

NETTAB 2012 Integrated Bio-Search

November 14-16, 2012, Como, Italy



Dipartimento di Elettronica e Informazione



Ranking-Aware Integration and Explorative Search of Distributed Bio-Data

Marco Masseroli, Matteo Picozzi, Giorgio Ghisalberti marco.masseroli@polimi.it



In the <u>Life Sciences</u>:

- <u>Numerous data</u>, sparsely <u>distributed</u> in <u>many</u> heterogeneous <u>sources</u>
 - Many are <u>ranked data</u> (or partially ranked) of various types, representing different phenomena, e.g.:
 - physical ordering, e.g. within a genome
 - <u>Analytical</u> order through <u>algorithmically</u> assigned <u>scores</u>,
 e.g. representing levels of <u>sequence similarity</u>
 - <u>experimentally</u> measured values, such as <u>gene expression</u> <u>levels</u>
 - The <u>ordering</u> may represent a range of different notions, such as <u>quantity</u>, <u>confidence</u>, or <u>location</u>

Life Sciences computational and data access web services

14



<u>BLAST</u> search result for the sequence "Human asparagine synthetase mRNA"

2) "5	-hvdroxvtry	optamine (seroton	in) rece	otor 2A" in UniProtKB - Mozilla Firefox					"5-hydroxytryptamine (serotonin) receptor 2A" in UniProtKB - Mozilla Firefox							
Ele	le Edit Vjew History Bookmarks Iools Help															
<	🗘 💽 🕫 🗶 🏠 💭 http://www.uniprot.org/uniprot/?query="5-hydroxytryptamine+(serotonin)+receptor+24"&sort=s 🗟 🏠 🔹 🔔 🔍 Yahool Search 🔎															
0) "5-hydroxytryptamine (serotonin) r +															
Uni	Prot 5 U	niProtKB	-			Downloads · Cor	tact · Docu	mentation/Help	^							
	Search Blast Alion Retrieve ID Manning *															
54	Search Blast Align Retrieve ID Mapping "															
F	Search in Query Protein Knowledgebase [UnProtKB] > "5-hydroxytryptamine (serotonin) receptor 2A" Search Clear Fields >>															
21 n	11 results for "5-lydroxyfryptamine (serotonin) receptor 2A" 🛛 in UniProtKB sorted by score descending 🖾															
> Re	strict term "5	hydroxytryptamine s	erotonin	Protein name	Gono namos 🊔	Organiem 🏛	Longth 🖨	Page 1 of 1								
	Accession	Entry name *	status	5 Invitronationation (Secotopin) recentor 24	Gene names *	Organism *	Lengun *	Score								
	Q543D4	Q543D4_MOUSE	*	(Putative uncharacterized protein) (5-hydroxytryptamine (Serotonin) receptor 2 A)	Htr2a (mCG_48994)	Mus musculus (Mouse)	471	7.060 [1.100×6.417]								
	Q9P2Q9	Q9P2Q9_HUMAN	*	5-hydroxytryptamine (Serotonin) receptor 2A (Fragment)	HTR2A	Homo sapiens (Human)	245	6.719 [1.523 × 4.413]								
	Q9N2F4	Q9N2F4_PANTR	*	5-hydroxytryptamine (Serotonin) receptor 2A (Fragment)	HTR2A	Pan troglodytes (Chimpanzee)	245	5.964 [1.436 × 4.153]								
	B3VRB5	B3VRB5_HUMAN	*	5-hydroxytryptamine (Serotonin) receptor 2A (Fragment)	HTR2A	Homo sapiens (Human)	66	5.690 [1.488 × 3.825]								
	B3VRB0	B3VRB0_HUMAN	*	5-hydroxytryptamine (Serotonin) receptor 2A (Fragment)	HTR2A	Homo sapiens (Human)	266	5.690 [1.443×3.943]								
	B3VRC0	B3VRC0_HUMAN	*	5-hydroxytryptamine (Serotonin) receptor 2A (Fragment)	HTR2A	Homo sapiens (Human)	137	5.690 [1.518 × 3.748]	~							
Done																

