges
- Labels show most frequent term with highest IC for each similarity range



### Semantic Similarity Analysis



### EPIDEMIC AND VPH INFORMATION RETRIEVAL

Luis Filipe Lopes, Fabricio A.B. Silva, Francisco Couto, João Zambir, Hugo Ferreira, Carla Sousa, Mário J. Silva, Epidemic Marketplace: An Information Management System for Epidemiological Data. Proceedings of the ITBAM - DEXA 2010 2010

## Epidemiological and Virtual Physiological Human

- Epidemic
  - Diseases
  - Symptoms
  - Anatomy
  - Phenotypes
  - Chemical substances
  - Proteins/Genes
  - Clinical Records
  - Geospatial
  - Therapeutics
  - Vaccine
  - Transmission modes
  - Organisms
  - ...
- VPH
  - Diseases
  - Symptoms
  - Anatomy
  - Phenotypes
  - Chemical substances
  - Proteins/Genes
  - Clinical Records
  - Cellular Component
  - Gene Expression
  - Cell Type
  - Pathways
  - Radiology
  - ....


## Information Retrieval

- Input
  - search keywords
- Output:
  - Ranked list of datasets and/or models
- Ranking Method
  - Semantic Similarity
  - Using multi-domain annotations
  - Mapping search keywords to ontology terms

Search Results
<input checked="" type="checkbox"/> Flu in Mexico Dataset 2010
<input checked="" type="checkbox"/> HIV in Latin America Model 2010
<input checked="" type="checkbox"/> Dengue in Brasil Dataset 2006-10
<input checked="" type="checkbox"/> Flu in Latin America Dataset 2005-09


## Challenges

- Dataset and Model annotation
  - Ontology Selection
  - Manual or/and semi-automated
- Semantic Similarity
  - Extend to another ontologies
    - Nowadays: GO, CHEBI, HPO
    - Next: FMA
- Integrate multi-domain measures
  - Weights



## OUR WEB TOOLS

<http://xldb.di.fc.ul.pt/wiki/BOA>



ProteinON  
Protein Interactions and Ontology

Home About Sources Team

Step 1: Query Step 2: Options Step 3: Input

compute protein semantic similarity

Measure: semGC  
GO type: Molecular Function  
 ignore IEA

Reset Back

Proteins

Select All		
Protein 1	Protein 2	Score
<<prev 14/36 next>>		
<input type="checkbox"/> Q5Z8B1	<input type="checkbox"/> Q9XG78	45.6%
<input type="checkbox"/> Q5Z8B1	<input type="checkbox"/> P29417	45.8%
<input type="checkbox"/> Q5Z8B1	<input type="checkbox"/> Q2PQV8	45.8%
<input type="checkbox"/> Q5Z8B1	<input type="checkbox"/> P93183	45.8%
<input type="checkbox"/> Q5XQ46	<input type="checkbox"/> P50308	100%
<input type="checkbox"/> Q4B8P5	<input type="checkbox"/> Q6P8W6	100%
<input type="checkbox"/> Q4B8P5	<input type="checkbox"/> P93185	100%
<input type="checkbox"/> Q4B8P5	<input type="checkbox"/> Q0H904	100%
<input type="checkbox"/> Q4B8P5	<input type="checkbox"/> Q5XQ46	100%
<input type="checkbox"/> Q4B8P5	<input type="checkbox"/> Q9P031	100%



ProteinON  
Protein Interactions and Ontology

Home About Sources Team

Step 1: Query Step 2: Options Step 3: Input

compute protein semantic similarity

Measure: semGC  
GO type: Molecular Function  
 ignore IEA

Proteins

P93186, Q9P184, Q9P171, Q1111 =  
P4\_242903, P91371, Q72343, Q1111  
P93186, Q9P184, Q9P171, Q1111, Q1111  
L218061, Q48615, Q9P079, Q21 =  
P93186, P93186, Q48615

Reset

The query **compute protein semantic similarity** returns the semantic similarity scores between all protein defined in multi-Reset.


The option **Measure** allows you to choose one of several semantic similarity measures: **Resnik**, **Lin** or **Song & Corbridge** measures will or without the **SCA** approach, plus the graph-based **sim3** and **simIC** measures. These measures are listed by order of performance, as calculated with protein sequence similarity.

The option **GO type** allows you to choose one of the aspects of GO: **molecular function**, **biological process** and **cellular component**.

The option **ignore IEA** lists the query to non-evidence annotations, actually evidence types: IEA, NAS, ND, NR.