<u>UniProt</u> search result for protein "5-hydroxytryptamine (serotonin) receptor 2A"

ignment	DB:ID	Source	<u>Length</u>	<u>Score</u>	ldentity%	Positives%	<u>E0</u>
	EM_PAT:DD130059	Diagnosis and Prognosis of Breast Cancer Patients.	1992	9960	100	100	0.
	EM_PAT:DD208683	Expression Profile of Prostate Cancer.	1992	9960	100	100	0.
	EM_PAT:DD415310	Diagnosis and Prognosis of Breast Cancer Patients.	1992	9960	100	100	0.
	EM_PAT:GM974767	Sequence 120 from Patent EP2003213.	1992	9960	100	100	0.
	EM_PAT:AR274918	Sequence 55 from patent US 6506607.	1992	9960	100	100	0.
	EM_PAT:EA062820	Sequence 645 from patent US 7171311.	1992	9960	100	100	0.
	EM_PAT:EA248485	Sequence 120 from patent US 7229774.	1992	9960	100	100	0.
	EM_PAT:EA427947	Sequence 120 from patent US 7332290.	1992	9960	100	100	0.
	EM_PAT:GP320972	Sequence 645 from patent US 7514209.	1992	9960	100	100	0.
	EM_HUM:M27396	Human asparagine synthetase mRNA, complete cds.	1992	9960	100	100	0.
	EM_PAT:CQ875273	Sequence 16 from Patent WO2004076613.	1994	9895	99	99	0.
2	EM_PAT:CS063065	Sequence 49 from Patent EP1522594.	1994	9895	99	99	0.
:	EM_PAT:CS080846	Sequence 49 from Patent WO2005040414.	1994	9895	99	99	0.
i 🗖	EM_PAT:DD387278	COMPOSITIONS AND METHODS FOR THE DIAGNOSIS AND TREATMENT OF TUMOR.	1994	9895	99	99	0.
i 🗖	EM_PAT:DL464877	COMPOSITIONS, KITS, AND METHODS FOR IDENTIFICATION, ASSESSMENT, PREVENTION, AND THERAPY OF CANCER.	1994	9895	99	99	0.
	EM_PAT:FB671589	Sequence 49 from Patent EP1892306	1994	9895	99	99	Ω

🥹 Gene Expression Atlas Search Results - Gene Expression Atlas - Mozilla Firefox	
Eile Edit View History Bookmarks Yahoo! Tools Help	
🕜 💽 🗸 🏡 📧 http://www.ebi.ac.uk/gxa/qrs?gprop_0=&gval_0=&fexp	>_0=UP_DOWN&fact_(🔂 ☆ +) 😹 🔍 • Yahoo! Search 🛛 🔎
EMBL-EBI 🔢 🚦 EB-eye All Databases 🔽 Enter Text Here	Go Reset @ Give us Advanced Search GeodDack
Databases Tools EBI Groups Training Industry About Us	Help Site Index 🔂 🚭
ATLAS home about t	he project faq feedback blog das api new help
Genes Organism (all genes) up/down in 💌 Saccharomyces cerevisiae 💌	Conditions View View O Heatmap Search Atlas Image: Search Atlas Image: Search Atlas Image: Search Atlas Image: Search Atlas
e.g. ASPM, "p33 binding" e	i.g. liver, cancer, diabetes advanced search

1 2 3 4 5 ... 464 > Genes 1-10 of 4638 total found (you can refine your query) • Download all results • REST AP

Gene	φ.	Organism	Experimental Factor	Factor Value	φ	φ.	P-value¢
∃ DTD1		Saccharomyces cerevisiae	Growth condition	rehydration		4	4.42E-8
E FAS1		Saccharomyces cerevisiae	Growth condition	rehydration		1 2	1.06E-8
EMP27		Saccharomyces cerevisiae	Growth condition	rehydration		1 2	5.72E-7
■ YPR117W		Saccharomyces cerevisiae	Growth condition	rehydration		1/2	1.01E-6
∃ PDR5		Saccharomyces cerevisiae	Growth condition	rehydration		2	3.26E-9
E CHL1		Saccharomyces cerevisiae	Growth condition	rehydration		2	8.66E-6
IRA2		Saccharomyces cerevisiae	Growth condition	rehydration		2	1.59E-6
∃ TUS1		Saccharomyces cerevisiae	Growth condition	rehydration		2	1.35E-5
E POL2		Saccharomyces cerevisiae	Growth condition	rehydration		2	7.96E-8
E NCR1		Saccharomyces cerevisiae	Growth condition	rehydration		² 1	3.18E-5

Gene expression data result from <u>Array Express</u>

© Marco Masseroli, PhD

GPDW: Genomic and Proteomic Data Warehouse http://www.bioinformatics.dei.polimi.it/GPKB/



Several integrated databanks, including:

- Entrez Gene, Ensembl
- Homologene
- IPI, UniProt/Swiss-Prot
- Gene Ontology, GOA
- BioCyc, KEGG, Reactome
- InterPro, Pfam
- OMIM, eVOC, ...

Numerous integrated data, including:

- 8,085,152 genes of 8,410 organisms
- 31,347,655 proteins of 367,853 specie
- 33,252 Gene Ontology terms and 61,899 relations (is a, part of)
- 27,667 biochemical pathways
- 14,163 protein domains; 7,215 OMIM genetic disorders; ...



Data Warehouse





- Several Life Science questions:
 - are <u>complex</u>
 - to be answered <u>require</u> <u>integration</u> and <u>comprehensive</u> evaluation of different data
 - often distributed, many of which <u>ranked</u>

Answering <u>complex questions</u> requires <u>integration</u> of vertical search services to create <u>multi-topic searches</u>

• where the different topic searches either <u>refine</u> or <u>augment</u> previous search results

Bioinformatics data <u>integration platforms</u> exist

Ordered data are poorly served or no supported at all by current data integration platforms





- 1. "Which genes encode proteins in different organisms with high sequence similarity to a protein *X* and have some biomedical features in common e.g. up/down significantly co-expressed in the same biological tissue or condition *Y* and involved in the biological function *Z*?"
- 2. "Which proteins of a given biochemical pathway are encoded by co-expressed genes and are likely to interact?"
- 3. "Which proteins in different organisms are most structurally and functionally similar to a given protein?"
- 4. "Which drugs treat diseases that are likely to be associated with a given genetic mutation?"

<u>Information</u> to answer such queries is <u>available</u> on the Internet, but <u>no available software system</u> is capable of <u>computing</u> the <u>answer</u>





Common Aspects:

- Multi-topic queries (e.g. sequence similarity, gene expression)
- Ranking composition (e.g. similarity score, diff. expression p-value)
- The answers are on the Web
- A <u>knowledgeable user</u> would do the <u>query</u> <u>step-by-step</u>:
 - Search proteins similar to a given protein and get their ID
 - Search genes that codify such proteins and get their symbol
 - Search a gene expression DB and find the differential expression of such genes in the given biological condition / tissue
 - Order results by best similarity and differential expression values

After hours of painful search the user might actually succeed!

• Can this be <u>done better</u>?



<u>Search Computing</u> (SeCo) is a 5 year project funded in November 2008 by the European Research Council (ERC) Advanced Grant program It <u>aims</u>:

- 1. Develop the informatics framework required for computing <u>multi-topic searches</u> by combing single topic <u>search results</u> from search engines, which are often <u>ranked</u>, with other data and computational resources
 - directly supporting <u>multi-topic</u> ordered data
 - taking into account <u>order</u> when the <u>results</u> of several requests are <u>combined</u>
 - enabling exploration and expansion of search results
- 2. Apply SeCo technology in different fields, including <u>Life Sciences</u>
 => Bio-SeCo: <u>Support answering complex bioinformatics queries</u>



Life Science example query:

"Which genes encode proteins in different organisms with high sequence similarity to a protein X and have some biomedical features in common, e.g. up/down significantly co-expressed in the same biological tissue or condition Y and involved in a biological function Z?"

This multi-topic case study question can be <u>decomposed</u> into the following <u>four single topic sub-queries</u>, each of these sub-queries can be <u>mapped</u> to an available <u>search service</u>:

Bio-SeCo: SeCo technologies to answer Life Science questions



- "Which proteins in different organisms have high sequence similarity to a protein *X* ?"
 - → <u>BLAST</u>, a <u>sequence similarity search program</u>, in one of its many implementations, e.g. WU-BLAST (http://www.ebi.ac.uk/blast2/) or NCBI-Blast (http://blast.ncbi.nlm.nih.gov/Blast)
- "Which genes encode which proteins ?"
 - → <u>GPDW</u> (Genomic and Proteomic Data Warehouse), a <u>query</u> <u>service</u> to a database of genomic and proteomic data (<u>GPDW_protein2gene</u>)