Enter 2 to 1000 protein (UniProt) separated by a blank space.  
(eg: Q12345, P12322)

note: entries beyond the first 1000 will be ignored



ProteinON  
Protein Interactions and Ontology

Home About Sources Team

Step 1: Query Step 2: Options Step 3: Input

find GO terms representativity

Measure: semGC  
GO type: Molecular Function  
 ignore IEA

Reset Back

Proteins

Select All		
Protein 1	Protein 2	Score
<<prev 1/6 next>>		
<input checked="" type="checkbox"/> Q74717	<input checked="" type="checkbox"/> Q6P8W6	100%
<input checked="" type="checkbox"/> P93186	<input checked="" type="checkbox"/> Q94163	100%
<input checked="" type="checkbox"/> P93186	<input checked="" type="checkbox"/> Q721V6	100%
<input checked="" type="checkbox"/> P93186	<input checked="" type="checkbox"/> Q8M808	100%
<input checked="" type="checkbox"/> P93186	<input checked="" type="checkbox"/> P93187	100%
<input checked="" type="checkbox"/> P93186	<input checked="" type="checkbox"/> Q9SVR6	100%
<input checked="" type="checkbox"/> P93186	<input checked="" type="checkbox"/> A2Q1V7	100%
<input checked="" type="checkbox"/> P93186	<input checked="" type="checkbox"/> Q74717	100%
<input checked="" type="checkbox"/> P93186	<input checked="" type="checkbox"/> Q9HGX1	100%
<input checked="" type="checkbox"/> P93186	<input checked="" type="checkbox"/> Q4H018	100%

ReBIL Protein Interactions and Ontology

Home About Sources Team

Step 1: Query Step 2: Options Step 3: Input

compute protein semantic similarity

Measure:  GO type:  ignore NAs:

View Chart Save

Select All

Term ID	Term Name	# Proteins	% Occurrence	enrichment	Info Content
-	all	11			0.0000
-	GO:0001114	7	100.0%	0.0000	0.0000
-	GO:0004916	7	100.0%	0.0000	0.0000
-	GO:0030448	2	7.3%	2.0000	0.0002

Compare compounds

14245,15377,00011,17234,28  
157,16646

Measure:

Run

CMPSim

Compare compounds

14245,15377,00011,17234,28  
157,16646

Measure:

Run

Search

Tools

BOA BIOINFORMATICS ONTOLOGY APPLICATIONS

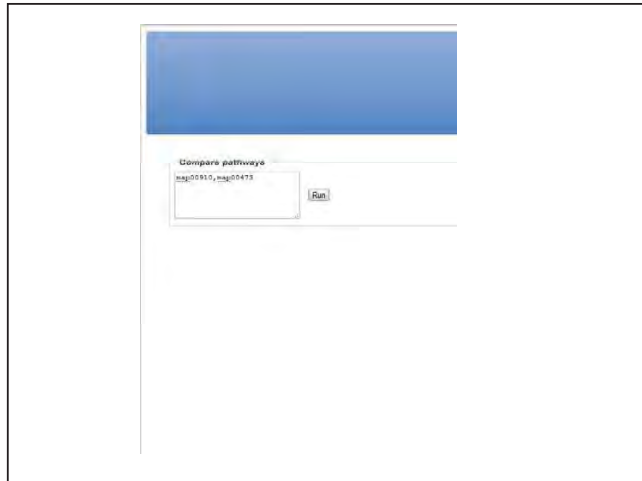
Compare compounds

14245,15377,00011,17234,28  
157,16646

Measure:

Run

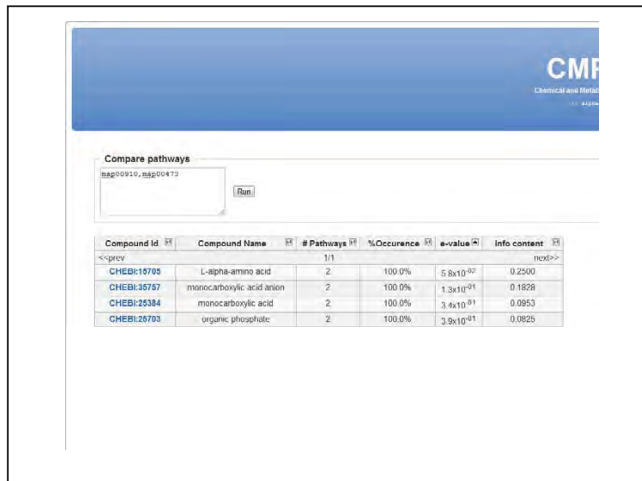
Compound 1	Compound 2	Score
enrich		
CHEBI:17234	CHEBI:28767	33.50%
CHEBI:17234	CHEBI:16646	7.98%
CHEBI:28767	CHEBI:16646	7.37%
CHEBI:46245	CHEBI:15377	6.76%
CHEBI:15377	CHEBI:16526	2.15%
CHEBI:16526	CHEBI:16646	2.05%
CHEBI:46245	CHEBI:16526	1.15%
CHEBI:16526	CHEBI:17234	1.08%
CHEBI:16526	CHEBI:28767	1.04%
CHEBI:15377	CHEBI:16646	0.75%