- "Which genes are up/down significantly co-expressed in the same biological condition / tissue *Y* ?"
 - →Array Express Gene Expression Atlas, a search engine of gene expression data (http://www.ebi.ac.uk/gxa/)
- "Which genes are involved in a biological function *Z*?"
 - → <u>GPDW</u> (Genomic and Proteomic Data Warehouse), a <u>query</u> <u>service</u> to a database of genomic and proteomic data (<u>GPDW_gene2biologicalFunctionFeature</u>)





- According to the Search Computing framework each search service has to be **modelled** in order to allow an organic connection with the other services
- SeCo modelling is performed at 3 different levels: Conceptual, Logical and Physical





- For modelling each **service** they are realized:
 - a <u>Service Mart</u> (SM)
 - one or more <u>Access Patterns</u> (AP)
 - a <u>Service Interface</u> (SI)

WU-BLAST

- 1 Service Mart
- 2 Access Patterns
- 1 Service Interface

GPDW_Gene2Protein

- 1 Service Mart
- 1 Access Patterns
- 1 Service Interface

ArrayExpress

- 1 Service Mart
- 2 Access Patterns
- 1 Service Interface

DUITECNICO



Service mart

sequenceAlignmentSearch(sequenceAlignmentProgram, searchedDB, querySequence, querySequenceID, querySequenceIDName, foundSequenceSymbol, foundSequenceID, foundSequenceIDName, foundSequenceDescription, foundSequenceOrganism, bestAlignmentScore, bestAlignmentExpectation, bestAlignmentProbability, alignments(score, expectation, probability, matchQuerySequence, matchFoundSequence, matchPattern))

Ex. Access pattern

sequenceAlignmentSearch_byID(sequenceAlignmentProgram^I, searchedDB^I, querySequenceID^I, querySequenceIDName^I, foundSequenceSymbol^O, foundSequenceID^O, foundSequenceIDName^O, foundSequenceDescription^O, foundSequenceOrganism^O, bestAlignmentScore^R, bestAlignmentExpectation^R, bestAlignmentProbability^R)





Service interface

WU_BLAST_byID("Washington University BLAST", sequenceAlignmentSearch_byID, http://www.ebi.ac.uk/Tools/webservices/wsdl/WSWUBlast.wsdl)

Input example:

- seaquenchAlignmentProgram: BLASTP
- searchedDB:
- querySequenceID:

Output example:

- foundSequenceSymbol:
- foundSequenceID:
- foundSequenceOrganism:
- foundSequenceDescription:
- bestAlignmentScore:
- bestAlignmentExpectation:

- uniprotKB
- 014543 querySequenceIDName: uniprot

SOCS3_MOUSE

- O35718 foundSequenceIDName: uniprot
- Mus musculus
- Suppressor of citokine signaling 3
- 990

2.99 e⁻⁹⁸



Their <u>pair-wise</u> coupling *connection patterns* useful for computing the answer to the considered case study question are as follows:

existsCodingGene_byProteinID(sequenceAlignmentSearch, protein2gene):
 [(sequenceAlignmentSearch.foundSequenceID = protein2gene.proteinID
 AND sequenceAlignmentSearch.foundSequenceIDName =
 protein2gene.proteinIDName)]

existsExpressedGene_byGeneID(protein2gene, geneExpressionSearch): [(protein2gene.geneIDName = "ensembl" AND geneExpressionSearch.queryEnsemblGeneID = protein2gene.geneID)]



<u>Services registered</u> in the framework are <u>pair-wise related</u> each other through connection patterns that define the available resource network





POLITECNICO DI MILANO Query interface for multi-topic search http://www.bioinformatics.dei.polimi.it/bio-seco/seco/



Bio Search Computing

New Query Session

Query interface for multi-topic search





Results of sequence alignment search on NCBI-BLAST



Session () Bio Search Computing

"Which proteins in different organisms have high sequence similarity to the protein with <u>UniProt ID: P26367</u>?"

Using <u>BLAST</u>, a <u>sequence</u> <u>similarity search program</u>, in one of its implementations, e.g. NCBI-BLAST



Results of sequence alignment search on NCBI-BLAST



"Which proteins in different organisms have high sequence similarity to the protein with <u>UniProt ID: P26367</u>?"

Using <u>**BLAST**</u>, a <u>sequence similarity search program</u>, in one of its implementations, e.g. **NCBI-BLAST**

Global Tuple Data		e NCBI Blast Sequence Alignment Search by Protein ID								
	Global Score ≎	Tuple Score ≎	Found Sequence Symbol ≎	Found Sequence ID ≎	Found Sequence ID Name ≎	Found Sequence Length ≎	Found Sequence Description \$	Best Alignment Expectation ≎	Concept OID ⇔	
	1.00000	1.00000	PAX6_HUMAN	P26367	uniprot	422	paired box protein Pax-6 isoform a [Homo sapiens]	1.0E-300	2	
	1.00000	1.00000	PAX6_BOVIN	Q1LZF1	uniprot	422	RecName: Full=Paired box protein Pax-6; AltName: Full=Oculorhombin	1.0E-300	2	
	1.00000	1.00000	PAX6_RAT	P63016	uniprot	422	RecName: Full=Paired box protein Pax-6; AltName: Full=Oculorhombin	1.0E-300	2	
	1.00000	1.00000	РАХБ_СОТЈА	P47238	uniprot	416	RecName: Full=Paired box protein Pax-6; AltName: Full=Pax-QNR	1.0E-300	2	
	1.00000	1.00000	PAX6_XENLA	P55864	uniprot	422	RecName: Full=Paired box protein Pax-6	1.0E-300	2	
	1.00000	1.00000	PAX6_DANRE	P26630	uniprot	437	RecName: Full=Paired box protein Pax-6; AltName: Full=Pax[Zf-a]	1.0E-300	2	
	1.00000	1.00000	PAX6_ORYLA	073917	uniprot	437	RecName: Full=Paired box protein Pax-6	1.0E-300	2	
	0.50638	0.50638	PAX6_CHICK	P47237	uniprot	216	RecName: Full=Paired box protein Pax-6	1.22218E-152	2	
	©	Marco N	lasseroli, PhD							



"Which genes encode which proteins?"

Using a <u>query service</u> (**GPDW_protein2gene**) to our <u>**GPDW**</u> (Genomic and Proteomic Data Warehouse), e.g. for protein with UniProt ID: *P26367*

Sec d	Bio	o Searc	h Com	puting Visua	ization	Source	Session
Session O	0/	🙁 Session	1 828				
Menu	Glo	bal Tuple Data		GPD	W Gene by Prote	in ID	
0		Global Score ≎	Tuple Score ≎	Gene ID ≎	Gene ID Name ≎	Gene Concept OID ≎	Concept OID ≎
2		1.00000	1.00000	8620	hgnc	769478066	1
		1.00000	1.00000	5080	entrez_gene	760897928	1
		1.00000	1.00000	ENSG0000007372	ensembl	760897928	1
		tions					
	M	lore All More O	ne				
	Glo	obal Score	✓ ~~ ✓				
	Fi	ilter		,			
		© Marco Mas	seroli, PhD				



"Which genes are significantly up or down expressed in tumor?"

Using <u>Array Express</u> Gene Expression Atlas, a <u>search engine</u> of <u>gene expression data</u> (http://www.ebi.ac.uk/gxa/), e.g. for gene with Ensembl ID: *ENSG0000007372*

Sec.	Bio S	;earc	h Com	visualizati puting	ion Source	Session
Session O	0 2 8	Session 1	028	Session 2 🔍 🖉 🕄		
Menu 0	Sessio	n: Ses	sion 2			
	Global Tu	uple Data		Array Express Ger	e Expression Search	by Gene ID
	Globa	al Score ≎	Tuple Score	Found Gene Symi	ool ≎ ∣Best Expression	l Pvalue ≎ Concept OID ≎
	1.000	000	1.00000	РАХ6	1.0E-11	423838596
	Actions-					
	More All	More Or	ie			
	Global Sc	ore	► ~~	*		
	Filter					

© Marco Masseroli, PhD

Results of gene2biologicalFunctionFeature search on GPDW



"Which genes are involved in a biological process?