## CESSM

- Platform to evaluate GO semantic similarity measure
- Provides a list of 13,430 protein pairs
  - 1,039 distinct proteins
  - Blast E-value <math>< 10^{-4}</math>
  - At least one PFam and EC annotation
- Calculate their similarity with your measure
- Provides a quantitative comparison

Similarity	GI	JA	JB	JM	LA	LB	LM	RA	RB	RM	UI	tm
ECC	0.62	0.64	0.4	0.43	0.6	0.42	0.45	0.64	0.34	0.36	0.56	0.53
PFam	0.64	0.62	0.44	0.18	0.57	0.45	0.18	0.56	0.33	0.13	0.49	0.55
SeqSim	0.72	0.58	0.5	0.12	0.67	0.47	0.12	0.61	0.3	0.1	0.55	0.7



## Web Services

Ontology/concept|instance/inputList/feature/parameter1/paramValue1/...

- Ontology:
  - GO: term (concept) and protein (instance)
  - CHEBI: compound (concept) and pathway (instance)
- Feature:
  - ssm (semantic similarity measures)
  - characterization (only available for GO\_protein)
- Examples:
  - <http://10.10.4.28/~btavares/BOA/GO/term/GO:0050681,GO:0030331/ssm/measure/simGIC/IEA/false/>
  - [http://10.10.4.28/~btavares/BOA/GO/protein/Q13263,P35222/ssm/IEA/false/goType/molecular\\_function/measure/simUI](http://10.10.4.28/~btavares/BOA/GO/protein/Q13263,P35222/ssm/IEA/false/goType/molecular_function/measure/simUI)
  - <http://10.10.4.28/~btavares/BOA/Chebi/compound/46245,15377,C00011,C00049/ssm/measure/simUI>
  - <http://10.10.4.28/~btavares/BOA/Chebi/pathway/map00910,map00473/ssm/>

```

<?xml version="1.0"?>
<!DOCTYPE
  xldbData["http://www.w3.org/1999/02/22-rdf-syntax-ns#"]
  xldbSchema["http://www.w3.org/2001/XMLSchema#"]
  xldbTypes["http://www.geneontology.org/d4d4#"]
  xldbClasses["http://www.uniprot.org/d4d4#"]
  xldbInstances["http://xldb.fc.ul.pt/xldb#"]
  xldbOntology ["rdf:type" "http://www.w3.org/2001/XMLSchema#string" "gene_ontology_unidentified_version_Biological_Process/xldb/ontology"]
  xldbInstance ["rdf:type" "http://www.w3.org/2001/XMLSchema#string" "xldb/instance"]
  xldbEntity
  xldbProtein ["rdf:type" "http://www.uniprot.org/uniprot/P38083"]
  xldbAccession ["P38083" "uniprot:accession"]
  xldbProtein
  xldbEntity
  xldbProtein ["rdf:type" "http://www.uniprot.org/uniprot/P54194"]
  xldbAccession ["P54194" "uniprot:accession"]
  xldbProtein
  xldbEntity
  xldbProtein ["rdf:type" "http://www.w3.org/2001/XMLSchema#float" "0.1384" "residue"]
  xldbValue
  xldbEntity
  xldbProtein ["rdf:type" "http://www.uniprot.org/uniprot/P27018"]
  xldbAccession ["P27018" "uniprot:accession"]
  xldbProtein
  xldbEntity
  xldbProtein ["rdf:type" "http://www.uniprot.org/uniprot/P55057"]
  xldbAccession ["P55057" "uniprot:accession"]
  xldbProtein
  xldbEntity
  xldbProtein ["rdf:type" "http://www.w3.org/2001/XMLSchema#float" "0.5886" "residue"]
  xldbValue
  xldbEntity
  xldbProtein ["rdf:type" "http://www.uniprot.org/uniprot/P27913"]
  xldbAccession ["P27913" "uniprot:accession"]
  xldbProtein
  xldbEntity
  xldbProtein ["rdf:type" "http://www.uniprot.org/uniprot/P54194"]
  xldbAccession ["P54194" "uniprot:accession"]
  xldbProtein
  xldbEntity
  xldbProtein ["rdf:type" "http://www.w3.org/2001/XMLSchema#float" "0.8183" "residue"]
  xldbValue
  xldbEntity
  xldbProtein ["rdf:type" "http://www.uniprot.org/uniprot/P55057"]
  xldbAccession ["P55057" "uniprot:accession"]
  xldbProtein
  xldbEntity
  xldbProtein ["rdf:type" "http://www.uniprot.org/uniprot/P28043"]
  xldbAccession ["P28043" "uniprot:accession"]
  xldbProtein
  xldbEntity
  
```

Thanks for your attention!



Biomedical Informatics research line  
at University of Lisbon

[http://xldb.fc.ul.pt/wiki/Biomedical\\_Informatics](http://xldb.fc.ul.pt/wiki/Biomedical_Informatics)