Using a <u>query service</u> **GPDW_gene2biologicalFunctionFeature**) to our <u>GPDW</u> (Genomic and Proteomic Data Warehouse), e.g. for gene with Entrez Gene ID: 9021 and biological process *regulation of*

metabolic process

Globa	al Tuple Dat	a	GF	DW Biological Function F	eature by Gene ID and Featu	re Name
	Global Score <	Tuple Score ≎ ≎	Biological Function Feature ID ≎	Biological Function Feature ID Name ≎	Biological Function Feature Name ≎	Biological Function Feature Definition ≎
	1.00000	1.00000	GO:0019220	go	regulation of phosphate metabolic process	Any process that modulates the frequency, rate or extent of the chemical reactions and pathways involving phosphates.
	1.00000	1.00000	GO:0031323	go	regulation of cellular metabolic process	Any process that modulates the frequency, rate or extent of the chemical reactions and pathways by which individual cells transform chemical substances.
	1.00000	1.00000	GO:0032268	go	regulation of cellular protein metabolic process	Any process that modulates the frequency, rate or extent of the chemical reactions and pathways involving a protein, occurring at the level of an individual cell.
	1.00000	1.00000	GO:0051174	go	regulation of phosphorus metabolic process	Any process that modulates the frequency, rate or extent of the chemical reactions and pathways involving phosphorus or compounds containing phosphorus.
	1.00000	1.00000	GO:0051246	go	regulation of protein metabolic process	Any process that modulates the frequency, rate or extent of the chemical reactions and pathways involving a protein





The submitted final <u>global query</u> included as input:

- The human *Paired box protein Pax-6 isoform a* protein (UniProt ID *P26367*) as <u>amino acid sequence X</u>
- *tumor* as pathological <u>biological condition</u> <u>Y</u>
- *regulation of programmed cell death* as <u>biological process</u> Z

Unpredictably, on October 8th 2012, Bio-SeCo discovered the <u>human</u> <u>*PAX7*</u> and <u>*PAX2*</u>, <u>mouse *Pax8*</u> and <u>human *PAX8*</u> genes, ranked by their global score of 0.90661, 0.90407, 0.90354 and 0.90289, respectively (with 1.0 as best score).

The <u>global score</u> is computed according to a <u>score function</u> as a combination of partial scores of intermediate ranked results, e.g. of ranked *sequence alignment expectation* and *gene expression p-value*



Combined search results



Visualization

Source

Session

Bio Search Computing

Session 0 🔞 🖉 🔞

			GPDW Gene by F				
	Tuple Score	e ≎ Found Sequence Syn	1bol ≎ Found Sequend	ce ID ≎ Best Alignment E	Expectation ≎		Gene ID
	0.25289	PAX7_HUMAN	P23759	1.35413E-76	N		ENSG0000009709
	0.23255	PAX2_HUMAN	Q02962	1.72295E-70	K		ENSG0000075891
	0.22831	PAX8_MOUSE	Q00288	3.22281E-69			ENSMUSG000002697
	0.22311	PAX8_HUMAN	Q06710	1.16475E-67			ENSG00000125618
1							18510

GPDW Gene by Protein ID							
Gene ID 🗧	≎ Gene ID Name ≎						
ENSG0000009709	ensembl						
ENSG00000075891	ensembl						
ENSMUSG0000026976	ensembl						
ENSG00000125618	ensembl						
18510	entrez_gene						
10010	onilioz_gono						

Menu

Array Express	Gene E	xpression	Search	by	Gene ID
---------------	--------	-----------	--------	----	---------

Tuple Score :	C Found Gene Sy	/mbol ◇ Best Expression Pvalue ◇
1.00000	PAX7	1.0E-11
1.00000	PAX2	1.0E-11
1.00000	Pax8	1.0E-11
1.00000	PAX8	1.0E-11
0.19063	emb2444	0.0080
0.17146	AT5G18510	0.013
0.13016	AT3G18510	0.037

Menu

GPDW Biological Function Feature by Gene Concept

Biological Function Feature ID 🌣	Biological Function Feature Name 🌣
GO:0043067	regulation of programmed cell death
GO:0043067	regulation of programmed cell death
GO:0043067	regulation of programmed cell death
GO:0043067	regulation of programmed cell death

© Marco Masseroli, PhD



See Bio-SeCo online at

http://www.bioinformatics.dei.polimi.it/bio-seco/seco/

Tomorrow **DEMO**

Thank you for your attention!

Any question?